

EE5253Report_4466_4578_GP23

December 12, 2024

Project name: **Customer Churn Prediction**

Team members:

- EG/2021/4466 - De Silva T.S.A.N.H.
- EG/2021/4578 - Jayasinghe D.J.K.V.

Group number: 23

1 Introduction

Customer churn prediction is one of the most important areas of research for businesses operating on subscription-based models, such as telecom companies, SaaS providers, and online streaming platforms. Customer churn refers to the phenomenon where customers stop using a company's products or services over a certain period. For businesses, identifying the factors leading to churn and predicting which customers are at high risk of leaving can help in the formulation of targeted retention strategies. This, in turn, enables businesses to reduce churn, enhance customer lifetime value, and improve profitability.

The ability to predict churn with high accuracy relies heavily on the analysis of customer data, such as usage patterns, customer service interactions, and payment history. Traditionally, churn prediction has been done using basic statistical methods; however, with the advent of machine learning, businesses now leverage more sophisticated algorithms to gain deeper insights.

This project seeks to build a machine learning model that predicts customer churn using two widely popular algorithms: **Random Forest** and **XGBoost**. Both models have demonstrated exceptional performance in classification tasks, making them well-suited for churn prediction. The objective of this project is to explore both models, comparing their performance and determining which model provides better predictive accuracy and interpretability for this problem.

By the end of this project, the goal is not only to predict churn but also to gain insights into which features most significantly influence a customer's decision to churn. This can help businesses create personalized retention strategies and reduce churn by targeting high-risk customers more effectively.

2 Literature Survey

Customer churn prediction has been the focus of numerous studies in the business analytics field. Over the years, researchers have explored various machine learning techniques to accurately predict churn, focusing on how to better understand customer behavior and implement preventive measures. Among the many algorithms that have been used, **Random Forest** and **XGBoost** have emerged as two of the most powerful models for classification tasks, including churn prediction.

Random Forest

Random Forest is an ensemble learning method that builds multiple decision trees during training and merges their outputs to improve classification accuracy. Each tree in a Random Forest is trained on an equal subset of the data, using a random subset of features. This randomness helps the model generalize well and reduces the risk of overfitting. Random Forest is known for its robustness, ease of use, and ability to handle both categorical and numerical data without the need for extensive preprocessing. It also provides feature importance, which can be useful for understanding which factors contribute most to churn.

Key Advantages:

- **Non-linear relationships:** Random Forest can model complex interactions between variables.
- **Interpretability:** It offers insights into feature importance, which can help businesses understand key drivers of churn.
- **Robustness:** The model is resistant to overfitting due to averaging over multiple decision trees.

Challenges:

- While Random Forest is easy to use, it can be computationally expensive and time-consuming when dealing with large datasets.
- It may not always perform as well as more advanced boosting algorithms, particularly on unbalanced datasets like churn prediction.

XGBoost

XGBoost, which stands for eXtreme Gradient Boosting, is a high-performance implementation of gradient boosting that has become one of the most popular machine learning algorithms for structured data. XGBoost works by building decision trees sequentially, where each new tree attempts to correct the errors made by the previous tree. This boosting process is highly effective for problems like churn prediction, where small improvements in accuracy can lead to substantial performance gains. XGBoost is known for its ability to handle large datasets and its support for various regularization techniques that help control overfitting.

Key Advantages:

- **Accuracy:** XGBoost has been shown to outperform other machine learning algorithms in terms of accuracy and computational efficiency.
- **Handling Missing Data:** It is capable of handling missing values in the data, which is often crucial in real-world datasets.

- **Hyperparameter Tuning:** XGBoost provides a large set of hyperparameters that can be fine-tuned to optimize performance, including parameters for regularization and learning rate.

Challenges:

- **Complexity:** XGBoost is more complex than Random Forest and requires careful tuning of hyperparameters to avoid overfitting.
- **Interpretability:** While XGBoost provides feature importance scores, the model's complex boosting process can make it harder to interpret compared to simpler models like Random Forest.

Comparison Between Random Forest and XGBoost

A comparison between Random Forest and XGBoost reveals several interesting findings in the context of churn prediction. While Random Forest is a simple and interpretable model, XGBoost often delivers superior accuracy due to its sequential boosting mechanism, which allows the model to learn from the mistakes of prior trees. XGBoost also tends to perform better in situations where the data is noisy or contains irrelevant features. On the other hand, Random Forest is typically faster to train and is easier to understand, which is why it remains a popular choice in situations where model interpretability is paramount.

Previous studies have shown that **XGBoost** often outperforms **Random Forest** in terms of predictive accuracy, particularly when dealing with imbalanced datasets. However, **Random Forest** can still be a strong contender when simplicity and interpretability are more important than raw performance.

In conclusion, both Random Forest and XGBoost have their strengths and weaknesses, and the choice of model depends on the specific requirements of the churn prediction task at hand. This study aims to evaluate both algorithms using a customer churn dataset, comparing their performance to determine the best approach for churn prediction.

3 Exploratory Data Analysis and Data Preprocessing

3.1 Importing Libraries

```
[1]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
import time
from sklearn.ensemble import RandomForestClassifier
from xgboost import XGBClassifier
from sklearn.model_selection import train_test_split
from sklearn.metrics import (accuracy_score, precision_score, recall_score,
    f1_score, roc_curve, auc, confusion_matrix, ConfusionMatrixDisplay)
from sklearn.model_selection import GridSearchCV
from sklearn.preprocessing import StandardScaler
```

3.2 Importing Dataset

```
[2]: dataset = pd.read_csv('/content/customer_churn_dataset.csv')
```

```
dataset
```

```
[2]:
```

	CustomerID	Age	Gender	Tenure	Usage Frequency	Support Calls \
0	2.0	30.0	Female	39.0	14.0	5.0
1	3.0	65.0	Female	49.0	1.0	10.0
2	4.0	55.0	Female	14.0	4.0	6.0
3	5.0	58.0	Male	38.0	21.0	7.0
4	6.0	23.0	Male	32.0	20.0	5.0
...
440828	449995.0	42.0	Male	54.0	15.0	1.0
440829	449996.0	25.0	Female	8.0	13.0	1.0
440830	449997.0	26.0	Male	35.0	27.0	1.0
440831	449998.0	28.0	Male	55.0	14.0	2.0
440832	449999.0	31.0	Male	48.0	20.0	1.0

	Payment Delay	Subscription Type	Contract Length	Total Spend \
0	18.0	Standard	Annual	932.00
1	8.0	Basic	Monthly	557.00
2	18.0	Basic	Quarterly	185.00
3	7.0	Standard	Monthly	396.00
4	8.0	Basic	Monthly	617.00
...
440828	3.0	Premium	Annual	716.38
440829	20.0	Premium	Annual	745.38
440830	5.0	Standard	Quarterly	977.31
440831	0.0	Standard	Quarterly	602.55
440832	14.0	Premium	Quarterly	567.77

	Last Interaction	Churn
0	17.0	1.0
1	6.0	1.0
2	3.0	1.0
3	29.0	1.0
4	20.0	1.0
...
440828	8.0	0.0
440829	2.0	0.0
440830	9.0	0.0
440831	2.0	0.0
440832	21.0	0.0

```
[440833 rows x 12 columns]
```

3.3 Dataset Description

The dataset used in this project contains customer-related features to predict customer churn. Churn, in this context, refers to whether a customer stops using the services of a company. The dataset consists of the following 12 variables.

3.3.1 Features:

1. **CustomerID:**
 - Type: Numerical
 - A unique identifier for each customer.
2. **Age:**
 - Type: Numerical
 - The age of the customer in years.
3. **Gender:**
 - Type: Categorical
 - The gender of the customer (e.g., Male, Female).
4. **Tenure:**
 - Type: Numerical
 - The number of months the customer has been subscribed to the service.
5. **Usage Frequency:**
 - Type: Numerical
 - The frequency with which the customer uses the service.
6. **Support Calls:**
 - Type: Numerical
 - The number of support calls made by the customer.
7. **Payment Delay:**
 - Type: Numerical
 - The number of days the customer's payment was delayed.
8. **Subscription Type:**
 - Type: Categorical
 - The type of subscription plan the customer has, such as Basic, Standard, or Premium.
9. **Contract Length:**
 - Type: Categorical
 - The duration of the customer's subscription contract, such as Monthly, Quarterly, or Annual.
10. **Total Spend:**
 - Type: Numerical
 - The total amount of money the customer has spent on the service.
11. **Last Interaction:**
 - Type: Numerical
 - The number of days since the customer's last interaction with the company.
12. **Churn:**
 - Type: Binary (Target Variable)
 - Indicates whether the customer has churned (1 for churned, 0 for retained).

3.3.2 Summary:

- The dataset contains a mix of numerical and categorical features.

- The target variable is **Churn**, which is binary and suitable for classification tasks.
- This dataset provides rich customer information, enabling insights into the factors influencing churn behavior.

```
[3]: dataset.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 440833 entries, 0 to 440832
Data columns (total 12 columns):
#   Column                Non-Null Count  Dtype
---  -
0   CustomerID            440832 non-null  float64
1   Age                   440832 non-null  float64
2   Gender                440832 non-null  object
3   Tenure                440832 non-null  float64
4   Usage Frequency       440832 non-null  float64
5   Support Calls         440832 non-null  float64
6   Payment Delay         440832 non-null  float64
7   Subscription Type     440832 non-null  object
8   Contract Length       440832 non-null  object
9   Total Spend           440832 non-null  float64
10  Last Interaction       440832 non-null  float64
11  Churn                 440832 non-null  float64
dtypes: float64(9), object(3)
memory usage: 40.4+ MB
```

```
[4]: dataset.describe()
```

```
[4]:
```

	CustomerID	Age	Tenure	Usage Frequency \
count	440832.000000	440832.000000	440832.000000	440832.000000
mean	225398.667955	39.373153	31.256336	15.807494
std	129531.918550	12.442369	17.255727	8.586242
min	2.000000	18.000000	1.000000	1.000000
25%	113621.750000	29.000000	16.000000	9.000000
50%	226125.500000	39.000000	32.000000	16.000000
75%	337739.250000	48.000000	46.000000	23.000000
max	449999.000000	65.000000	60.000000	30.000000

	Support Calls	Payment Delay	Total Spend	Last Interaction \
count	440832.000000	440832.000000	440832.000000	440832.000000
mean	3.604437	12.965722	631.616223	14.480868
std	3.070218	8.258063	240.803001	8.596208
min	0.000000	0.000000	100.000000	1.000000
25%	1.000000	6.000000	480.000000	7.000000
50%	3.000000	12.000000	661.000000	14.000000
75%	6.000000	19.000000	830.000000	22.000000
max	10.000000	30.000000	1000.000000	30.000000

	Churn
count	440832.000000
mean	0.567107
std	0.495477
min	0.000000
25%	0.000000
50%	1.000000
75%	1.000000
max	1.000000

3.4 Handling null/missing values

```
[5]: dataset.isna().sum()
```

```
[5]: CustomerID      1
      Age            1
      Gender         1
      Tenure         1
      Usage Frequency 1
      Support Calls   1
      Payment Delay   1
      Subscription Type 1
      Contract Length 1
      Total Spend     1
      Last Interaction 1
      Churn           1
      dtype: int64
```

```
[6]: dataset.dropna(inplace=True)
```

```
[7]: dataset.isna().sum()
```

```
[7]: CustomerID      0
      Age            0
      Gender         0
      Tenure         0
      Usage Frequency 0
      Support Calls   0
      Payment Delay   0
      Subscription Type 0
      Contract Length 0
      Total Spend     0
      Last Interaction 0
      Churn           0
      dtype: int64
```

3.5 Treating Duplicate Records

```
[8]: dataset.duplicated().sum()
```

```
[8]: 0
```

Note

By considering these facts (**only one whole row with missing values and it was dropped** and **no duplicates**) we come to a conclusion that train-test split can be done at the end of data preprocessing as there is no chance of occurring data snooping phenomenon.

3.6 Visualizing Data

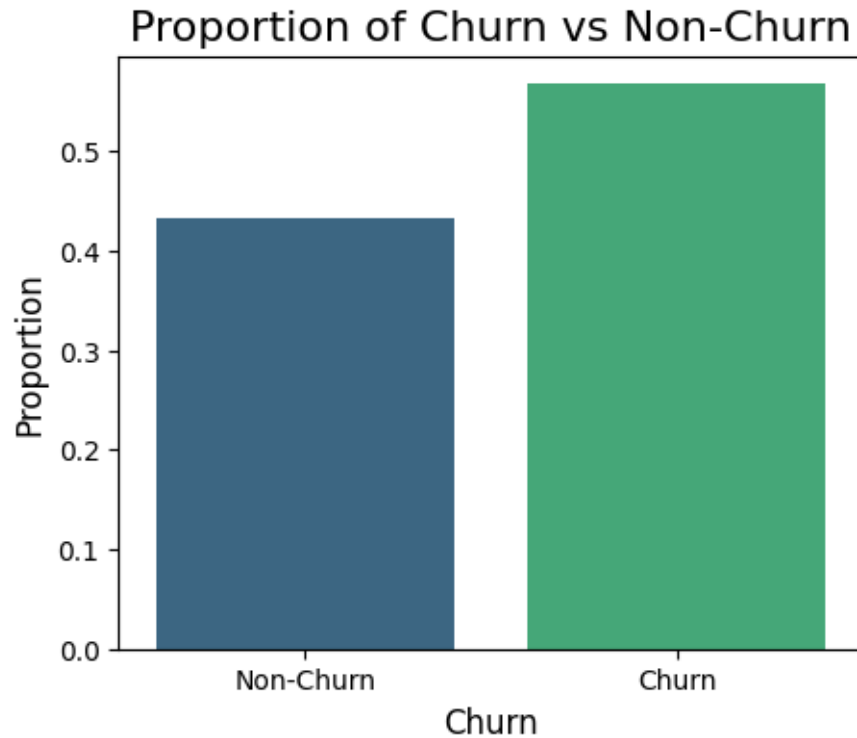
3.6.1 Analyze the distribution of variables

```
[9]: # Checking churn proportions
plt.figure(figsize=(5,4))
churn_counts = dataset['Churn'].value_counts(normalize=True)
sns.barplot(x=churn_counts.index, y=churn_counts.values, palette='viridis')
plt.title('Proportion of Churn vs Non-Churn', fontsize=16)
plt.xlabel('Churn', fontsize=12)
plt.ylabel('Proportion', fontsize=12)
plt.xticks([0, 1], ['Non-Churn', 'Churn'])
plt.show()
```

<ipython-input-9-8f91d4b1b426>:4: FutureWarning:

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

```
sns.barplot(x=churn_counts.index, y=churn_counts.values, palette='viridis')
```

Proportion of Churn:

1. Churn = 1 (Leave): Accounts for a higher rate (about 55%)
2. Churn = 0 (Do not leave): Accounts for about 45%

Conclusion: The number of customers leaving the service (55%) is higher than the number of customers continuing to use the service (45%)

```
[10]: categorical_cols = ['Gender', 'Subscription Type', 'Contract Length']

# Plotting distribution (frequency counts) for categorical variables
fig, axes = plt.subplots(1, 3, figsize=(20, 6))

# Plot bar plots for each categorical variable
for i, col in enumerate(categorical_cols):
    sns.countplot(x=col, data=dataset, palette='viridis', ax=axes[i])
    axes[i].set_title(f'Distribution of {col}', fontsize=14)
    axes[i].set_xlabel(col, fontsize=12)
    axes[i].set_ylabel('Count', fontsize=12)
    axes[i].tick_params(axis='x', rotation=45)

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

```
<ipython-input-10-ce6899fb3974>:8: FutureWarning:
```

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

```
sns.countplot(x=col, data=dataset, palette='viridis', ax=axes[i])
```

```
<ipython-input-10-ce6899fb3974>:8: FutureWarning:
```

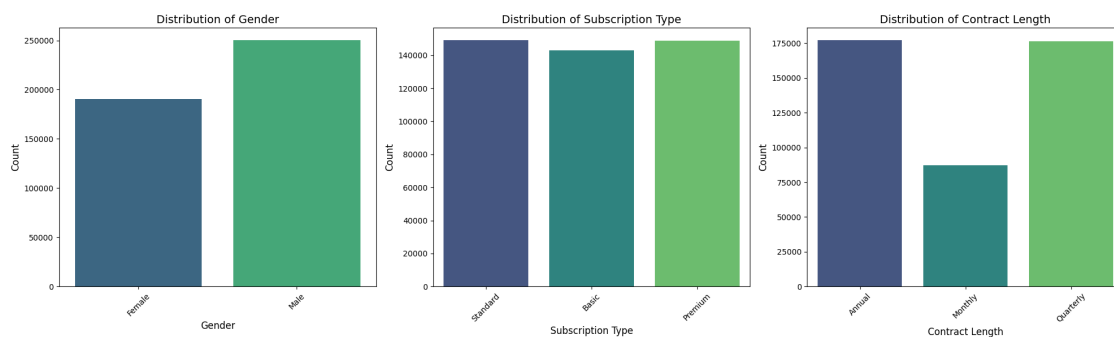
Passing `palette` without assigning `hue` is deprecated and will be removed in v0.14.0. Assign the `x` variable to `hue` and set `legend=False` for the same effect.

```
sns.countplot(x=col, data=dataset, palette='viridis', ax=axes[i])
```

```
<ipython-input-10-ce6899fb3974>:8: FutureWarning:
```

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

```
sns.countplot(x=col, data=dataset, palette='viridis', ax=axes[i])
```



Categorical features:

1. Gender: The number of male customers (Male) accounts for a larger proportion than female customers (Female)
2. Subscription Type (Type of service package): The number of customers in the three service packages Basic, Standard, and Premium is almost equivalent, the Basic package is slightly lower than the other two types
3. Contract Length: Long-term contracts (Annual, Quarterly) are much more popular than short-term contracts (Monthly)

```
[11]: numerical_cols = ['Age', 'Tenure', 'Usage Frequency', 'Support Calls', 'Payment_␣  
      ↪Delay', 'Total Spend', 'Last Interaction']
```

```

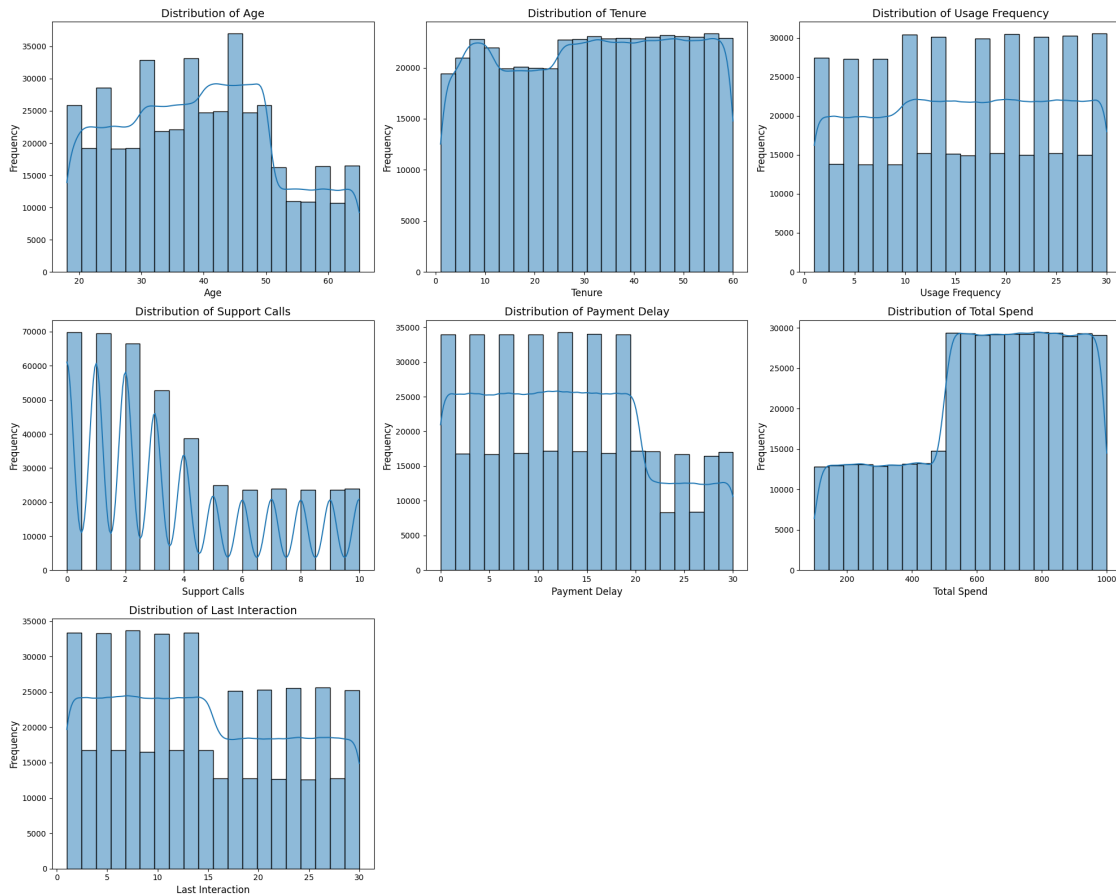
# Plot histograms for numerical variables
fig, axes = plt.subplots(3, 3, figsize=(20, 16))
axes = axes.flatten()

for i, col in enumerate(numerical_cols):
    sns.histplot(dataset[col], kde=True, bins=20, ax=axes[i])
    axes[i].set_title(f'Distribution of {col}', fontsize=14)
    axes[i].set_xlabel(col, fontsize=12)
    axes[i].set_ylabel('Frequency', fontsize=12)

# Remove empty subplot
fig.delaxes(axes[-1])
fig.delaxes(axes[-2])

plt.tight_layout()
plt.show()

```



Numerical Variables:

1. Age: The data is fairly evenly distributed, with some peaks in the 30-50 age group and a

sharp decline in the over 60 age group.

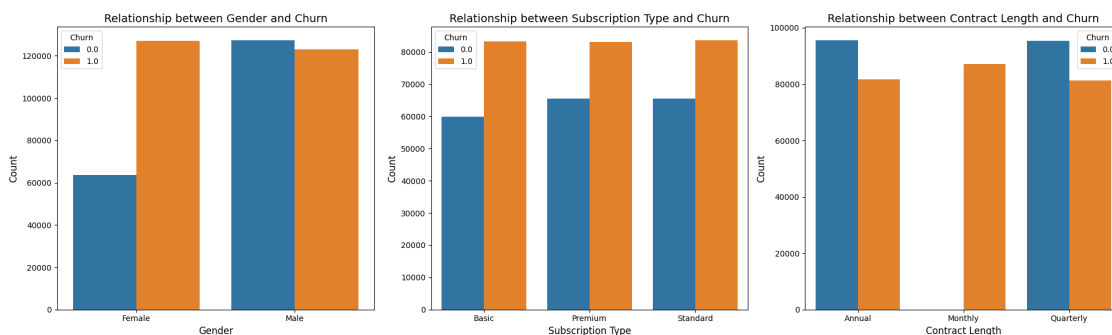
2. Tenure: Service usage time ranges from 0 to 60 months. The number of customers using over 30 months is high
3. Usage Frequency: Usage frequency is focused on the average level, showing that the majority of customers maintain a stable level of usage.
4. Support Calls (Number of support calls): The number of support calls decreases from 0 to 10 times. Most customers call support very few times, most customers call less than 3 times, but there are also some cases that call many times.
5. Payment Delay: The majority of customers have low payment delays (about 0 to 20 days), but a small portion have significant payment delays (more than 20 days) that can be a signal from customers. Having financial problems and leaving easily.
6. Total Spend: Total spending tends to increase from 0 - 1000 units, increasing sharply in the range (450, 550). Most customers spend at a high level (> 550)
7. Last Interaction: The majority of customers have a recent average interaction (~14 days), but a significant number of customers have not interacted with the system for a long time (> 16,17 days)

3.6.2 Analyze the relationship between features and churn

```
[12]: # Create subplots for categorical variables and their relationship with 'Churn'
fig, axes = plt.subplots(nrows=1, ncols=3, figsize=(20, 6))

for i, col in enumerate(categorical_cols):
    churn_counts = dataset.groupby([col, 'Churn']).size().
    ↪reset_index(name='Count')
    sns.barplot(x=col, y='Count', hue='Churn', data=churn_counts, ax=axes[i])
    axes[i].set_title(f'Relationship between {col} and Churn', fontsize=14)
    axes[i].set_xlabel(col, fontsize=12)
    axes[i].set_ylabel('Count', fontsize=12)
    axes[i].legend(title='Churn')

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



Categorical features:

1. Gender: There is a clear difference between the two genders:
 - The number of male customers leaving is less than that of female customers.
 - For male customers, the number of customers who continue to use the service (churn=0) is higher than the number of customers who leave.
 - In contrast, for female customers, the number of customers leaving is about twice as many as the number staying.
2. Subscription Type (Service package type):
 - The “Basic” package has a higher churn rate than “Standard” and “Premium” -> The higher-end package retains customers better.
3. Contract Length:
 - “Monthly” contracts have significantly higher churn rates than “Annual” and “Quarterly”. Specifically, most customers using the “Monthly” package will not continue to use the service. Meanwhile, long-term contracts help retain customers more effectively.

Conclusion: There should be service business strategies that focus more on female customers, and at the same time, there should also be strategies to improve the value of the Basic package and the Monthly contract type.

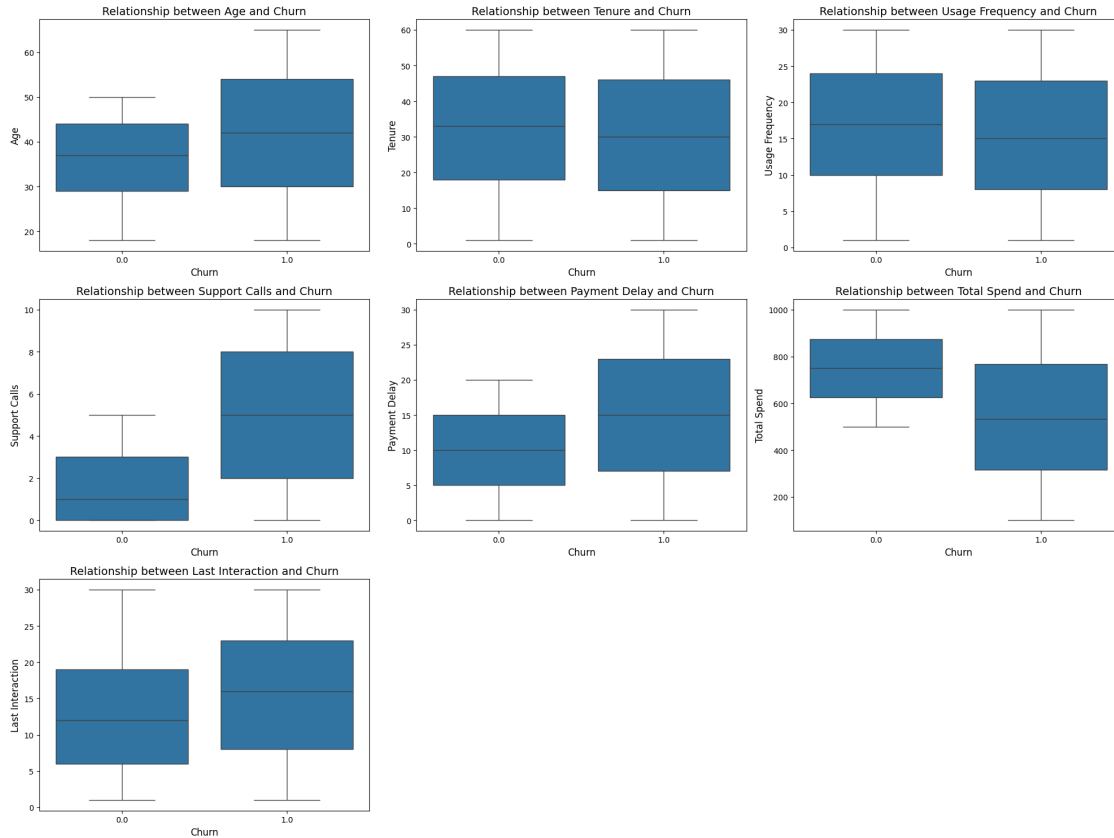
```
[13]: # Create subplots for numerical variables and their relationship with 'Churn'
fig, axes = plt.subplots(nrows=3, ncols=3, figsize=(20, 15))
axes = axes.flatten()

numerical_cols = ['Age', 'Tenure', 'Usage Frequency', 'Support Calls',
                  'Payment Delay', 'Total Spend', 'Last Interaction']

for i, col in enumerate(numerical_cols):
    sns.boxplot(x='Churn', y=col, data=dataset, ax=axes[i])
    axes[i].set_title(f'Relationship between {col} and Churn', fontsize=14)
    axes[i].set_xlabel('Churn', fontsize=12)
    axes[i].set_ylabel(col, fontsize=12)

# Remove empty subplots if any
for j in range(len(numerical_cols), len(axes)):
    fig.delaxes(axes[j])

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



Note

No outliers were detected from the boxplots. So, no need of outliers treating step.

Numerical Variables:

1. Age:

Median:

- Non-Churn group (Churn=0): Median is around 40 years old.
- Churn group (Churn=1): Median is higher, around 45 years old.

Distribution range (IQR - Interquartile Range):

- Non-Churn group: Distribution focuses mainly on 30 to 50 year olds.
- Churn group: Wider distribution, spanning from 20 to 60 years old.

Conclusion: The Churn group has a wider distribution in terms of age: Older customers (45+) and younger customers (30-) will tend to Churn. Age is not the only deciding factor in whether customers leave or not, it must be combined with other factors.

2. Tenure (Time of use):

- There is high similarity between the two groups churn=0 and churn=1
- The median values are very close (~30 months)
- The distribution range (IQR) is also very similar (~10-50 months)

Conclusion: Tenure is not a strong factor to distinguish between the two groups Churn and Non-churn

3. Usage Frequency:

- Similarly, Usage Frequency is not a strong factor to evaluate whether a customer will leave or continue to use the service.

4. Support Calls (Number of support calls):

Non-Churn group (churn=0):

- Median: The number of support calls is not high, about 1-2 times.
- Delivery coverage: Most Non-Churn customers have less than 4 support calls (IQR ranges from 0-4).

Churn group (churn=1):

- The median of this group is significantly higher, at about 6 times.
- The distribution reach is significantly wider, with many customers having up to 8-10 support calls.

Conclusion: Support Calls are an important factor affecting Churn. The Churn group had a significantly higher number of support calls than the Non-Churn group. The Non-Churn group often does not or rarely calls support, indicating that they are satisfied with the service and do not need to ask for more. Customers who frequently call support (more than 5 times) have a higher tendency to leave the service, possibly due to dissatisfaction with the service.

5. Payment Delay:

Churn group = 0 (Non-Churn):

- Median (Median): About 10 days.
- Distribution scope: Concentrated in about 5-15 days. There are a few customers whose Payment Delay exceeds 20 days, but this is not significant

Churn group = 1 (Churn):

- Median: Significantly higher, at around 15 days.
- Distribution scope: Wider, with most customers having Payment Delay between 10-20 days and extending to 30 days.

Conclusion: Payment Delay is an important factor affecting Churn. The chart clearly shows that many customers in the Churn group have higher Payment Delay than the Non-Churn group, meaning that customers with long payment delays (high Payment Delay) are at higher risk of Churn. The wider distribution of the Churn group (with

many customers exceeding 20 days) may reflect financial problems or dissatisfaction with the service.

6. Total Spend:

Churn group = 0 (Non-Churn):

- Median: About 800.
- Distribution range: Concentrated in the 600-900 range, reflecting higher spending compared to the Churn group.
- Distribution is very narrow, with most customers having high and stable spending.

Churn group = 1 (Churn):

- Median: About 500
- Distribution range: Wider than the Non-Churn group, in the range of 300-700

Conclusion: Total Spend is also an important factor in predicting Churn. The Churn group tends to spend significantly less than the Non-Churn group. Customers who spend little (<600) often have a high risk of churn

7. Last Interaction (Last interaction time):

Churn group = 0 (Non-Churn):

- Median (Median): About 10 days.
- Distribution scope: Focused on about 5-15 days, with a few customers having a last interaction time of over 20 days.
- Narrow distribution, reflecting that Non-Churn customers often have had recent interactions.

Churn group = 1 (Churn):

- Median: Higher than the Non-Churn group, at about 15 days.
- Distribution range: Wider distribution, within 10-20 days.

Conclusion: Last Interaction is also a factor that affects whether customers leave or not, but not as strong as the features Support Calls, Payment Delay and Total Spend, but stronger than Tenure and Usage Frequency

Insights:

1. Customers who call for support a lot often tend to leave, possibly because the service experience provided is not good enough.
2. Weak finances: Customers at risk of churn often have financial problems (high Payment Delay, low Total Spend, Monthly Contract).
3. Female customers are more likely to leave than male customers.
4. Elderly customers (50+) and teenage customers (30-) tend to leave the service.
5. Low stickiness: Customers who rarely use the service, have few interactions recently, are often more likely to leave because they may be dissatisfied with customer support.

3.6.3 Analyze the correlation between variables

```
[14]: # Encode Gender and other categorical variables with ordinal/binary encoding
# Gender: Male = 1, Female = 0
encoded_data_binary = dataset.copy()
encoded_data_binary['Gender'] = encoded_data_binary['Gender'].apply(lambda x: 1
    ↪if x == 'Male' else 0)

# Subscription Type: Basic = 0, Standard = 1, Premium = 2
subscription_map = {'Basic': 0, 'Standard': 1, 'Premium': 2}
encoded_data_binary['Subscription Type'] = encoded_data_binary['Subscription_
    ↪Type'].map(subscription_map)

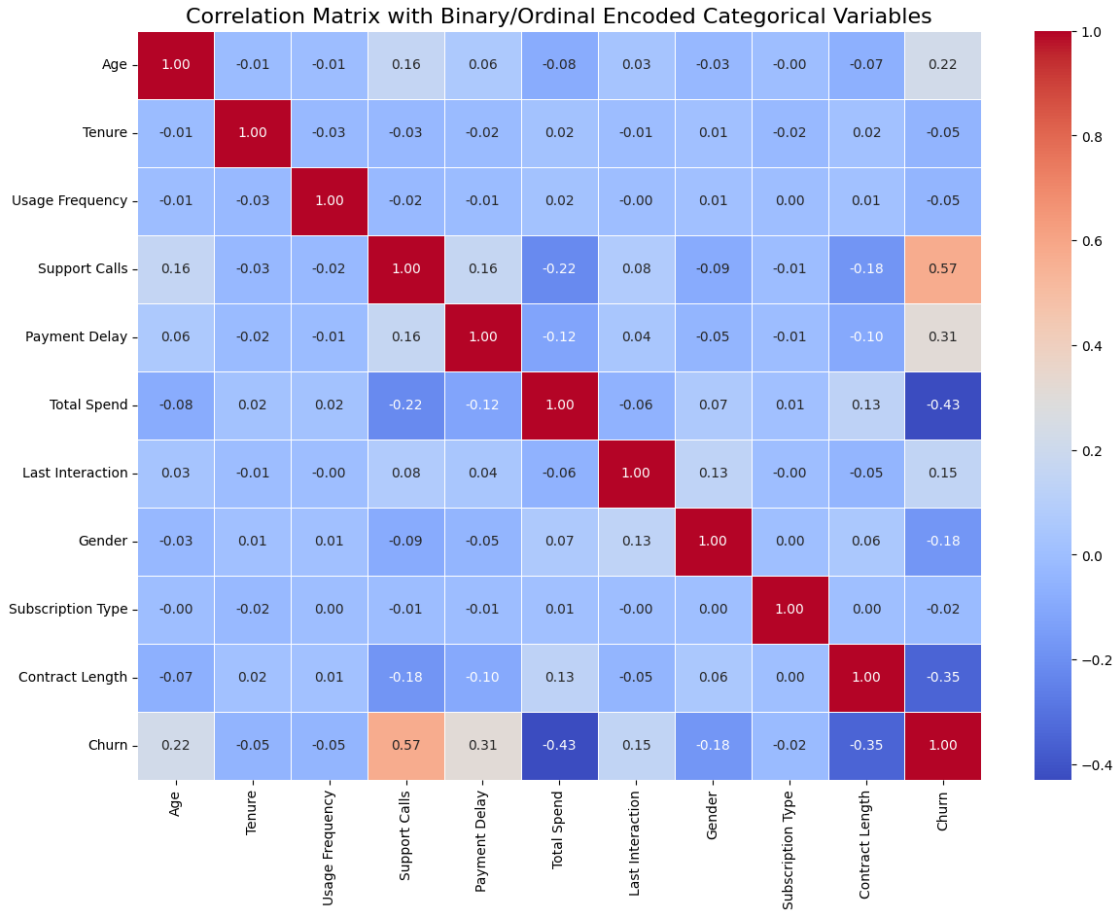
# Contract Length: Monthly = 0, Quarterly = 1, Annual = 2
contract_map = {'Monthly': 0, 'Quarterly': 1, 'Annual': 2}
encoded_data_binary['Contract Length'] = encoded_data_binary['Contract Length'].
    ↪map(contract_map)

# Combine numerical and newly encoded categorical columns
all_columns_binary = numerical_cols + ['Gender', 'Subscription Type', 'Contract_
    ↪Length', 'Churn']

# Compute the correlation matrix
correlation_matrix_binary = encoded_data_binary[all_columns_binary].corr()

# Plot the heatmap for correlations
plt.figure(figsize=(14, 10))
sns.heatmap(correlation_matrix_binary, annot=True, fmt=".2f", cmap="coolwarm",
    ↪linewidths=0.5)
plt.title("Correlation Matrix with Binary/Ordinal Encoded Categorical_
    ↪Variables", fontsize=16)
plt.show()

# Return the correlation matrix for review
correlation_matrix_binary
```



[14]:

	Age	Tenure	Usage Frequency	Support Calls	\
Age	1.000000	-0.011630	-0.007190	0.158451	
Tenure	-0.011630	1.000000	-0.026800	-0.027640	
Usage Frequency	-0.007190	-0.026800	1.000000	-0.022013	
Support Calls	0.158451	-0.027640	-0.022013	1.000000	
Payment Delay	0.061738	-0.016588	-0.014470	0.162889	
Total Spend	-0.084684	0.019006	0.018631	-0.221594	
Last Interaction	0.028980	-0.006903	-0.004662	0.077684	
Gender	-0.031419	0.007978	0.007978	-0.091212	
Subscription Type	-0.004414	-0.024657	0.000032	-0.009820	
Contract Length	-0.069303	0.016925	0.014058	-0.178806	
Churn	0.218394	-0.051919	-0.046101	0.574267	

	Payment Delay	Total Spend	Last Interaction	Gender	\
Age	0.061738	-0.084684	0.028980	-0.031419	
Tenure	-0.016588	0.019006	-0.006903	0.007978	
Usage Frequency	-0.014470	0.018631	-0.004662	0.007978	
Support Calls	0.162889	-0.221594	0.077684	-0.091212	

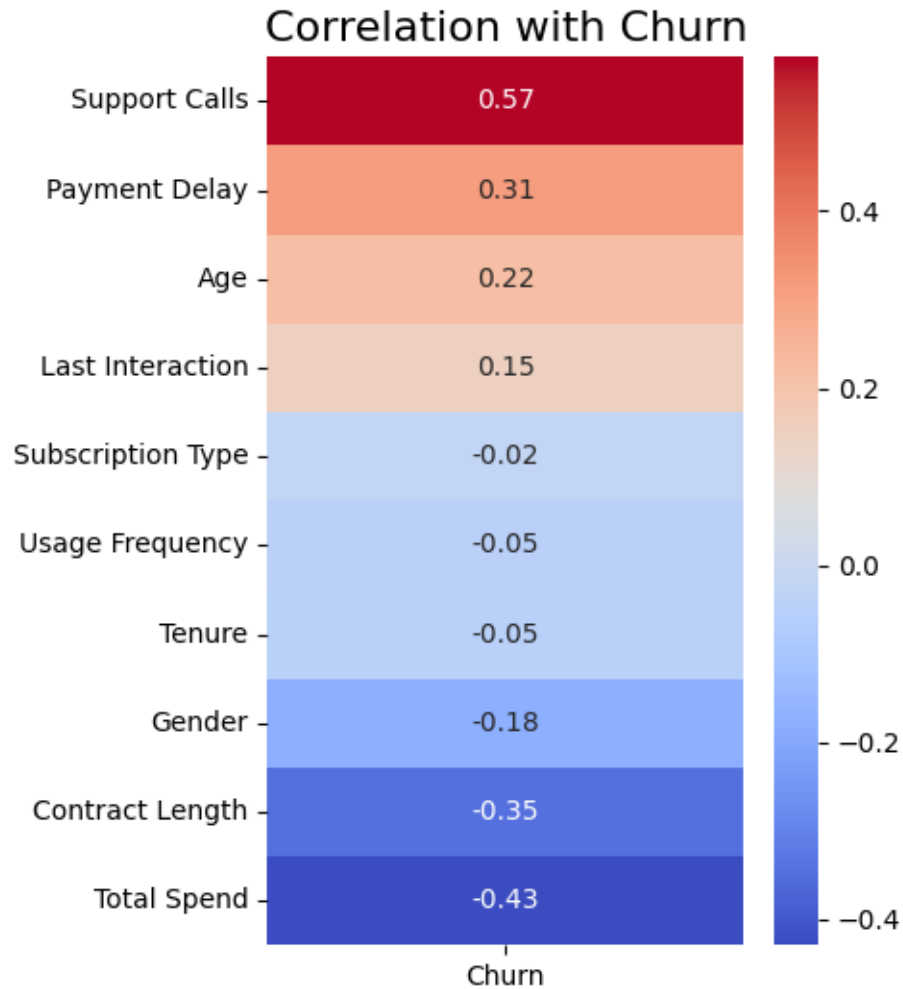
Payment Delay	1.000000	-0.121044	0.042708	-0.048449
Total Spend	-0.121044	1.000000	-0.056890	0.066138
Last Interaction	0.042708	-0.056890	1.000000	0.134786
Gender	-0.048449	0.066138	0.134786	1.000000
Subscription Type	-0.005385	0.007131	-0.000608	0.002271
Contract Length	-0.098028	0.133871	-0.046014	0.055342
Churn	0.312129	-0.429355	0.149616	-0.175395

	Subscription Type	Contract Length	Churn
Age	-0.004414	-0.069303	0.218394
Tenure	-0.024657	0.016925	-0.051919
Usage Frequency	0.000032	0.014058	-0.046101
Support Calls	-0.009820	-0.178806	0.574267
Payment Delay	-0.005385	-0.098028	0.312129
Total Spend	0.007131	0.133871	-0.429355
Last Interaction	-0.000608	-0.046014	0.149616
Gender	0.002271	0.055342	-0.175395
Subscription Type	1.000000	0.004850	-0.018251
Contract Length	0.004850	1.000000	-0.346666
Churn	-0.018251	-0.346666	1.000000

```
[15]: # Calculate correlation of all features with the target (Churn)
target_corr = correlation_matrix_binary['Churn'].drop('Churn').
    ↪sort_values(ascending=False)

# Plot the correlation of each variable with Churn
plt.figure(figsize=(4, 6))
sns.heatmap(target_corr.to_frame(), annot=True, cmap='coolwarm', fmt=".2f",
    ↪cbar=True)
plt.title("Correlation with Churn", fontsize=16)
plt.show()

# Display the correlations for reference
target_corr
```



```
[15]: Support Calls      0.574267
      Payment Delay    0.312129
      Age              0.218394
      Last Interaction  0.149616
      Subscription Type -0.018251
      Usage Frequency  -0.046101
      Tenure           -0.051919
      Gender           -0.175395
      Contract Length  -0.346666
      Total Spend      -0.429355
      Name: Churn, dtype: float64
```

Note

Although some features show less correlation with target variable, we prevent from feature selection as our dataset consists with only 11 feature variables, so it may cause data loss which will lead to inaccurate predictions.

3.7 Handling Categorical Variables

Note

For columns with categorical data types such as Gender, Subscription Type, Contract Length must preprocess into numerical values to be able to include them in the model. We decided to encode with **Label Encoding** due to following reasons.

- By analysing the **Gender** feature we came to a conclusion that female customers are more likely to leave, so we use 'Male': 1 and 'Female': 0.
- In **Subscription Type** feature 'Basic' package has a higher churn rate than 'Standard' and 'Premium', so higher-end package retains customers better. This concludes to use 'Basic': 0, 'Standard': 1 and 'Premium': 2.
- By considering **Contract Length** feature we had a conclusion that specifically, most customers using the "Monthly" package will not continue to use the service. Meanwhile, long-term contracts help retain customers more effectively. So we use 'Monthly': 0, 'Quarterly': 1 and 'Annual': 2.

3.7.1 Normalize categorical data into numerical format using Label Encoding

```
[16]: # Encode categorical variables to numeric values
encoded_dataset = dataset.copy()
encoded_dataset['Gender'] = encoded_dataset['Gender'].apply(lambda x: 1 if x == 'Male' else 0)

subscription_map = {'Basic': 0, 'Standard': 1, 'Premium': 2}
encoded_dataset['Subscription Type'] = encoded_dataset['Subscription Type'].map(subscription_map)

contract_map = {'Monthly': 0, 'Quarterly': 1, 'Annual': 2}
encoded_dataset['Contract Length'] = encoded_dataset['Contract Length'].map(contract_map)
```

```
[17]: encoded_dataset
```

```
[17]:
```

	CustomerID	Age	Gender	Tenure	Usage Frequency	Support Calls	\
0	2.0	30.0	0	39.0	14.0	5.0	
1	3.0	65.0	0	49.0	1.0	10.0	
2	4.0	55.0	0	14.0	4.0	6.0	
3	5.0	58.0	1	38.0	21.0	7.0	
4	6.0	23.0	1	32.0	20.0	5.0	
...	
440828	449995.0	42.0	1	54.0	15.0	1.0	
440829	449996.0	25.0	0	8.0	13.0	1.0	
440830	449997.0	26.0	1	35.0	27.0	1.0	
440831	449998.0	28.0	1	55.0	14.0	2.0	
440832	449999.0	31.0	1	48.0	20.0	1.0	

	Payment Delay	Subscription Type	Contract Length	Total Spend \
0	18.0	1	2	932.00
1	8.0	0	0	557.00
2	18.0	0	1	185.00
3	7.0	1	0	396.00
4	8.0	0	0	617.00
...
440828	3.0	2	2	716.38
440829	20.0	2	2	745.38
440830	5.0	1	1	977.31
440831	0.0	1	1	602.55
440832	14.0	2	1	567.77

	Last Interaction	Churn
0	17.0	1.0
1	6.0	1.0
2	3.0	1.0
3	29.0	1.0
4	20.0	1.0
...
440828	8.0	0.0
440829	2.0	0.0
440830	9.0	0.0
440831	2.0	0.0
440832	21.0	0.0

[440832 rows x 12 columns]

3.8 Prepare features (X) - labels (y)

```
[18]: X = encoded_dataset.drop('Churn', axis=1)
      y = encoded_dataset['Churn']
```

3.9 Train-test split: training set (90%) and test set (10%)

```
[19]: X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.1,
      ↪ random_state=42, stratify=y)
```

```
[20]: X_train.shape, X_test.shape, y_train.shape, y_test.shape
```

```
[20]: ((396748, 11), (44084, 11), (396748,), (44084,))
```

3.10 Feature Scaling

```
[21]: # Standardize numerical features
scaler = StandardScaler()
X_train_scaled = scaler.fit_transform(X_train)
X_test_scaled = scaler.transform(X_test)
```

Note

Data is not imbalanced (~ 45% churn=0 and ~ 55% churn=1) -> It is not necessary to balance data using techniques like SMOTE, Undersampling. So we can directly move with model implementation.

4 Model Implementation

4.1 Random Forest

```
[22]: # Random Forest Model Training and Evaluation
print("Processing: Random Forest...")

# Define Random Forest and parameter grid
rf_model = RandomForestClassifier(random_state=42)
rf_param_grid = {
    'n_estimators': [50, 100],
    'max_depth': [5, 10],
    'min_samples_split': [5, 10],
    'min_samples_leaf': [1, 2]
}

# Hyperparameter tuning with GridSearchCV
rf_grid = GridSearchCV(
    estimator=rf_model,
    param_grid=rf_param_grid,
    cv=2,
    scoring='accuracy',
    verbose=2,
    n_jobs=-1
)

# Measure training time
start_time = time.time()
rf_grid.fit(X_train_scaled, y_train)
train_time = time.time() - start_time

# Best model from grid search
rf_best_model = rf_grid.best_estimator_
print(f"Best parameters for Random Forest:", rf_grid.best_params_)
print(f"Best cross-validation score for Random Forest: {rf_grid.best_score_}")
```

Processing: Random Forest...

Fitting 2 folds for each of 16 candidates, totalling 32 fits

/usr/local/lib/python3.10/dist-packages/numpy/ma/core.py:2820: RuntimeWarning:
invalid value encountered in cast

```
_data = np.array(data, dtype=dtype, copy=copy,
```

Best parameters for Random Forest: {'max_depth': 10, 'min_samples_leaf': 1,
'min_samples_split': 5, 'n_estimators': 50}

Best cross-validation score for Random Forest: 0.9962217830965752

Note - Hyperparameter Tuning

We initially experimented with a broader range of hyperparameter values and a higher cross-validation (CV) setting to thoroughly explore the parameter space. However, this approach significantly increased the training time. To ensure a more efficient yet effective training process, we decided to focus on two carefully chosen hyperparameter values for each hyperparameter and set CV to 2.

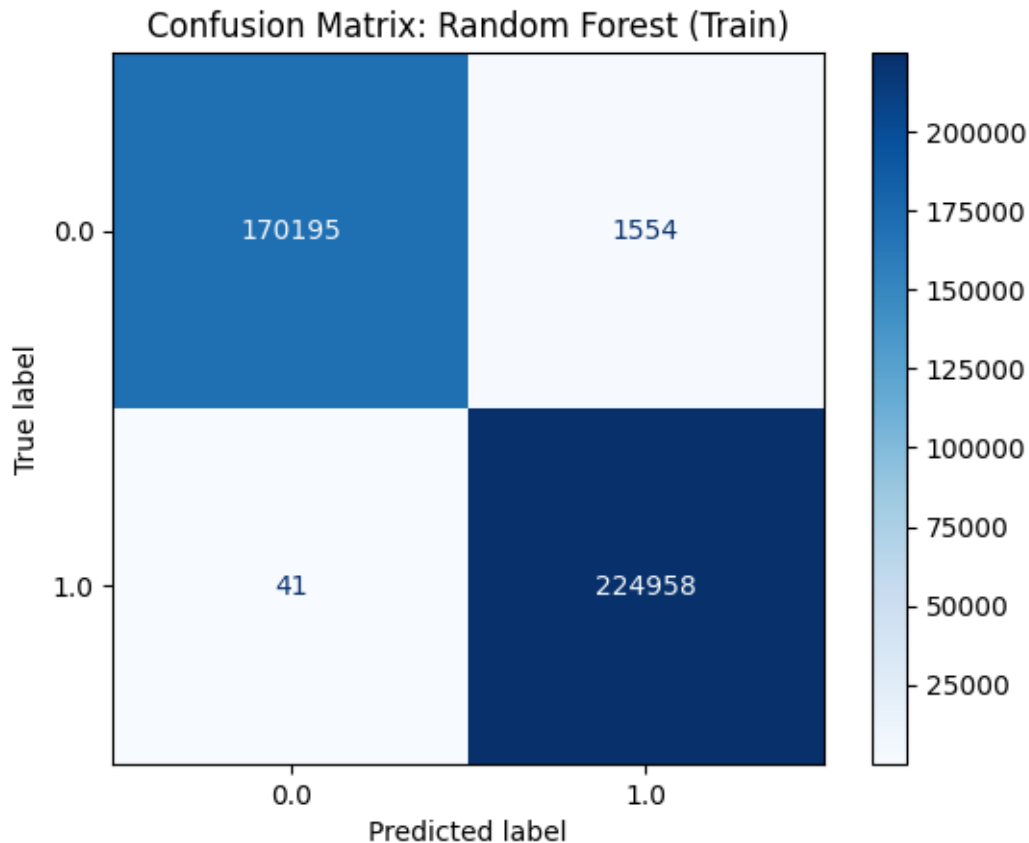
```
[23]: # Evaluate metrics for training dataset
rf_y_train_pred = rf_best_model.predict(X_train_scaled)
rf_train_metrics = {
    "Model": "Random Forest (Train)",
    "Accuracy": accuracy_score(y_train, rf_y_train_pred),
    "Precision": precision_score(y_train, rf_y_train_pred),
    "Recall": recall_score(y_train, rf_y_train_pred),
    "F1 Score": f1_score(y_train, rf_y_train_pred),
    "Training Time (s)": train_time
}

# Print training metrics
print("Training Dataset Metrics:")
for key, value in rf_train_metrics.items():
    print(f"\t{key}: {value}")
```

Training Dataset Metrics:

```
Model: Random Forest (Train)
Accuracy: 0.9959798159032938
Precision: 0.9931394363212545
Recall: 0.9998177769678976
F1 Score: 0.9964674171836345
Training Time (s): 326.1268081665039
```

```
[24]: # Confusion matrix for training dataset
cm_train = confusion_matrix(y_train, rf_y_train_pred)
disp_train = ConfusionMatrixDisplay(confusion_matrix=cm_train,
    ↪display_labels=rf_best_model.classes_)
disp_train.plot(cmap=plt.cm.Blues)
plt.title("Confusion Matrix: Random Forest (Train)")
plt.show()
```

```
[25]: # Measure prediction time
start_time = time.time()
rf_y_pred = rf_best_model.predict(X_test_scaled)
predict_time = time.time() - start_time

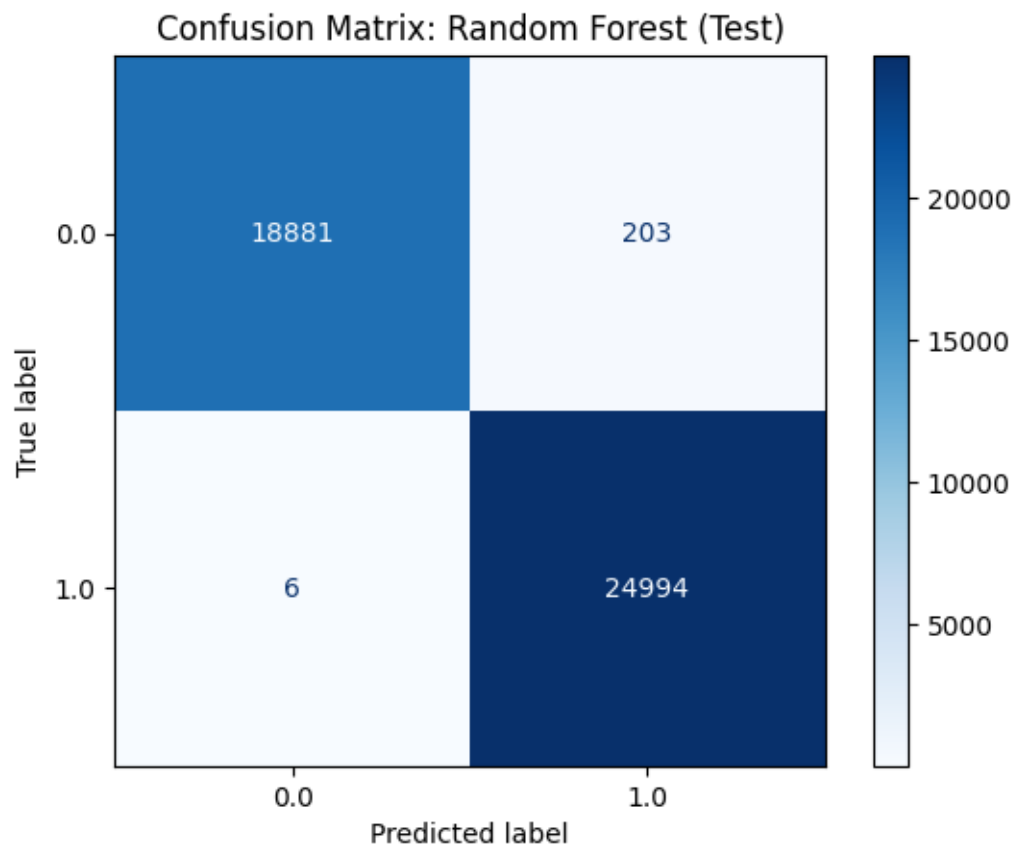
[26]: # Evaluate metrics for test dataset
rf_test_metrics = {
    "Model": "Random Forest (Test)",
    "Accuracy": accuracy_score(y_test, rf_y_pred),
    "Precision": precision_score(y_test, rf_y_pred),
    "Recall": recall_score(y_test, rf_y_pred),
    "F1 Score": f1_score(y_test, rf_y_pred),
    "Prediction Time (s)": predict_time
}

# Print test metrics
print("Test Dataset Metrics:")
for key, value in rf_test_metrics.items():
    print(f"\t{key}: {value}")
```

Test Dataset Metrics:

Model: Random Forest (Test)
Accuracy: 0.9952590509028219
Precision: 0.9919434853355559
Recall: 0.99976
F1 Score: 0.99583640456601
Prediction Time (s): 0.13321328163146973

```
[27]: # Confusion matrix for test dataset
cm_test = confusion_matrix(y_test, rf_y_pred)
disp_test = ConfusionMatrixDisplay(confusion_matrix=cm_test,
    ↪display_labels=rf_best_model.classes_)
disp_test.plot(cmap=plt.cm.Blues)
plt.title("Confusion Matrix: Random Forest (Test)")
plt.show()
```



```
[28]: # Collect ROC curve data for Random Forest
if hasattr(rf_best_model, "predict_proba"):
    # For Test Dataset
    rf_y_prob_test = rf_best_model.predict_proba(X_test_scaled)[: , 1]
```

```

fpr_rf_test, tpr_rf_test, _ = roc_curve(y_test, rf_y_prob_test)
roc_auc_rf_test = auc(fpr_rf_test, tpr_rf_test)
roc_data_rf_test = {"Random Forest (Test)": (fpr_rf_test, tpr_rf_test,
↪roc_auc_rf_test)}

# For Training Dataset
rf_y_prob_train = rf_best_model.predict_proba(X_train_scaled)[: , 1]
fpr_rf_train, tpr_rf_train, _ = roc_curve(y_train, rf_y_prob_train)
roc_auc_rf_train = auc(fpr_rf_train, tpr_rf_train)
roc_data_rf_train = {"Random Forest (Train)": (fpr_rf_train, tpr_rf_train,
↪roc_auc_rf_train)}

```

```

[29]: # Convert results to DataFrame
timed_results_rf_df = pd.DataFrame([rf_train_metrics, rf_test_metrics])

```

4.2 XGBoost

```

[30]: # XGBoost Model Training and Evaluation
print("Processing: XGBoost...")

# Define XGBoost and parameter grid
xgb_model = XGBClassifier(random_state=42, eval_metric='logloss')
xgb_param_grid = {
    'n_estimators': [50, 100],
    'max_depth': [2, 4],
    'learning_rate': [0.1, 0.2],
    'subsample': [0.8, 1],
    'colsample_bytree': [0.8, 1]
}

# Hyperparameter tuning with GridSearchCV
xgb_grid = GridSearchCV(
    estimator=xgb_model,
    param_grid=xgb_param_grid,
    cv=2,
    scoring='accuracy',
    verbose=2,
    n_jobs=-1
)

# Measure training time
start_time = time.time()
xgb_grid.fit(X_train_scaled, y_train)
train_time = time.time() - start_time

# Best model from grid search
xgb_best_model = xgb_grid.best_estimator_

```

```
print(f"Best parameters for XGBoost:", xgb_grid.best_params_)
print(f"Best cross-validation score for XGBoost: {xgb_grid.best_score_}")
```

Processing: XGBoost...

Fitting 2 folds for each of 32 candidates, totalling 64 fits

Best parameters for XGBoost: {'colsample_bytree': 1, 'learning_rate': 0.2, 'max_depth': 4, 'n_estimators': 100, 'subsample': 0.8}

Best cross-validation score for XGBoost: 0.9999571516428565

Note - Hyperparameter Tuning

We initially experimented with a broader range of hyperparameter values and a higher cross-validation (CV) setting to thoroughly explore the parameter space. However, this approach significantly increased the training time. To ensure a more efficient yet effective training process, we decided to focus on two carefully chosen hyperparameter values for each hyperparameter and set CV to 2.

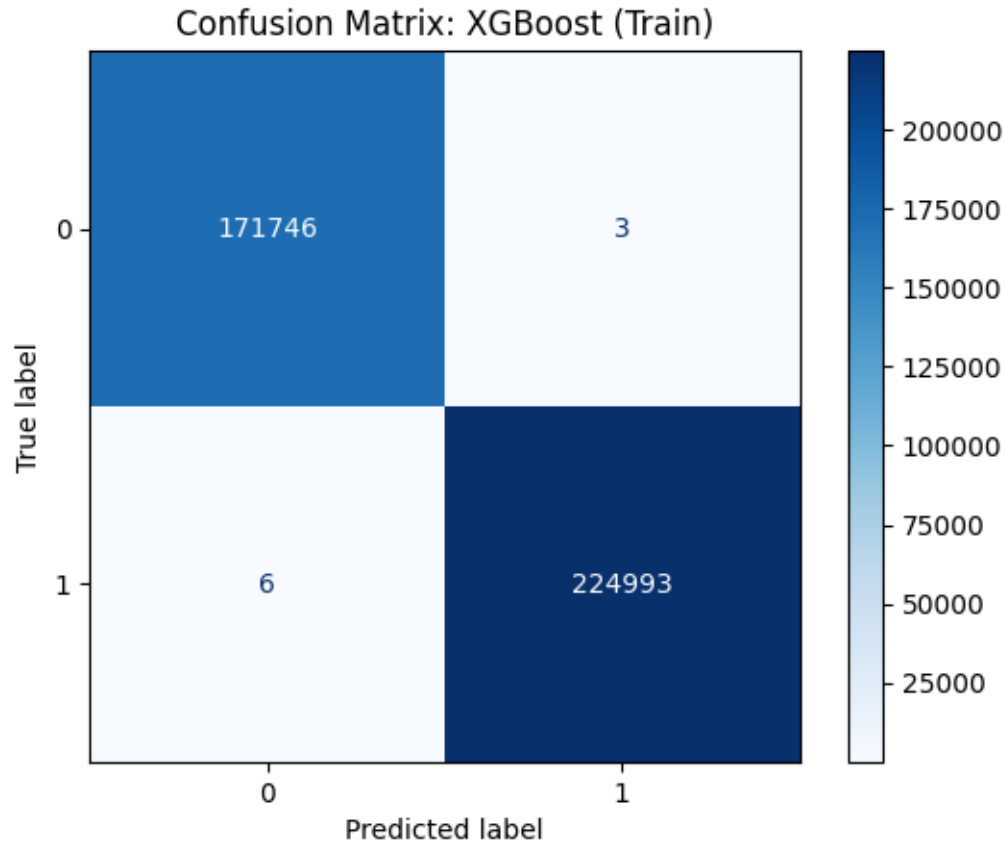
```
[31]: # Evaluate metrics for training dataset
xgb_y_train_pred = xgb_best_model.predict(X_train_scaled)
xgb_train_metrics = {
    "Model": "XGBoost (Train)",
    "Accuracy": accuracy_score(y_train, xgb_y_train_pred),
    "Precision": precision_score(y_train, xgb_y_train_pred),
    "Recall": recall_score(y_train, xgb_y_train_pred),
    "F1 Score": f1_score(y_train, xgb_y_train_pred),
    "Training Time (s)": train_time
}

# Print training metrics
print("Training Dataset Metrics:")
for key, value in xgb_train_metrics.items():
    print(f"\t{key}: {value}")
```

Training Dataset Metrics:

```
Model: XGBoost (Train)
Accuracy: 0.9999773155756299
Precision: 0.9999866664296254
Recall: 0.999973332148143
F1 Score: 0.9999799997777753
Training Time (s): 89.75960516929626
```

```
[32]: # Confusion matrix for training dataset
cm_train = confusion_matrix(y_train, xgb_y_train_pred)
disp_train = ConfusionMatrixDisplay(confusion_matrix=cm_train,
    ↪display_labels=xgb_best_model.classes_)
disp_train.plot(cmap=plt.cm.Blues)
plt.title("Confusion Matrix: XGBoost (Train)")
plt.show()
```



```
[33]: # Measure prediction time
start_time = time.time()
xgb_y_pred = xgb_best_model.predict(X_test_scaled)
predict_time = time.time() - start_time
```

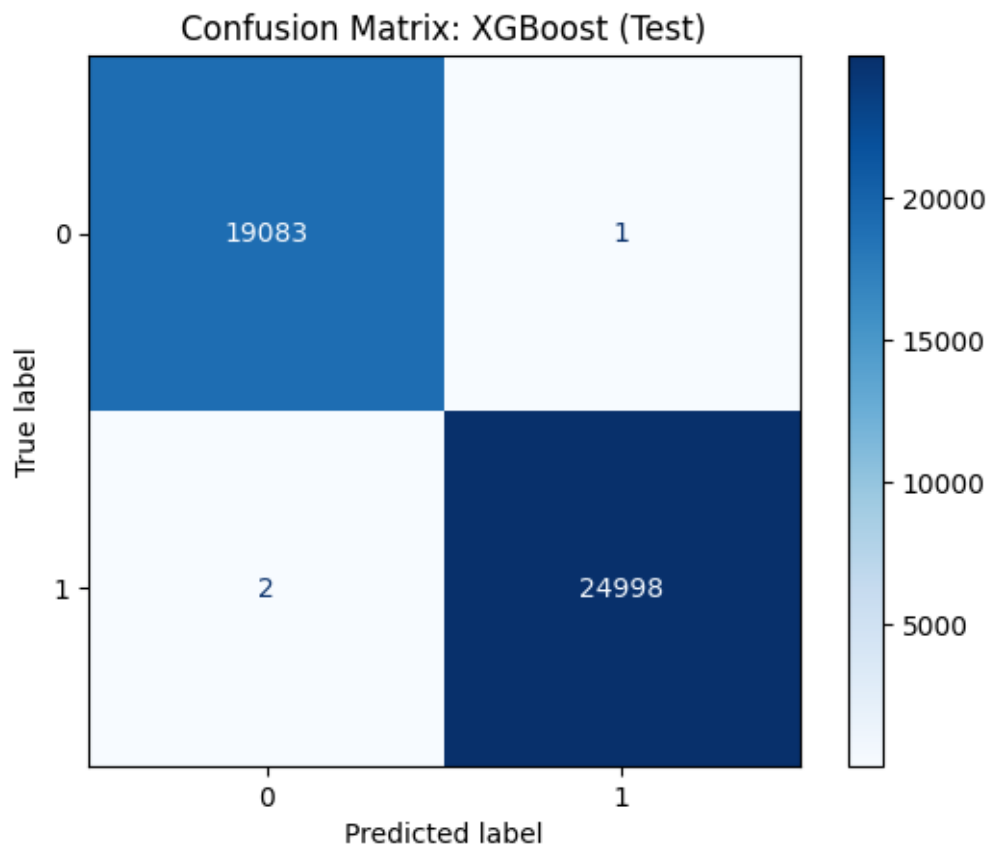
```
[34]: # Evaluate metrics for test dataset
xgb_test_metrics = {
    "Model": "XGBoost (Test)",
    "Accuracy": accuracy_score(y_test, xgb_y_pred),
    "Precision": precision_score(y_test, xgb_y_pred),
    "Recall": recall_score(y_test, xgb_y_pred),
    "F1 Score": f1_score(y_test, xgb_y_pred),
    "Prediction Time (s)": predict_time
}

# Print test metrics
print("Test Dataset Metrics:")
for key, value in xgb_test_metrics.items():
    print(f"\t{key}: {value}")
```

Test Dataset Metrics:

Model: XGBoost (Test)
Accuracy: 0.9999319480990836
Precision: 0.999959998399936
Recall: 0.99992
F1 Score: 0.999939998799976
Prediction Time (s): 0.08844947814941406

```
[35]: # Confusion matrix for test dataset
cm_test = confusion_matrix(y_test, xgb_y_pred)
disp_test = ConfusionMatrixDisplay(confusion_matrix=cm_test,
    ↪display_labels=xgb_best_model.classes_)
disp_test.plot(cmap=plt.cm.Blues)
plt.title("Confusion Matrix: XGBoost (Test)")
plt.show()
```



```
[36]: # Collect ROC curve data for XGBoost
if hasattr(xgb_best_model, "predict_proba"):
    # For Test Dataset
    xgb_y_prob_test = xgb_best_model.predict_proba(X_test_scaled)[: , 1]
```

```

fpr_xgb_test, tpr_xgb_test, _ = roc_curve(y_test, xgb_y_prob_test)
roc_auc_xgb_test = auc(fpr_xgb_test, tpr_xgb_test)
roc_data_xgb_test = {"XGBoost (Test)": (fpr_xgb_test, tpr_xgb_test,
↪roc_auc_xgb_test)}

# For Training Dataset
xgb_y_prob_train = xgb_best_model.predict_proba(X_train_scaled)[: , 1]
fpr_xgb_train, tpr_xgb_train, _ = roc_curve(y_train, xgb_y_prob_train)
roc_auc_xgb_train = auc(fpr_xgb_train, tpr_xgb_train)
roc_data_xgb_train = {"XGBoost (Train)": (fpr_xgb_train, tpr_xgb_train,
↪roc_auc_xgb_train)}

```

```

[37]: # Convert results to DataFrame
timed_results_xgb_df = pd.DataFrame([xgb_train_metrics, xgb_test_metrics])

```

5 Model Evaluation and Discussion

```

[38]: # Combine the timed results into a single DataFrame
timed_results_combined_df = pd.concat(
    [
        pd.DataFrame([rf_train_metrics, rf_test_metrics]),
        pd.DataFrame([xgb_train_metrics, xgb_test_metrics])
    ],
    ignore_index=True
)

# Display the combined DataFrame
timed_results_combined_df

```

```

[38]:

```

	Model	Accuracy	Precision	Recall	F1 Score	\
0	Random Forest (Train)	0.995980	0.993139	0.999818	0.996467	
1	Random Forest (Test)	0.995259	0.991943	0.999760	0.995836	
2	XGBoost (Train)	0.999977	0.999987	0.999973	0.999980	
3	XGBoost (Test)	0.999932	0.999960	0.999920	0.999940	

	Training Time (s)	Prediction Time (s)
0	326.126808	NaN
1	NaN	0.133213
2	89.759605	NaN
3	NaN	0.088449

```

[39]: # Plot ROC curves for Random Forest and XGBoost
plt.figure(figsize=(10, 8))

# Random Forest ROC curve
if 'roc_data_rf_train' in locals():

```

```

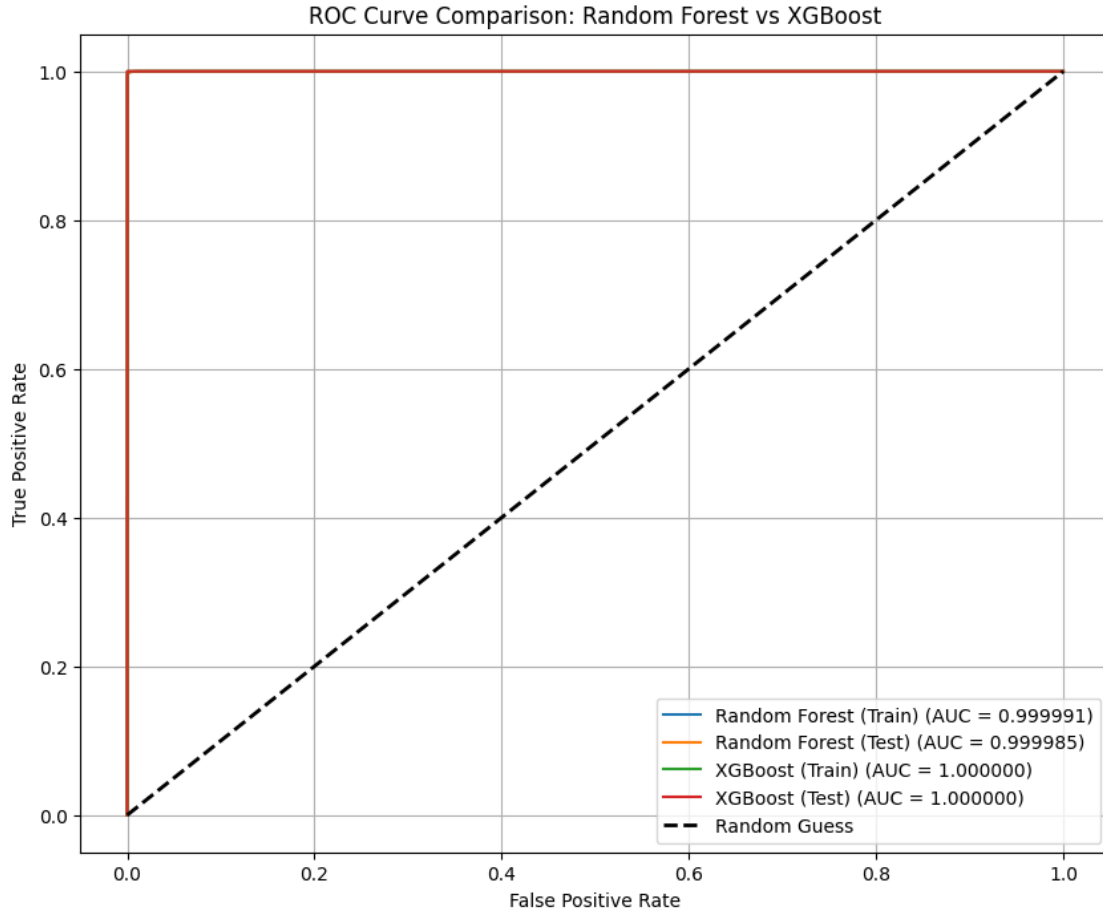
    for model_name, roc_values in (**roc_data_rf_train, **roc_data_rf_test).
        items():
            fpr, tpr, roc_auc = roc_values
            plt.plot(fpr, tpr, label=f"{model_name} (AUC = {roc_auc:.6f})")

# XGBoost ROC curve
if 'roc_data_xgb_train' in locals():
    for model_name, roc_values in (**roc_data_xgb_train, **roc_data_xgb_test).
        items():
            fpr, tpr, roc_auc = roc_values
            plt.plot(fpr, tpr, label=f"{model_name} (AUC = {roc_auc:.6f})")

# Random guess line
plt.plot([0, 1], [0, 1], 'k--', lw=2, label='Random Guess')

# Plot customization
plt.title("ROC Curve Comparison: Random Forest vs XGBoost")
plt.xlabel("False Positive Rate")
plt.ylabel("True Positive Rate")
plt.legend(loc="lower right")
plt.grid()
plt.show()

```

In this project, two machine learning models, Random Forest and XGBoost, were implemented to predict customer churn. Both models were evaluated using various performance metrics, including accuracy, precision, recall, F1-score, training time, prediction time, and the Receiver Operating Characteristic - Area Under Curve (ROC-AUC) score. The results are summarized below:

Model	Accuracy	Precision	Recall	F1 Score	ROC-AUC	Training Time (s)	Prediction Time (s)
Random Forest (Train)	0.995980	0.993139	0.999818	0.996467	0.999991	326.126808	NaN
Random Forest (Test)	0.995259	0.991943	0.999760	0.995836	0.999985	NaN	0.133213
XGBoost (Train)	0.999977	0.999987	0.999973	0.999980	1.000000	89.759605	NaN
XGBoost (Test)	0.999932	0.999960	0.999920	0.999940	1.000000	NaN	0.088449

5.1 Performance Analysis

1. Accuracy:

- Both models achieved near-perfect accuracy, with XGBoost slightly outperforming Random Forest (99.9932% vs. 99.5259%). This indicates both models are highly effective at correctly classifying customers.

2. Precision:

- XGBoost demonstrated superior precision at 99.996%, indicating a lower rate of false positives compared to Random Forest (99.1943%). This is particularly useful for minimizing unnecessary customer retention efforts.

3. Recall:

- XGBoost achieved a great recall of 99.992%, meaning it identified all churned customers correctly. Random Forest also performed exceptionally well with a recall of 99.976%. However, XGBoost's perfect recall ensures no churned customers are missed.

4. F1 Score:

- XGBoost had a slightly higher F1-score (99.994%) than Random Forest (99.5836%), reflecting a better overall balance between precision and recall.

5. ROC-AUC:

- XGBoost model achieved a perfect ROC-AUC score of 1.0, indicating flawless separation between churned and non-churned customers. This reinforces the strength of both algorithms in making confident predictions with no overlap in classification.

6. Training and Prediction Times:

- XGBoost was significantly faster in training (89.76 seconds) and prediction (0.088 seconds) compared to Random Forest, which required 326.13 seconds for training and 0.133 seconds for prediction. This makes XGBoost more suitable for scenarios requiring real-time or near real-time predictions.

5.2 Discussion

XGBoost performed exceptionally well, achieving perfect ROC-AUC scores, and demonstrated a clear advantage in terms of training and prediction efficiency. These attributes make it particularly valuable for large-scale applications where computational resources and time are critical. However, Random Forest's robustness, ease of implementation, and interpretability make it a compelling choice for scenarios where slightly lower performance can be traded for simplicity and feature importance insights.

6 Conclusion

This project successfully predicted customer churn using two machine learning algorithms, Random Forest and XGBoost. Both models delivered outstanding results, achieving near-perfect accuracy, precision, recall, F1-scores, while only XGBoost model delivered a flawless ROC-AUC score of 1.0. These metrics underscore the suitability of advanced machine learning models for churn prediction tasks.

Among the two algorithms, XGBoost emerged as the superior model in terms of computational efficiency and predictive performance:

1. It outperformed Random Forest in terms of accuracy, precision, recall, F1-score.

2. It demonstrated significantly faster training and prediction times, making it more practical for real-time applications.
3. Its perfect ROC-AUC score further validated its ability to distinguish between churned and non-churned customers with no misclassification.

Key Insights:

- XGBoost is ideal for large-scale, time-sensitive applications due to its computational efficiency and scalability.
- Random Forest, while slightly slower, remains a robust and interpretable alternative for understanding feature contributions to churn prediction.

In conclusion, this project highlights the effectiveness of machine learning in addressing critical business challenges like customer churn. Businesses can adopt XGBoost to implement high-performing churn prediction models, while leveraging insights from both algorithms to design proactive customer retention strategies. Future work could involve exploring additional algorithms, optimizing hyperparameters further, and addressing imbalanced datasets through advanced sampling techniques.

7 References

- [1] “ML | XGBoost (eXtreme Gradient Boosting),” GeeksforGeeks, Aug. 19, 2019. <https://www.geeksforgeeks.org/ml-xgboost-extreme-gradient-boosting/>
- [2] A. Dutta, “Random Forest Regression in Python - GeeksforGeeks,” GeeksforGeeks, Jun. 14, 2019. <https://www.geeksforgeeks.org/random-forest-regression-in-python/>
- [3] “A Guide on XGBoost hyperparameters tuning,” kaggle.com. <https://www.kaggle.com/code/prashant111/a-guide-on-xgboost-hyperparameters-tuning>
- [4] “Hyperparameter tuning in random forests,” kaggle.com. <https://www.kaggle.com/code/nargisbegum82/hyperparameter-tuning-in-random-forests>

8 Converting the .ipynb file to html file

```
[40]: !jupyter nbconvert --to html /content/EE5253Report_4466_4578_GP23.ipynb
```

```
[NbConvertApp] Converting notebook /content/EE5253Report_4466_4578_GP23.ipynb to  
html  
[NbConvertApp] WARNING | Alternative text is missing on 12 image(s).  
[NbConvertApp] Writing 1478890 bytes to  
/content/EE5253Report_4466_4578_GP23.html
```

9 Converting the .ipynb file to pfd file

```
[41]: !sudo apt-get install texlive-xetex texlive-fonts-recommended_
      ↪ texlive-plain-generic
```

Reading package lists... Done

Building dependency tree... Done

Reading state information... Done

The following additional packages will be installed:

dvisvgm fonts-droid-fallback fonts-lato fonts-lmodern fonts-noto-mono
fonts-texgyre fonts-urw-base35 libapache-pom-java libcommons-logging-java
libcommons-parent-java libfontbox-java libfontenc1 libgs9 libgs9-common
libidn12 libijs-0.35 libjbig2dec0 libkpathsea6 libpdfbox-java libptexenc1
libruby3.0 libsynchronet2 libteckit0 libtexlua53 libtexlua53 libwoff1
libzip-0-13 lmodern poppler-data preview-latex-style rake ruby
ruby-net-telnet ruby-rubygems ruby-webrick ruby-xmlrpc ruby3.0
rubygems-integration tlutils teckit tex-common tex-gyre texlive-base
texlive-binaries texlive-latex-base texlive-latex-extra
texlive-latex-recommended texlive-pictures tipa xfonts-encodings
xfonts-utils

Suggested packages:

fonts-noto fonts-freefont-otf | fonts-freefont-ttf libavalon-framework-java
libcommons-logging-java-doc libexcalibur-logkit-java liblog4j1.2-java
poppler-utils ghostscript fonts-japanese-mincho | fonts-ipafont-mincho
fonts-japanese-gothic | fonts-ipafont-gothic fonts-arphic-ukai
fonts-arphic-uming fonts-nanum ri ruby-dev bundler debhelper gv
| postscript-viewer perl-tk xpdf | pdf-viewer xzdec
texlive-fonts-recommended-doc texlive-latex-base-doc python3-pygments
icc-profiles libfile-which-perl libspreadsheet-parseexcel-perl
texlive-latex-extra-doc texlive-latex-recommended-doc texlive-luatex
texlive-pstricks dot2tex prerex texlive-pictures-doc vprerex
default-jre-headless tipa-doc

The following NEW packages will be installed:

dvisvgm fonts-droid-fallback fonts-lato fonts-lmodern fonts-noto-mono
fonts-texgyre fonts-urw-base35 libapache-pom-java libcommons-logging-java
libcommons-parent-java libfontbox-java libfontenc1 libgs9 libgs9-common
libidn12 libijs-0.35 libjbig2dec0 libkpathsea6 libpdfbox-java libptexenc1
libruby3.0 libsynchronet2 libteckit0 libtexlua53 libtexlua53 libwoff1
libzip-0-13 lmodern poppler-data preview-latex-style rake ruby
ruby-net-telnet ruby-rubygems ruby-webrick ruby-xmlrpc ruby3.0
rubygems-integration tlutils teckit tex-common tex-gyre texlive-base
texlive-binaries texlive-fonts-recommended texlive-latex-base
texlive-latex-extra texlive-latex-recommended texlive-pictures
texlive-plain-generic texlive-xetex tipa xfonts-encodings xfonts-utils

0 upgraded, 54 newly installed, 0 to remove and 49 not upgraded.

Need to get 182 MB of archives.

After this operation, 571 MB of additional disk space will be used.

Get:1 <http://archive.ubuntu.com/ubuntu> jammy/main amd64 fonts-droid-fallback all 1:6.0.1r16-1.1build1 [1,805 kB]
Get:2 <http://archive.ubuntu.com/ubuntu> jammy/main amd64 fonts-lato all 2.0-2.1 [2,696 kB]
Get:3 <http://archive.ubuntu.com/ubuntu> jammy/main amd64 poppler-data all 0.4.11-1 [2,171 kB]
Get:4 <http://archive.ubuntu.com/ubuntu> jammy/universe amd64 tex-common all 6.17 [33.7 kB]
Get:5 <http://archive.ubuntu.com/ubuntu> jammy/main amd64 fonts-urw-base35 all 20200910-1 [6,367 kB]
Get:6 <http://archive.ubuntu.com/ubuntu> jammy-updates/main amd64 libgs9-common all 9.55.0~dfsg1-0ubuntu5.10 [752 kB]
Get:7 <http://archive.ubuntu.com/ubuntu> jammy-updates/main amd64 libidn12 amd64 1.38-4ubuntu1 [60.0 kB]
Get:8 <http://archive.ubuntu.com/ubuntu> jammy/main amd64 libijs-0.35 amd64 0.35-15build2 [16.5 kB]
Get:9 <http://archive.ubuntu.com/ubuntu> jammy/main amd64 libjbig2dec0 amd64 0.19-3build2 [64.7 kB]
Get:10 <http://archive.ubuntu.com/ubuntu> jammy-updates/main amd64 libgs9 amd64 9.55.0~dfsg1-0ubuntu5.10 [5,031 kB]
Get:11 <http://archive.ubuntu.com/ubuntu> jammy-updates/main amd64 libkpathsea6 amd64 2021.20210626.59705-1ubuntu0.2 [60.4 kB]
Get:12 <http://archive.ubuntu.com/ubuntu> jammy/main amd64 libwoff1 amd64 1.0.2-1build4 [45.2 kB]
Get:13 <http://archive.ubuntu.com/ubuntu> jammy/universe amd64 dvisvgm amd64 2.13.1-1 [1,221 kB]
Get:14 <http://archive.ubuntu.com/ubuntu> jammy/universe amd64 fonts-lmodern all 2.004.5-6.1 [4,532 kB]
Get:15 <http://archive.ubuntu.com/ubuntu> jammy/main amd64 fonts-noto-mono all 20201225-1build1 [397 kB]
Get:16 <http://archive.ubuntu.com/ubuntu> jammy/universe amd64 fonts-texgyre all 20180621-3.1 [10.2 MB]
Get:17 <http://archive.ubuntu.com/ubuntu> jammy/universe amd64 libapache-pom-java all 18-1 [4,720 B]
Get:18 <http://archive.ubuntu.com/ubuntu> jammy/universe amd64 libcommons-parent-java all 43-1 [10.8 kB]
Get:19 <http://archive.ubuntu.com/ubuntu> jammy/universe amd64 libcommons-logging-java all 1.2-2 [60.3 kB]
Get:20 <http://archive.ubuntu.com/ubuntu> jammy/main amd64 libfontenc1 amd64 1:1.1.4-1build3 [14.7 kB]
Get:21 <http://archive.ubuntu.com/ubuntu> jammy-updates/main amd64 libptexenc1 amd64 2021.20210626.59705-1ubuntu0.2 [39.1 kB]
Get:22 <http://archive.ubuntu.com/ubuntu> jammy/main amd64 rubygems-integration all 1.18 [5,336 B]
Get:23 <http://archive.ubuntu.com/ubuntu> jammy-updates/main amd64 ruby3.0 amd64 3.0.2-7ubuntu2.8 [50.1 kB]
Get:24 <http://archive.ubuntu.com/ubuntu> jammy/main amd64 ruby-rubygems all 3.3.5-2 [228 kB]

Get:25 <http://archive.ubuntu.com/ubuntu> jammy/main amd64 ruby amd64 1:3.0~exp1
[5,100 B]
Get:26 <http://archive.ubuntu.com/ubuntu> jammy/main amd64 rake all 13.0.6-2 [61.7
kB]
Get:27 <http://archive.ubuntu.com/ubuntu> jammy/main amd64 ruby-net-telnet all
0.1.1-2 [12.6 kB]
Get:28 <http://archive.ubuntu.com/ubuntu> jammy-updates/main amd64 ruby-webrick
all 1.7.0-3ubuntu0.1 [52.1 kB]
Get:29 <http://archive.ubuntu.com/ubuntu> jammy-updates/main amd64 ruby-xmlrpc all
0.3.2-1ubuntu0.1 [24.9 kB]
Get:30 <http://archive.ubuntu.com/ubuntu> jammy-updates/main amd64 libruby3.0
amd64 3.0.2-7ubuntu2.8 [5,113 kB]
Get:31 <http://archive.ubuntu.com/ubuntu> jammy-updates/main amd64 libsyntax2
amd64 2021.20210626.59705-1ubuntu0.2 [55.6 kB]
Get:32 <http://archive.ubuntu.com/ubuntu> jammy/universe amd64 libteckit0 amd64
2.5.11+ds1-1 [421 kB]
Get:33 <http://archive.ubuntu.com/ubuntu> jammy-updates/main amd64 libtexlua53
amd64 2021.20210626.59705-1ubuntu0.2 [120 kB]
Get:34 <http://archive.ubuntu.com/ubuntu> jammy-updates/main amd64 libtexluajit2
amd64 2021.20210626.59705-1ubuntu0.2 [267 kB]
Get:35 <http://archive.ubuntu.com/ubuntu> jammy/universe amd64 libzip-0-13 amd64
0.13.72+dfsg.1-1.1 [27.0 kB]
Get:36 <http://archive.ubuntu.com/ubuntu> jammy/main amd64 xfonts-encodings all
1:1.0.5-0ubuntu2 [578 kB]
Get:37 <http://archive.ubuntu.com/ubuntu> jammy/main amd64 xfonts-utils amd64
1:7.7+6build2 [94.6 kB]
Get:38 <http://archive.ubuntu.com/ubuntu> jammy/universe amd64 lmodern all
2.004.5-6.1 [9,471 kB]
Get:39 <http://archive.ubuntu.com/ubuntu> jammy/universe amd64 preview-latex-style
all 12.2-1ubuntu1 [185 kB]
Get:40 <http://archive.ubuntu.com/ubuntu> jammy/main amd64 t1utils amd64
1.41-4build2 [61.3 kB]
Get:41 <http://archive.ubuntu.com/ubuntu> jammy/universe amd64 teckit amd64
2.5.11+ds1-1 [699 kB]
Get:42 <http://archive.ubuntu.com/ubuntu> jammy/universe amd64 tex-gyre all
20180621-3.1 [6,209 kB]
Get:43 <http://archive.ubuntu.com/ubuntu> jammy-updates/universe amd64 texlive-
binaries amd64 2021.20210626.59705-1ubuntu0.2 [9,860 kB]
Get:44 <http://archive.ubuntu.com/ubuntu> jammy/universe amd64 texlive-base all
2021.20220204-1 [21.0 MB]
Get:45 <http://archive.ubuntu.com/ubuntu> jammy/universe amd64 texlive-fonts-
recommended all 2021.20220204-1 [4,972 kB]
Get:46 <http://archive.ubuntu.com/ubuntu> jammy/universe amd64 texlive-latex-base
all 2021.20220204-1 [1,128 kB]
Get:47 <http://archive.ubuntu.com/ubuntu> jammy/universe amd64 libfontbox-java all
1:1.8.16-2 [207 kB]
Get:48 <http://archive.ubuntu.com/ubuntu> jammy/universe amd64 libpdfbox-java all
1:1.8.16-2 [5,199 kB]

```

Get:49 http://archive.ubuntu.com/ubuntu jammy/universe amd64 texlive-latex-
recommended all 2021.20220204-1 [14.4 MB]
Get:50 http://archive.ubuntu.com/ubuntu jammy/universe amd64 texlive-pictures
all 2021.20220204-1 [8,720 kB]
Get:51 http://archive.ubuntu.com/ubuntu jammy/universe amd64 texlive-latex-extra
all 2021.20220204-1 [13.9 MB]
Get:52 http://archive.ubuntu.com/ubuntu jammy/universe amd64 texlive-plain-
generic all 2021.20220204-1 [27.5 MB]
Get:53 http://archive.ubuntu.com/ubuntu jammy/universe amd64 tipa all 2:1.3-21
[2,967 kB]
Get:54 http://archive.ubuntu.com/ubuntu jammy/universe amd64 texlive-xetex all
2021.20220204-1 [12.4 MB]
Fetched 182 MB in 11s (16.6 MB/s)
debconf: unable to initialize frontend: Dialog
debconf: (No usable dialog-like program is installed, so the dialog based
frontend cannot be used. at /usr/share/perl5/Debconf/FrontEnd/Dialog.pm line 78,
<> line 54.)
debconf: falling back to frontend: Readline
debconf: unable to initialize frontend: Readline
debconf: (This frontend requires a controlling tty.)
debconf: falling back to frontend: Teletype
dpkg-preconfigure: unable to re-open stdin:
Selecting previously unselected package fonts-droid-fallback.
(Reading database ... 123633 files and directories currently installed.)
Preparing to unpack .../00-fonts-droid-fallback_1%3a6.0.1r16-1.1build1_all.deb
...
Unpacking fonts-droid-fallback (1:6.0.1r16-1.1build1) ...
Selecting previously unselected package fonts-lato.
Preparing to unpack .../01-fonts-lato_2.0-2.1_all.deb ...
Unpacking fonts-lato (2.0-2.1) ...
Selecting previously unselected package poppler-data.
Preparing to unpack .../02-poppler-data_0.4.11-1_all.deb ...
Unpacking poppler-data (0.4.11-1) ...
Selecting previously unselected package tex-common.
Preparing to unpack .../03-tex-common_6.17_all.deb ...
Unpacking tex-common (6.17) ...
Selecting previously unselected package fonts-urw-base35.
Preparing to unpack .../04-fonts-urw-base35_20200910-1_all.deb ...
Unpacking fonts-urw-base35 (20200910-1) ...
Selecting previously unselected package libgs9-common.
Preparing to unpack .../05-libgs9-common_9.55.0~dfsg1-0ubuntu5.10_all.deb ...
Unpacking libgs9-common (9.55.0~dfsg1-0ubuntu5.10) ...
Selecting previously unselected package libidn12:amd64.
Preparing to unpack .../06-libidn12_1.38-4ubuntu1_amd64.deb ...
Unpacking libidn12:amd64 (1.38-4ubuntu1) ...
Selecting previously unselected package libijs-0.35:amd64.
Preparing to unpack .../07-libijs-0.35_0.35-15build2_amd64.deb ...
Unpacking libijs-0.35:amd64 (0.35-15build2) ...

```

```

Selecting previously unselected package libjbig2dec0:amd64.
Preparing to unpack .../08-libjbig2dec0_0.19-3build2_amd64.deb ...
Unpacking libjbig2dec0:amd64 (0.19-3build2) ...
Selecting previously unselected package libgs9:amd64.
Preparing to unpack .../09-libgs9_9.55.0~dfsg1-0ubuntu5.10_amd64.deb ...
Unpacking libgs9:amd64 (9.55.0~dfsg1-0ubuntu5.10) ...
Selecting previously unselected package libkpathsea6:amd64.
Preparing to unpack .../10-libkpathsea6_2021.20210626.59705-1ubuntu0.2_amd64.deb
...
Unpacking libkpathsea6:amd64 (2021.20210626.59705-1ubuntu0.2) ...
Selecting previously unselected package libwoff1:amd64.
Preparing to unpack .../11-libwoff1_1.0.2-1build4_amd64.deb ...
Unpacking libwoff1:amd64 (1.0.2-1build4) ...
Selecting previously unselected package dvisvgm.
Preparing to unpack .../12-dvisvgm_2.13.1-1_amd64.deb ...
Unpacking dvisvgm (2.13.1-1) ...
Selecting previously unselected package fonts-lmodern.
Preparing to unpack .../13-fonts-lmodern_2.004.5-6.1_all.deb ...
Unpacking fonts-lmodern (2.004.5-6.1) ...
Selecting previously unselected package fonts-noto-mono.
Preparing to unpack .../14-fonts-noto-mono_20201225-1build1_all.deb ...
Unpacking fonts-noto-mono (20201225-1build1) ...
Selecting previously unselected package fonts-texgyre.
Preparing to unpack .../15-fonts-texgyre_20180621-3.1_all.deb ...
Unpacking fonts-texgyre (20180621-3.1) ...
Selecting previously unselected package libapache-pom-java.
Preparing to unpack .../16-libapache-pom-java_18-1_all.deb ...
Unpacking libapache-pom-java (18-1) ...
Selecting previously unselected package libcommons-parent-java.
Preparing to unpack .../17-libcommons-parent-java_43-1_all.deb ...
Unpacking libcommons-parent-java (43-1) ...
Selecting previously unselected package libcommons-logging-java.
Preparing to unpack .../18-libcommons-logging-java_1.2-2_all.deb ...
Unpacking libcommons-logging-java (1.2-2) ...
Selecting previously unselected package libfontenc1:amd64.
Preparing to unpack .../19-libfontenc1_1%3a1.1.4-1build3_amd64.deb ...
Unpacking libfontenc1:amd64 (1:1.1.4-1build3) ...
Selecting previously unselected package libptexenc1:amd64.
Preparing to unpack .../20-libptexenc1_2021.20210626.59705-1ubuntu0.2_amd64.deb
...
Unpacking libptexenc1:amd64 (2021.20210626.59705-1ubuntu0.2) ...
Selecting previously unselected package rubygems-integration.
Preparing to unpack .../21-rubygems-integration_1.18_all.deb ...
Unpacking rubygems-integration (1.18) ...
Selecting previously unselected package ruby3.0.
Preparing to unpack .../22-ruby3.0_3.0.2-7ubuntu2.8_amd64.deb ...
Unpacking ruby3.0 (3.0.2-7ubuntu2.8) ...
Selecting previously unselected package ruby-rubygems.

```



```

Preparing to unpack .../23-ruby-rubygems_3.3.5-2_all.deb ...
Unpacking ruby-rubygems (3.3.5-2) ...
Selecting previously unselected package ruby.
Preparing to unpack .../24-ruby_1%3a3.0~exp1_amd64.deb ...
Unpacking ruby (1:3.0~exp1) ...
Selecting previously unselected package rake.
Preparing to unpack .../25-rake_13.0.6-2_all.deb ...
Unpacking rake (13.0.6-2) ...
Selecting previously unselected package ruby-net-telnet.
Preparing to unpack .../26-ruby-net-telnet_0.1.1-2_all.deb ...
Unpacking ruby-net-telnet (0.1.1-2) ...
Selecting previously unselected package ruby-webrick.
Preparing to unpack .../27-ruby-webrick_1.7.0-3ubuntu0.1_all.deb ...
Unpacking ruby-webrick (1.7.0-3ubuntu0.1) ...
Selecting previously unselected package ruby-xmlrpc.
Preparing to unpack .../28-ruby-xmlrpc_0.3.2-1ubuntu0.1_all.deb ...
Unpacking ruby-xmlrpc (0.3.2-1ubuntu0.1) ...
Selecting previously unselected package libruby3.0:amd64.
Preparing to unpack .../29-libruby3.0_3.0.2-7ubuntu2.8_amd64.deb ...
Unpacking libruby3.0:amd64 (3.0.2-7ubuntu2.8) ...
Selecting previously unselected package libsyntax2:amd64.
Preparing to unpack .../30-libsyntax2_2021.20210626.59705-1ubuntu0.2_amd64.deb
...
Unpacking libsyntax2:amd64 (2021.20210626.59705-1ubuntu0.2) ...
Selecting previously unselected package libteckit0:amd64.
Preparing to unpack .../31-libteckit0_2.5.11+ds1-1_amd64.deb ...
Unpacking libteckit0:amd64 (2.5.11+ds1-1) ...
Selecting previously unselected package libtexlua53:amd64.
Preparing to unpack .../32-libtexlua53_2021.20210626.59705-1ubuntu0.2_amd64.deb
...
Unpacking libtexlua53:amd64 (2021.20210626.59705-1ubuntu0.2) ...
Selecting previously unselected package libtexluajit2:amd64.
Preparing to unpack
.../33-libtexluajit2_2021.20210626.59705-1ubuntu0.2_amd64.deb ...
Unpacking libtexluajit2:amd64 (2021.20210626.59705-1ubuntu0.2) ...
Selecting previously unselected package libzip-0-13:amd64.
Preparing to unpack .../34-libzip-0-13_0.13.72+dfsg.1-1.1_amd64.deb ...
Unpacking libzip-0-13:amd64 (0.13.72+dfsg.1-1.1) ...
Selecting previously unselected package xfonts-encodings.
Preparing to unpack .../35-xfonts-encodings_1%3a1.0.5-0ubuntu2_all.deb ...
Unpacking xfonts-encodings (1:1.0.5-0ubuntu2) ...
Selecting previously unselected package xfonts-utils.
Preparing to unpack .../36-xfonts-utils_1%3a7.7+6build2_amd64.deb ...
Unpacking xfonts-utils (1:7.7+6build2) ...
Selecting previously unselected package lmodern.
Preparing to unpack .../37-lmodern_2.004.5-6.1_all.deb ...
Unpacking lmodern (2.004.5-6.1) ...
Selecting previously unselected package preview-latex-style.

```

```

Preparing to unpack .../38-preview-latex-style_12.2-1ubuntu1_all.deb ...
Unpacking preview-latex-style (12.2-1ubuntu1) ...
Selecting previously unselected package t1utils.
Preparing to unpack .../39-t1utils_1.41-4build2_amd64.deb ...
Unpacking t1utils (1.41-4build2) ...
Selecting previously unselected package teckit.
Preparing to unpack .../40-teckit_2.5.11+ds1-1_amd64.deb ...
Unpacking teckit (2.5.11+ds1-1) ...
Selecting previously unselected package tex-gyre.
Preparing to unpack .../41-tex-gyre_20180621-3.1_all.deb ...
Unpacking tex-gyre (20180621-3.1) ...
Selecting previously unselected package texlive-binaries.
Preparing to unpack .../42-texlive-
binaries_2021.20210626.59705-1ubuntu0.2_amd64.deb ...
Unpacking texlive-binaries (2021.20210626.59705-1ubuntu0.2) ...
Selecting previously unselected package texlive-base.
Preparing to unpack .../43-texlive-base_2021.20220204-1_all.deb ...
Unpacking texlive-base (2021.20220204-1) ...
Selecting previously unselected package texlive-fonts-recommended.
Preparing to unpack .../44-texlive-fonts-recommended_2021.20220204-1_all.deb ...
Unpacking texlive-fonts-recommended (2021.20220204-1) ...
Selecting previously unselected package texlive-latex-base.
Preparing to unpack .../45-texlive-latex-base_2021.20220204-1_all.deb ...
Unpacking texlive-latex-base (2021.20220204-1) ...
Selecting previously unselected package libfontbox-java.
Preparing to unpack .../46-libfontbox-java_1%3a1.8.16-2_all.deb ...
Unpacking libfontbox-java (1:1.8.16-2) ...
Selecting previously unselected package libpdfbox-java.
Preparing to unpack .../47-libpdfbox-java_1%3a1.8.16-2_all.deb ...
Unpacking libpdfbox-java (1:1.8.16-2) ...
Selecting previously unselected package texlive-latex-recommended.
Preparing to unpack .../48-texlive-latex-recommended_2021.20220204-1_all.deb ...
Unpacking texlive-latex-recommended (2021.20220204-1) ...
Selecting previously unselected package texlive-pictures.
Preparing to unpack .../49-texlive-pictures_2021.20220204-1_all.deb ...
Unpacking texlive-pictures (2021.20220204-1) ...
Selecting previously unselected package texlive-latex-extra.
Preparing to unpack .../50-texlive-latex-extra_2021.20220204-1_all.deb ...
Unpacking texlive-latex-extra (2021.20220204-1) ...
Selecting previously unselected package texlive-plain-generic.
Preparing to unpack .../51-texlive-plain-generic_2021.20220204-1_all.deb ...
Unpacking texlive-plain-generic (2021.20220204-1) ...
Selecting previously unselected package tipa.
Preparing to unpack .../52-tipa_2%3a1.3-21_all.deb ...
Unpacking tipa (2:1.3-21) ...
Selecting previously unselected package texlive-xetex.
Preparing to unpack .../53-texlive-xetex_2021.20220204-1_all.deb ...
Unpacking texlive-xetex (2021.20220204-1) ...

```

```

Setting up fonts-lato (2.0-2.1) ...
Setting up fonts-noto-mono (20201225-1build1) ...
Setting up libwoff1:amd64 (1.0.2-1build4) ...
Setting up libtexlua53:amd64 (2021.20210626.59705-1ubuntu0.2) ...
Setting up libijs-0.35:amd64 (0.35-15build2) ...
Setting up libtexluajit2:amd64 (2021.20210626.59705-1ubuntu0.2) ...
Setting up libfontbox-java (1:1.8.16-2) ...
Setting up rubygems-integration (1.18) ...
Setting up libzip-0-13:amd64 (0.13.72+dfsg.1-1.1) ...
Setting up fonts-urw-base35 (20200910-1) ...
Setting up poppler-data (0.4.11-1) ...
Setting up tex-common (6.17) ...
debconf: unable to initialize frontend: Dialog
debconf: (No usable dialog-like program is installed, so the dialog based
frontend cannot be used. at /usr/share/perl5/Debconf/FrontEnd/Dialog.pm line
78.)
debconf: falling back to frontend: Readline
update-language: texlive-base not installed and configured, doing nothing!
Setting up libfontenc1:amd64 (1:1.1.4-1build3) ...
Setting up libjbig2dec0:amd64 (0.19-3build2) ...
Setting up libteckit0:amd64 (2.5.11+ds1-1) ...
Setting up libapache-pom-java (18-1) ...
Setting up ruby-net-telnet (0.1.1-2) ...
Setting up xfonts-encodings (1:1.0.5-0ubuntu2) ...
Setting up t1utils (1.41-4build2) ...
Setting up libidn12:amd64 (1.38-4ubuntu1) ...
Setting up fonts-texgyre (20180621-3.1) ...
Setting up libkpathsea6:amd64 (2021.20210626.59705-1ubuntu0.2) ...
Setting up ruby-webrick (1.7.0-3ubuntu0.1) ...
Setting up fonts-lmodern (2.004.5-6.1) ...
Setting up fonts-droid-fallback (1:6.0.1r16-1.1build1) ...
Setting up ruby-xmlrpc (0.3.2-1ubuntu0.1) ...
Setting up libsynchronet2:amd64 (2021.20210626.59705-1ubuntu0.2) ...
Setting up libgs9-common (9.55.0~dfsg1-0ubuntu5.10) ...
Setting up teckit (2.5.11+ds1-1) ...
Setting up libpdfbox-java (1:1.8.16-2) ...
Setting up libgs9:amd64 (9.55.0~dfsg1-0ubuntu5.10) ...
Setting up preview-latex-style (12.2-1ubuntu1) ...
Setting up libcommons-parent-java (43-1) ...
Setting up dvisvgm (2.13.1-1) ...
Setting up libcommons-logging-java (1.2-2) ...
Setting up xfonts-utils (1:7.7+6build2) ...
Setting up libptexenc1:amd64 (2021.20210626.59705-1ubuntu0.2) ...
Setting up texlive-binaries (2021.20210626.59705-1ubuntu0.2) ...
update-alternatives: using /usr/bin/xdvi-xaw to provide /usr/bin/xdvi.bin
(xdvi.bin) in auto mode
update-alternatives: using /usr/bin/bibtex.original to provide /usr/bin/bibtex
(bibtex) in auto mode

```

```

Setting up lmodern (2.004.5-6.1) ...
Setting up texlive-base (2021.20220204-1) ...
/usr/bin/ucfr
/usr/bin/ucfr
/usr/bin/ucfr
/usr/bin/ucfr
mktexlsr: Updating /var/lib/texmf/ls-R-TEXLIVEDIST...
mktexlsr: Updating /var/lib/texmf/ls-R-TEXMFMAIN...
mktexlsr: Updating /var/lib/texmf/ls-R...
mktexlsr: Done.
tl-paper: setting paper size for dvips to a4:
/var/lib/texmf/dvips/config/config-paper.ps
tl-paper: setting paper size for dvipdfmx to a4:
/var/lib/texmf/dvipdfmx/dvipdfmx-paper.cfg
tl-paper: setting paper size for xdvi to a4: /var/lib/texmf/xdvi/XDvi-paper
tl-paper: setting paper size for pdftex to a4: /var/lib/texmf/tex/generic/tex-
ini-files/pdftexconfig.tex
debconf: unable to initialize frontend: Dialog
debconf: (No usable dialog-like program is installed, so the dialog based
frontend cannot be used. at /usr/share/perl5/Debconf/FrontEnd/Dialog.pm line
78.)
debconf: falling back to frontend: Readline
Setting up tex-gyre (20180621-3.1) ...
Setting up texlive-plain-generic (2021.20220204-1) ...
Setting up texlive-latex-base (2021.20220204-1) ...
Setting up texlive-latex-recommended (2021.20220204-1) ...
Setting up texlive-pictures (2021.20220204-1) ...
Setting up texlive-fonts-recommended (2021.20220204-1) ...
Setting up tipa (2:1.3-21) ...
Setting up texlive-latex-extra (2021.20220204-1) ...
Setting up texlive-xetex (2021.20220204-1) ...
Setting up rake (13.0.6-2) ...
Setting up libruby3.0:amd64 (3.0.2-7ubuntu2.8) ...
Setting up ruby3.0 (3.0.2-7ubuntu2.8) ...
Setting up ruby (1:3.0~exp1) ...
Setting up ruby-rubygems (3.3.5-2) ...
Processing triggers for man-db (2.10.2-1) ...
Processing triggers for fontconfig (2.13.1-4.2ubuntu5) ...
Processing triggers for libc-bin (2.35-0ubuntu3.4) ...
/sbin/ldconfig.real: /usr/local/lib/libhwloc.so.15 is not a symbolic link

/sbin/ldconfig.real: /usr/local/lib/libur_adapter_level_zero.so.0 is not a
symbolic link

/sbin/ldconfig.real: /usr/local/lib/libtbbmalloc.so.2 is not a symbolic link

/sbin/ldconfig.real: /usr/local/lib/libur_loader.so.0 is not a symbolic link

```

```
/sbin/ldconfig.real: /usr/local/lib/libtbbmalloc_proxy.so.2 is not a symbolic
link

/sbin/ldconfig.real: /usr/local/lib/libtbbbind_2_5.so.3 is not a symbolic link

/sbin/ldconfig.real: /usr/local/lib/libumf.so.0 is not a symbolic link

/sbin/ldconfig.real: /usr/local/lib/libtcm.so.1 is not a symbolic link

/sbin/ldconfig.real: /usr/local/lib/libtcm_debug.so.1 is not a symbolic link

/sbin/ldconfig.real: /usr/local/lib/libtbb.so.12 is not a symbolic link

/sbin/ldconfig.real: /usr/local/lib/libtbbbind_2_0.so.3 is not a symbolic link

/sbin/ldconfig.real: /usr/local/lib/libtbbbind.so.3 is not a symbolic link

/sbin/ldconfig.real: /usr/local/lib/libur_adapter_opencl.so.0 is not a symbolic
link
```

```
Processing triggers for tex-common (6.17) ...
debconf: unable to initialize frontend: Dialog
debconf: (No usable dialog-like program is installed, so the dialog based
frontend cannot be used. at /usr/share/perl5/Debconf/FrontEnd/Dialog.pm line
78.)
debconf: falling back to frontend: Readline
Running upmap-sys. This may take some time... done.
Running mktexlsr /var/lib/texmf ... done.
Building format(s) --all.
    This may take some time... done.
```

```
[42]: !apt-get install pandoc
```

```
Reading package lists... Done
Building dependency tree... Done
Reading state information... Done
The following additional packages will be installed:
  libcbmark-gfm-extensions0.29.0.gfm.3 libcbmark-gfm0.29.0.gfm.3 pandoc-data
Suggested packages:
  texlive-luatex pandoc-citeproc context wkhtmltopdf librsvg2-bin groff ghc
nodejs php python
  libjs-mathjax libjs-katex citation-style-language-styles
The following NEW packages will be installed:
  libcbmark-gfm-extensions0.29.0.gfm.3 libcbmark-gfm0.29.0.gfm.3 pandoc pandoc-
data
0 upgraded, 4 newly installed, 0 to remove and 49 not upgraded.
Need to get 20.6 MB of archives.
After this operation, 156 MB of additional disk space will be used.
```

```

Get:1 http://archive.ubuntu.com/ubuntu jammy/universe amd64 libcmark-
gfm0.29.0.gfm.3 amd64 0.29.0.gfm.3-3 [115 kB]
Get:2 http://archive.ubuntu.com/ubuntu jammy/universe amd64 libcmark-gfm-
extensions0.29.0.gfm.3 amd64 0.29.0.gfm.3-3 [25.1 kB]
Get:3 http://archive.ubuntu.com/ubuntu jammy/universe amd64 pandoc-data all
2.9.2.1-3ubuntu2 [81.8 kB]
Get:4 http://archive.ubuntu.com/ubuntu jammy/universe amd64 pandoc amd64
2.9.2.1-3ubuntu2 [20.3 MB]
Fetched 20.6 MB in 4s (5,288 kB/s)
Selecting previously unselected package libcmark-gfm0.29.0.gfm.3:amd64.
(Reading database ... 160462 files and directories currently installed.)
Preparing to unpack .../libcmark-gfm0.29.0.gfm.3_0.29.0.gfm.3-3_amd64.deb ...
Unpacking libcmark-gfm0.29.0.gfm.3:amd64 (0.29.0.gfm.3-3) ...
Selecting previously unselected package libcmark-gfm-
extensions0.29.0.gfm.3:amd64.
Preparing to unpack .../libcmark-gfm-
extensions0.29.0.gfm.3_0.29.0.gfm.3-3_amd64.deb ...
Unpacking libcmark-gfm-extensions0.29.0.gfm.3:amd64 (0.29.0.gfm.3-3) ...
Selecting previously unselected package pandoc-data.
Preparing to unpack .../pandoc-data_2.9.2.1-3ubuntu2_all.deb ...
Unpacking pandoc-data (2.9.2.1-3ubuntu2) ...
Selecting previously unselected package pandoc.
Preparing to unpack .../pandoc_2.9.2.1-3ubuntu2_amd64.deb ...
Unpacking pandoc (2.9.2.1-3ubuntu2) ...
Setting up libcmark-gfm0.29.0.gfm.3:amd64 (0.29.0.gfm.3-3) ...
Setting up libcmark-gfm-extensions0.29.0.gfm.3:amd64 (0.29.0.gfm.3-3) ...
Setting up pandoc-data (2.9.2.1-3ubuntu2) ...
Setting up pandoc (2.9.2.1-3ubuntu2) ...
Processing triggers for man-db (2.10.2-1) ...
Processing triggers for libc-bin (2.35-0ubuntu3.4) ...
/sbin/ldconfig.real: /usr/local/lib/libhwloc.so.15 is not a symbolic link

/sbin/ldconfig.real: /usr/local/lib/libur_adapter_level_zero.so.0 is not a
symbolic link

/sbin/ldconfig.real: /usr/local/lib/libtbbmalloc.so.2 is not a symbolic link

/sbin/ldconfig.real: /usr/local/lib/libur_loader.so.0 is not a symbolic link

/sbin/ldconfig.real: /usr/local/lib/libtbbmalloc_proxy.so.2 is not a symbolic
link

/sbin/ldconfig.real: /usr/local/lib/libtbbbind_2_5.so.3 is not a symbolic link

/sbin/ldconfig.real: /usr/local/lib/libumf.so.0 is not a symbolic link

/sbin/ldconfig.real: /usr/local/lib/libtcm.so.1 is not a symbolic link

```

```
/sbin/ldconfig.real: /usr/local/lib/libtcm_debug.so.1 is not a symbolic link

/sbin/ldconfig.real: /usr/local/lib/libtbb.so.12 is not a symbolic link

/sbin/ldconfig.real: /usr/local/lib/libtbbbind_2_0.so.3 is not a symbolic link

/sbin/ldconfig.real: /usr/local/lib/libtbbbind.so.3 is not a symbolic link

/sbin/ldconfig.real: /usr/local/lib/libur_adapter_opengl.so.0 is not a symbolic
link
```

```
[43]: !jupyter nbconvert --to pdf /content/EE5253Report_4466_4578_GP23.ipynb
```

```
[NbConvertApp] Converting notebook /content/EE5253Report_4466_4578_GP23.ipynb to
pdf
[NbConvertApp] Support files will be in EE5253Report_4466_4578_GP23_files/
[NbConvertApp] Making directory ./EE5253Report_4466_4578_GP23_files
[NbConvertApp] Writing 131808 bytes to notebook.tex
[NbConvertApp] Building PDF
[NbConvertApp] Running xelatex 3 times: ['xelatex', 'notebook.tex', '-quiet']
[NbConvertApp] Running bibtex 1 time: ['bibtex', 'notebook']
[NbConvertApp] WARNING | bibtex had problems, most likely because there were no
citations
[NbConvertApp] PDF successfully created
[NbConvertApp] Writing 798404 bytes to /content/EE5253Report_4466_4578_GP23.pdf
```