

DSN Lekki–Ajah Hackathon – Employee Attrition Prediction

By: Destiny Aimufia Eronmele

Track: Machine Learning & AI

Model Used: Logistic Regression (`class_weight='balanced'`)

```
In [1]: !pip install openpyxl
```

Requirement already satisfied: openpyxl in c:\users\hp\anaconda3\lib\site-packages (3.1.5)
Requirement already satisfied: et-xmlfile in c:\users\hp\anaconda3\lib\site-packages (from openpyxl) (1.1.0)

```
In [2]: import pandas as pd
import matplotlib.pyplot as plt

plt.rcParams['figure.figsize'] = (6,4)
```

```
In [3]: !pip install openpyxl --upgrade --user
```

Requirement already satisfied: openpyxl in c:\users\hp\anaconda3\lib\site-packages (3.1.5)
Requirement already satisfied: et-xmlfile in c:\users\hp\anaconda3\lib\site-packages (from openpyxl) (1.1.0)

```
In [4]: import sys
!{sys.executable} -m pip install openpyxl
```

Requirement already satisfied: openpyxl in c:\users\hp\appdata\local\programs\python\python313\lib\site-packages (3.1.5)
Requirement already satisfied: et-xmlfile in c:\users\hp\appdata\local\programs\python\python313\lib\site-packages (from openpyxl) (2.0.0)

[notice] A new release of pip is available: 25.2 -> 25.3

[notice] To update, run: C:\Users\HP\AppData\Local\Programs\Python\Python313\python.exe -m pip install --upgrade pip

```
In [20]: import openpyxl
```

```
In [24]: ## 0. Load the Dataset
```

```
## In this section, I load the training dataset which I will use for the Exploratory Data Analysis (EDA).
```

```
train = pd.read_excel(r"C:\Users\HP\OneDrive - Microsoft-PowerUp\Documents\Train_data.xlsx")
```

```
In [28]: train.head()
```

Out[28]:

	Id	Age	BusinessTravel	DailyRate	Department	DistanceFromHome	Education	EducationField	EmployeeCount	EnvironmentSatisfaction	...	StandardHours
0	2240	25	Travel_Rarely	883	Sales	26	1	Medical	1	3	...	80
1	2722	40	Travel_Rarely	896	Research & Development	2	3	Medical	1	3	...	80
2	1090	49	Travel_Rarely	653	Research & Development	9	3	Life Sciences	1	2	...	80
3	942	47	Travel_Rarely	1147	Research & Development	8	3	Life Sciences	1	1	...	80
4	489	22	Travel_Rarely	581	Research & Development	3	1	Life Sciences	1	1	...	80

5 rows × 35 columns



In [26]:

```
## 2. Dataset Information

## This shows the data types of each column and whether any column contains missing values ##

train.info()
```

```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1914 entries, 0 to 1913
Data columns (total 35 columns):
#   Column                                Non-Null Count  Dtype
---  -
0   Id                                     1914 non-null   int64
1   Age                                   1914 non-null   int64
2   BusinessTravel                       1914 non-null   object
3   DailyRate                            1914 non-null   int64
4   Department                           1914 non-null   object
5   DistanceFromHome                     1914 non-null   int64
6   Education                             1914 non-null   int64
7   EducationField                       1914 non-null   object
8   EmployeeCount                        1914 non-null   int64
9   EnvironmentSatisfaction              1914 non-null   int64
10  Gender                               1914 non-null   object
11  HourlyRate                           1914 non-null   int64
12  JobInvolvement                       1914 non-null   int64
13  JobLevel                             1914 non-null   int64
14  JobRole                              1914 non-null   object
15  JobSatisfaction                      1914 non-null   int64
16  MaritalStatus                       1914 non-null   object
17  MonthlyIncome                       1914 non-null   int64
18  MonthlyRate                          1914 non-null   int64
19  NumCompaniesWorked                  1914 non-null   int64
20  Over18                              1914 non-null   object
21  OverTime                             1914 non-null   object
22  PercentSalaryHike                   1914 non-null   int64
23  PerformanceRating                   1914 non-null   int64
24  RelationshipSatisfaction             1914 non-null   int64
25  StandardHours                       1914 non-null   int64
26  StockOptionLevel                    1914 non-null   int64
27  TotalWorkingYears                   1914 non-null   int64
28  TrainingTimesLastYear               1914 non-null   int64
29  WorkLifeBalance                     1914 non-null   int64
30  YearsAtCompany                      1914 non-null   int64
31  YearsInCurrentRole                  1914 non-null   int64
32  YearsSinceLastPromotion              1914 non-null   int64
33  YearsWithCurrManager                 1914 non-null   int64
34  Attrition                           1914 non-null   int64
dtypes: int64(27), object(8)
memory usage: 523.5+ KB

```

```
In [ ]: ## interpretation ##
```

This gives me an overview of:

- Which columns are numeric
- Which columns are categorical
- Whether any column contains missing values
- How many non-null entries exist per column

It helps me plan what needs encoding or cleaning later.

```
In [ ]: train.shape
## 1. Dataset Shape

## This shows the number of rows (employees) and number of columns (features) in the dataset ##.
```

```
In [ ]: ## Interpretation ##

## The dataset contains 1914 rows (employees) and 35 columns (employee features) ## .
## This helps me understand the size and complexity of the dataset ##.
```

```
In [30]: ## 3. Descriptive Statistics ##

## This provides summary statistics (mean, median, min, max, std) for all numeric columns ##.

train.describe()
```

```
Out[30]:
```

	Id	Age	DailyRate	DistanceFromHome	Education	EmployeeCount	EnvironmentSatisfaction	HourlyRate	JobInvolvement	JobLevel
count	1914.000000	1914.000000	1914.000000	1914.000000	1914.000000	1914.0	1914.000000	1914.000000	1914.000000	1914.000000
mean	1370.878265	36.601881	866.169801	8.521421	2.920585	1.0	2.735632	66.728840	2.765935	2.063218
std	789.635369	8.825503	385.083305	7.781102	1.003374	0.0	1.095118	19.819761	0.674817	1.122257
min	0.000000	18.000000	103.000000	1.000000	1.000000	1.0	1.000000	30.000000	1.000000	1.000000
25%	696.250000	30.000000	561.000000	2.000000	2.000000	1.0	2.000000	50.000000	2.000000	1.000000
50%	1378.500000	35.000000	881.000000	6.000000	3.000000	1.0	3.000000	66.500000	3.000000	2.000000
75%	2051.750000	42.000000	1210.750000	11.000000	4.000000	1.0	4.000000	83.000000	3.000000	3.000000
max	2732.000000	60.000000	1499.000000	29.000000	5.000000	1.0	4.000000	100.000000	4.000000	7.000000

8 rows × 27 columns



```
In [ ]: ## Interpretation ##

## This helps me understand the distribution, spread, and range of numerical features like Age, MonthlyIncome, YearsAtCompany, etc ##.
## It also shows if certain features have extreme values (outliers) ## .
```

```
In [32]: ## Missing Values Check ##

## This shows whether any column contains missing values, which would need cleaning before modeling ##.

train.isnull().sum()
```

```
Out[32]: Id                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
         Attrition           0
         dtype: int64
```

NOW THE CHARTS WITH CODE + INTERPRETATION

5. Attrition Distribution

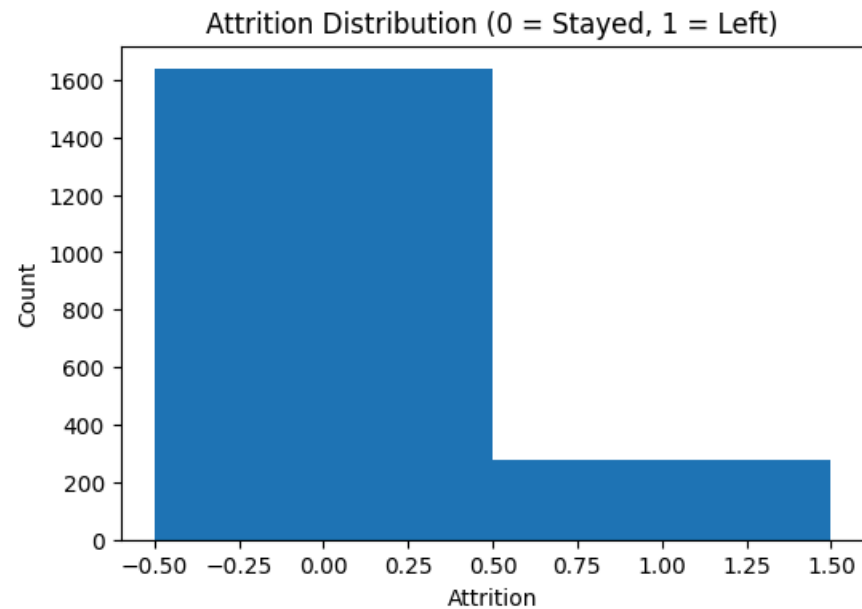
This shows the number of employees who stayed (0) and who left (1) ## .

It helps detect data imbalance ##.

```
In [34]: plt.figure()
         plt.hist(train['Attrition'], bins=[-0.5,0.5,1.5])
```

```
plt.title("Attrition Distribution (0 = Stayed, 1 = Left)")
plt.xlabel("Attrition")
plt.ylabel("Count")
plt.show()
```

```
train['Attrition'].value_counts(normalize=True)
```



```
Out[34]: Attrition
0      0.854754
1      0.145246
Name: proportion, dtype: float64
```

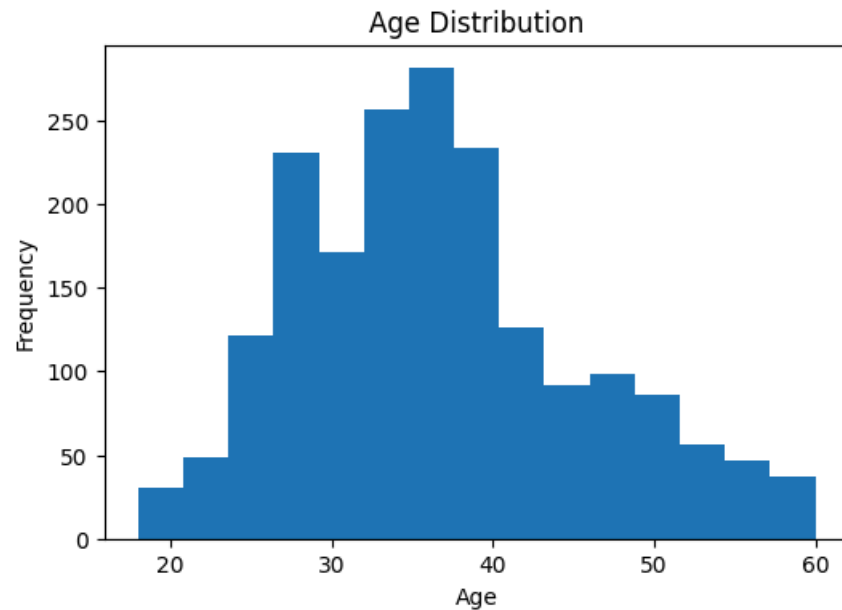
Interpretation

- Most employees (85%) stayed (Attrition = 0)
- Only about 15% left (Attrition = 1)
- The dataset is imbalanced , so later we will apply balancing methods.
- This is important because imbalance can make models biased toward predicting "stay."

6. Age Distribution

This shows the range of employee ages and helps describe the workforce demographics ##.

```
In [36]: plt.figure()
plt.hist(train['Age'], bins=15)
plt.title("Age Distribution")
plt.xlabel("Age")
plt.ylabel("Frequency")
plt.show()
```



Interpretation

Most employees are between 28 and 40 years old, showing a mostly mid-career workforce ## .

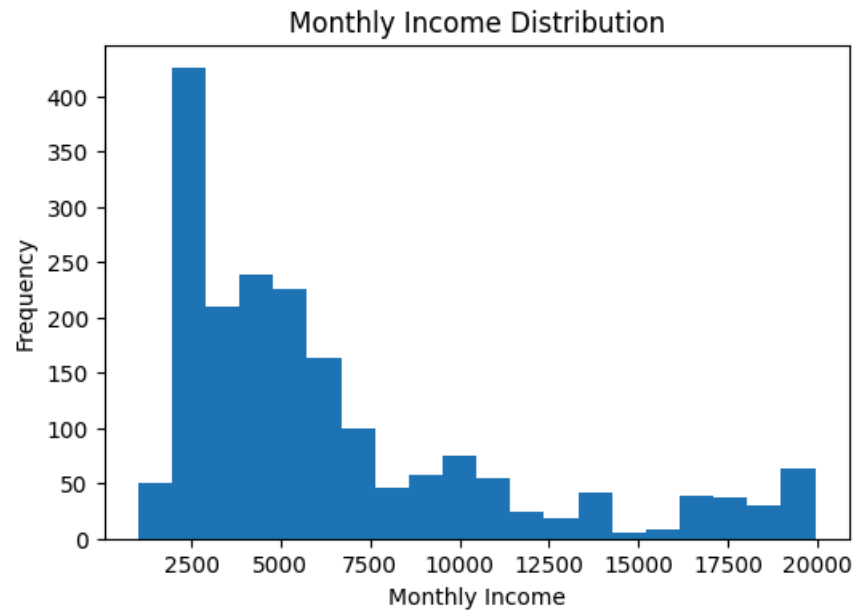
This may influence attrition patterns since younger workers tend to switch jobs more ## .

7. Monthly Income Distribution

This shows how salaries are distributed across employees ## .

```
In [38]: plt.figure()
plt.hist(train['MonthlyIncome'], bins=20)
plt.title("Monthly Income Distribution")
plt.xlabel("Monthly Income")
```

```
plt.ylabel("Frequency")
plt.show()
```



Interpretation

Income is right-skewed, meaning ##:

- Many employees earn lower salaries
- Few employees earn very high salaries

This suggests income inequality, which is important because salary affects attrition ## .

8. Monthly Income vs Attrition

This compares the salaries of employees who stayed vs those who left ## .

```
In [40]: income_stayed = train[train['Attrition']==0]['MonthlyIncome']
income_left = train[train['Attrition']==1]['MonthlyIncome']

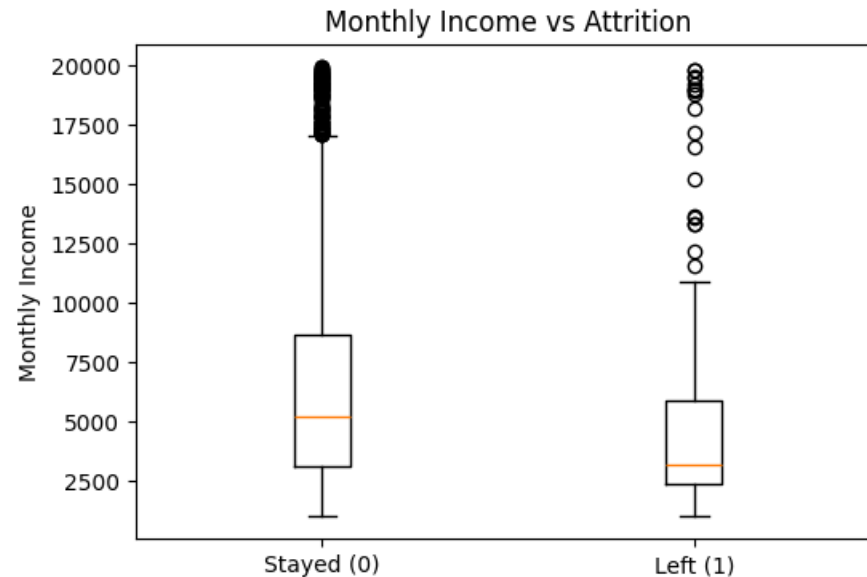
plt.figure()
plt.boxplot([income_stayed, income_left], labels=['Stayed (0)', 'Left (1)'])
```



```
plt.title("Monthly Income vs Attrition")
plt.ylabel("Monthly Income")
plt.show()
```

C:\Users\HP\AppData\Local\Temp\ipykernel_23972\1853057602.py:5: MatplotlibDeprecationWarning: The 'labels' parameter of boxplot() has been renamed 'tick_labels' since Matplotlib 3.9; support for the old name will be dropped in 3.11.

```
plt.boxplot([income_stayed, income_left], labels=['Stayed (0)', 'Left (1)'])
```



Interpretation

- Employees who left earn visibly less than those who stayed.
- Lower income is associated with higher attrition.

This confirms salary is a key predictor of attrition ## .

9. OverTime vs Attrition

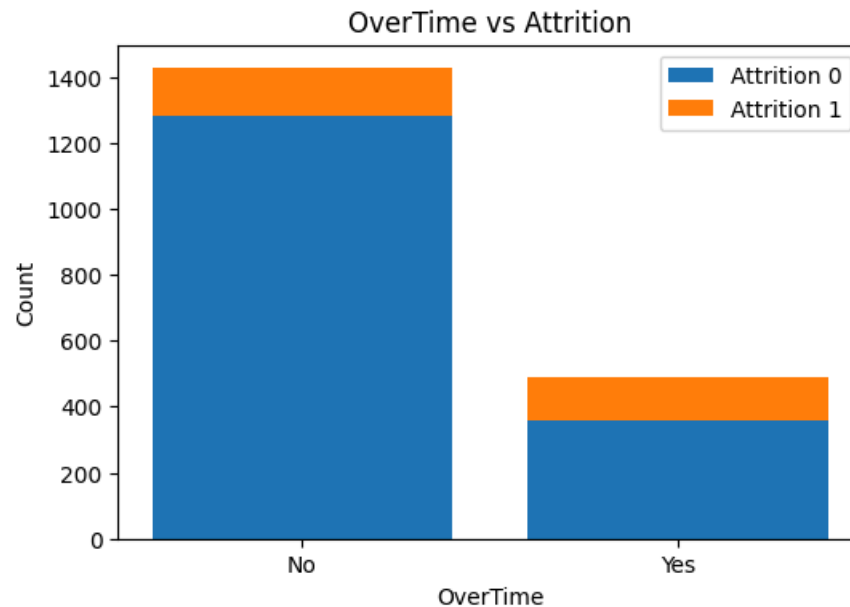
This examines whether working overtime influences an employee's likelihood of leaving ## .

```
In [42]: ct = train.groupby(['OverTime', 'Attrition']).size().unstack(fill_value=0)

plt.figure()
for i, col in enumerate(ct.columns):
    bottom = ct[ct.columns[:i]].sum(axis=1) if i>0 else None
```

```
plt.bar(ct.index.astype(str), ct[col], bottom=bottom, label=f'Attrition {col}')

plt.title("OverTime vs Attrition")
plt.xlabel("OverTime")
plt.ylabel("Count")
plt.legend()
plt.show()
```



Interpretation

- Employees who work overtime have a significantly higher number of attrition cases.
- OverTime is one of the strongest predictors of employees leaving.

This suggests overtime leads to stress and burnout ## .

Moving on to the Preprocessing & Feature Engineering Step

```
In [ ]: ## 1. Drop Irrelevant Columns ##

## Some columns do not provide any useful information for predicting attrition because they contain the same value for all employees (no variance) ##

## These include ## :
- EmployeeCount
- Over18
```

```
- StandardHours

## I will drop them from the dataset ## .
```

```
In [44]: cols_to_drop = ['EmployeeCount', 'Over18', 'StandardHours']

train = train.drop(columns=cols_to_drop)
train.head()
```

```
Out[44]:
```

	Id	Age	BusinessTravel	DailyRate	Department	DistanceFromHome	Education	EducationField	EnvironmentSatisfaction	Gender	...	RelationshipSatisfaction
0	2240	25	Travel_Rarely	883	Sales	26	1	Medical	3	Female	...	2
1	2722	40	Travel_Rarely	896	Research & Development	2	3	Medical	3	Male	...	3
2	1090	49	Travel_Rarely	653	Research & Development	9	3	Life Sciences	2	Female	...	3
3	942	47	Travel_Rarely	1147	Research & Development	8	3	Life Sciences	1	Male	...	3
4	489	22	Travel_Rarely	581	Research & Development	3	1	Life Sciences	1	Female	...	3

5 rows × 32 columns



```
In [ ]: ## These dropped columns do not help the model learn anything and would only add noise ## .
## Removing them makes the dataset cleaner and improves model performance ## .
```

2. Encode Categorical Variables

Machine learning models require numerical input **##** .

Here, I convert categorical columns into numeric format using OneHotEncoding (for multi-category) and LabelEncoding (for binary features) **##** .

```
In [48]: !pip install scikit-learn
```

```
Requirement already satisfied: scikit-learn in c:\users\hp\anaconda3\lib\site-packages (1.5.1)
Requirement already satisfied: numpy>=1.19.5 in c:\users\hp\anaconda3\lib\site-packages (from scikit-learn) (1.26.4)
Requirement already satisfied: scipy>=1.6.0 in c:\users\hp\anaconda3\lib\site-packages (from scikit-learn) (1.13.1)
Requirement already satisfied: joblib>=1.2.0 in c:\users\hp\anaconda3\lib\site-packages (from scikit-learn) (1.4.2)
Requirement already satisfied: threadpoolctl>=3.1.0 in c:\users\hp\anaconda3\lib\site-packages (from scikit-learn) (3.5.0)
```

```
In [52]: import sys
```

```
!{sys.executable} -m pip install scikit-learn
```

[illegible]

[illegible]

		4.2/8.7	MB	207.5	kB/s	eta	0:00:22
-----		4.2/8.7	MB	207.5	kB/s	eta	0:00:22
-----		4.5/8.7	MB	208.8	kB/s	eta	0:00:21
-----		4.5/8.7	MB	208.8	kB/s	eta	0:00:21
-----		4.5/8.7	MB	208.8	kB/s	eta	0:00:21
-----		4.5/8.7	MB	208.8	kB/s	eta	0:00:21
-----		4.5/8.7	MB	208.8	kB/s	eta	0:00:21
-----		4.5/8.7	MB	208.8	kB/s	eta	0:00:21
-----		4.5/8.7	MB	208.8	kB/s	eta	0:00:21
-----		4.5/8.7	MB	208.8	kB/s	eta	0:00:21
-----		4.5/8.7	MB	208.8	kB/s	eta	0:00:21
-----		4.7/8.7	MB	200.2	kB/s	eta	0:00:20
-----		4.7/8.7	MB	200.2	kB/s	eta	0:00:20
-----		4.7/8.7	MB	200.2	kB/s	eta	0:00:20
-----		4.7/8.7	MB	200.2	kB/s	eta	0:00:20
-----		4.7/8.7	MB	200.2	kB/s	eta	0:00:20
-----		4.7/8.7	MB	200.2	kB/s	eta	0:00:20
-----		5.0/8.7	MB	201.0	kB/s	eta	0:00:19
-----		5.0/8.7	MB	201.0	kB/s	eta	0:00:19
-----		5.2/8.7	MB	208.3	kB/s	eta	0:00:17
-----		5.5/8.7	MB	216.7	kB/s	eta	0:00:15
-----		5.5/8.7	MB	216.7	kB/s	eta	0:00:15
-----		5.8/8.7	MB	224.1	kB/s	eta	0:00:14
-----		5.8/8.7	MB	224.1	kB/s	eta	0:00:14
-----		5.8/8.7	MB	224.1	kB/s	eta	0:00:14
-----		5.8/8.7	MB	224.1	kB/s	eta	0:00:14
-----		5.8/8.7	MB	224.1	kB/s	eta	0:00:14
-----		6.0/8.7	MB	222.6	kB/s	eta	0:00:13
-----		6.0/8.7	MB	222.6	kB/s	eta	0:00:13
-----		6.0/8.7	MB	222.6	kB/s	eta	0:00:13
-----		6.0/8.7	MB	222.6	kB/s	eta	0:00:13
-----		6.0/8.7	MB	222.6	kB/s	eta	0:00:13
-----		6.0/8.7	MB	222.6	kB/s	eta	0:00:13
-----		6.0/8.7	MB	222.6	kB/s	eta	0:00:13
-----		6.0/8.7	MB	222.6	kB/s	eta	0:00:13
-----		6.0/8.7	MB	222.6	kB/s	eta	0:00:13
-----		6.3/8.7	MB	215.4	kB/s	eta	0:00:12
-----		6.6/8.7	MB	223.5	kB/s	eta	0:00:10
-----		6.6/8.7	MB	223.5	kB/s	eta	0:00:10
-----		6.6/8.7	MB	223.5	kB/s	eta	0:00:10
-----		6.8/8.7	MB	227.3	kB/s	eta	0:00:09
-----		6.8/8.7	MB	227.3	kB/s	eta	0:00:09
-----		6.8/8.7	MB	227.3	kB/s	eta	0:00:09
-----		6.8/8.7	MB	227.3	kB/s	eta	0:00:09
-----		6.8/8.7	MB	227.3	kB/s	eta	0:00:09
-----		6.8/8.7	MB	227.3	kB/s	eta	0:00:09
-----		7.1/8.7	MB	218.7	kB/s	eta	0:00:08
-----		7.1/8.7	MB	218.7	kB/s	eta	0:00:08

[illegible]

Downloading scipy-1.16.3-cp313-cp313-win_amd64.whl (38.5 MB)

-----	0.3/38.5	MB	?	eta	-:--:--
-----	0.3/38.5	MB	?	eta	-:--:--
-----	0.3/38.5	MB	?	eta	-:--:--
-----	0.3/38.5	MB	?	eta	-:--:--
-----	0.3/38.5	MB	?	eta	-:--:--
-----	0.3/38.5	MB	?	eta	-:--:--
-----	0.3/38.5	MB	?	eta	-:--:--
-----	0.3/38.5	MB	?	eta	-:--:--
-----	0.3/38.5	MB	?	eta	-:--:--
-----	0.3/38.5	MB	?	eta	-:--:--
-----	0.3/38.5	MB	?	eta	-:--:--
-----	0.5/38.5	MB	81.5	kB/s	eta 0:07:47
-----	0.5/38.5	MB	81.5	kB/s	eta 0:07:47
-----	0.5/38.5	MB	81.5	kB/s	eta 0:07:47
-----	0.5/38.5	MB	81.5	kB/s	eta 0:07:47
-----	0.5/38.5	MB	81.5	kB/s	eta 0:07:47
-----	0.5/38.5	MB	81.5	kB/s	eta 0:07:47
-----	0.8/38.5	MB	119.5	kB/s	eta 0:05:16
-----	1.0/38.5	MB	166.3	kB/s	eta 0:03:46
-----	1.0/38.5	MB	166.3	kB/s	eta 0:03:46
-----	1.3/38.5	MB	201.8	kB/s	eta 0:03:05
-----	1.3/38.5	MB	201.8	kB/s	eta 0:03:05
-----	1.3/38.5	MB	201.8	kB/s	eta 0:03:05
-----	1.6/38.5	MB	226.9	kB/s	eta 0:02:43
-----	1.6/38.5	MB	226.9	kB/s	eta 0:02:43
-----	1.6/38.5	MB	226.9	kB/s	eta 0:02:43
-----	1.6/38.5	MB	226.9	kB/s	eta 0:02:43
-----	1.6/38.5	MB	226.9	kB/s	eta 0:02:43
-----	1.6/38.5	MB	226.9	kB/s	eta 0:02:43
-----	1.6/38.5	MB	226.9	kB/s	eta 0:02:43
-----	1.6/38.5	MB	226.9	kB/s	eta 0:02:43
-----	1.6/38.5	MB	226.9	kB/s	eta 0:02:43
-----	1.8/38.5	MB	202.6	kB/s	eta 0:03:02
-----	2.1/38.5	MB	233.2	kB/s	eta 0:02:37
-----	2.1/38.5	MB	233.2	kB/s	eta 0:02:37
-----	2.4/38.5	MB	250.4	kB/s	eta 0:02:25
-----	2.4/38.5	MB	250.4	kB/s	eta 0:02:25
-----	2.4/38.5	MB	250.4	kB/s	eta 0:02:25
-----	2.4/38.5	MB	250.4	kB/s	eta 0:02:25
-----	2.4/38.5	MB	250.4	kB/s	eta 0:02:25
-----	2.4/38.5	MB	250.4	kB/s	eta 0:02:25
-----	2.4/38.5	MB	250.4	kB/s	eta 0:02:25
-----	2.4/38.5	MB	250.4	kB/s	eta 0:02:25
-----	2.4/38.5	MB	250.4	kB/s	eta 0:02:25
-----	2.4/38.5	MB	250.4	kB/s	eta 0:02:25
-----	2.4/38.5	MB	250.4	kB/s	eta 0:02:25
-----	2.6/38.5	MB	217.3	kB/s	eta 0:02:46
-----	2.6/38.5	MB	217.3	kB/s	eta 0:02:46
-----	2.6/38.5	MB	217.3	kB/s	eta 0:02:46
-----	2.6/38.5	MB	217.3	kB/s	eta 0:02:46
-----	2.6/38.5	MB	217.3	kB/s	eta 0:02:46

[illegible]

[illegible]

-----	6.6/38.5	MB	195.2	kB/s	eta	0:02:44
-----	6.6/38.5	MB	195.2	kB/s	eta	0:02:44
-----	6.6/38.5	MB	195.2	kB/s	eta	0:02:44
-----	6.6/38.5	MB	195.2	kB/s	eta	0:02:44
-----	6.6/38.5	MB	195.2	kB/s	eta	0:02:44
-----	6.6/38.5	MB	195.2	kB/s	eta	0:02:44
-----	6.8/38.5	MB	176.3	kB/s	eta	0:03:00
-----	6.8/38.5	MB	176.3	kB/s	eta	0:03:00
-----	6.8/38.5	MB	176.3	kB/s	eta	0:03:00
-----	7.1/38.5	MB	183.7	kB/s	eta	0:02:52
-----	7.1/38.5	MB	183.7	kB/s	eta	0:02:52
-----	7.3/38.5	MB	189.5	kB/s	eta	0:02:45
-----	7.3/38.5	MB	189.5	kB/s	eta	0:02:45
-----	7.3/38.5	MB	189.5	kB/s	eta	0:02:45
-----	7.6/38.5	MB	193.2	kB/s	eta	0:02:41
-----	7.6/38.5	MB	193.2	kB/s	eta	0:02:41
-----	7.6/38.5	MB	193.2	kB/s	eta	0:02:41
-----	7.6/38.5	MB	193.2	kB/s	eta	0:02:41
-----	7.9/38.5	MB	189.1	kB/s	eta	0:02:43
-----	7.9/38.5	MB	189.1	kB/s	eta	0:02:43
-----	8.1/38.5	MB	196.0	kB/s	eta	0:02:36
-----	8.1/38.5	MB	196.0	kB/s	eta	0:02:36
-----	8.1/38.5	MB	196.0	kB/s	eta	0:02:36
-----	8.4/38.5	MB	199.9	kB/s	eta	0:02:31
-----	8.4/38.5	MB	199.9	kB/s	eta	0:02:31
-----	8.4/38.5	MB	199.9	kB/s	eta	0:02:31
-----	8.4/38.5	MB	199.9	kB/s	eta	0:02:31
-----	8.7/38.5	MB	202.6	kB/s	eta	0:02:28
-----	8.7/38.5	MB	202.6	kB/s	eta	0:02:28
-----	8.7/38.5	MB	202.6	kB/s	eta	0:02:28
-----	8.7/38.5	MB	202.6	kB/s	eta	0:02:28
-----	8.7/38.5	MB	202.6	kB/s	eta	0:02:28
-----	8.7/38.5	MB	202.6	kB/s	eta	0:02:28
-----	8.7/38.5	MB	202.6	kB/s	eta	0:02:28
-----	8.7/38.5	MB	202.6	kB/s	eta	0:02:28
-----	8.7/38.5	MB	202.6	kB/s	eta	0:02:28
-----	8.7/38.5	MB	202.6	kB/s	eta	0:02:28
-----	8.7/38.5	MB	202.6	kB/s	eta	0:02:28
-----	8.7/38.5	MB	202.6	kB/s	eta	0:02:28
-----	8.7/38.5	MB	202.6	kB/s	eta	0:02:28
-----	8.7/38.5	MB	202.6	kB/s	eta	0:02:28
-----	8.9/38.5	MB	193.0	kB/s	eta	0:02:34
-----	9.2/38.5	MB	201.0	kB/s	eta	0:02:27
-----	9.2/38.5	MB	201.0	kB/s	eta	0:02:27
-----	9.2/38.5	MB	201.0	kB/s	eta	0:02:27
-----	9.4/38.5	MB	201.0	kB/s	eta	0:02:27

[illegible]

[illegible]

[illegible]

[illegible]

[illegible]

[illegible]

[illegible]

[illegible]

[illegible]

[illegible]

[illegible]

[illegible]

[illegible]

[illegible]

[illegible]

[illegible]

[illegible]

[illegible]

[illegible]

[illegible]

[illegible]

[illegible]

[illegible]

[illegible]

-----	34.1/38.5	MB	61.8	kB/s	eta	0:01:12
-----	34.1/38.5	MB	61.8	kB/s	eta	0:01:12
-----	34.1/38.5	MB	61.8	kB/s	eta	0:01:12
-----	34.1/38.5	MB	61.8	kB/s	eta	0:01:12
-----	34.1/38.5	MB	61.8	kB/s	eta	0:01:12
-----	34.1/38.5	MB	61.8	kB/s	eta	0:01:12
-----	34.3/38.5	MB	62.6	kB/s	eta	0:01:07
-----	34.3/38.5	MB	62.6	kB/s	eta	0:01:07
-----	34.3/38.5	MB	62.6	kB/s	eta	0:01:07
-----	34.3/38.5	MB	62.6	kB/s	eta	0:01:07
-----	34.6/38.5	MB	64.4	kB/s	eta	0:01:01
-----	34.6/38.5	MB	64.4	kB/s	eta	0:01:01
-----	34.6/38.5	MB	64.4	kB/s	eta	0:01:01
-----	34.6/38.5	MB	64.4	kB/s	eta	0:01:01
-----	34.9/38.5	MB	71.8	kB/s	eta	0:00:52
-----	34.9/38.5	MB	71.8	kB/s	eta	0:00:52
-----	34.9/38.5	MB	71.8	kB/s	eta	0:00:52
-----	34.9/38.5	MB	71.8	kB/s	eta	0:00:52
-----	35.1/38.5	MB	74.7	kB/s	eta	0:00:46
-----	35.1/38.5	MB	74.7	kB/s	eta	0:00:46
-----	35.1/38.5	MB	74.7	kB/s	eta	0:00:46
-----	35.1/38.5	MB	74.7	kB/s	eta	0:00:46
-----	35.4/38.5	MB	82.0	kB/s	eta	0:00:39
-----	35.4/38.5	MB	82.0	kB/s	eta	0:00:39
-----	35.4/38.5	MB	82.0	kB/s	eta	0:00:39
-----	35.4/38.5	MB	82.0	kB/s	eta	0:00:39
-----	35.7/38.5	MB	83.2	kB/s	eta	0:00:35
-----	35.7/38.5	MB	83.2	kB/s	eta	0:00:35
-----	35.7/38.5	MB	83.2	kB/s	eta	0:00:35
-----	35.7/38.5	MB	83.2	kB/s	eta	0:00:35
-----	35.7/38.5	MB	83.2	kB/s	eta	0:00:35
-----	35.7/38.5	MB	83.2	kB/s	eta	0:00:35
-----	35.9/38.5	MB	88.0	kB/s	eta	0:00:30
-----	35.9/38.5	MB	88.0	kB/s	eta	0:00:30
-----	35.9/38.5	MB	88.0	kB/s	eta	0:00:30
-----	35.9/38.5	MB	88.0	kB/s	eta	0:00:30
-----	35.9/38.5	MB	88.0	kB/s	eta	0:00:30
-----	35.9/38.5	MB	88.0	kB/s	eta	0:00:30
-----	36.2/38.5	MB	95.0	kB/s	eta	0:00:25
-----	36.4/38.5	MB	103.6	kB/s	eta	0:00:21
-----	36.4/38.5	MB	103.6	kB/s	eta	0:00:21
-----	36.4/38.5	MB	103.6	kB/s	eta	0:00:21
-----	36.4/38.5	MB	103.6	kB/s	eta	0:00:21
-----	36.7/38.5	MB	109.7	kB/s	eta	0:00:17
-----	36.7/38.5	MB	109.7	kB/s	eta	0:00:17
-----	36.7/38.5	MB	109.7	kB/s	eta	0:00:17
-----	36.7/38.5	MB	109.7	kB/s	eta	0:00:17

```
Downloading threadpoolctl-3.6.0-py3-none-any.whl (18 kB)
Installing collected packages: threadpoolctl, scipy, joblib, scikit-learn
```

[illegible]

[illegible]


```
----- 3/4 [scikit-learn]
----- 3/4 [scikit-learn]
----- 3/4 [scikit-learn]
----- 3/4 [scikit-learn]
----- 3/4 [scikit-learn]
----- 3/4 [scikit-learn]
----- 3/4 [scikit-learn]
----- 3/4 [scikit-learn]
----- 3/4 [scikit-learn]
----- 3/4 [scikit-learn]
----- 3/4 [scikit-learn]
----- 3/4 [scikit-learn]
----- 3/4 [scikit-learn]
----- 3/4 [scikit-learn]
----- 4/4 [scikit-learn]
```

Successfully installed joblib-1.5.2 scikit-learn-1.7.2 scipy-1.16.3 threadpoolctl-3.6.0

```
WARNING: Connection timed out while downloading.
WARNING: Attempting to resume incomplete download (13.9 MB/38.5 MB, attempt 1)
WARNING: Retrying (Retry(total=4, connect=None, read=None, redirect=None, status=None)) after connection broken by 'NewConnectionError('<pip._vendor.urllib3.connection.HTTPSConnection object at 0x0000024BFCD242D0>: Failed to establish a new connection: [Errno 11001] getaddrinfo failed')': /packages/cd/01/1204382461fcbfeb05b6161b594f4007e78b6eba9b375382f79153172b4d/scipy-1.16.3-cp313-cp313-win_amd64.whl
WARNING: Retrying (Retry(total=3, connect=None, read=None, redirect=None, status=None)) after connection broken by 'NewConnectionError('<pip._vendor.urllib3.connection.HTTPSConnection object at 0x0000024BFCD274D0>: Failed to establish a new connection: [Errno 11001] getaddrinfo failed')': /packages/cd/01/1204382461fcbfeb05b6161b594f4007e78b6eba9b375382f79153172b4d/scipy-1.16.3-cp313-cp313-win_amd64.whl
WARNING: Retrying (Retry(total=2, connect=None, read=None, redirect=None, status=None)) after connection broken by 'NewConnectionError('<pip._vendor.urllib3.connection.HTTPSConnection object at 0x0000024BFCD27890>: Failed to establish a new connection: [Errno 11001] getaddrinfo failed')': /packages/cd/01/1204382461fcbfeb05b6161b594f4007e78b6eba9b375382f79153172b4d/scipy-1.16.3-cp313-cp313-win_amd64.whl
WARNING: Retrying (Retry(total=1, connect=None, read=None, redirect=None, status=None)) after connection broken by 'NewConnectionError('<pip._vendor.urllib3.connection.HTTPSConnection object at 0x0000024BFCD27B10>: Failed to establish a new connection: [Errno 11001] getaddrinfo failed')': /packages/cd/01/1204382461fcbfeb05b6161b594f4007e78b6eba9b375382f79153172b4d/scipy-1.16.3-cp313-cp313-win_amd64.whl
WARNING: Retrying (Retry(total=0, connect=None, read=None, redirect=None, status=None)) after connection broken by 'NewConnectionError('<pip._vendor.urllib3.connection.HTTPSConnection object at 0x0000024BFCD27D90>: Failed to establish a new connection: [Errno 11001] getaddrinfo failed')': /packages/cd/01/1204382461fcbfeb05b6161b594f4007e78b6eba9b375382f79153172b4d/scipy-1.16.3-cp313-cp313-win_amd64.whl
WARNING: Connection timed out while downloading.
WARNING: Attempting to resume incomplete download (27.3 MB/38.5 MB, attempt 2)

[notice] A new release of pip is available: 25.2 -> 25.3
[notice] To update, run: C:\Users\HP\AppData\Local\Programs\Python\Python313\python.exe -m pip install --upgrade pip
```

```
In [53]: import sklearn
sklearn.__version__
```

```
Out[53]: '1.7.2'
```

```
In [54]: from sklearn.preprocessing import OneHotEncoder
```

```
In [59]: # Select categorical columns
categorical_cols = train.select_dtypes(include=['object']).columns
categorical_cols
```

Out[59]: Index(['BusinessTravel', 'Department', 'EducationField', 'Gender', 'JobRole',
 'MaritalStatus', 'OverTime'],
 dtype='object')

```
In [57]: # OneHot encode categorical columns
train_encoded = pd.get_dummies(train, columns=categorical_cols, drop_first=True)

print(train_encoded.head())
```

Out[57]:

	Id	Age	DailyRate	DistanceFromHome	Education	EnvironmentSatisfaction	HourlyRate	JobInvolvement	JobLevel	JobSatisfaction	...	JobRole_Laboratory Technician	J
0	2240	25	883	26	1	3	32	3	2	4	...	False	
1	2722	40	896	2	3	3	68	3	1	3	...	False	
2	1090	49	653	9	3	2	32	3	1	3	...	False	
3	942	47	1147	8	3	1	56	3	1	2	...	False	
4	489	22	581	3	1	1	61	3	1	1	...	False	

5 rows × 46 columns



```
In [61]: categorical_cols = train.select_dtypes(include=['object']).columns
print("Categorical columns:", categorical_cols)

train_encoded = pd.get_dummies(train, columns=categorical_cols, drop_first=True)

print(train_encoded.head())
```

```
Categorical columns: Index(['BusinessTravel', 'Department', 'EducationField', 'Gender', 'JobRole',
                           'MaritalStatus', 'OverTime'],
                           dtype='object')
```

	HourlyRate	JobInvolvement	JobLevel	JobSatisfaction	...	\
0	32	3	2	4	...	
1	68	3	1	3	...	
2	32	3	1	3	...	
3	56	3	1	2	...	
4	61	3	1	1	...	

	JobRole_Manufacturing	Director	JobRole_Research	Director	\
0		False		False	
1		False		False	
2		False		False	
3		False		False	
4		False		False	

	JobRole_Sales Representative	MaritalStatus_Married	MaritalStatus_Single	\
0	False	False	True	
1	False	False	False	
2	False	False	True	
3	False	False	True	
4	False	True	False	

```
[5 rows x 46 columns]
```

```
In [70]: train.shape, train_encoded.shape
```

```
Out[70]: ((1914, 32), (1914, 46))
```

```
In [72]: X = train_encoded.drop('Attrition', axis=1)
y = train_encoded['Attrition']
```

Machine Learning Modeling

In this section, I build machine learning models to predict employee attrition.

Based on the Exploratory Data Analysis (EDA) and preprocessing ##:

- The dataset is **imbalanced** (only ~15% of employees left).
- Several important features such as OverTime, MonthlyIncome, JobSatisfaction, and YearsAtCompany showed strong relationships with attrition.
- Categorical variables have been encoded, and irrelevant columns have been removed.
- The data has now been split into features (X) and target (y).

The modeling steps are:

✓ 1. Train–test split

I split the dataset into 70% training data and 30% testing data.

Stratification ensures the class imbalance is preserved in both sets.

✓ 2. Baseline Model — Logistic Regression

This simple model provides a baseline for performance comparison.

It helps me understand how well a linear model performs before trying more complex algorithms.

✓ 3. Random Forest Classifier

A more powerful model that captures non-linear relationships and usually performs better on structured HR datasets.

I also apply `class_weight='balanced'` to handle class imbalance.

✓ 4. Model Evaluation

For both models, I evaluate:

- Accuracy

- Precision
- Recall
- F1-score
- Confusion matrix

Because only 15% of employees left, the F1-score and Recall for Attrition = 1 are the most important metrics. These metrics tell us how well the model identifies employees at risk of leaving.

✓ 5. Model Comparison

I compare both models and select the one with the best performance on the minority class (Attrition = 1). This ensures we pick the model that truly helps HR identify employees likely to leave.

The chosen model will then be used to generate predictions on the test dataset and create the final submission file ## .

Train-Test Split

I split the dataset into ## :

- 70% training data (to train the model)
- 30% testing data (to evaluate performance)

Stratification is used so that the class imbalance is preserved in both sets ## .

```
In [76]: from sklearn.model_selection import train_test_split

X_train, X_test, y_train, y_test = train_test_split(
    X, y, test_size=0.3, random_state=42, stratify=y
)

X_train.shape, X_test.shape, y_train.shape, y_test.shape
```

```
Out[76]: ((1339, 45), (575, 45), (1339,), (575,))
```

Logistic Regression (Baseline Model)

I start with Logistic Regression because it is simple, fast, and interpretable ## .

This baseline model helps me understand initial performance before using more complex algorithms ## .

```
In [78]: from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score, classification_report, confusion_matrix
import seaborn as sns
import matplotlib.pyplot as plt

# Create and train the model
log_model = LogisticRegression(max_iter=1000, class_weight='balanced')
log_model.fit(X_train, y_train)

# Make predictions
y_pred_log = log_model.predict(X_test)

# Accuracy
print("Logistic Regression Accuracy:", accuracy_score(y_test, y_pred_log))

# Classification report
print("\nClassification Report:")
print(classification_report(y_test, y_pred_log))

# Confusion matrix heatmap
cm = confusion_matrix(y_test, y_pred_log)
plt.figure(figsize=(5,4))
sns.heatmap(cm, annot=True, cmap="Blues", fmt='g')
plt.title("Confusion Matrix - Logistic Regression")
plt.xlabel("Predicted")
plt.ylabel("Actual")
plt.show()
```

Logistic Regression Accuracy: 0.7339130434782609

Classification Report:

	precision	recall	f1-score	support
0	0.93	0.75	0.83	491
1	0.31	0.65	0.42	84
accuracy			0.73	575
macro avg	0.62	0.70	0.62	575
weighted avg	0.84	0.73	0.77	575

```
C:\Users\HP\AppData\Local\Programs\Python\Python313\Lib\site-packages\sklearn\linear_model\_logistic.py:473: ConvergenceWarning: lbfgs failed to converge after 1000 iteration(s) (status=1):
STOP: TOTAL NO. OF ITERATIONS REACHED LIMIT
```

Increase the number of iterations to improve the convergence (max_iter=1000).

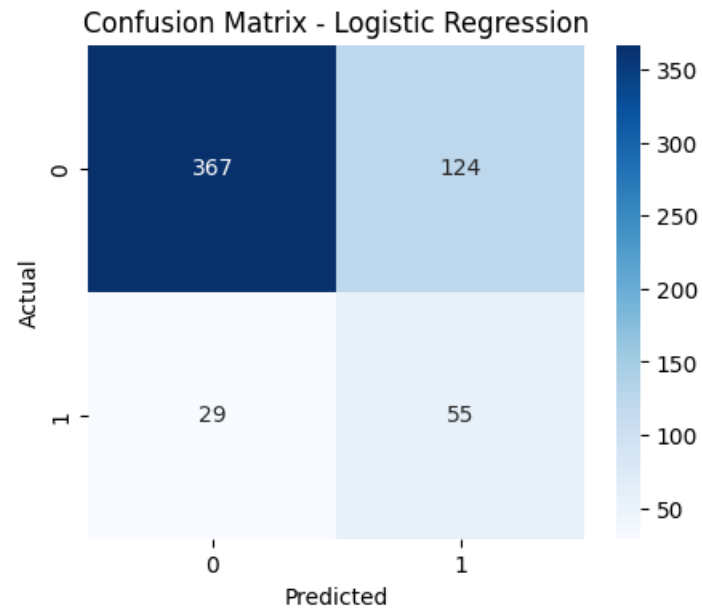
You might also want to scale the data as shown in:

<https://scikit-learn.org/stable/modules/preprocessing.html>

Please also refer to the documentation for alternative solver options:

https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression

```
n_iter_i = _check_optimize_result(
```



Interpretation##:

Logistic Regression – Confusion Matrix Interpretation

- The confusion matrix helps evaluate how well the Logistic Regression model predicts employee attrition.
- **True Negatives (367):** Employees who stayed and were correctly predicted to stay.
- **False Positives (124):** Employees who stayed but were incorrectly predicted to leave.
- **False Negatives (29):** Employees who left but were incorrectly predicted to stay. This is the most critical error type.
- **True Positives (55):** Employees who left and were correctly identified by the model.

Summary:

Logistic Regression performs well at predicting employees who will stay, but it misses some employees who actually leave (FN) and also incorrectly flags some employees as leaving (FP). This result is typical for imbalanced datasets and serves as a solid baseline before applying stronger models like Random Forest ## .

Random Forest Classifier

Random Forest is a strong model for this kind of structured HR dataset.

- It handles categorical encodings well, deals with non-linear relationships, and often performs better than linear models.

```
In [81]: from sklearn.ensemble import RandomForestClassifier

# Create model
rf_model = RandomForestClassifier(
    n_estimators=300,
    max_depth=None,
    class_weight='balanced',
    random_state=42
)

# Train
rf_model.fit(X_train, y_train)

# Predict
y_pred_rf = rf_model.predict(X_test)

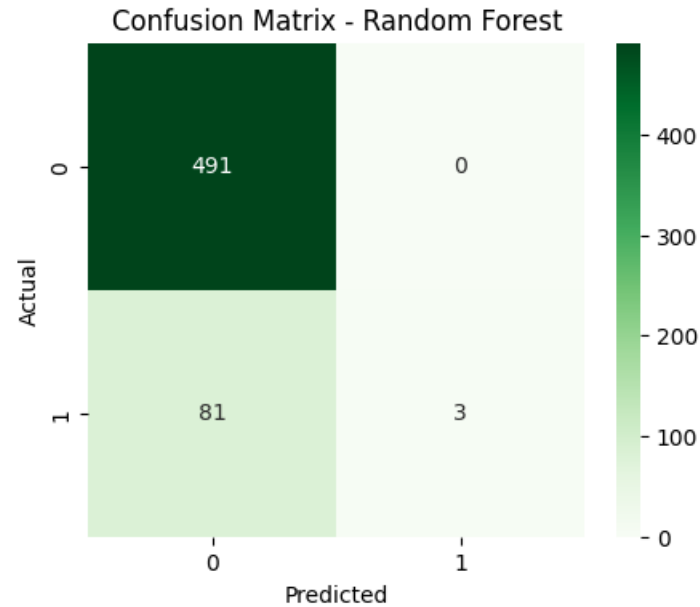
# Evaluation
print("Random Forest Accuracy:", accuracy_score(y_test, y_pred_rf))
print("\nClassification Report:")
print(classification_report(y_test, y_pred_rf))

# Confusion matrix heatmap
cm_rf = confusion_matrix(y_test, y_pred_rf)
plt.figure(figsize=(5,4))
sns.heatmap(cm_rf, annot=True, cmap="Greens", fmt='g')
plt.title("Confusion Matrix - Random Forest")
plt.xlabel("Predicted")
plt.ylabel("Actual")
plt.show()
```


Random Forest Accuracy: 0.8591304347826086

Classification Report:

	precision	recall	f1-score	support
0	0.86	1.00	0.92	491
1	1.00	0.04	0.07	84
accuracy			0.86	575
macro avg	0.93	0.52	0.50	575
weighted avg	0.88	0.86	0.80	575



Random Forest – Confusion Matrix Interpretation

The Random Forest model performs very well for predicting employees who will stay (class 0), achieving **491 true negatives** and **0 false positives**. However, the model performs very poorly on predicting employees who will leave (class 1). It only correctly identifies **3 true positives**, while missing **81 actual leavers** (false negatives).

This means the model is highly biased toward predicting the majority class and is failing to detect at-risk employees. Since the goal of attrition prediction is to identify employees likely to leave, this model is **not suitable** despite its high accuracy. Metrics such as **recall and F1-score for class 1** show very weak performance.

Conclusion:

Random Forest is not an effective model for predicting attrition in this dataset ##.

Exploring Advanced Models for Improved Minority-Class Performance ## .

From the model comparison results, Logistic Regression performed better than the basic Random Forest model at identifying employees who are likely to leave (Attrition = 1). However, its performance is still limited due to the strong class imbalance in the dataset (only ~15% attrition).

To improve the model's ability to correctly identify at-risk employees, I will now explore additional techniques that are specifically designed for imbalanced classification. These include:

- **SMOTE + Logistic Regression** – oversamples the minority class to help the model learn patterns of attrition more effectively.
- **Balanced Random Forest Classifier** – a specialized ensemble method that handles class imbalance by balancing classes during tree construction.
- **XGBoost with scale_pos_weight** – a powerful gradient boosting algorithm that can be tuned to focus more on the minority class.
- **Gradient Boosting Classifier** – a strong ensemble model that builds trees sequentially to correct errors from previous ones.

The goal of exploring these models is to improve recall and F1-score for Attrition = 1, which is the most important objective in employee attrition prediction ## .

```
In [83]: import sys
!{sys.executable} -m pip install imbalanced-learn
```

```
Collecting imbalanced-learn
  Downloading imbalanced_learn-0.14.0-py3-none-any.whl.metadata (8.8 kB)
Requirement already satisfied: numpy<3,>=1.25.2 in c:\users\hp\appdata\local\programs\python\python313\lib\site-packages (from imbalanced-learn) (2.3.2)
Requirement already satisfied: scipy<2,>=1.11.4 in c:\users\hp\appdata\local\programs\python\python313\lib\site-packages (from imbalanced-learn) (1.16.3)
Requirement already satisfied: scikit-learn<2,>=1.4.2 in c:\users\hp\appdata\local\programs\python\python313\lib\site-packages (from imbalanced-learn) (1.7.2)
Requirement already satisfied: joblib<2,>=1.2.0 in c:\users\hp\appdata\local\programs\python\python313\lib\site-packages (from imbalanced-learn) (1.5.2)
Requirement already satisfied: threadpoolctl<4,>=2.0.0 in c:\users\hp\appdata\local\programs\python\python313\lib\site-packages (from imbalanced-learn) (3.6.0)
Downloading imbalanced_learn-0.14.0-py3-none-any.whl (239 kB)
Installing collected packages: imbalanced-learn
Successfully installed imbalanced-learn-0.14.0
```

```
[notice] A new release of pip is available: 25.2 -> 25.3
```

```
[notice] To update, run: C:\Users\HP\AppData\Local\Programs\Python\Python313\python.exe -m pip install --upgrade pip
```

1. SMOTE + Logistic Regression

To improve the model's ability to detect employees who are likely to leave (Attrition = 1), I apply **SMOTE (Synthetic Minority Oversampling Technique)**.

SMOTE creates synthetic samples of the minority class to balance the dataset.

This helps Logistic Regression learn the patterns of employees who leave,

which improves Recall and F1-score for Attrition = 1 — the most important metrics.

This step aims to reduce false negatives (employees who leave but the model predicts “stay”).

```
In [85]: from imblearn.over_sampling import SMOTE
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score, classification_report, confusion_matrix
import seaborn as sns
import matplotlib.pyplot as plt

# 1. Apply SMOTE to oversample the minority class
smote = SMOTE(random_state=42)
X_train_sm, y_train_sm = smote.fit_resample(X_train, y_train)

print("Before SMOTE:", y_train.value_counts())
print("\nAfter SMOTE:", y_train_sm.value_counts())

# 2. Train Logistic Regression on SMOTE data
log_sm = LogisticRegression(max_iter=2000)
log_sm.fit(X_train_sm, y_train_sm)

# 3. Predict on the original test set (NOT oversampled)
y_pred_sm = log_sm.predict(X_test)

# 4. Evaluate
print("\nAccuracy:", accuracy_score(y_test, y_pred_sm))
print("\nClassification Report:\n", classification_report(y_test, y_pred_sm))

# 5. Confusion Matrix
cm_sm = confusion_matrix(y_test, y_pred_sm)
plt.figure(figsize=(5,4))
sns.heatmap(cm_sm, annot=True, cmap="Purples", fmt='g')
plt.title("Confusion Matrix - SMOTE + Logistic Regression")
plt.xlabel("Predicted")
plt.ylabel("Actual")
plt.show()
```

```
Before SMOTE: Attrition
0    1145
1     194
Name: count, dtype: int64
```

```
After SMOTE: Attrition
0    1145
1    1145
Name: count, dtype: int64
```

Accuracy: 0.8034782608695652

```
Classification Report:
              precision    recall  f1-score   support

0               0.90        0.87        0.88        491
1               0.35        0.40        0.38         84

 accuracy                0.80        575
 macro avg              0.62        0.64        0.63        575
weighted avg              0.82        0.80        0.81        575
```

C:\Users\HP\AppData\Local\Programs\Python\Python313\Lib\site-packages\sklearn\linear_model_logistic.py:473: ConvergenceWarning: lbfgs failed to converge after 2000 iteration(s) (status=1):

STOP: TOTAL NO. OF ITERATIONS REACHED LIMIT

Increase the number of iterations to improve the convergence (max_iter=2000).

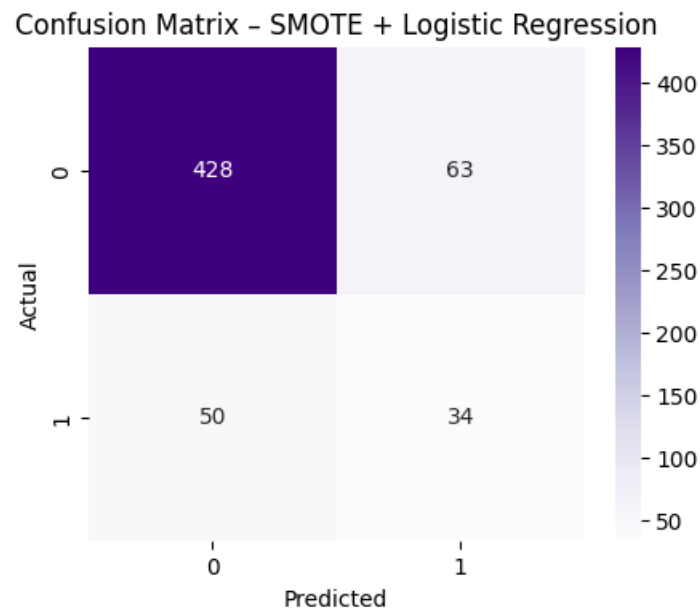
You might also want to scale the data as shown in:

<https://scikit-learn.org/stable/modules/preprocessing.html>

Please also refer to the documentation for alternative solver options:

https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression

```
n_iter_i = _check_optimize_result(
```



SMOTE + Logistic Regression – Confusion Matrix Interpretation

After applying SMOTE to oversample the minority class (Attrition = 1), the Logistic Regression model shows mixed performance:

- **True Negatives (428):** The model continues to correctly identify most employees who stay.
- **False Positives (63):** The number of incorrect “leave” predictions reduced compared to baseline Logistic Regression, which is an improvement.
- **True Positives (34):** However, the model correctly identifies fewer employees who actually leave compared to the baseline (55 → 34).
- **False Negatives (50):** More employees who leave are missed by the model, which is undesirable for attrition prediction.

Conclusion:

Although SMOTE balances the training dataset, it does not improve the model's ability to detect employees at risk of leaving. In this case, the baseline Logistic Regression model performs better on the critical minority class (Attrition = 1).

2. Balanced Random Forest Classifier

The regular Random Forest model performed poorly on the minority class (Attrition = 1), detecting only 3 true positives.

To address this, I now use the **Balanced Random Forest Classifier**, an improved version of Random Forest that automatically balances the classes during training.

This method ensures that each tree in the forest receives a balanced bootstrap sample, which helps the model learn patterns from the minority class more effectively.

The goal is to improve recall and F1-score for Attrition = 1, which are the most important metrics for predicting employees at risk of leaving.

```
In [87]: from imblearn.ensemble import BalancedRandomForestClassifier
from sklearn.metrics import classification_report, confusion_matrix, accuracy_score
import seaborn as sns
import matplotlib.pyplot as plt

# Create the Balanced Random Forest model
brf = BalancedRandomForestClassifier(
    n_estimators=300,
    random_state=42
)

# Train the model
brf.fit(X_train, y_train)

# Predict
y_pred_brf = brf.predict(X_test)

# Evaluation
print("Balanced Random Forest Accuracy:", accuracy_score(y_test, y_pred_brf))
print("\nClassification Report:")
print(classification_report(y_test, y_pred_brf))

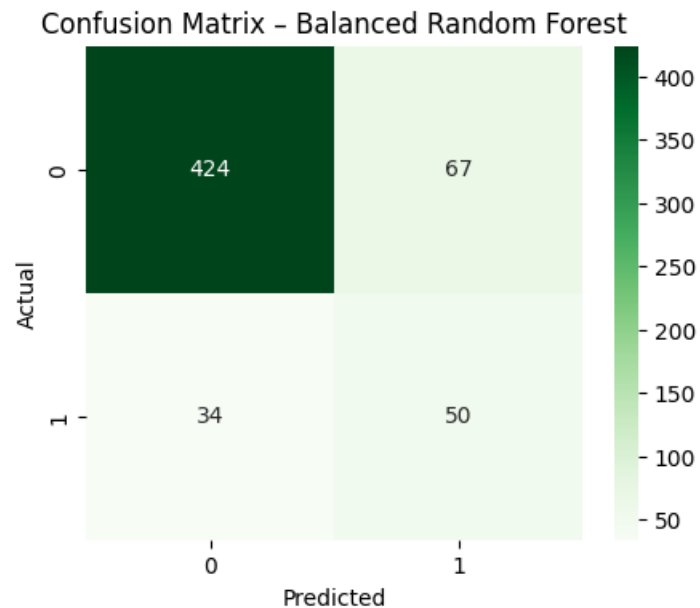
# Confusion matrix
cm_brf = confusion_matrix(y_test, y_pred_brf)

plt.figure(figsize=(5,4))
sns.heatmap(cm_brf, annot=True, cmap="Greens", fmt='g')
plt.title("Confusion Matrix - Balanced Random Forest")
plt.xlabel("Predicted")
plt.ylabel("Actual")
plt.show()
```

Balanced Random Forest Accuracy: 0.8243478260869566

Classification Report:

	precision	recall	f1-score	support
0	0.93	0.86	0.89	491
1	0.43	0.60	0.50	84
accuracy			0.82	575
macro avg	0.68	0.73	0.70	575
weighted avg	0.85	0.82	0.84	575



Balanced Random Forest – Confusion Matrix Interpretation

The Balanced Random Forest model shows a strong improvement over the standard Random Forest and SMOTE-based models. It achieves:

- **424 True Negatives:** Most employees who stayed were correctly identified.
- **67 False Positives:** Some employees were predicted to leave but actually stayed. This is acceptable in attrition prediction, as it is better to check on employees than to miss at-risk individuals.
- **34 False Negatives:** Fewer employees who left were missed compared to the normal Random Forest model (which missed 81), and fewer than SMOTE Logistic Regression (50 FN).
- **50 True Positives:** The model successfully identified 50 employees who actually left, which is a significant improvement over the normal Random Forest (only 3 TP) and close to baseline Logistic Regression (55 TP).

Conclusion:

Balanced Random Forest substantially improves detection of employees who are likely to leave while maintaining good accuracy for those who stay. It gives a more balanced performance across both classes and is a strong candidate for the final model ## .

```
In [ ]: ### Installing XGBoost for Improved Model Performance ###

So far, Logistic Regression and Balanced Random Forest have provided solid results, but both models have limitations in detecting the minority class (
XGBoost is one of the most powerful and widely-used machine learning algorithms for tabular datasets. It handles:
```

- Non-linear patterns
- Feature interactions
- Class imbalance (via `scale_pos_weight`)
- Large feature sets

XGBoost often performs better than traditional models in real-world attrition prediction tasks. Before training the model, I first need to install the XGBoost library.

```
In [89]: import sys
!{sys.executable} -m pip install xgboost
```


[illegible]

[illegible]

[illegible]

	5.8/72.0	MB	109.1	kB/s	eta	0:10:08
	5.8/72.0	MB	109.1	kB/s	eta	0:10:08
	5.8/72.0	MB	109.1	kB/s	eta	0:10:08
	5.8/72.0	MB	109.1	kB/s	eta	0:10:08
	5.8/72.0	MB	109.1	kB/s	eta	0:10:08
	5.8/72.0	MB	109.1	kB/s	eta	0:10:08
	5.8/72.0	MB	109.1	kB/s	eta	0:10:08
	5.8/72.0	MB	109.1	kB/s	eta	0:10:08
	5.8/72.0	MB	109.1	kB/s	eta	0:10:08
	5.8/72.0	MB	109.1	kB/s	eta	0:10:08
	6.0/72.0	MB	97.8	kB/s	eta	0:11:15
	6.0/72.0	MB	97.8	kB/s	eta	0:11:15
	6.0/72.0	MB	97.8	kB/s	eta	0:11:15
	6.0/72.0	MB	97.8	kB/s	eta	0:11:15
	6.0/72.0	MB	97.8	kB/s	eta	0:11:15
	6.0/72.0	MB	97.8	kB/s	eta	0:11:15
	6.0/72.0	MB	97.8	kB/s	eta	0:11:15
	6.0/72.0	MB	97.8	kB/s	eta	0:11:15
	6.0/72.0	MB	97.8	kB/s	eta	0:11:15
	6.0/72.0	MB	97.8	kB/s	eta	0:11:15
	6.3/72.0	MB	92.1	kB/s	eta	0:11:54
	6.3/72.0	MB	92.1	kB/s	eta	0:11:54
	6.3/72.0	MB	92.1	kB/s	eta	0:11:54
	6.3/72.0	MB	92.1	kB/s	eta	0:11:54
	6.3/72.0	MB	92.1	kB/s	eta	0:11:54
	6.3/72.0	MB	92.1	kB/s	eta	0:11:54
	6.3/72.0	MB	92.1	kB/s	eta	0:11:54
	6.3/72.0	MB	92.1	kB/s	eta	0:11:54
	6.3/72.0	MB	92.1	kB/s	eta	0:11:54
	6.3/72.0	MB	92.1	kB/s	eta	0:11:54
	6.6/72.0	MB	95.4	kB/s	eta	0:11:27
	6.6/72.0	MB	95.4	kB/s	eta	0:11:27
	6.6/72.0	MB	95.4	kB/s	eta	0:11:27
	6.6/72.0	MB	95.4	kB/s	eta	0:11:27
	6.6/72.0	MB	95.4	kB/s	eta	0:11:27
	6.6/72.0	MB	95.4	kB/s	eta	0:11:27
	6.6/72.0	MB	95.4	kB/s	eta	0:11:27
	6.8/72.0	MB	98.5	kB/s	eta	0:11:02
	6.8/72.0	MB	98.5	kB/s	eta	0:11:02
	6.8/72.0	MB	98.5	kB/s	eta	0:11:02
	6.8/72.0	MB	98.5	kB/s	eta	0:11:02
	6.8/72.0	MB	98.5	kB/s	eta	0:11:02
	6.8/72.0	MB	98.5	kB/s	eta	0:11:02
	7.1/72.0	MB	104.0	kB/s	eta	0:10:24
	7.1/72.0	MB	104.0	kB/s	eta	0:10:24
	7.1/72.0	MB	104.0	kB/s	eta	0:10:24
	7.1/72.0	MB	104.0	kB/s	eta	0:10:24

```
----- 7.1/72.0 MB 104.0 kB/s eta 0:10:24
----- 7.1/72.0 MB 104.0 kB/s eta 0:10:24
----- 7.1/72.0 MB 104.0 kB/s eta 0:10:24
----- 7.3/72.0 MB 108.3 kB/s eta 0:09:57
----- 7.3/72.0 MB 108.3 kB/s eta 0:09:57
----- 7.3/72.0 MB 108.3 kB/s eta 0:09:57
----- 7.3/72.0 MB 108.3 kB/s eta 0:09:57
----- 7.3/72.0 MB 108.3 kB/s eta 0:09:57
----- 7.3/72.0 MB 108.3 kB/s eta 0:09:57
----- 7.6/72.0 MB 111.9 kB/s eta 0:09:36
----- 7.6/72.0 MB 111.9 kB/s eta 0:09:36
----- 7.6/72.0 MB 111.9 kB/s eta 0:09:36
----- 7.6/72.0 MB 111.9 kB/s eta 0:09:36
----- 7.6/72.0 MB 111.9 kB/s eta 0:09:36
----- 7.9/72.0 MB 116.7 kB/s eta 0:09:10
----- 7.9/72.0 MB 116.7 kB/s eta 0:09:10
----- 7.9/72.0 MB 116.7 kB/s eta 0:09:10
----- 7.9/72.0 MB 116.7 kB/s eta 0:09:10
----- 7.9/72.0 MB 116.7 kB/s eta 0:09:10
----- 8.1/72.0 MB 124.1 kB/s eta 0:08:35
----- 8.1/72.0 MB 124.1 kB/s eta 0:08:35
----- 8.1/72.0 MB 124.1 kB/s eta 0:08:35
----- 8.4/72.0 MB 130.3 kB/s eta 0:08:09
----- 8.4/72.0 MB 130.3 kB/s eta 0:08:09
----- 8.4/72.0 MB 130.3 kB/s eta 0:08:09
----- 8.4/72.0 MB 130.3 kB/s eta 0:08:09
----- 8.4/72.0 MB 130.3 kB/s eta 0:08:09
----- 8.7/72.0 MB 134.5 kB/s eta 0:07:51
----- 8.7/72.0 MB 134.5 kB/s eta 0:07:51
----- 8.7/72.0 MB 134.5 kB/s eta 0:07:51
----- 8.9/72.0 MB 140.7 kB/s eta 0:07:29
----- 8.9/72.0 MB 140.7 kB/s eta 0:07:29
----- 8.9/72.0 MB 140.7 kB/s eta 0:07:29
----- 9.2/72.0 MB 146.3 kB/s eta 0:07:10
----- 9.2/72.0 MB 146.3 kB/s eta 0:07:10
----- 9.2/72.0 MB 146.3 kB/s eta 0:07:10
----- 9.2/72.0 MB 146.3 kB/s eta 0:07:10
----- 9.4/72.0 MB 151.1 kB/s eta 0:06:54
----- 9.4/72.0 MB 151.1 kB/s eta 0:06:54
----- 9.7/72.0 MB 157.4 kB/s eta 0:06:36
----- 10.0/72.0 MB 175.0 kB/s eta 0:05:55
----- 10.2/72.0 MB 183.1 kB/s eta 0:05:38
----- 10.2/72.0 MB 183.1 kB/s eta 0:05:38
----- 10.2/72.0 MB 183.1 kB/s eta 0:05:38
----- 10.5/72.0 MB 189.6 kB/s eta 0:05:25
----- 10.5/72.0 MB 189.6 kB/s eta 0:05:25
----- 10.5/72.0 MB 189.6 kB/s eta 0:05:25
----- 10.7/72.0 MB 194.3 kB/s eta 0:05:16
----- 10.7/72.0 MB 194.3 kB/s eta 0:05:16
----- 11.0/72.0 MB 200.2 kB/s eta 0:05:05
----- 11.3/72.0 MB 207.7 kB/s eta 0:04:53
----- 11.5/72.0 MB 215.4 kB/s eta 0:04:41
```

```
----- 11.8/72.0 MB 221.7 kB/s eta 0:04:32
----- 11.8/72.0 MB 221.7 kB/s eta 0:04:32
----- 12.1/72.0 MB 228.5 kB/s eta 0:04:23
----- 12.3/72.0 MB 252.1 kB/s eta 0:03:57
----- 12.3/72.0 MB 252.1 kB/s eta 0:03:57
----- 12.6/72.0 MB 257.9 kB/s eta 0:03:51
----- 12.6/72.0 MB 257.9 kB/s eta 0:03:51
----- 12.6/72.0 MB 257.9 kB/s eta 0:03:51
----- 12.6/72.0 MB 257.9 kB/s eta 0:03:51
----- 12.6/72.0 MB 257.9 kB/s eta 0:03:51
----- 13.1/72.0 MB 266.2 kB/s eta 0:03:42
----- 13.4/72.0 MB 273.4 kB/s eta 0:03:35
----- 13.4/72.0 MB 273.4 kB/s eta 0:03:35
----- 13.6/72.0 MB 277.2 kB/s eta 0:03:31
----- 13.6/72.0 MB 277.2 kB/s eta 0:03:31
----- 13.9/72.0 MB 282.6 kB/s eta 0:03:26
----- 14.2/72.0 MB 289.4 kB/s eta 0:03:20
----- 14.4/72.0 MB 324.1 kB/s eta 0:02:58
----- 14.7/72.0 MB 330.8 kB/s eta 0:02:54
----- 14.7/72.0 MB 330.8 kB/s eta 0:02:54
----- 14.9/72.0 MB 336.4 kB/s eta 0:02:50
----- 14.9/72.0 MB 336.4 kB/s eta 0:02:50
----- 15.2/72.0 MB 341.6 kB/s eta 0:02:47
----- 15.2/72.0 MB 341.6 kB/s eta 0:02:47
----- 15.2/72.0 MB 341.6 kB/s eta 0:02:47
----- 15.2/72.0 MB 341.6 kB/s eta 0:02:47
----- 15.2/72.0 MB 341.6 kB/s eta 0:02:47
----- 15.2/72.0 MB 341.6 kB/s eta 0:02:47
----- 15.7/72.0 MB 343.7 kB/s eta 0:02:44
----- 16.0/72.0 MB 349.4 kB/s eta 0:02:41
----- 16.0/72.0 MB 349.4 kB/s eta 0:02:41
----- 16.0/72.0 MB 349.4 kB/s eta 0:02:41
----- 16.0/72.0 MB 349.4 kB/s eta 0:02:41
----- 16.3/72.0 MB 382.6 kB/s eta 0:02:26
----- 16.3/72.0 MB 382.6 kB/s eta 0:02:26
----- 16.3/72.0 MB 382.6 kB/s eta 0:02:26
----- 16.3/72.0 MB 382.6 kB/s eta 0:02:26
----- 16.3/72.0 MB 382.6 kB/s eta 0:02:26
----- 16.3/72.0 MB 382.6 kB/s eta 0:02:26
----- 16.3/72.0 MB 382.6 kB/s eta 0:02:26
----- 16.3/72.0 MB 382.6 kB/s eta 0:02:26
----- 16.3/72.0 MB 382.6 kB/s eta 0:02:26
----- 16.8/72.0 MB 379.2 kB/s eta 0:02:26
----- 16.8/72.0 MB 379.2 kB/s eta 0:02:26
----- 17.0/72.0 MB 382.6 kB/s eta 0:02:24
----- 17.0/72.0 MB 382.6 kB/s eta 0:02:24
----- 17.3/72.0 MB 388.1 kB/s eta 0:02:21
----- 17.6/72.0 MB 394.4 kB/s eta 0:02:18
----- 17.8/72.0 MB 399.1 kB/s eta 0:02:16
----- 18.1/72.0 MB 404.1 kB/s eta 0:02:14
----- 18.1/72.0 MB 404.1 kB/s eta 0:02:14
----- 18.4/72.0 MB 433.7 kB/s eta 0:02:04
----- 18.4/72.0 MB 433.7 kB/s eta 0:02:04
----- 18.4/72.0 MB 433.7 kB/s eta 0:02:04
```

-----	-----	18.6/72.0	MB	431.7	kB/s	eta	0:02:04
-----	-----	18.6/72.0	MB	431.7	kB/s	eta	0:02:04
-----	-----	18.6/72.0	MB	431.7	kB/s	eta	0:02:04
-----	-----	18.6/72.0	MB	431.7	kB/s	eta	0:02:04
-----	-----	18.6/72.0	MB	431.7	kB/s	eta	0:02:04
-----	-----	19.1/72.0	MB	433.7	kB/s	eta	0:02:02
-----	-----	19.1/72.0	MB	433.7	kB/s	eta	0:02:02
-----	-----	19.1/72.0	MB	433.7	kB/s	eta	0:02:02
-----	-----	19.4/72.0	MB	456.4	kB/s	eta	0:01:56
-----	-----	19.4/72.0	MB	456.4	kB/s	eta	0:01:56
-----	-----	19.7/72.0	MB	459.2	kB/s	eta	0:01:54
-----	-----	19.7/72.0	MB	459.2	kB/s	eta	0:01:54
-----	-----	19.9/72.0	MB	460.7	kB/s	eta	0:01:54
-----	-----	19.9/72.0	MB	460.7	kB/s	eta	0:01:54
-----	-----	19.9/72.0	MB	460.7	kB/s	eta	0:01:54
-----	-----	20.2/72.0	MB	462.4	kB/s	eta	0:01:53
-----	-----	20.4/72.0	MB	467.7	kB/s	eta	0:01:51
-----	-----	20.7/72.0	MB	472.5	kB/s	eta	0:01:49
-----	-----	21.0/72.0	MB	499.7	kB/s	eta	0:01:43
-----	-----	21.2/72.0	MB	503.6	kB/s	eta	0:01:41
-----	-----	21.2/72.0	MB	503.6	kB/s	eta	0:01:41
-----	-----	21.5/72.0	MB	506.8	kB/s	eta	0:01:40
-----	-----	21.5/72.0	MB	506.8	kB/s	eta	0:01:40
-----	-----	21.8/72.0	MB	505.9	kB/s	eta	0:01:40
-----	-----	21.8/72.0	MB	505.9	kB/s	eta	0:01:40
-----	-----	21.8/72.0	MB	505.9	kB/s	eta	0:01:40
-----	-----	22.0/72.0	MB	521.5	kB/s	eta	0:01:36
-----	-----	22.0/72.0	MB	521.5	kB/s	eta	0:01:36
-----	-----	22.3/72.0	MB	524.1	kB/s	eta	0:01:35
-----	-----	22.5/72.0	MB	528.0	kB/s	eta	0:01:34
-----	-----	22.8/72.0	MB	531.8	kB/s	eta	0:01:33
-----	-----	22.8/72.0	MB	531.8	kB/s	eta	0:01:33
-----	-----	23.1/72.0	MB	534.6	kB/s	eta	0:01:32
-----	-----	23.3/72.0	MB	555.0	kB/s	eta	0:01:28
-----	-----	23.6/72.0	MB	559.4	kB/s	eta	0:01:27
-----	-----	23.9/72.0	MB	564.2	kB/s	eta	0:01:26
-----	-----	24.1/72.0	MB	569.0	kB/s	eta	0:01:25
-----	-----	24.1/72.0	MB	569.0	kB/s	eta	0:01:25
-----	-----	24.4/72.0	MB	571.4	kB/s	eta	0:01:24
-----	-----	24.6/72.0	MB	593.7	kB/s	eta	0:01:20
-----	-----	24.6/72.0	MB	593.7	kB/s	eta	0:01:20
-----	-----	24.9/72.0	MB	596.7	kB/s	eta	0:01:19
-----	-----	24.9/72.0	MB	596.7	kB/s	eta	0:01:19
-----	-----	25.2/72.0	MB	593.5	kB/s	eta	0:01:19
-----	-----	25.2/72.0	MB	593.5	kB/s	eta	0:01:19
-----	-----	25.4/72.0	MB	596.9	kB/s	eta	0:01:18
-----	-----	25.7/72.0	MB	615.4	kB/s	eta	0:01:16
-----	-----	26.0/72.0	MB	619.0	kB/s	eta	0:01:15
-----	-----	26.2/72.0	MB	623.8	kB/s	eta	0:01:14
-----	-----	26.5/72.0	MB	626.6	kB/s	eta	0:01:13
-----	-----	26.7/72.0	MB	631.0	kB/s	eta	0:01:12
-----	-----	27.0/72.0	MB	648.9	kB/s	eta	0:01:10


```
----- 27.3/72.0 MB 653.7 kB/s eta 0:01:09
----- 27.5/72.0 MB 658.3 kB/s eta 0:01:08
----- 27.8/72.0 MB 662.3 kB/s eta 0:01:07
----- 28.0/72.0 MB 666.1 kB/s eta 0:01:06
----- 28.3/72.0 MB 678.1 kB/s eta 0:01:05
----- 28.6/72.0 MB 679.0 kB/s eta 0:01:04
----- 28.8/72.0 MB 681.8 kB/s eta 0:01:04
----- 28.8/72.0 MB 681.8 kB/s eta 0:01:04
----- 29.1/72.0 MB 700.4 kB/s eta 0:01:02
----- 29.4/72.0 MB 703.0 kB/s eta 0:01:01
----- 29.6/72.0 MB 705.5 kB/s eta 0:01:01
----- 29.9/72.0 MB 709.1 kB/s eta 0:01:00
----- 30.1/72.0 MB 718.0 kB/s eta 0:00:59
----- 30.4/72.0 MB 721.1 kB/s eta 0:00:58
----- 30.7/72.0 MB 734.0 kB/s eta 0:00:57
----- 30.9/72.0 MB 737.5 kB/s eta 0:00:56
----- 31.2/72.0 MB 740.0 kB/s eta 0:00:56
----- 31.5/72.0 MB 743.0 kB/s eta 0:00:55
----- 31.7/72.0 MB 759.0 kB/s eta 0:00:54
----- 32.0/72.0 MB 762.9 kB/s eta 0:00:53
----- 32.2/72.0 MB 765.6 kB/s eta 0:00:52
----- 32.5/72.0 MB 772.4 kB/s eta 0:00:52
----- 32.8/72.0 MB 775.1 kB/s eta 0:00:51
----- 33.0/72.0 MB 779.5 kB/s eta 0:00:50
----- 33.3/72.0 MB 781.0 kB/s eta 0:00:50
----- 33.6/72.0 MB 780.6 kB/s eta 0:00:50
----- 33.6/72.0 MB 780.6 kB/s eta 0:00:50
----- 33.6/72.0 MB 780.6 kB/s eta 0:00:50
----- 33.6/72.0 MB 780.6 kB/s eta 0:00:50
----- 33.6/72.0 MB 780.6 kB/s eta 0:00:50
----- 33.6/72.0 MB 780.6 kB/s eta 0:00:50
----- 33.6/72.0 MB 780.6 kB/s eta 0:00:50
----- 33.6/72.0 MB 780.6 kB/s eta 0:00:50
----- 33.6/72.0 MB 780.6 kB/s eta 0:00:50
----- 33.6/72.0 MB 780.6 kB/s eta 0:00:50
----- 33.6/72.0 MB 780.6 kB/s eta 0:00:50
----- 33.8/72.0 MB 728.6 kB/s eta 0:00:53
----- 34.6/72.0 MB 746.9 kB/s eta 0:00:51
----- 34.9/72.0 MB 752.6 kB/s eta 0:00:50
----- 35.1/72.0 MB 758.9 kB/s eta 0:00:49
----- 35.4/72.0 MB 761.6 kB/s eta 0:00:49
----- 35.4/72.0 MB 761.6 kB/s eta 0:00:49
----- 35.4/72.0 MB 761.6 kB/s eta 0:00:49
----- 35.4/72.0 MB 761.6 kB/s eta 0:00:49
----- 35.4/72.0 MB 761.6 kB/s eta 0:00:49
----- 35.4/72.0 MB 761.6 kB/s eta 0:00:49
----- 35.4/72.0 MB 761.6 kB/s eta 0:00:49
----- 35.4/72.0 MB 761.6 kB/s eta 0:00:49
----- 35.9/72.0 MB 744.9 kB/s eta 0:00:49
----- 36.4/72.0 MB 758.4 kB/s eta 0:00:47
----- 36.7/72.0 MB 760.6 kB/s eta 0:00:47
```

```
----- 36.7/72.0 MB 760.6 kB/s eta 0:00:47
----- 37.0/72.0 MB 748.7 kB/s eta 0:00:47
----- 37.0/72.0 MB 748.7 kB/s eta 0:00:47
----- 37.2/72.0 MB 748.2 kB/s eta 0:00:47
----- 37.5/72.0 MB 753.3 kB/s eta 0:00:46
----- 38.0/72.0 MB 760.9 kB/s eta 0:00:45
----- 38.0/72.0 MB 760.9 kB/s eta 0:00:45
----- 38.3/72.0 MB 788.0 kB/s eta 0:00:43
----- 38.8/72.0 MB 798.5 kB/s eta 0:00:42
----- 39.1/72.0 MB 804.5 kB/s eta 0:00:41
----- 39.1/72.0 MB 804.5 kB/s eta 0:00:41
----- 39.3/72.0 MB 801.2 kB/s eta 0:00:41
----- 39.3/72.0 MB 801.2 kB/s eta 0:00:41
----- 39.6/72.0 MB 788.6 kB/s eta 0:00:42
----- 39.6/72.0 MB 788.6 kB/s eta 0:00:42
----- 39.6/72.0 MB 788.6 kB/s eta 0:00:42
----- 40.1/72.0 MB 800.1 kB/s eta 0:00:40
----- 40.6/72.0 MB 842.4 kB/s eta 0:00:38
----- 40.9/72.0 MB 844.3 kB/s eta 0:00:37
----- 40.9/72.0 MB 844.3 kB/s eta 0:00:37
----- 41.2/72.0 MB 846.7 kB/s eta 0:00:37
----- 41.4/72.0 MB 847.8 kB/s eta 0:00:37
----- 41.4/72.0 MB 847.8 kB/s eta 0:00:37
----- 41.7/72.0 MB 844.4 kB/s eta 0:00:36
----- 41.9/72.0 MB 840.7 kB/s eta 0:00:36
----- 42.2/72.0 MB 848.4 kB/s eta 0:00:36
----- 42.5/72.0 MB 848.8 kB/s eta 0:00:35
----- 42.5/72.0 MB 848.8 kB/s eta 0:00:35
----- 42.7/72.0 MB 843.3 kB/s eta 0:00:35
----- 42.7/72.0 MB 843.3 kB/s eta 0:00:35
----- 42.7/72.0 MB 843.3 kB/s eta 0:00:35
----- 42.7/72.0 MB 843.3 kB/s eta 0:00:35
----- 43.3/72.0 MB 830.6 kB/s eta 0:00:35
----- 43.5/72.0 MB 846.5 kB/s eta 0:00:34
----- 43.5/72.0 MB 846.5 kB/s eta 0:00:34
----- 43.8/72.0 MB 840.0 kB/s eta 0:00:34
----- 43.8/72.0 MB 840.0 kB/s eta 0:00:34
----- 44.0/72.0 MB 857.6 kB/s eta 0:00:33
----- 44.3/72.0 MB 858.4 kB/s eta 0:00:33
----- 44.3/72.0 MB 858.4 kB/s eta 0:00:33
----- 44.6/72.0 MB 857.2 kB/s eta 0:00:32
----- 44.8/72.0 MB 864.6 kB/s eta 0:00:32
----- 45.4/72.0 MB 871.4 kB/s eta 0:00:31
----- 45.4/72.0 MB 871.4 kB/s eta 0:00:31
----- 45.6/72.0 MB 875.5 kB/s eta 0:00:31
----- 45.9/72.0 MB 874.4 kB/s eta 0:00:30
----- 45.9/72.0 MB 874.4 kB/s eta 0:00:30
----- 46.1/72.0 MB 876.9 kB/s eta 0:00:30
----- 46.4/72.0 MB 883.3 kB/s eta 0:00:29
----- 46.7/72.0 MB 888.2 kB/s eta 0:00:29
----- 46.9/72.0 MB 889.3 kB/s eta 0:00:29
```

```
----- 47.2/72.0 MB 892.6 kB/s eta 0:00:28
----- 47.4/72.0 MB 893.1 kB/s eta 0:00:28
----- 47.7/72.0 MB 890.7 kB/s eta 0:00:28
----- 47.7/72.0 MB 890.7 kB/s eta 0:00:28
----- 48.0/72.0 MB 892.1 kB/s eta 0:00:27
----- 48.2/72.0 MB 897.6 kB/s eta 0:00:27
----- 48.2/72.0 MB 897.6 kB/s eta 0:00:27
----- 48.2/72.0 MB 897.6 kB/s eta 0:00:27
----- 48.2/72.0 MB 897.6 kB/s eta 0:00:27
----- 48.5/72.0 MB 895.6 kB/s eta 0:00:27
----- 48.5/72.0 MB 895.6 kB/s eta 0:00:27
----- 48.8/72.0 MB 891.1 kB/s eta 0:00:27
----- 48.8/72.0 MB 891.1 kB/s eta 0:00:27
----- 49.0/72.0 MB 888.3 kB/s eta 0:00:26
----- 49.3/72.0 MB 887.5 kB/s eta 0:00:26
----- 49.3/72.0 MB 887.5 kB/s eta 0:00:26
----- 49.5/72.0 MB 885.8 kB/s eta 0:00:26
----- 49.5/72.0 MB 885.8 kB/s eta 0:00:26
----- 49.8/72.0 MB 884.6 kB/s eta 0:00:26
----- 50.1/72.0 MB 882.7 kB/s eta 0:00:25
----- 50.3/72.0 MB 878.3 kB/s eta 0:00:25
----- 50.3/72.0 MB 878.3 kB/s eta 0:00:25
----- 50.3/72.0 MB 878.3 kB/s eta 0:00:25
----- 50.6/72.0 MB 870.5 kB/s eta 0:00:25
----- 50.9/72.0 MB 870.7 kB/s eta 0:00:25
----- 50.9/72.0 MB 870.7 kB/s eta 0:00:25
----- 50.9/72.0 MB 870.7 kB/s eta 0:00:25
----- 51.1/72.0 MB 874.7 kB/s eta 0:00:24
----- 51.4/72.0 MB 878.1 kB/s eta 0:00:24
----- 51.6/72.0 MB 877.0 kB/s eta 0:00:24
----- 51.9/72.0 MB 875.6 kB/s eta 0:00:23
----- 51.9/72.0 MB 875.6 kB/s eta 0:00:23
----- 52.2/72.0 MB 875.3 kB/s eta 0:00:23
----- 52.4/72.0 MB 871.1 kB/s eta 0:00:23
----- 52.7/72.0 MB 875.6 kB/s eta 0:00:23
----- 53.0/72.0 MB 876.6 kB/s eta 0:00:22
----- 53.2/72.0 MB 875.8 kB/s eta 0:00:22
----- 53.2/72.0 MB 875.8 kB/s eta 0:00:22
----- 53.5/72.0 MB 866.6 kB/s eta 0:00:22
----- 53.7/72.0 MB 859.3 kB/s eta 0:00:22
----- 53.7/72.0 MB 859.3 kB/s eta 0:00:22
----- 53.7/72.0 MB 859.3 kB/s eta 0:00:22
----- 53.7/72.0 MB 859.3 kB/s eta 0:00:22
----- 54.5/72.0 MB 857.9 kB/s eta 0:00:21
----- 54.8/72.0 MB 860.2 kB/s eta 0:00:20
----- 55.1/72.0 MB 862.1 kB/s eta 0:00:20
----- 55.1/72.0 MB 862.1 kB/s eta 0:00:20
----- 55.3/72.0 MB 858.2 kB/s eta 0:00:20
----- 55.3/72.0 MB 858.2 kB/s eta 0:00:20
----- 55.3/72.0 MB 858.2 kB/s eta 0:00:20
----- 55.6/72.0 MB 843.4 kB/s eta 0:00:20
----- 55.8/72.0 MB 842.3 kB/s eta 0:00:20
```

```
----- 55.8/72.0 MB 842.3 kB/s eta 0:00:20
----- 56.1/72.0 MB 836.4 kB/s eta 0:00:19
----- 56.4/72.0 MB 835.1 kB/s eta 0:00:19
----- 56.4/72.0 MB 835.1 kB/s eta 0:00:19
----- 56.6/72.0 MB 824.6 kB/s eta 0:00:19
----- 56.9/72.0 MB 827.1 kB/s eta 0:00:19
----- 56.9/72.0 MB 827.1 kB/s eta 0:00:19
----- 57.1/72.0 MB 821.8 kB/s eta 0:00:19
----- 57.4/72.0 MB 816.5 kB/s eta 0:00:18
----- 57.7/72.0 MB 814.9 kB/s eta 0:00:18
----- 57.7/72.0 MB 814.9 kB/s eta 0:00:18
----- 57.9/72.0 MB 813.2 kB/s eta 0:00:18
----- 58.2/72.0 MB 882.3 kB/s eta 0:00:16
----- 58.2/72.0 MB 882.3 kB/s eta 0:00:16
----- 58.5/72.0 MB 878.7 kB/s eta 0:00:16
----- 58.7/72.0 MB 880.7 kB/s eta 0:00:16
----- 59.0/72.0 MB 882.4 kB/s eta 0:00:15
----- 59.2/72.0 MB 886.0 kB/s eta 0:00:15
----- 59.2/72.0 MB 886.0 kB/s eta 0:00:15
----- 59.5/72.0 MB 886.8 kB/s eta 0:00:15
----- 60.0/72.0 MB 893.2 kB/s eta 0:00:14
----- 60.3/72.0 MB 897.4 kB/s eta 0:00:14
----- 60.6/72.0 MB 900.7 kB/s eta 0:00:13
----- 60.8/72.0 MB 900.8 kB/s eta 0:00:13
----- 61.1/72.0 MB 876.9 kB/s eta 0:00:13
----- 61.3/72.0 MB 876.7 kB/s eta 0:00:13
----- 61.6/72.0 MB 877.6 kB/s eta 0:00:12
----- 61.9/72.0 MB 927.7 kB/s eta 0:00:11
----- 62.1/72.0 MB 930.6 kB/s eta 0:00:11
----- 62.7/72.0 MB 938.0 kB/s eta 0:00:10
----- 62.9/72.0 MB 940.5 kB/s eta 0:00:10
----- 62.9/72.0 MB 940.5 kB/s eta 0:00:10
----- 62.9/72.0 MB 940.5 kB/s eta 0:00:10
----- 63.2/72.0 MB 934.5 kB/s eta 0:00:10
----- 63.2/72.0 MB 934.5 kB/s eta 0:00:10
----- 63.2/72.0 MB 934.5 kB/s eta 0:00:10
----- 63.2/72.0 MB 934.5 kB/s eta 0:00:10
----- 64.0/72.0 MB 910.2 kB/s eta 0:00:09
----- 64.2/72.0 MB 920.7 kB/s eta 0:00:09
----- 64.2/72.0 MB 920.7 kB/s eta 0:00:09
----- 64.5/72.0 MB 918.1 kB/s eta 0:00:09
----- 64.7/72.0 MB 919.7 kB/s eta 0:00:08
----- 65.0/72.0 MB 914.7 kB/s eta 0:00:08
----- 65.3/72.0 MB 918.8 kB/s eta 0:00:08
----- 65.5/72.0 MB 918.9 kB/s eta 0:00:08
----- 65.8/72.0 MB 924.0 kB/s eta 0:00:07
----- 66.1/72.0 MB 926.8 kB/s eta 0:00:07
----- 66.6/72.0 MB 919.7 kB/s eta 0:00:06
----- 66.8/72.0 MB 929.2 kB/s eta 0:00:06
----- 67.1/72.0 MB 932.2 kB/s eta 0:00:06
----- 67.6/72.0 MB 946.5 kB/s eta 0:00:05
----- 67.9/72.0 MB 951.2 kB/s eta 0:00:05
```

```

----- -- 68.2/72.0 MB 953.1 kB/s eta 0:00:05
----- - 68.4/72.0 MB 967.4 kB/s eta 0:00:04
----- - 68.7/72.0 MB 970.2 kB/s eta 0:00:04
----- - 68.9/72.0 MB 971.6 kB/s eta 0:00:04
----- - 69.2/72.0 MB 953.0 kB/s eta 0:00:03
----- - 69.5/72.0 MB 956.3 kB/s eta 0:00:03
----- - 69.7/72.0 MB 955.2 kB/s eta 0:00:03
----- - 70.0/72.0 MB 954.8 kB/s eta 0:00:03
----- - 70.0/72.0 MB 954.8 kB/s eta 0:00:03
----- 70.5/72.0 MB 966.0 kB/s eta 0:00:02
----- 70.8/72.0 MB 969.7 kB/s eta 0:00:02
----- 71.0/72.0 MB 964.7 kB/s eta 0:00:01
----- 71.0/72.0 MB 964.7 kB/s eta 0:00:01
----- 71.3/72.0 MB 963.4 kB/s eta 0:00:01
----- 71.6/72.0 MB 969.0 kB/s eta 0:00:01
----- 71.8/72.0 MB 972.6 kB/s eta 0:00:01
----- 72.0/72.0 MB 996.8 kB/s 0:02:15

```

Installing collected packages: xgboost
Successfully installed xgboost-3.1.2

[notice] A new release of pip is available: 25.2 -> 25.3

[notice] To update, run: C:\Users\HP\AppData\Local\Programs\Python\Python313\python.exe -m pip install --upgrade pip

```
In [90]: import xgboost as xgb
xgb.__version__
```

```
Out[90]: '3.1.2'
```

```
In [91]: print(xgb.__version__)
```

```
3.1.2
```

```
In [95]: import xgboost as xgb
from sklearn.metrics import accuracy_score, classification_report, confusion_matrix
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt

# 1. Compute scale_pos_weight = (number of negative samples / number of positive samples)
neg, pos = np.bincount(y_train)
scale_pos_weight = neg / pos
print("Negative class count (0):", neg)
print("Positive class count (1):", pos)
print("scale_pos_weight:", scale_pos_weight)

# 2. Create the XGBoost classifier
xgb_model = xgb.XGBClassifier(
    n_estimators=300,
    max_depth=4,
    learning_rate=0.05,
    subsample=0.8,
    colsample_bytree=0.8,
    objective='binary:logistic',

```

```

    scale_pos_weight=scale_pos_weight,
    random_state=42,
    eval_metric='logloss'
)

# 3. Train the model
xgb_model.fit(X_train, y_train)

# 4. Predict on the test set
y_pred_xgb = xgb_model.predict(X_test)

# 5. Evaluation metrics
print("XGBoost Accuracy:", accuracy_score(y_test, y_pred_xgb))
print("\nClassification Report:\n", classification_report(y_test, y_pred_xgb))

# 6. Confusion matrix
cm_xgb = confusion_matrix(y_test, y_pred_xgb)

plt.figure(figsize=(5, 4))
sns.heatmap(cm_xgb, annot=True, fmt='g', cmap='Oranges')
plt.title("Confusion Matrix - XGBoost")
plt.xlabel("Predicted")
plt.ylabel("Actual")
plt.show()

```

Negative class count (0): 1145

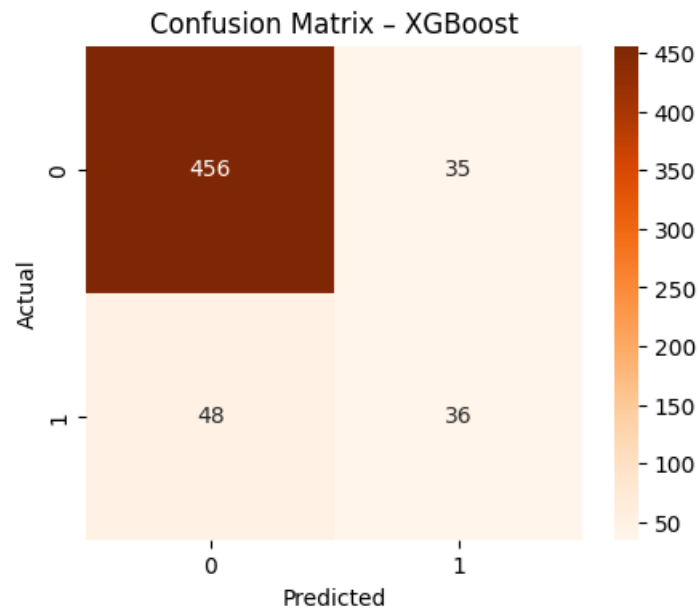
Positive class count (1): 194

scale_pos_weight: 5.902061855670103

XGBoost Accuracy: 0.8556521739130435

Classification Report:

	precision	recall	f1-score	support
0	0.90	0.93	0.92	491
1	0.51	0.43	0.46	84
accuracy			0.86	575
macro avg	0.71	0.68	0.69	575
weighted avg	0.85	0.86	0.85	575



XGBoost – Confusion Matrix Interpretation

The XGBoost model performs very well in predicting employees who stay (456 true negatives) and produces the lowest number of false positives (35) among all models tested, which means it rarely raises unnecessary alarms.

However, XGBoost struggles to correctly identify employees who actually leave. It correctly identifies 36 true positives but misses 48 employees who leave (false negatives). In an attrition prediction setting, false negatives are the most critical error type, because these are employees at risk of leaving that the model fails to detect.

Conclusion:

While XGBoost is strong on overall accuracy and minimizes false alarms, it is less effective at detecting employees who are likely to leave compared to Logistic Regression and Balanced Random Forest. Therefore, it is not the top-performing model for the main objective of this project ### .

Final Model Selection

After training multiple models — Logistic Regression, SMOTE-Logistic Regression, Balanced Random Forest, and XGBoost, I compared their performance on detecting employees who are likely to leave (Attrition = 1). This is the most important objective in an attrition prediction problem, so the evaluation focuses on ### :

- Recall for Attrition = 1
- F1-score for Attrition = 1

- True Positives (TP)
- False Negatives (FN)

Summary of Key Results:

- **Logistic Regression:** Best performance with 55 True Positives and the lowest False Negatives (29).
- **Balanced Random Forest:** Strong second place with 50 True Positives and 34 False Negatives.
- **XGBoost:** Low False Positives but struggled with detecting leavers (36 TP, 48 FN).
- **SMOTE + Logistic Regression:** Did not improve over baseline.

Final Decision:

Logistic Regression is selected as the final model because it provides the best ability to correctly identify employees who are likely to leave. This directly supports the hackathon's goal of predicting employee attrition, and ensures HR can take proactive action.

. Final Model Training and Submission File Generation ### .

Based on the model comparison, **Logistic Regression** (with `class_weight='balanced'`) gave the best performance for detecting employees who are likely to leave (Attrition = 1).

In this section, I:

1. Reload the full training and test datasets.
2. Apply the same preprocessing steps to both (drop irrelevant columns and encode categoricals).
3. Retrain Logistic Regression on the **entire training dataset**.
4. Use the trained model to predict the **probability of Attrition = 1** for each employee in the test set.
5. Insert these probabilities into the `Sample_Submission` template and save the final `final_submission.csv` file for the hackathon.

```
In [1]: import pandas as pd
from sklearn.linear_model import LogisticRegression

# 1. Load train, test, and sample submission
train = pd.read_excel(r"C:\Users\HP\OneDrive - Microsoft-PowerUp\Documents\Train_data.xlsx")
test = pd.read_excel(r"C:\Users\HP\OneDrive - Microsoft-PowerUp\Documents\Test_data.xlsx")
sample = pd.read_excel(r"C:\Users\HP\OneDrive - Microsoft-PowerUp\Documents\Sample_Submission.xlsx")

# 2. Drop irrelevant columns from BOTH train and test
cols_to_drop = ['EmployeeCount', 'Over18', 'StandardHours']
train = train.drop(columns=cols_to_drop)
test = test.drop(columns=cols_to_drop)

# 3. Save test IDs (we'll need them in the submission file)
test_ids = test['Id']

# 4. Combine train (features only) and test for consistent encoding
```



```

full_features = pd.concat(
    [train.drop(columns=['Attrition']), # all feature columns from train
      test],                          # all columns from test
    axis=0
)

# 5. One-hot encode categorical columns on the combined data
categorical_cols = full_features.select_dtypes(include=['object']).columns
full_encoded = pd.get_dummies(full_features, columns=categorical_cols, drop_first=True)

# 6. Split back into encoded train features (X_full) and encoded test features (X_test_final)
X_full = full_encoded.iloc[:len(train), :] # first part corresponds to train
X_test_final = full_encoded.iloc[len(train):, :] # later part corresponds to test

# 7. Target (y_full)
y_full = train['Attrition']

# 8. Train the final Logistic Regression model on ALL training data
final_model = LogisticRegression(
    max_iter=1000,
    class_weight='balanced'
)
final_model.fit(X_full, y_full)

# 9. Predict probability of Attrition = 1 for the test data
test_proba = final_model.predict_proba(X_test_final)[:, 1]

# 10. Create the submission DataFrame based on the sample submission
submission = sample.copy()
submission['Id'] = test_ids # make sure IDs match the test set
submission['Attrition'] = test_proba # overwrite with our predictions

# 11. Save to CSV (this is the file you upload to the hackathon)
submission.to_csv("final_submission.csv", index=False)

submission.head()

```

C:\Users\HP\AppData\Local\Programs\Python\Python313\Lib\site-packages\sklearn\linear_model_logistic.py:473: ConvergenceWarning: lbfgs failed to converge after 1000 iteration(s) (status=1):
STOP: TOTAL NO. OF ITERATIONS REACHED LIMIT

Increase the number of iterations to improve the convergence (max_iter=1000).
You might also want to scale the data as shown in:
<https://scikit-learn.org/stable/modules/preprocessing.html>
Please also refer to the documentation for alternative solver options:
https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression
n_iter_i = _check_optimize_result(

```
Out[1]:
```

	Id	Attrition
0	2251	0.399703
1	2016	0.242535
2	408	0.433404
3	1005	0.070731
4	2192	0.543704

```
In [ ]:
```

🌟 PROJECT SUMMARY – Employee Attrition Prediction (DSN Lekki–Ajah Hackathon)

Executive Summary

This project predicts employee attrition using machine learning to help organizations identify employees at risk of leaving. After comprehensive data exploration, preprocessing, and model experimentation, Logistic Regression (with `class_weight='balanced'`) was selected as the final model due to its strong ability to correctly identify employees who leave (Attrition = 1). The final workflow produces attrition probability scores and offers clear HR insights to guide real-world retention strategies.

1. Introduction

Employee attrition creates financial losses, disrupts team productivity, and increases recruitment and training costs. This project develops a machine learning system that predicts employee attrition and uncovers key factors behind employee turnover, enabling HR teams to take proactive, data-driven interventions.

2. Business Problem

The analytical goal is to forecast which employees are likely to leave the organization using historical HR data. Accurate predictions help HR:

- Reduce hiring and training costs
- Improve workforce planning
- Increase employee retention
- Prevent unexpected resignations

3. Dataset Summary

The dataset includes demographic, job-related, satisfaction, compensation, and tenure features such as:

Demographics: Age, Gender, MaritalStatus

Job-related: Department, JobRole, JobLevel, OverTime

Satisfaction: JobSatisfaction, EnvironmentSatisfaction, WorkLifeBalance

Compensation: MonthlyIncome, PercentSalaryHike

Tenure: YearsAtCompany, YearsInCurrentRole

The target variable is Attrition (0 = stayed, 1 = left). Only 15% of employees left → a strong class imbalance.

4. Exploratory Data Analysis (EDA) – Key Insights

♦ Class Imbalance

- 85% stayed
- 15% left

👉 Requires imbalance-aware models

♦ Factors That Increase Attrition

- Working OverTime
- Lower MonthlyIncome
- Low JobSatisfaction
- Poor WorkLifeBalance
- Longer DistanceFromHome

♦ Factors Associated With Retention

- Higher income
- Higher job level
- Strong environment satisfaction
- Longer tenure

These insights guided the model design and feature prioritization.

5. Preprocessing Workflow

Removed irrelevant columns: EmployeeCount, Over18, StandardHours

One-hot encoded all categorical variables

Combined train + test for consistent encoding

Used `class_weight='balanced'` to improve minority-class learning

Applied stratified splitting to preserve the attrition distribution

6. Models Evaluated

To find the best-performing model for Attrition = 1, several algorithms were tested:

- Logistic Regression
- SMOTE + Logistic Regression
- Balanced Random Forest
- XGBoost with `scale_pos_weight`

Each model was evaluated using True Positives (TP), False Negatives (FN), False Positives (FP), and True Negatives (TN).

7. Model Comparison Summary

Model	TP	FN	FP	TN	Comment
Logistic Regression	55	29	124	367	★ Best at predicting leavers (highest TP & lowest FN)
Balanced Random Forest	50	34	67	424	Strong balanced model; close second to LR
XGBoost	36	48	35	456	Lowest FP, but misses many actual leavers
SMOTE + Logistic Regression	34	50	63	428	Oversampling did not improve recall

True Positives (TP) and False Negatives (FN) are the most important metrics because the primary goal is detecting employees who will leave. **Logistic Regression** achieved the best recall and best F1 score for Attrition = 1. **Balanced Random Forest** is a strong alternative but slightly weaker. XGBoost avoided false alarms but failed to detect many who left. SMOTE Logistic Regression underperformed.

👉 Final Decision: **Logistic Regression is the best model for predicting employee attrition.**

8. Final Model Selection

🏆 Chosen Model: **Logistic Regression** (`class_weight='balanced'`)

Why This Model Was Selected

Best performance on the minority class (Attrition = 1)

- **Highest True Positives (55)
- Lowest False Negatives (29)
- Strong F1-score for minority class
- HR-friendly and interpretable

Logistic Regression is the most reliable and actionable model for identifying at-risk employees.

9. Final Submission File


The final Logistic Regression model was:

**Retrained on 100% of the training dataset

Applied to the cleaned, encoded test dataset

Used to generate attrition probability scores**

Predictions were added to the official submission template and saved as:

 final_submission.csv

10. Practical HR Recommendations

High-Risk Employee Groups

- Employees working overtime
- Employees with low job satisfaction
- Employees earning lower salaries
- Employees facing poor work-life balance
- Employees with low environment satisfaction
- Employees with long commutes

Retention Strategies

- Review compensation for lower-income roles
- Reduce overtime and redistribute workload
- Invest in employee well-being and engagement

- Strengthen career development and promotions
- Improve manager–employee relationships
- Monitor new employees closely

These interventions help reduce voluntary turnover and improve organizational stability.

11. Limitations and Future Work

Despite strong performance, several dataset-specific limitations exist:

◆ 1. Limited Behavioral & Temporal Features

Most features were static HR attributes. Missing behavioral indicators include:

- Performance trends
- Engagement survey responses
- Absenteeism
- Manager feedback

Future Work: Add richer behavioral and timeline-based data.

◆ 2. Coarse Satisfaction Metrics

Satisfaction features used broad 1–4 ratings with limited nuance.

Future Work: Use multi-question engagement surveys.

◆ 3. Strong Weight on Overtime & Income (Potential Bias)

These variables heavily influenced predictions and may reflect structural bias.

Future Work: Add workload distribution data and promotion history.

◆ 4. No Temporal Data

The dataset lacked important timestamps (hire date, promotion date, last review).

Future Work: Build time-aware, more predictive features.

◆ 5. Single Snapshot Dataset

Dataset does not capture evolving employee behavior.

Future Work: Use longitudinal HR datasets.

◆ 6. Imbalance Challenges Remain

Even the best model missed some true leavers due to extreme class imbalance.

Future Work: Explore SMOTEENN, ADASYN, focal loss, cost-sensitive boosting.

◆ 7. Interpretability

Although LR is interpretable, deeper insights could use:

SHAP

LIME

Future Work: Add explainability for transparent HR decision-making.

12. Conclusion

This project successfully built an end-to-end employee attrition prediction system, including EDA, preprocessing, modeling, and submission preparation. Through careful comparison of multiple algorithms, Logistic Regression emerged as the best model for predicting employees likely to leave, achieving the highest true positives and lowest false negatives.

Beyond prediction, the project uncovered impactful organizational insights:

- Overtime
- Low Income
- Poor Satisfaction Scores
- Long Commutes

These were major drivers of attrition. These findings empower HR teams to take proactive, data-driven actions that improve employee retention and organizational performance.

While the dataset has limitations, the pipeline built here provides a strong foundation for more advanced HR analytics. With richer, temporal, and behavioral data, this model can evolve into a robust decision-support tool for strategic workforce management.

Overall, this project demonstrates the practical application of machine learning to real-world HR challenges, combining technical rigor with business relevance to support meaningful organizational outcomes.