

▼ Advanced Certification Programme in AI and MLOps

A programme by IISc and TalentSprint

Mini Project Notebook: Employee Attrition Prediction

▼ Problem Statement

To predict employee attrition using CatBoost and XgBoost

▼ Learning Objectives

At the end of the experiment, you will be able to

- explore the employee attrition dataset
- apply CatBoost and XgBoost on the dataset
- tune the model hyperparameters to improve accuracy
- evaluate the model using suitable metrics

▼ Introduction

Employee attrition is the gradual reduction in employee numbers. Employee attrition happens when the size of your workforce diminishes over time. This means that employees are leaving faster than they are hired. Employee attrition happens when employees retire, resign, or simply aren't replaced. Although employee attrition can be company-wide, it may also be confined to specific parts of a business.

Employee attrition can happen for several reasons. These include unhappiness about employee benefits or the pay structure, a lack of employee development opportunities, and even poor conditions in the workplace.

To know more about the factors that lead to employee attrition, refer [here](#).

Gradient Boosted Decision Trees

- Gradient boosted decision trees (GBDTs) are one of the most important machine learning models.
- GBDTs originate from AdaBoost, an algorithm that ensembles weak learners and uses the majority vote, weighted by their individual accuracy, to solve binary classification problems. The weak learners in this case are decision trees with a single split, called decision stumps.
- Some of the widely used gradient boosted decision trees are XgBoost, CatBoost and LightGBM.

▼ Dataset

The dataset used for this mini-project is [HR Employee Attrition dataset](#). This dataset is synthetically created by IBM data scientists. There are 35 features and 1470 records.

There are numerical features such as:

- Age
- DistanceFromHome
- EmployeeNumber
- PerformanceRating

There are several categorical features such as:

- JobRole
- EducationField
- Department
- BusinessTravel

Dependent or target feature is 'attrition' which has values as Yes/No.

▼ Download the data

```
1 #@title Download the data
2 !wget -qq https://cdn.iisc.talentsprint.com/CDS/Datasets/wa_fn_usec_hr_employee_attrition_tsv.csv
```

<https://colab.research.google.com/drive/14tvm8H3gbGn3O78BX0-O1v2hZkJIt7U8#scrollTo=qFGMmqS-CEVu&uniqifier=3&printMode=true>

```
3 print("Data Downloaded Successfully!!")
```

```
Data Downloaded Successfully!!
```

▼ Grading = 10 Points

▼ Install CatBoost

```
1 !pip -qq install catboost
```

```
98.7/98.7 MB 7.7 MB/s eta 0:00:00
```

▼ Import Required Packages

```
1 import numpy as np
2 import pandas as pd
3 import seaborn as sns
4 from matplotlib import pyplot as plt
5 from sklearn.metrics import roc_auc_score, accuracy_score, confusion_matrix, f1_score
6 from sklearn.model_selection import train_test_split
7 from lightgbm import LGBMClassifier
8 from xgboost import XGBClassifier
9 from catboost import CatBoostClassifier, metrics
10 import warnings
11 warnings.filterwarnings("ignore")
12 plt.style.use('fivethirtyeight')
13 pd.set_option('display.max_columns', 100)
14 %matplotlib inline
```

Please refer to the [ReadMe](#) before proceeding further.

▼ Part-A

▼ Load the Dataset

Exercise 1: Read the dataset [0.5 Mark]

Hint: `pd.read_csv()`

```
1 # read the dataset
2 # YOUR CODE HERE
3 import pandas as pd
4
5 # The file is named 'wa_fn_usec_hr_employee_attrition_tsv.csv' and is in the same directory
6 #file_name = 'wa_fn_usec_hr_employee_attrition_tsv.csv'
7 #dataset = pd.read_csv(file_name, delimiter='\t') # TSV format uses tab as delimiter
8
9 # Display the first few rows to ensure it's loaded correctly
10 #print(dataset.head(10))
11 # Reading the dataset as CSV instead of TSV
12 dataset = pd.read_csv(file_name, delimiter=',') # Using comma as delimiter
13
14 # Check the first few rows of the dataset to ensure it's loaded properly:
15 print("Dataset head:")
16 print(dataset.head())
17 print("\n-----\n")
18
19
```



Dataset head:

	age	attrition	business	travel	dailyrate	department	\
0	41	Yes	Travel_Rarely	1102		Sales	
1	49	No	Travel_Frequently	279	Research & Development		
2	37	Yes	Travel_Rarely	1373	Research & Development		
3	33	No	Travel_Frequently	1392	Research & Development		
4	27	No	Travel_Rarely	591	Research & Development		

	distancefromhome	education	educationfield	employee	count	employee	number	\
0	1	2	Life Sciences	1		1		
1	8	1	Life Sciences	1		2		
2	2	2	Other	1		4		

3	3	4	Life Sciences	1	5
4	2	1	Medical	1	7

	environmentsatisfaction	gender	hourlyrate	jobinvolvement	joblevel	\
0		2 Female	94	3	2	
1		3 Male	61	2	2	
2		4 Male	92	2	1	
3		4 Female	56	3	1	
4		1 Male	40	3	1	

	jobrole	jobsatisfaction	maritalstatus	monthlyincome	\
0	Sales Executive	4	Single	5993	
1	Research Scientist	2	Married	5130	
2	Laboratory Technician	3	Single	2090	
3	Research Scientist	3	Married	2909	
4	Laboratory Technician	2	Married	3468	

	monthlyrate	numcompaniesworked	over18	overtime	percentsalaryhike	\
0	19479	8	Y	Yes	11	
1	24907	1	Y	No	23	
2	2396	6	Y	Yes	15	
3	23159	1	Y	Yes	11	
4	16632	9	Y	No	12	

	performancerating	relationshipsatisfaction	standardhours	\
0	3	1	80	
1	4	4	80	
2	3	2	80	
3	3	3	80	
4	3	4	80	

	stockoptionlevel	totalworkingyears	trainingtimeslastyear	\
0	0	8	0	
1	1	10	3	
2	0	7	3	
3	0	8	3	
4	1	6	3	

	worklifebalance	yearsatcompany	yearsincurrentrole	\
0	1	6	4	
1	3	10	7	
2	3	0	0	
3	3	8	7	
4	3	2	2	

	yearssincelastpromotion	yearswithcurrmanager
0		
1		
2		
3		
4		

```

1 # Check the shape of dataframe.
2 # YOUR CODE HERE
3 #print(dataset.shape)
4 dataset.shape
5

```

```
(1470, 35)
```

There can be more than one file to read as this is introduced as a competition, dataset has one file for training the model. Their can be other files as one containing the test features and the other can be the true labels.

▼ Data Exploration

- Check for missing values
- Check for consistent data type across a feature
- Check for outliers or inconsistencies in data columns
- Check for correlated features
- Do we have a target label imbalance
- How our independent variables are distributed relative to our target label
- Are there features that have strong linear or monotonic relationships? Making correlation heatmaps makes it easy to identify possible collinearity

Exercise 2: Create a list of numerical and categorical columns. Display a statistical description of the dataset. Remove missing values (if any) [0.5 Mark]

Hint: Use `for` to iterate through each column.

```

1 # YOUR CODE HERE
2 import pandas as pd
3
4 # Assuming you've already read the dataset:

```

```

5 # dataset = pd.read_csv(file_name, delimiter='\t')
6
7 # Create empty lists for numerical and categorical columns
8 numerical_columns = []
9 categorical_columns = []
10
11 # Iterate through columns to classify them as numerical or categorical using a for loop
12 for column in dataset.columns:
13     if dataset[column].dtype in ['int64', 'float64']:
14         numerical_columns.append(column)
15     else:
16         categorical_columns.append(column)
17
18 print("Numerical Columns:", numerical_columns)
19 print("Categorical Columns:", categorical_columns)
20
21 # Display a statistical description of the dataset
22 print("\nStatistical Description of the Dataset:\n")
23 print(dataset.describe(include='all'))
24
25 # Remove rows with missing values
26 missing_before = dataset.isnull().sum().sum()
27 dataset.dropna(inplace=True)
28 missing_after = dataset.isnull().sum().sum()
29
30 print(f"\nNumber of missing values before removal: {missing_before}")
31 print(f"Number of missing values after removal: {missing_after}")
32
33
34
35
36
37

```

Numerical Columns: ['age', 'dailyrate', 'distancefromhome', 'education', 'employeecount', 'employeenumber', 'environmentsatisfaction', 'gender', 'joblevel', 'jobrole', 'maritalstatus', 'monthlyincome', 'yearsexperience']
 Categorical Columns: ['attrition', 'businesstravel', 'department', 'educationfield', 'gender', 'jobrole', 'maritalstatus']

Statistical Description of the Dataset:

	age	attrition	businesstravel	dailyrate	\
count	1470.000000	1470	1470	1470.000000	
unique	NaN	2	3	NaN	
top	NaN	No	Travel_Rarely	NaN	
freq	NaN	1233	1043	NaN	
mean	36.923810	NaN	NaN	802.485714	
std	9.135373	NaN	NaN	403.509100	
min	18.000000	NaN	NaN	102.000000	
25%	30.000000	NaN	NaN	465.000000	
50%	36.000000	NaN	NaN	802.000000	
75%	43.000000	NaN	NaN	1157.000000	
max	60.000000	NaN	NaN	1499.000000	

	department	distancefromhome	education	educationfield	\
count	1470	1470.000000	1470.000000	1470	
unique	3	NaN	NaN	6	
top	Research & Development	NaN	NaN	Life Sciences	
freq	961	NaN	NaN	606	
mean	NaN	9.192517	2.912925	NaN	
std	NaN	8.106864	1.024165	NaN	
min	NaN	1.000000	1.000000	NaN	
25%	NaN	2.000000	2.000000	NaN	
50%	NaN	7.000000	3.000000	NaN	
75%	NaN	14.000000	4.000000	NaN	
max	NaN	29.000000	5.000000	NaN	

	employeecount	employeenumber	environmentsatisfaction	gender	\
count	1470.0	1470.000000	1470.000000	1470	
unique	NaN	NaN	NaN	2	
top	NaN	NaN	NaN	Male	
freq	NaN	NaN	NaN	882	
mean	1.0	1024.865306	2.721769	NaN	
std	0.0	602.024335	1.093082	NaN	
min	1.0	1.000000	1.000000	NaN	
25%	1.0	491.250000	2.000000	NaN	
50%	1.0	1020.500000	3.000000	NaN	
75%	1.0	1555.750000	4.000000	NaN	
max	1.0	2068.000000	4.000000	NaN	

	hourlyrate	jobinvolvement	joblevel	jobrole	\
count	1470.000000	1470.000000	1470.000000	1470	
unique	NaN	NaN	NaN	9	
top	NaN	NaN	NaN	Sales Executive	
freq	NaN	NaN	NaN	326	
mean	65.891156	2.729932	2.063946	NaN	
std	20.329428	0.711561	1.106940	NaN	
min	30.000000	1.000000	1.000000	NaN	
25%	48.000000	2.000000	1.000000	NaN	

50%	66.000000	3.000000	2.000000	NaN
75%	83.750000	3.000000	3.000000	NaN
max	100.000000	4.000000	5.000000	NaN

```
jobsatisfaction maritalstatus monthlvincome monthlvrate \
```

First, we want to get a sense of our data:

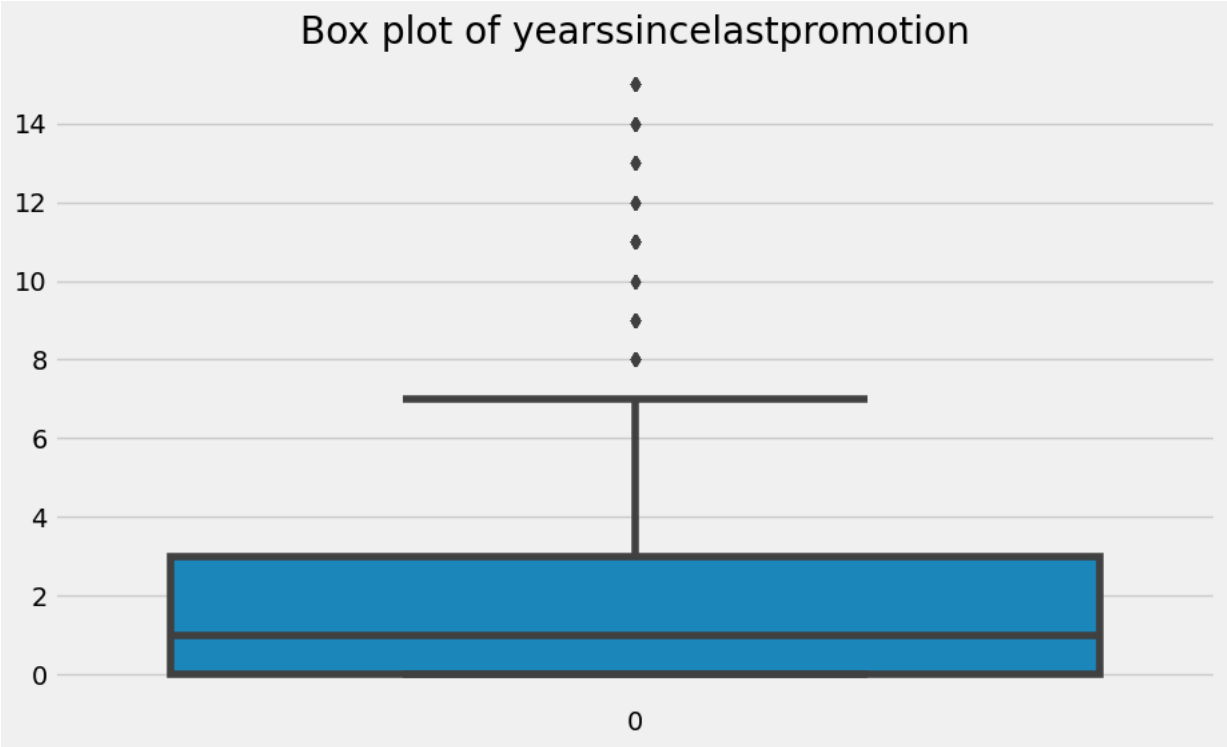
- What features have the most divergent distributions based on target class
- Do we have a target label imbalance
- How our independent variables are distributed relative to our target label
- Are there features that have strong linear or monotonic relationships, making correlation heatmaps makes it easy to identify possible colinearity

▼ Check for outliers

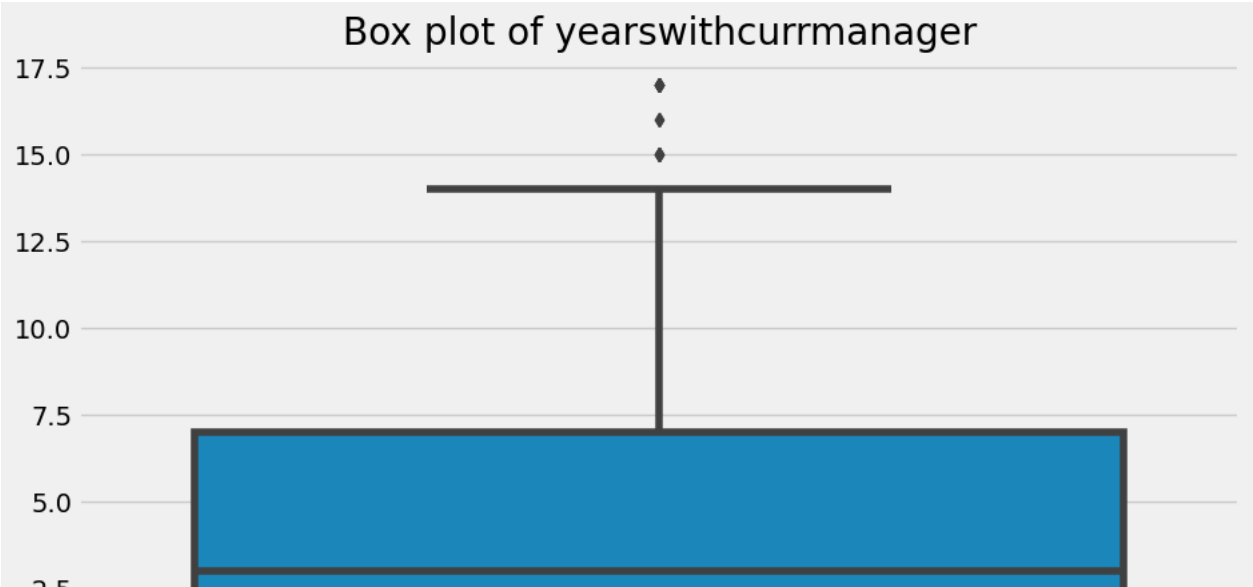
Exercise 3: Create a box plot to check for outliers [0.5 Mark]

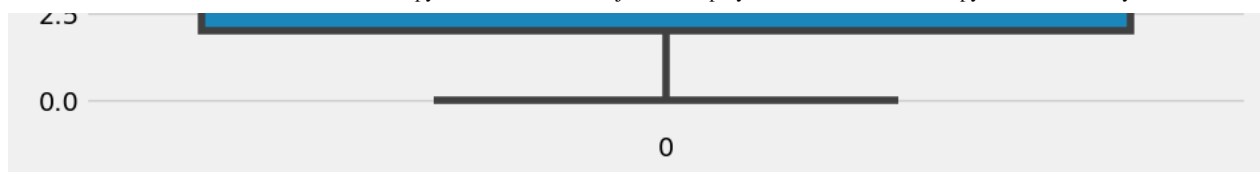
```
1 # Check for outliers
2 # YOUR CODE HERE
3 import pandas as pd
4 import matplotlib.pyplot as plt
5 import seaborn as sns
6
7 # Assuming you've already read your dataset into a DataFrame named 'dataset'
8
9 # Extract numerical columns if you haven't already
10 numerical_columns = [column for column in dataset.columns if dataset[column].dtype in ['int64', 'float64']]
11
12 # Create box plots for each numerical column and print the outliers
13 for column in numerical_columns:
14     plt.figure(figsize=(10, 6))
15     sns.boxplot(dataset[column])
16     plt.title(f"Box plot of {column}")
17     plt.show()
18
19     # Calculate IQR for the column
20     Q1 = dataset[column].quantile(0.25)
21     Q3 = dataset[column].quantile(0.75)
22     IQR = Q3 - Q1
23
24     # Determine bounds for outliers
25     lower_bound = Q1 - 1.5 * IQR
26     upper_bound = Q3 + 1.5 * IQR
27
28     # Extract outliers
29     outliers = dataset[(dataset[column] < lower_bound) | (dataset[column] > upper_bound)][column]
30
31     # Print outliers if they exist
32     if outliers.empty:
33         print(f"No outliers found for {column}.\n")
34     else:
35         print(f"Outliers for {column}:")
36         print(outliers, "\n")
37
38
```

```
190    18
231    17
281    16
417    15
466    16
595    15
716    16
746    16
861    15
976    16
1024   17
1150   15
1156   15
1221   15
1327   17
1351   17
1430   16
Name: yearsincurrentrole, dtype: int64
```



```
Outliers for yearssincelastpromotion:
15      8
45     15
46      8
55      8
61      9
..
1414   12
1425    8
1444    9
1447   11
1462    9
Name: yearssincelastpromotion, Length: 107, dtype: int64
```





Outliers for yearswithcurrmanager:

28	17
123	15
153	15
187	15
231	15
386	17
561	16
616	17
635	15
686	17
875	17
926	17
1078	17
1348	16

Name: yearswithcurrmanager, dtype: int64

▼ Handling outliers

Exercise 4: Use lower bound as 25% and upper bound as 75% to handle the outliers [0.5 Mark]

```
1 # YOUR CODE HERE
2 import pandas as pd
3
4 # Assuming you've already read your dataset into a DataFrame named 'dataset'
5
6 # Check the first few rows of the dataset to ensure it's loaded properly:
7 print("Dataset head:")
8 print(dataset.head())
9 print("\n-----\n")
10
11 # Extract numerical columns if you haven't already
12 numerical_columns = [column for column in dataset.columns if dataset[column].dtype in ['int64', 'float64']]
13
14 # Display the numerical columns to verify:
15 print("Numerical columns:")
16 print(numerical_columns)
17 print("\n-----\n")
18
19 # For each numerical column, cap the data with the 25% and 75% quantiles
20 for column in numerical_columns:
21     Q1 = dataset[column].quantile(0.25)
22     Q3 = dataset[column].quantile(0.75)
23
24     # Replacing values below Q1 with Q1 and values above Q3 with Q3
25     dataset[column] = dataset[column].apply(lambda x: Q1 if x < Q1 else (Q3 if x > Q3 else x))
26
27 # Check the revised data
28 print("Revised Data Description:")
```

```
29 print(dataset[numerical_columns].describe())
```

```
30
```

Dataset head:

	age	attrition	businesstravel	dailyrate	department	\
0	41	Yes	Travel_Rarely	1102	Sales	
1	49	No	Travel_Frequently	279	Research & Development	
2	37	Yes	Travel_Rarely	1373	Research & Development	
3	33	No	Travel_Frequently	1392	Research & Development	
4	27	No	Travel_Rarely	591	Research & Development	

	distancefromhome	education	educationfield	employeeecount	employeenumber	\
0	1	2	Life Sciences	1	1	
1	8	1	Life Sciences	1	2	
2	2	2	Other	1	4	
3	3	4	Life Sciences	1	5	
4	2	1	Medical	1	7	

	environmentsatisfaction	gender	hourlyrate	jobinvolvement	joblevel	\
0	2	Female	94	3	2	
1	3	Male	61	2	2	
2	4	Male	92	2	1	
3	4	Female	56	3	1	
4	1	Male	40	3	1	

	jobrole	jobsatisfaction	maritalstatus	monthlyincome	\
0	Sales Executive	4	Single	5993	
1	Research Scientist	2	Married	5130	
2	Laboratory Technician	3	Single	2090	
3	Research Scientist	3	Married	2909	
4	Laboratory Technician	2	Married	3468	

	monthlyrate	numcompaniesworked	over18	overtime	percentsalaryhike	\
0	19479	8	Y	Yes	11	
1	24907	1	Y	No	23	
2	2396	6	Y	Yes	15	
3	23159	1	Y	Yes	11	
4	16632	9	Y	No	12	

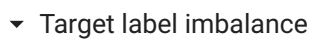
	performancerating	relationshipsatisfaction	standardhours	\
0	3	1	80	
1	4	4	80	
2	3	2	80	
3	3	3	80	
4	3	4	80	

	stockoptionlevel	totalworkingyears	trainingtimeslastyear	\
0	0	8	0	
1	1	10	3	
2	0	7	3	
3	0	8	3	
4	1	6	3	

	worklifebalance	yearsatcompany	yearsincurrentrole	\
0	1	6	4	
1	3	10	7	
2	3	0	0	
3	3	8	7	
4	3	2	2	

yearssincelastpromotion yearswithcurrmanager

```
1 # Recheck for outliers
2 # YOUR CODE HERE
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5
6 # Setting the size for the plots
7 plt.figure(figsize=(20, 10))
8
9 # Plotting boxplots for all numerical columns
10 for i, column in enumerate(numerical_columns, 1):
11     plt.subplot(2, len(numerical_columns)//2, i)
12     sns.boxplot(y=dataset[column])
13     plt.title(column)
14     plt.ylabel('Value')
15     plt.tight_layout()
16
17 plt.show()
18
```



```
1 # Count of unique values in Attrition column
2 # YOUR CODE HERE
3 # Checking the distribution of the target label
4 print(dataset['attrition'].value_counts())
5
6 # Plotting the distribution
7 plt.figure(figsize=(8, 5))
8 sns.countplot(data=dataset, x='attrition')
9 plt.title("Distribution of Target Label - Attrition")
10 plt.xlabel('Attrition')
11 plt.ylabel('Count')
12 plt.show()
13
```

Distribution of Target Label - Attrition

A bar chart titled 'Distribution of Target Label - Attrition'. The x-axis is labeled 'Attrition' and has two categories: 'Yes' and 'No'. The y-axis is labeled 'Count' and ranges from 0 to 1200 with major grid lines every 200 units. The 'Yes' bar is blue and reaches a count of approximately 240. The 'No' bar is orange and reaches a count of approximately 1240.

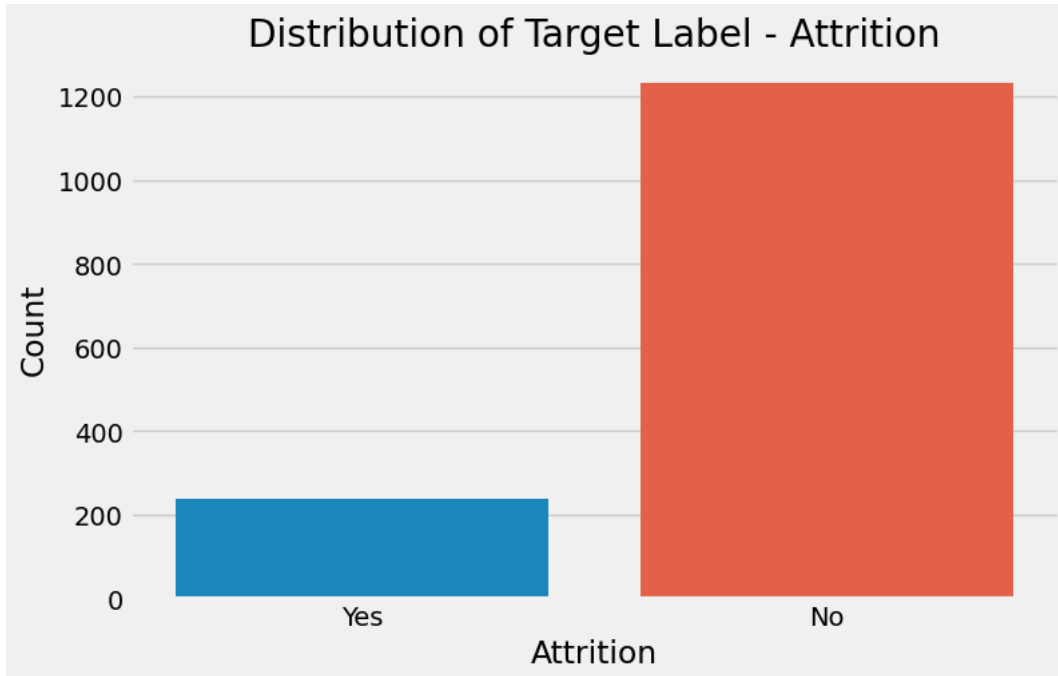
Attrition	Count
Yes	240
No	1240

19/21

```

3 import matplotlib.pyplot as plt
4 import seaborn as sns
5
6 # Plotting the distribution
7 plt.figure(figsize=(8, 5))
8 sns.countplot(data=dataset, x='attrition')
9 plt.title("Distribution of Target Label - Attrition")
10 plt.xlabel('Attrition')
11 plt.ylabel('Count')
12 plt.show()
13

```



If there is any imbalance in the dataset then a few techniques can be utilised (optional):

1. SMOTE
2. Cross Validation
3. Regularizing the model's parameters

▼ Plot pairplot

Exercise 6: Visualize the relationships between the predictor variables and the target variable using a pairplot [0.5 Mark]

Hint: Use sns.pairplot

```

1 # Visualize a pairplot with relevant features
2 # YOUR CODE HERE
3 import seaborn as sns
4
5 # Selecting a few features for demonstration purposes. Adjust this list as needed.
6 selected_features = ['age', 'dailyrate', 'distancefromhome', 'monthlyincome']
7
8 # Adding the target label to visualize the relationships based on the target classes
9 selected_features.append('attrition')
10
11 # Creating the pairplot
12 sns.pairplot(dataset[selected_features], hue='attrition', plot_kws={'alpha':0.5})
13 plt.show()
14

```



▼ Explore Correlation

- Plotting the Heatmap



Exercise 7: Visualize the correlation among IBM employee attrition numerical features using a heatmap [0.5 Mark]

```

1 # Visualize heatmap
2 # YOUR CODE HERE
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5
6 # Compute the correlation matrix
7 corr_matrix = dataset[numerical_columns].corr()
8
9 # Create a heatmap to visualize the correlations
10 plt.figure(figsize=(15, 10))
11 sns.heatmap(corr_matrix, annot=True, cmap='coolwarm', vmin=-1, vmax=1)
12 plt.title("Correlation Heatmap of Numerical Features")
13 plt.show()
14

```