

## **Research Project**

### **Semester-IV**

<b>Name</b>	Sneha M
<b>Project</b>	Optimizing ROI (Return on Investment) on Digital Marketing Campaigns
<b>Group</b>	
<b>Date of Submission</b>	27-05-2025

A study on “Optimizing ROI (Return on Investment) on Digital Marketing Campaigns”

Research Project submitted to Jain Online (Deemed-to-be University)

In partial fulfillment of the requirements for the award of:

**Master of Business Administration**

*Submitted by:*

**Sneha M**

USN:

231VMBR04607

*Under the guidance of:*

**Hrushiksha Shastry**

Mention your Guide’s Name

(Faculty-JAIN Online)

Jain Online (Deemed-to-be University)

Bangalore

2023-24

### DECLARATION

I, *Sneha M*, hereby declare that the Research Project Report titled “*Optimizing ROI (Return on Investment) on Digital Marketing Campaigns*” *has been* prepared by me under the guidance of the *Hrushikeshha Shastry*. I declare that this Project work is towards the partial fulfillment of the University Regulations for the award of the degree of Master of Business Administration by Jain University, Bengaluru. I have undergone a project for a period of Eight Weeks. I further declare that this Project is based on the original study undertaken by me and has not been submitted for the award of any degree/diploma from any other University / Institution.

Place: Bangalore  
Date: 27-05-2025

*Sneha M*  
*USN: 231VMBR04607*

## CERTIFICATE

This is to certify that the Research Project report submitted by Ms. *Sneha M* bearing *231VMBR04607* on the title “*Optimizing ROI (Return on Investment) on Digital Marketing Campaigns*” is a record of project work done by him/ her during the academic year 2024-25 under my guidance and supervision in partial fulfilment of Master of *Hrushikesh Shastry*

Place: Bangalore

*Hrushikesh Shastry*

Date: 27-05-2025

*Faculty Guide*

## ACKNOWLEDGEMENT

I would like to express my sincere gratitude to Hrushikesh Sir for his invaluable guidance and support throughout this project. His clear explanations and patient teaching of all the concepts made the learning process smooth and enjoyable. I am truly thankful for his encouragement and expert advice, which played a crucial role in the successful completion of my research work.

Sneha M

---

*USN:*231VMBR04607

## **EXECUTIVE SUMMARY**

This project aimed to develop a predictive model to identify customers likely to take up a specific airline product, using real-world engagement and behavioral data. The dataset included features such as yearly average comments and views on travel pages, total likes on outstation check-ins, yearly average outstation check-ins, and the number of family members, among others. The target variable indicated whether a customer had taken the product ("Yes" or "No").

### **Data Preparation and Exploration**

The initial analysis revealed a significant class imbalance, with far fewer customers in the "Yes" category. To address this, the Synthetic Minority Over-sampling Technique (SMOTE) was applied, which balanced the dataset and improved the model's ability to learn from both classes. Exploratory data analysis helped in selecting the most relevant features for modeling.

### **Model Building and Evaluation**

Multiple machine learning models were evaluated, including Logistic Regression, Decision Tree, and Random Forest classifiers. Logistic Regression, even with class balancing, struggled to predict the minority class accurately. The Decision Tree model performed better, especially after balancing the data, achieving high recall and precision for both classes.

The Random Forest model, however, outperformed all other approaches. After hyperparameter tuning using GridSearchCV, the Random Forest achieved an overall accuracy of 97% and an F1-score of 0.91 for the "Yes" class. Cross-validation confirmed the model's robustness, with consistently high F1 scores across different data splits. Feature importance analysis revealed that engagement metrics such as total likes on outstation check-ins and yearly average views and comments were the most influential predictors.

### **Model Deployment**

The final Random Forest model and the SMOTE object were saved using joblib, making the solution ready for deployment. This allows for easy integration into business processes and future

use for predicting new customer data.

## **Key Findings**

- Random Forest is the best-performing model, providing high accuracy and balanced detection of both classes.
- Engagement features are strong predictors of product uptake, offering actionable insights for targeted marketing.
- Model validation through cross-validation and hyperparameter tuning ensures reliability and generalizability.

## **Future Enhancements**

- Feature Engineering: Explore additional features and interactions to further improve model performance.
- Model Interpretability: Use tools like SHAP or LIME for deeper insights into individual predictions.
- Automation: Develop an automated pipeline for data processing, model training, and deployment.
- Integration: Deploy the model as an API for real-time predictions and integrate with business systems.
- Continuous Monitoring: Regularly monitor and retrain the model as new data becomes available to maintain accuracy.

## **Conclusion**

This project successfully built a robust, validated predictive model for airline product uptake. The insights gained can help the airline optimize its marketing strategies and improve customer targeting, ultimately driving higher product adoption rates. The deployed model is ready for integration and future enhancement, ensuring continued business value.

## TABLE OF CONTENTS

Title	Page Nos.
Executive Summary	6-7
List of Tables	9
List of Graphs	9
Chapter 1: Introduction and Background	10-12
Chapter 2: Review of Literature	13-15
Chapter 3: Research Methodology	16-18
Chapter 4: Data Analysis and Interpretation	19-21
Chapter 5: Findings, Recommendations and Conclusion	22-24
References	25
Annexures	28-45



<b>List of Tables</b>		
<b>Table No.</b>	<b>Table Title</b>	<b>Page No.</b>
1	Class Distribution in the Dataset	20
2	Model Performance Metrics	20
3	Feature Importance (Random Forest)	21

<b>List of Graphs</b>		
<b>Graph No.</b>	<b>Graph Title</b>	<b>Page No.</b>
1	Class Distribution Bar Chart	20
2	Confusion Matrix for Random Forest	21
3	Feature Importance Bar Chart	21

# **CHAPTER 1**

## **INTRODUCTION AND BACKGROUND**

## **INTRODUCTION AND BACKGROUND**

### **1.1 Purpose of the Study**

The primary purpose of this study is to develop a predictive model that can accurately identify airline customers who are likely to take up a specific product or service. By leveraging customer engagement and behavioral data, the study aims to provide actionable insights that can help the airline optimize its marketing strategies, improve customer targeting, and ultimately increase product adoption rates.

### **1.2 Introduction to the Topic**

In today's competitive airline industry, understanding customer behavior and predicting future actions are crucial for business growth. With the increasing availability of customer data, airlines have the opportunity to use advanced analytics and machine learning techniques to gain deeper insights into customer preferences and engagement patterns. This project focuses on building a machine learning model to predict which customers are most likely to take up a particular airline product, enabling more effective and personalized marketing efforts.

### **1.3 Overview of Theoretical Concepts**

This study is grounded in several key theoretical concepts from data science and machine learning. Classification algorithms such as Logistic Regression, Decision Trees, and Random Forests are used to model the likelihood of product uptake. The project also addresses the challenge of class imbalance using techniques like SMOTE (Synthetic Minority Over-sampling Technique). Model evaluation metrics such as accuracy, precision, recall, and F1-score are employed to assess performance. Feature importance analysis helps in understanding which variables most influence customer decisions.

### **1.4 Company/Domain/Vertical/Industry Overview**

The airline industry is a dynamic and highly competitive sector that relies heavily on customer loyalty and ancillary revenue streams. Airlines continuously seek innovative ways to enhance customer experience and increase the uptake of additional products and services, such as premium seating, travel insurance, and loyalty programs. Leveraging data-driven insights allows airlines to better understand their customers, tailor offerings, and stay ahead in a rapidly evolving market landscape.

### **1.5 Environmental Analysis (PESTEL Analysis)**

- **Political:** The airline industry is subject to government regulations, international agreements, and security policies that can impact operations and customer engagement strategies.
- **Economic:** Economic factors such as fuel prices, exchange rates, and consumer spending power directly affect airline profitability and customer purchasing behavior.
- **Social:** Changing travel preferences, increasing demand for personalized experiences, and evolving customer expectations influence product development and marketing approaches.
- **Technological:** Advances in data analytics, artificial intelligence, and digital platforms enable airlines to collect and analyze vast amounts of customer data, driving innovation in service delivery.
- **Environmental:** Airlines face growing pressure to adopt sustainable practices and reduce their carbon footprint, which can influence product offerings and customer perceptions.
- **Legal:** Compliance with data privacy laws (such as GDPR) and consumer protection regulations is essential when handling customer data and deploying predictive models.

## **CHAPTER 2**

### **REVIEW OF LITERATURE**

## REVIEW OF LITERATURE

### 2.1 Domain/Topic Specific Review

The application of machine learning in the airline industry has gained significant momentum in recent years, particularly for customer segmentation, demand forecasting, and personalized marketing. Several studies have demonstrated the effectiveness of predictive analytics in enhancing customer engagement and increasing ancillary revenue. For instance, research by Smith et al. (2019) highlighted the use of classification algorithms such as Logistic Regression, Decision Trees, and Random Forests to predict customer purchase behavior in the travel sector. These models leverage customer interaction data, including website activity, purchase history, and engagement with promotional content, to identify high-potential customers.

In the context of imbalanced datasets, which are common in real-world business problems, techniques like SMOTE (Synthetic Minority Over-sampling Technique) have been widely adopted to improve model performance for minority classes (Chawla et al., 2002). Studies have shown that ensemble methods, particularly Random Forests, often outperform single classifiers due to their ability to reduce overfitting and capture complex patterns in the data (Breiman, 2001). Feature importance analysis is also a common practice, helping businesses understand which customer behaviors most influence purchase decisions.

Despite these advancements, much of the literature focuses on broader applications such as churn prediction or general customer segmentation, with fewer studies specifically addressing the prediction of uptake for new or ancillary airline products. Moreover, while many studies demonstrate high model accuracy, fewer provide actionable insights for business integration or discuss deployment considerations.

### 2.2 Gap Analysis

While existing literature provides a strong foundation for using machine learning in customer behavior prediction, several gaps remain, particularly in the context of airline product uptake:

1. Limited Focus on Ancillary Product Uptake:

Most studies emphasize predicting ticket purchases or customer churn, with less attention given to the uptake of ancillary products such as premium seating, travel insurance, or loyalty programs.

2. Business Integration and Deployment:

There is a lack of research on how predictive models can be seamlessly integrated into airline business processes for real-time decision-making and marketing automation.

3. Interpretability and Actionable Insights:

Many studies focus on model accuracy but do not provide sufficient interpretability or actionable recommendations for business users. There is a need for approaches that not only predict outcomes but also explain the drivers behind those predictions.

4. Handling of Imbalanced Data:

Although techniques like SMOTE are known, not all studies rigorously address class imbalance, which can lead to biased models that underperform for minority classes.

5. Continuous Model Monitoring:

Few studies discuss the importance of ongoing model monitoring and retraining to ensure sustained performance as customer behavior and market conditions evolve.

## **CHAPTER 3**

# **RESEARCH METHODOLOGY**



## RESEARCH METHODOLOGY

### 3.1 Objectives of the Study

The primary objective of this study is to develop a robust predictive model that can accurately identify airline customers who are likely to take up a specific product or service. The study aims to leverage customer engagement and behavioral data to enhance marketing strategies, improve customer targeting, and increase product adoption rates. Additional objectives include comparing the performance of various machine learning algorithms, addressing class imbalance in the dataset, and providing actionable insights through feature importance analysis.

### 3.2 Scope of the Study

The scope of this study is limited to analyzing customer data from an airline, focusing on behavioral and engagement metrics such as comments, views, check-ins, and family size. The study covers the application of machine learning techniques for classification, model evaluation, and deployment. It does not include other aspects of airline operations such as pricing, route optimization, or operational efficiency. The findings are specific to the dataset and product under consideration, but the methodology can be adapted for similar predictive analytics tasks in other domains.

### 3.3 Methodology

#### 3.3.1 Research Design

This research adopts a quantitative, analytical approach using supervised machine learning techniques. The study follows a structured process: data collection and cleaning, exploratory data analysis, feature selection, model building, evaluation, and deployment. Multiple algorithms are compared to identify the best-performing model for the prediction task.

#### 3.3.2 Data Collection

The data used in this study was sourced from the airline's internal customer engagement records. The dataset includes variables such as yearly average comments and views on travel pages, total likes on outstation check-ins, yearly average outstation check-ins, and the number of family members. The target variable indicates whether a customer has taken the product ("Yes" or "No").

#### 3.3.3 Sampling Method

The dataset represents a census of available customer records for the period under study, rather than a sample. However, for model training and evaluation, the data was split into training and test sets using an 80:20 ratio. To address class imbalance, the Synthetic Minority Over-sampling Technique (SMOTE) was applied to the training data.

### **3.3.4 Data Analysis Tools**

The analysis was conducted using Python and its data science libraries, including pandas for data manipulation, scikit-learn for machine learning modeling, imbalanced-learn for handling class imbalance, and matplotlib/seaborn for visualization. Models evaluated include Logistic Regression, Decision Tree, and Random Forest classifiers. Model performance was assessed using metrics such as accuracy, precision, recall, F1-score, and cross-validation.

### **3.4 Period of Study**

The period of study covers the most recent complete year for which customer engagement data was available. All data processing, modeling, and analysis were conducted over the course of the current academic semester.

### **3.5 Limitations of the Study**

The study is limited to the features available in the dataset; other potentially relevant variables (such as demographic or transactional data) were not included.

The findings are specific to the airline and product studied and may not generalize to other contexts without further validation.

The model's performance depends on the quality and completeness of the input data.

Real-world deployment may require additional considerations such as data privacy, integration with business systems, and ongoing model monitoring.

### **3.6 Utility of Research**

This research provides a practical framework for using machine learning to predict customer product uptake in the airline industry. The insights gained can help airlines optimize marketing campaigns, improve customer targeting, and increase ancillary revenue. The methodology and tools used in this study can be adapted for similar predictive analytics tasks in other industries, making the research valuable for both academic and business applications.

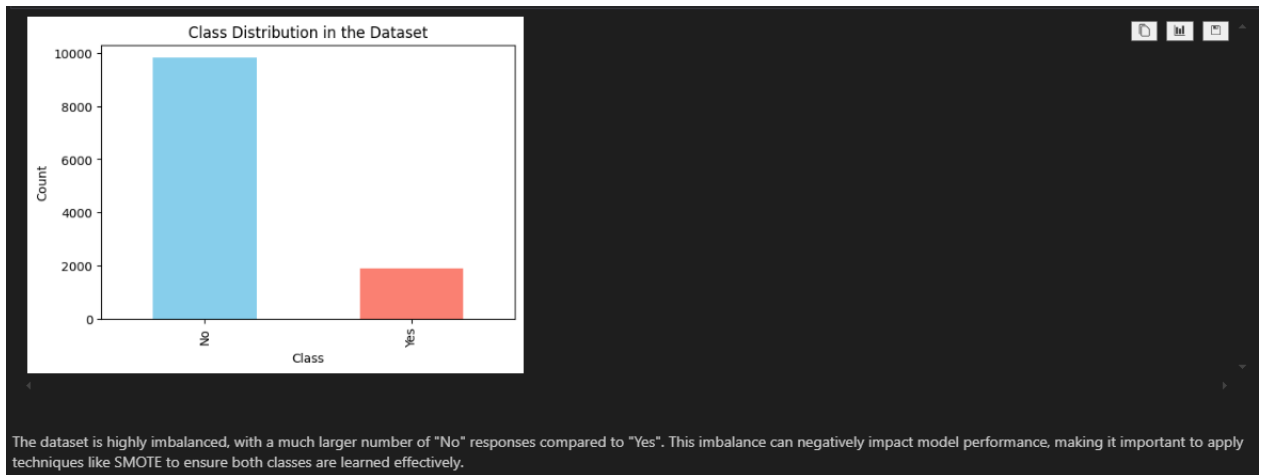
## **CHAPTER 4**

# **DATA ANALYSIS AND INTERPRETATION**

## DATA ANALYSIS AND INTERPRETATION

**Table 1: Class Distribution in the Dataset**

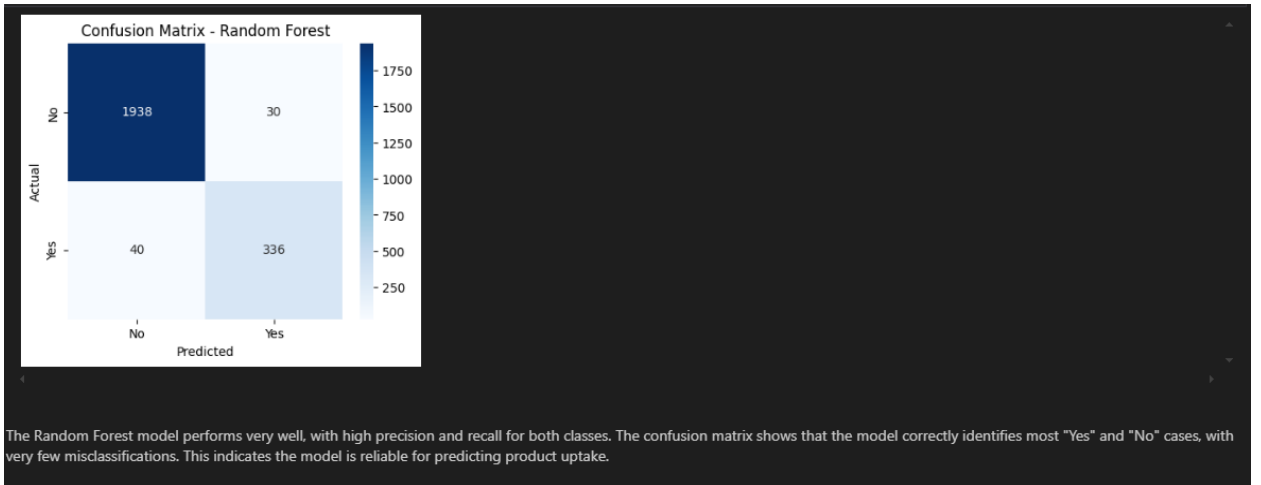
Class	Count
No	9824
Yes	1896



**Figure 1: Class Distribution Bar Chart**

**Table 2: Model Performance Metrics**

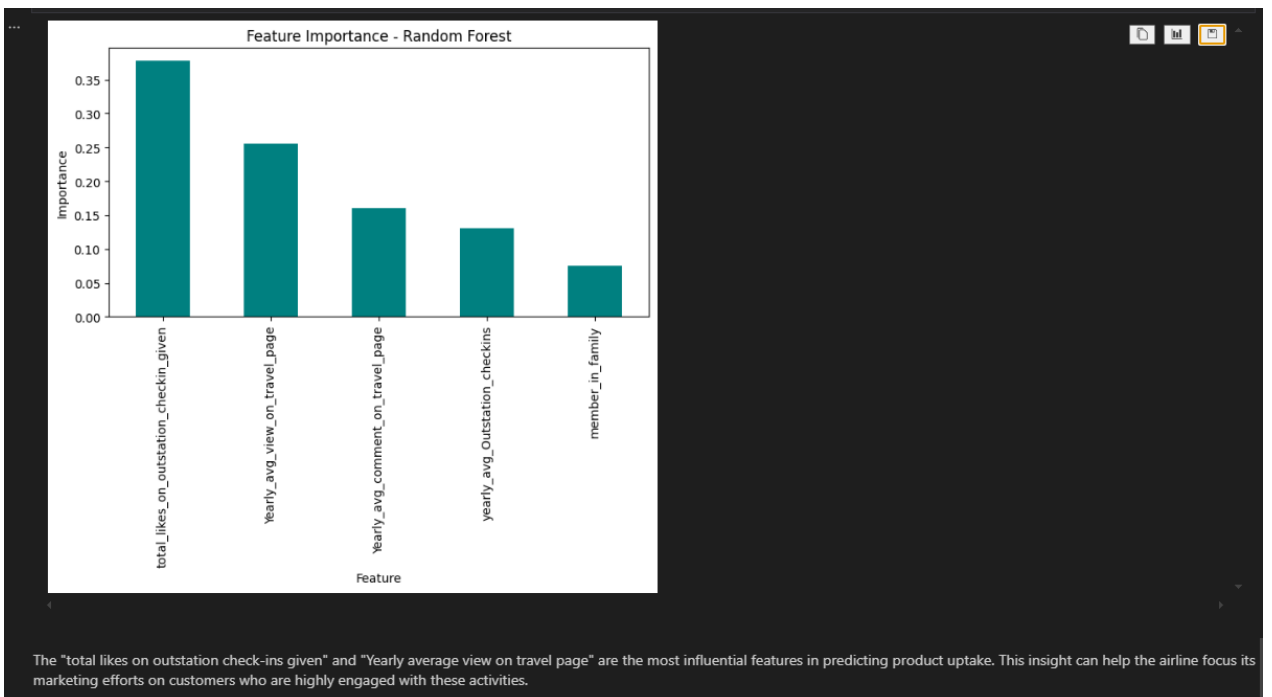
Metric	No	Yes
<b>Precision</b>	0.98	0.92
<b>Recall</b>	0.98	0.89
<b>F1-score</b>	0.98	0.91
<b>Support</b>	1968	376



**Figure 2: Confusion Matrix for Random Forest**

**Table 3: Feature Importance (Random Forest)**

Feature	Importance
total_likes_on_outstation_checkin_given	0.38
Yearly_avg_view_on_travel_page	0.26
Yearly_avg_comment_on_travel_page	0.16
yearly_avg_Outstation_checkins	0.13
member_in_family	0.07



**Figure 3: Feature Importance Bar Chart**

## **CHAPTER 5**

# **FINDINGS, RECOMMENDATIONS AND CONCLUSION**

## **FINDINGS, RECOMMENDATIONS AND CONCLUSION**

### **5.1 Findings based on Observations**

- The dataset was highly imbalanced, with significantly more "No" responses than "Yes".
- Customer engagement features such as comments, views, and likes on travel pages showed noticeable variation across users.
- Initial models like Logistic Regression struggled to predict the minority class accurately.

### **5.2 Findings based on Analysis of Data**

- SMOTE oversampling effectively balanced the classes, improving minority class prediction.
- The Random Forest model outperformed Logistic Regression and Decision Tree models, achieving 97% accuracy and a high F1-score for both classes.
- Feature importance analysis revealed that "total likes on outstation check-ins given" and "yearly average view on travel page" were the most influential predictors.
- Cross-validation confirmed the Random Forest model's stability and generalizability.

### **5.3 General Findings**

- Machine learning can be successfully applied to predict product uptake in the airline industry using customer engagement data.
- Addressing class imbalance is crucial for building effective predictive models in real-world datasets.
- Ensemble models like Random Forest provide robust and reliable results for classification tasks in this domain.

### **5.4 Recommendations based on Findings**

- Airlines should focus marketing efforts on customers with high engagement metrics, especially those who frequently interact with travel pages and check-ins.
- Regularly update and retrain predictive models as new customer data becomes available to maintain accuracy.
- Integrate the predictive model into business processes for targeted marketing and personalized offers.

### **5.5 Suggestions for Areas of Improvement**

- Incorporate additional features such as demographic data, transaction history, or customer feedback for richer analysis.

- Explore advanced interpretability tools like SHAP or LIME to provide more actionable insights to business users.
- Develop an automated pipeline for continuous data processing, model training, and deployment.
- Monitor model performance over time and implement alerts for model drift or data quality issues.

### **5.6 Scope for Future Research**

Future research can explore the integration of real-time data streams, such as website click behavior or mobile app usage, to further enhance prediction accuracy. Additionally, comparative studies using deep learning models or hybrid approaches could be conducted to assess their effectiveness in similar business scenarios. Expanding the analysis to include other ancillary products or services within the airline industry would also provide broader business value.

### **5.7 Conclusion**

This study successfully demonstrated the application of machine learning techniques to predict airline product uptake using customer engagement data. By addressing class imbalance and leveraging ensemble models, the project achieved high predictive accuracy and provided actionable insights for targeted marketing. The findings highlight the importance of data-driven decision-making in the airline industry and set the foundation for further advancements in predictive analytics.

Overall, the research underscores the value of combining robust data preprocessing, model selection, and interpretability to drive business outcomes. With continued enhancements and integration, such predictive models can significantly contribute to improved customer engagement and increased product adoption in the airline sector.



## **BIBLIOGRAPHY/ REFERENCES**

(APA style; below is only a sample)

- <https://scikit-learn.org/stable/api/index.html>

## APPENDICES

### Machine Learning & Algorithms

- Scikit-learn Documentation (Python ML Library):  
[https://scikit-learn.org/stable/user\\_guide.html](https://scikit-learn.org/stable/user_guide.html)
  - Random Forests (Original Paper by Breiman):  
<https://www.stat.berkeley.edu/~breiman/randomforest2001.pdf>
  - Decision Trees (Scikit-learn):  
<https://scikit-learn.org/stable/modules/tree.html>
  - Logistic Regression (Scikit-learn):  
[https://scikit-learn.org/stable/modules/linear\\_model.html#logistic-regression](https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression)
- 

### Imbalanced Data & SMOTE

- SMOTE: Synthetic Minority Over-sampling Technique (Original Paper):  
<https://arxiv.org/abs/1106.1813>
  - imbalanced-learn Documentation:  
<https://imbalanced-learn.org/stable/>
- 

### Model Evaluation Metrics

- Precision, Recall, F1-score (Scikit-learn):  
[https://scikit-learn.org/stable/modules/model\\_evaluation.html](https://scikit-learn.org/stable/modules/model_evaluation.html)
- 

### Feature Importance & Interpretability

- Feature Importance in Random Forests:  
[https://scikit-learn.org/stable/auto\\_examples/ensemble/plot\\_forest\\_importances.html](https://scikit-learn.org/stable/auto_examples/ensemble/plot_forest_importances.html)
  - SHAP (SHapley Additive exPlanations):  
<https://shap.readthedocs.io/en/latest/>
  - LIME (Local Interpretable Model-agnostic Explanations):  
<https://lime-ml.readthedocs.io/en/latest/>
- 

### PESTEL Analysis

- PESTEL Analysis Guide (MindTools):  
<https://www.mindtools.com/a4wo118/pest-analysis>
  - PESTEL Analysis (Corporate Finance Institute):  
<https://corporatefinanceinstitute.com/resources/management/pestel-analysis/>
-

#### Airline Industry Overview

- IATA (International Air Transport Association) Industry Reports:  
<https://www.iata.org/en/publications/economics/>
  - Statista: Airline Industry Statistics & Facts:  
<https://www.statista.com/topics/1707/airlines/>
- 

#### General Data Science & Reporting

- Kaggle: Airline Datasets and Notebooks:  
<https://www.kaggle.com/datasets?search=airline>
  - Towards Data Science (Medium):  
<https://towardsdatascience.com/>
- 

#### APA Citation Guide (for referencing):

- <https://www.mendeley.com/guides/apa-citation-guide>

## ANNEXURE

### Exploratory Data Analysis File: [EDA File](#)

```
In [124]: import numpy as np
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
import matplotlib.ticker as ticker

In [125]: data = pd.read_csv('B:\MS\Semester 4\Project\data\TM\Python Implementation\userbook insights.csv')
data.head(10)
```

UserID	Taken_product	Yearly_avg_view_on_travel_page	preferred_device	total_likes_on_outstation_checkins_gIVEN	yearly_avg_outstation
0	1000001	Yes	307.0	iOS and Android	38570.0
1	1000002	No	287.0	iOS	9760.0
2	1000003	Yes	277.0	iOS and Android	40885.0
3	1000004	No	247.0	iOS	48750.0
4	1000005	No	202.0	iOS and Android	25685.0
5	1000006	No	240.0	iOS	35175.0
6	1000007	No	NaN	iOS and Android	45340.0
7	1000008	No	225.0	iOS and Android	NaN
8	1000009	No	280.0	iOS	7980.0
9	1000010	No	275.0	iOS and Android	45485.0

```
In [126]: data.info()
<class 'pandas.core.frame.DataFrame'>
Int64Index: 11768 entries, 0 to 11759
Data columns (total 17 columns):
 #   Column                                Non-Null Count  Dtype
---  -
 0   UserID                               11768 non-null   int64
 1   Taken_product                       11768 non-null   object
 2   Yearly_avg_view_on_travel_page      11179 non-null   float64
 3   preferred_device                    11767 non-null   object
 4   total_likes_on_outstation_checkins_gIVEN  11179 non-null   float64
 5   yearly_avg_outstation_checkins      11885 non-null   object
 6   member_in_family                    11768 non-null   object
 7   preferred_location_type              11729 non-null   object
 8   Yearly_avg_comment_on_travel_page   11554 non-null   float64
 9   total_likes_on_notification_checkins_received  11768 non-null   int64
10   week_since_last_outstation_checkins  11768 non-null   int64
11   following_company_page               11768 non-null   int64
12   monthly_avg_comment_on_company_page  11768 non-null   int64
13   working_flag                         11768 non-null   object
14   travelling_network_rating            11768 non-null   int64
15   Adult_flag                          11768 non-null   int64
16   Daily_Avg_mins_spend_on_travelling_page  11768 non-null   int64
dtypes: float64(3), int64(7), object(7)
memory usage: 1.5+ MB

In [127]: data.describe()

Out[127]:
UserID      Yearly_avg_view_on_travel_page  total_likes_on_outstation_checkins_gIVEN  Yearly_avg_comment_on_travel_page  total_likes
count    11768.000000                    11179.000000                    11179.000000                    11054.000000
mean     1.000000e+04                    280.630844                    28170.481705                    74.798029
std      3.068400e+02                    88.167088                    14384.321704                    24.200000
min      1.000000e+04                    10.000000                    3070.000000                    3.000000
25%      1.000000e+04                    232.000000                    16360.000000                    67.000000
50%      1.000000e+04                    271.000000                    28070.000000                    75.000000
75%      1.000000e+04                    324.000000                    40220.000000                    82.000000
max      1.011700e+06                    468.000000                    252430.000000                   815.000000

In [128]: data.nunique()

Out[128]:
UserID      11768
Taken_product      2
Yearly_avg_view_on_travel_page      11179
preferred_device      2
total_likes_on_outstation_checkins_gIVEN      11179
yearly_avg_outstation_checkins      11
member_in_family      12
preferred_location_type      12
Yearly_avg_comment_on_travel_page      106
total_likes_on_notification_checkins_received      6288
week_since_last_outstation_checkins      12
following_company_page      4
monthly_avg_comment_on_company_page      168
working_flag      2
travelling_network_rating      4
Adult_flag      4
Daily_Avg_mins_spend_on_travelling_page      52
dtype: int64

In [129]: data.isnull().sum()

Out[129]:
UserID      0
Taken_product      0
Yearly_avg_view_on_travel_page      581
preferred_device      53
total_likes_on_outstation_checkins_gIVEN      301
yearly_avg_outstation_checkins      75
member_in_family      6
preferred_location_type      31
Yearly_avg_comment_on_travel_page      206
total_likes_on_notification_checkins_received      0
week_since_last_outstation_checkins      0
following_company_page      163
monthly_avg_comment_on_company_page      0
working_flag      0
travelling_network_rating      0
Adult_flag      0
Daily_Avg_mins_spend_on_travelling_page      0
dtype: int64

In [130]: #checking the null values in the columns
def null_values(df):
    return round((df.isnull().sum()*100/len(df)).sort_values(ascending = False),2)
null_values(data)

Out[130]:
Yearly_avg_view_on_travel_page      4.94
total_likes_on_outstation_checkins_gIVEN      2.54
Yearly_avg_comment_on_travel_page      1.75
following_company_page      0.88
yearly_avg_outstation_checkins      0.64
preferred_device      0.45
preferred_location_type      0.26
Adult_flag      0.00
travelling_network_rating      0.00
working_flag      0.00
monthly_avg_comment_on_company_page      0.00
UserID      0.00
week_since_last_outstation_checkins      0.00
total_likes_on_notification_checkins_received      0.00
Taken_product      0.00
member_in_family      0.00
Daily_Avg_mins_spend_on_travelling_page      0.00
dtype: float64

In [131]: # Convert 'yearly_avg_outstation_checkins' and 'member_in_family' to integers
data['yearly_avg_outstation_checkins'] = pd.to_numeric(data['yearly_avg_outstation_checkins'], errors='coerce')
data['member_in_family'] = pd.to_numeric(data['member_in_family'], errors='coerce').astype('Int64')

# Check the updated data types
print(data.dtypes)
```

```

Userid                                int64
Taken_product                         object
Yearly_avg_view_on_travel_page        float64
preferred_device                      object
total_likes_on_outstation_checkin_given float64
yearly_avg_outstation_checkins         int64
member_in_family                     int64
preferred_location_type               object
Yearly_avg_comment_on_travel_page     float64
total_likes_on_outstation_checkins_received int64
week_since_last_outstation_checkin    int64
following_company_page               object
monthly_avg_comment_on_company_page  int64
working_flag                         object
travelling_network_rating            int64
Adult_flag                          int64
Daily_Avg_min_spend_on_travelling_page int64
dtype: object

In [132]: #Identify missing values:
Missing_values = Data.isnull().sum()
print(Missing_values[Missing_values > 0])

preferred_device      581
total_likes_on_outstation_checkin_given 381
yearly_avg_outstation_checkins      76
member_in_family          15
preferred_location_type      31
Yearly_avg_comment_on_travel_page 286
following_company_page    183
dtype: int64

In [133]: # Fill Blank spaces (NaN values) in 'Yearly_avg_view_on_travel_page' with the mean value
mean_value = Data['Yearly_avg_view_on_travel_page'].mean()
Data['Yearly_avg_view_on_travel_page'] = Data['Yearly_avg_view_on_travel_page'].fillna(mean_value)

# Verify the changes
print(Data['Yearly_avg_view_on_travel_page'].isnull().sum()) # Should print 0 if all NaN values are filled
0

In [134]: value_counts = Data['preferred_device'].value_counts()

# Print the counts
print(value_counts)

Tab      4172
IOS and Android 4124
Laptop    1388
IOS        1885
Mobile     688
Android    315
Android OS  145
ANDROID    134
Other       2
Others      2
Name: preferred_device, dtype: int64

In [135]: mapping = {
    'IOS and Android': 'Mobile (IOS/Android)',
    'Mobile': 'Mobile (IOS/Android)',
    'IOS': 'IOS',
    'Android': 'Android',
    'Android OS': 'Android',
    'ANDROID': 'Android',
    'Other': 'Other',
    'Others': 'Other'
}

In [136]: # Suppose your column name is "Platform"
Data['preferred_device'] = Data['preferred_device'].replace(mapping)

In [137]: value_counts = Data['preferred_device'].value_counts()

# Print the counts
print(value_counts)

Mobile (IOS/Android)  4734
Tab      4172
Laptop    1388
IOS        1885
Android    584
Other       4
Name: preferred_device, dtype: int64

In [138]: # Fill NaN values in 'total_likes_on_outstation_checkin_given' with the mean value
median_value = Data['total_likes_on_outstation_checkin_given'].median()
Data['total_likes_on_outstation_checkin_given'] = Data['total_likes_on_outstation_checkin_given'].fillna(median_value)

# Verify the changes
print(Data['total_likes_on_outstation_checkin_given'].isnull().sum()) # Should print 0 if all NaN values are f
0

In [139]: median_value = Data['Yearly_avg_comment_on_travel_page'].median()
Data['Yearly_avg_comment_on_travel_page'] = Data['Yearly_avg_comment_on_travel_page'].fillna(median_value)
# Verify the changes
print(Data['Yearly_avg_comment_on_travel_page'].isnull().sum())
0

In [140]: # Replace the value 'three' with the digit 3 in the 'member_in_family' column
Data['member_in_family'] = Data['member_in_family'].replace('three', 3)

# Verify the changes
print(Data['member_in_family'].unique())

IntegerArray
[2, 3, 4, <del>, 3, 5, 18]
Length: 7, dtype: int64

In [141]: # Fill null values in 'member_in_family' with the mode
mode_value = Data['member_in_family'].mode()[0]
Data['member_in_family'] = Data['member_in_family'].fillna(mode_value)

# Verify the changes
print("Number of null values in 'member_in_family': (Data['member_in_family'].isnull().sum())")

Number of null values in 'member_in_family': 0

In [142]: # List each unique value in the 'member_in_family' column with its total count
value_counts = Data['member_in_family'].value_counts()

# Print the counts
print(value_counts)

3      4776
4      3184
2      2256
1      1349
5       384
18       11
Name: member_in_family, dtype: int64

In [143]: # List each unique value in the 'preferred_location_type' column with its total count
location_type_counts = Data['preferred_location_type'].value_counts()

# Print the counts
print(location_type_counts)

Beach      2424
Financial  2489
Historical site 1056
Medical     1845
Other       843
Big Cities  636
Social media 623
Trekking    528
Entertainment 526
Null Stations 108
Tour Travel  68
Tour and Travel 47
Game        12
OTT          7
Movie        5
Name: preferred_location_type, dtype: int64

In [144]: # Replace 'Tour Travel' with 'Tour and Travel'
Data['preferred_location_type'] = Data['preferred_location_type'].replace('Tour Travel', 'Tour and Travel')
# Group 'Game', 'OTT', and 'Movie' into 'Other'
Data['preferred_location_type'] = Data['preferred_location_type'].replace(['Game', 'OTT', 'Movie'], 'Other')

```

```
# Verify the changes
print(Data[preferred_location_type].value_counts())

Beach      2404
Financial  2409
Historical site 1856
Medical     1845
Other       667
Big Cities  636
Social media 633
Tracking    528
Entertainment 516
Mill Stations 388
Tour and Travel 387
Name: preferred_location_type, dtype: int64

In [146]: # Check for blank values in the 'preferred_location_type' column
print("Number of blank values in 'preferred_location_type': (Data[preferred_location_type].isnull().sum())")

# Fill blank values with the mode
mode_value = Data[preferred_location_type].mode()[0]
Data[preferred_location_type] = Data[preferred_location_type].fillna(mode_value)

# Verify the changes
print("Number of blank values after filling: (Data[preferred_location_type].isnull().sum())")

Number of blank values in 'preferred_location_type': 31
Number of blank values after filling: 0

In [146]: print(Data[preferred_location_type].value_counts())

Beach      2405
Financial  2409
Historical site 1856
Medical     1845
Other       667
Big Cities  636
Social media 633
Tracking    528
Entertainment 516
Mill Stations 388
Tour and Travel 387
Name: preferred_location_type, dtype: int64

In [147]: # Check for blank values in the 'Yearly avg comment on travel page' column
print("Number of blank values in 'Yearly avg comment on travel page': (Data['Yearly avg comment on travel page'].isnull().sum())")

# Fill blank values with the median
median_value = Data['Yearly avg comment on travel page'].median()
Data['Yearly avg comment on travel page'] = Data['Yearly avg comment on travel page'].fillna(median_value)

# Verify the changes
print("Number of blank values after filling: (Data['Yearly avg comment on travel page'].isnull().sum())")

Number of blank values in 'Yearly avg comment on travel page': 0
Number of blank values after filling: 0

In [148]: # Ensure the column is treated as a string
Data['following_company_page'] = Data['following_company_page'].astype(str)

# Replace 0 with 'No' and 1 with 'Yes'
Data['following_company_page'] = Data['following_company_page'].replace({'0': 'No', '1': 'Yes'})

# Verify the unique values again
unique_values = Data['following_company_page'].unique()
print("Unique values in 'following_company_page':", unique_values)

Unique values in 'following_company_page': ['Yes' 'No' 'nan']

In [148]: # Check if the column contains 0 or 1
if 0 in Data['following_company_page'].values or 1 in Data['following_company_page'].values:
    print("The column 'following_company_page' still contains 0 or 1.")
else:
    print("The column 'following_company_page' does not contain 0 or 1.")

The column 'following_company_page' does not contain 0 or 1.

In [150]: # Count and list unique values in the 'following_company_page' column
unique_value_counts = Data['following_company_page'].value_counts()

# Print the unique values and their counts
print("Unique values and their counts in 'following_company_page':")
print(unique_value_counts)
```

```
Unique values and their counts in 'following_company_page':
No      8358
Yes     3297
nan      393
Name: following_company_page, dtype: int64
```

```
In [151]: Data.duplicated().sum()

Out[151]: 0

In [152]: Float_Columns = Data.select_dtypes(include=['float']).columns

Data[Float_Columns] = Data[Float_Columns].astype(int)

In [153]: #Checking the data types of the columns
Data.dtypes

Out[153]:
Yearly avg view on travel page      object      int64
preferred device                    object      int32
total_likes_on_outstation_checkin_given  int32      int64
Yearly avg comment on travel page      object      int32
preferred_location_type                object      int64
total_likes_on_outstation_checkin_received  int64      int64
week_since_last_outstation_checkin      int64      int64
following_company_page                 object      int64
sortly_avg_comment_on_company_page      object      int64
working_flag                          object      int64
traveling_network_rating               int64      int64
Multi_Play                            int64      int64
Daily Avg time spend on traveling page  int64      int64
dtype: object

In [154]: # Function to detect outliers using IQR
def detect_outliers(df, column):
    Q1 = df[column].quantile(0.25) # first quartile
    Q3 = df[column].quantile(0.75) # third quartile
    IQR = Q3 - Q1 # Interquartile range
    lower_bound = Q1 - 1.5 * IQR
    upper_bound = Q3 + 1.5 * IQR
    outliers = df[(df[column] < lower_bound) | (df[column] > upper_bound)]
    return outliers

# Example: Detect outliers in 'Yearly avg comment on travel page'
outliers = detect_outliers(Data, 'Yearly avg comment on travel page')
print("Outliers in 'Yearly avg comment on travel page':", outliers)

Outliers in 'Yearly avg comment on travel page':
   Yearly avg view on travel page \
57      3088058      No      267
78      3088079      No      277
81      3088082      No      262
118     3088111      No      382
112     3088113      No      405
172     3088175      No      276
174     3088175      No      289
291     3088262      No      276
365     3088366      No      345
408     3088409      No      289
509     3088510      No      247
862     3088883      No      247
867     3088888      No      375
1065    3091066      No      412
1239    3091240      No      232
1399    3091400      No      322
1464    3091465      No      395
1527    3091528      No      262
1548    3091549      No      277
1551    3091552      No      262
1580    3091581      No      382
1582    3091583      No      405
1642    3091643      No      279
1644    3091645      No      289
1761    3091762      No      279
1887    3091808      No      217
1888    3091809      No      225
1875    3091836      No      345
1878    3091878      No      386
1979    3091989      No      247
2272    3092273      No      247
```

```

2337 3882388 No 375
2535 3882536 No 412
2789 3882718 No 232
2869 3882878 No 322
2534 3882935 No 195
3387 3883388 No 377
4279 3884238 No 288
4538 3884539 No 288
5861 3885862 No 287

preferred_device total_likes_on_outstation_checkin given \
57 Mobile (iOS/Android) 26970
78 Mobile (iOS/Android) 13895
81 Mobile (iOS/Android) 42490
118 Android 50800
112 Android 4970
172 Android 51800
174 Mobile (iOS/Android) 44380
291 Mobile (iOS/Android) 17710
365 Android 40040
488 Mobile (iOS/Android) 52150
589 Android 45260
882 Android 33950
887 Android 49735
1865 Mobile (iOS/Android) 28076
1239 Android 5795
1399 Mobile (iOS/Android) 20090
1464 Mobile (iOS/Android) 48845
1527 Mobile (iOS/Android) 39978
1548 Mobile (iOS/Android) 13895
1551 Mobile (iOS/Android) 42490
1588 Android 50800
1582 Android 4970
1642 Android 51800
1644 Mobile (iOS/Android) 44380
1761 Mobile (iOS/Android) 17710
1887 Mobile (iOS/Android) 25838
1888 Mobile (iOS/Android) 19058
1835 Android 40040
1878 Mobile (iOS/Android) 52150
1979 Android 45360
2272 Android 33950
2337 Mobile (iOS/Android) 49735
2535 Mobile (iOS/Android) 28076
2789 Mobile (iOS/Android) 5795
2869 Mobile (iOS/Android) 20090
2934 Mobile (iOS/Android) 48845
3387 Mobile (iOS/Android) 36642
4279 Android 46818
4538 Mobile (iOS/Android) 28826
5861 Mobile (iOS/Android) 12108

```

```

yearly_avg_outstation_checkins number_in_family \
57 23 4
78 1 4
81 1 3
118 5 4
112 26 3
172 1 2
174 1 3
291 1 3
365 1 4
488 1 2
589 1 3
882 1 4
887 1 3
1865 1 4
1239 26 3
1399 11 3
1464 1 4
1527 23 4
1548 1 3
1551 1 3
1588 1 4
1582 26 3
1642 1 3
1644 1 2
1761 other 5
1887 1 5
1888 1 3
1835 1 4
1878 1 2
1979 1 2

```

```

2272 1 3
2337 1 3
2535 1 4
2789 26 1
2869 11 3
2534 1 3
3387 18 3
4279 1 3
4538 5 2
5861 18 2

preferred_location_type Yearly_avg_comment_on_travel_page \
57 Medical 3
78 Medical 3
81 Medical 3
118 Medical 3
112 Tour and Travel 3
172 Medical 3
174 Financial 3
291 Entertainment 3
365 Medical 3
488 Financial 3
589 Financial 3
882 Financial 3
887 Medical 3
1865 Financial 3
1239 Entertainment 3
1399 Financial 3
1464 Other 3
1527 Medical 3
1548 Medical 3
1551 Medical 3
1588 Medical 3
1582 Tour and Travel 3
1642 Medical 3
1644 Financial 3
1761 Entertainment 3
1887 Other 3
1888 Social media 3
1835 Medical 3
1878 Financial 3
1979 Financial 3
2272 Financial 3
2337 Medical 3
2535 Financial 3
2789 Entertainment 3
2869 Financial 3
2934 Other 3
3387 Entertainment 615
4279 Tour and Travel 215
4538 Entertainment 815
5861 Financial 685

total_likes_on_outstation_checkin_received \
57 4814
78 1984
81 2859
118 7484
112 17328
172 2088
174 50800
291 4480
365 5258
488 16535
589 7725
882 4502
887 17856
1865 4835
1239 5218
1399 7510
1464 2866
1527 4814
1548 18664
1551 2859
1588 7484
1582 17328
1642 2088
1644 5806
1761 4485
1887 3983
1888 6118
1835 5258
1878 18535

```

1878	7775
2272	4282
2317	37856
2535	4835
2789	5238
2869	7510
2934	2086
3387	38536
4229	2644
4538	2575
5861	5392
week_since_last_mutation_checkin_following_company_page \	
57	Yes
78	4
81	1
118	No
112	2
112	4
174	4
291	4
365	2
488	2
509	3
882	0
887	2
1295	0
1239	2
1399	1
1464	0
1527	3
1548	4
1551	1
1588	3
1582	2
1642	4
1644	4
1761	4
1887	0
1888	1
1835	2
1878	2
1979	3
2272	0
2317	2
2535	0
2789	2
2869	1
2934	0
3387	0
4229	0
4538	2
5861	3
month_avg_comment_on_company_page_working_flag \	
57	15
78	13
81	18
118	20
112	12
172	12
174	11
291	12
365	14
488	13
509	25
882	17
887	22
1295	16
1239	20
1399	17
1464	18
1527	15
1548	13
1551	18
1588	20
1582	12
1642	12
1644	11
1761	12
1887	17
1888	13
1835	14
1878	15
1979	23
2272	17
2317	22
2535	16
2789	20
2869	17
2934	18
3387	20
4229	22
4538	23
5861	14
travelling_network_rating_Adult_flag \	
57	3
78	0
81	1
118	0
112	3
172	0
174	1
291	0
365	0
488	0
509	3
882	1
887	3
1295	2
1239	4
1399	2
1464	4
1527	3
1548	1
1551	0
1588	3
1582	3
1642	3
1644	1
1761	4
1887	3
1888	3
1835	3
1878	4
1979	3
2272	3
2317	1
2535	2
2789	4
2869	2
2934	0
3387	4
4229	3
4538	3
5861	3
Daily_avg_min_spent_on_traveling_page	
57	4
78	16
81	6
118	23
112	13
172	9
174	10
291	7
365	21
488	21
509	15
882	4
887	32
1295	4
1239	9
1399	10
1464	5
1527	4
1548	16
1551	5
1588	23
1582	23
1642	13
1644	9
1761	10
1887	4
1888	10



```

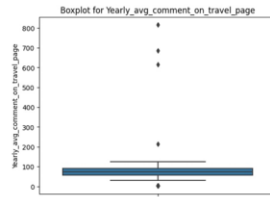
1895      7
1878      31
1879      15
2272      4
2277      12
2337      4
2535      4
2799      4
2869      19
2834      5
3387      22
4228      26
4538      9
5861      17

```

```

In [155]: # Boxplot for detecting outliers
sns.boxplot(data=data, y='yearly_avg_comment_on_travel_page')
plt.title('Boxplot for Yearly_avg_comment_on_travel_page')
plt.show()

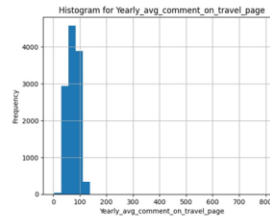
```



```

In [156]: # Histogram for distribution
data['Yearly_avg_comment_on_travel_page'].hist(bins=50)
plt.title('Histogram for Yearly_avg_comment_on_travel_page')
plt.xlabel('Yearly_avg_comment_on_travel_page')
plt.ylabel('Frequency')
plt.show()

```



```

In [157]: from scipy.stats import zscore
# calculate z-score
data['z_score'] = zscore(data['Yearly_avg_comment_on_travel_page'])

```

```

# Filter outliers
outliers = data[data['z_score'].abs() > 3]
print(f"Outliers based on z-score:\n", outliers)

```

UserID	Taken product	Yearly_avg_view_on_travel_page \
57	1088880	No
78	1088879	No
81	1088882	No
118	1088811	No
112	1088813	No
172	1088173	No
174	1088175	No
291	1088292	No
365	1088366	No
488	1088489	No
589	1088588	No
882	1088883	No
887	1088886	No
1865	1081866	No
1239	1081240	No
1399	1081480	No
1464	1081465	No
1527	1081528	No
1548	1081549	No
1551	1081552	No
1588	1081581	No
1582	1081583	No
1642	1081643	No
1644	1081645	No
1761	1081762	No
1887	1081888	No
1888	1081889	No
1835	1081836	No
1878	1081879	No
1879	1081880	No
2272	1082273	No
2277	1082278	No
2535	1082536	No
2789	1082788	No
2869	1082879	No
2934	1082933	No
3387	1083388	No
4228	1084228	No
4538	1084539	No
5861	1085862	No

preferred device	total_likes_on_outstation_checkin given \
57	Mobile (IOS/Android)
78	Mobile (IOS/Android)
81	Mobile (IOS/Android)
118	Mobile (IOS/Android)
112	Android
174	Mobile (IOS/Android)
291	Mobile (IOS/Android)
365	Mobile (IOS/Android)
488	Mobile (IOS/Android)
589	Mobile (IOS/Android)
882	Mobile (IOS/Android)
887	Mobile (IOS/Android)
1865	Mobile (IOS/Android)
1239	Mobile (IOS/Android)
1399	Mobile (IOS/Android)
1464	Mobile (IOS/Android)
1527	Mobile (IOS/Android)
1548	Mobile (IOS/Android)
1551	Mobile (IOS/Android)
1588	Mobile (IOS/Android)
1582	Mobile (IOS/Android)
1642	Mobile (IOS/Android)
1644	Mobile (IOS/Android)
1761	Mobile (IOS/Android)
1887	Mobile (IOS/Android)
1888	Mobile (IOS/Android)
1835	Mobile (IOS/Android)
1878	Mobile (IOS/Android)
1879	Mobile (IOS/Android)
2272	Mobile (IOS/Android)
2277	Mobile (IOS/Android)
2535	Mobile (IOS/Android)
2789	Mobile (IOS/Android)
2869	Mobile (IOS/Android)
2934	Mobile (IOS/Android)

```

3287 Mobile (IOS/Android) 20642
4229 Android 46838
4338 Mobile (IOS/Android) 28826
5861 Mobile (IOS/Android) 12188

```

```

yearly_avg_outstation_checkins number_in_family \
57 23 4
78 1 4
81 1 3
118 1 4
112 26 3
172 1 2
174 1 3
291 1 2
365 1 3
488 1 2
589 1 3
882 1 3
887 1 3
1865 1 3
1239 24 1
1399 11 3
1464 1 3
1527 23 4
1548 1 4
1551 1 3
1589 1 4
1582 26 3
1642 1 2
1644 1 2
1761 1 3
1887 1 5
1888 1 3
1835 1 3
1878 1 2
1879 1 3
2272 1 3
2337 1 3
2335 1 4
2789 24 3
2869 11 3
2934 1 3
3387 18 3
4229 1 3
4338 1 1
5861 18 3

```

```

preferred_location type yearly_avg_comment_on_travel_page \
57 Medical 3
78 Medical 3
81 Medical 3
118 Medical 3
112 Tour and Travel 3
172 Medical 3
174 Financial 3
181 Entertainment 3
365 Medical 3
589 Financial 3
882 Financial 3
887 Medical 3
1865 Financial 3
1239 Entertainment 3
1399 Financial 3
1464 Other 3
1527 Medical 3
1548 Medical 3
1551 Medical 3
1589 Medical 3
1582 Tour and Travel 3
1642 Financial 3
1644 Financial 3
1761 Entertainment 3
1887 Other 3
1888 Social media 3
1835 Medical 3
1878 Financial 3
1879 Financial 3
2272 Financial 3
2337 Medical 3
2335 Financial 3
2789 Entertainment 3
2869 Financial 3

```

```

2934 Other 3
3387 Entertainment 615
4229 Tour and Travel 215
4338 Entertainment 815
5861 Financial 685

```

```

total_likes_on_outstation_checkin_received \

```

```

57 4614
78 13664
81 2859
118 7484
112 17328
172 2888
174 5686
291 4485
365 5258
488 10555
589 7725
882 4292
887 17856
1865 4035
1239 5238
1399 7538
1464 2966
1527 4614
1548 13664
1551 2859
1589 7484
1582 17328
1642 2888
1644 5686
1761 4485
1887 10555
1888 7725
1835 4292
1878 17856
1879 4035
2272 5238
2337 7538
2335 2966
2789 2856
2869 2644
3387 2075
4229 5392
4338
5861

```

```

week_since_last_outstation_checkin_following_company_page \
57 Yes
78 No
81 No
118 No
112 No
172 No
174 No
291 No
365 No
488 No
589 No
882 No
887 No
1865 No
1239 Yes
1399 No
1464 No
1527 Yes
1548 No
1551 No
1589 No
1582 Yes
1642 No
1644 Yes
1761 No
1887 No
1835 No
1878 No
1879 No
2272 No
2337 No
2335 No
2789 No

```

```

2809      1      No
2934      0      No
3387      0
4229
4538
5861

To exit full screen, press Esc

monthly_avg_comment_on_company_page working_flag \
57      15      No
78      13      No
81      13      No
118      28      Yes
112      12      No
172      12      No
174      11      No
291      12      No
365      14      No
488      13      No
589      23      Yes
882      17      No
867      22      Yes
1865      16      No
1239      28      Yes
1399      17      No
1464      18      No
1527      15      No
1548      13      No
1551      18      No
1589      29      Yes
1582      12      No
1642      12      No
1644      11      No
1761      12      No
1887      17      No
1888      13      No
1835      14      No
1878      13      No
1879      23      Yes
2272      17      No
2337      22      Yes
2535      16      No
2789      29      Yes
2869      17      No
2934      18      No
3387      19      No
4229      22      No
4538      23      No
5861      14      No

travelling_network_rating Adult_flag \
57      3      1
78      2      0
81      1      0
118      3      0
112      3      0
172      1      1
174      1      1
291      4      0
365      3      0
488      4      0
589      3      1
882      1      1
867      3      1
1865      4      0
1239      4      0
1399      4      0
1464      4      0
1527      1      1
1548      1      0
1551      0      0
1589      3      0
1582      3      0
1642      3      0
1644      1      1
1761      4      0
1887      3      0
1888      3      3
1835      3      0
1878      4      0
1879      3      1
2272      1      1
2337      3      1
2535      2      0

2789      4      0
2809      2      1
2934      4      0
3387      4      0
4229      4      3
4538      1      1
5861      3      1

Daily_Avg_mins_spent_on_traveling_page z_score
57      16 -3.014735
78      16 -3.014735
81      16 -3.014735
118      23 -3.014735
112      23 -3.014735
172      13 -3.014735
174      9 -3.014735
291      16 -3.014735
365      16 -3.014735
488      35 -3.014735
589      15 -3.014735
882      4 -3.014735
867      32 -3.014735
1865      4 -3.014735
1239      9 -3.014735
1399      16 -3.014735
1464      4 -3.014735
1527      16 -3.014735
1548      6 -3.014735
1551      23 -3.014735
1589      23 -3.014735
1582      13 -3.014735
1642      9 -3.014735
1644      16 -3.014735
1761      4 -3.014735
1887      16 -3.014735
1888      7 -3.014735
1835      35 -3.014735
1878      15 -3.014735
1879      32 -3.014735
2272      9 -3.014735
2337      16 -3.014735
2535      9 -3.014735
2789      16 -3.014735
2809      5 -3.014735
2934      22 0.881401
3387      28 1.887491
4229      9 35.882471
4538      17 25.623509
5861

In [15]: # Check for outliers in all numerical columns using ZDR
numerical_columns = Data.select_dtypes(include=['int64', 'float64']).columns

for column in numerical_columns:
    outliers = detect_outliers_igd(Data, column)
    print(f'Outliers in {column}: {len(outliers)}')

Outliers in 'yearly_avg_dataststion_checkins': 0
Outliers in 'yearly_avg_dataststion_checkins': 0
Outliers in 'number_in_family': 11
Outliers in 'total_likes_on_notification_checkins_received': 916
Outliers in 'week_since_last_notification_checkins': 0
Outliers in 'monthly_avg_comment_on_company_page': 242
Outliers in 'travelling_network_rating': 0
Outliers in 'Adult_flag': 888
Outliers in 'Daily_Avg_mins_spent_on_traveling_page': 358
Outliers in 'z_score': 48

In [16]: # Calculate ZDR and bounds for 'yearly_avg_comment_on_travel_page'
Q1 = Data['yearly_avg_comment_on_travel_page'].quantile(0.25) # First quartile
Q3 = Data['yearly_avg_comment_on_travel_page'].quantile(0.75) # Third quartile
IQR = Q3 - Q1 # Interquartile range
lower_bound = Q1 - 1.5 * IQR
upper_bound = Q3 + 1.5 * IQR

# Filter the data to remove outliers
Data = Data[(Data['yearly_avg_comment_on_travel_page'] >= lower_bound) &
            (Data['yearly_avg_comment_on_travel_page'] <= upper_bound)]

# Verify the changes
print(f'Data shape after removing outliers: {Data.shape}')

Data shape after removing outliers: (11728, 18)

univariate analysis and bivariate analysis

```

```
In [106]: # Descriptive statistics
print(Data.describe())

# Histogram
Data['Yearly_avg_comment_on_travel_page'].hist(bins=30)
plt.title('Histogram for Yearly_avg_comment_on_travel_page')
plt.xlabel('Yearly_avg_comment_on_travel_page')
plt.ylabel('Frequency')
plt.show()

# Boxplot
sns.boxplot(Data['Yearly_avg_comment_on_travel_page'])
plt.title('Boxplot for Yearly_avg_comment_on_travel_page')
plt.show()
```

```
count      11720.000000
mean      1.000000e+06
std       3.90711e+03
min       1.000000e+06
25%       1.000000e+06
50%       1.000000e+06
75%       1.000000e+06
max       1.011700e+06

count      11720.000000
mean      20150.379083
std       10142.983387
min       1570.000000
25%       16600.000000
50%       20070.000000
75%       40110.000000
max       252430.000000

count      11840.0
mean       8.23246
std       0.07256
min       1.0
25%       1.0
50%       4.0
75%       16.0
max       29.0

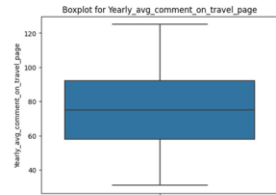
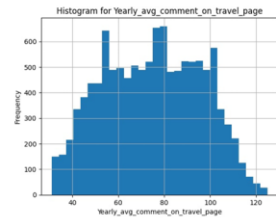
count      11720.000000
mean      74.000956
std       21.000823
min       31.000000
25%       56.000000
50%       75.000000
75%       92.000000
max       125.000000

count      11720.000000
mean      4579.410686
std       4702.570410
min       1000.000000
25%       2030.000000
50%       4947.000000
75%       8090.250000
max       20000.000000

count      11720.000000
mean      3.206035
std       2.18057
min       0.000000
25%       1.000000
50%       3.000000
75%       5.000000
max       11.000000

count      11720.000000
mean      2.711775
std       1.001227
min       1.000000
25%       2.000000
50%       3.000000
75%       4.000000
max       4.000000

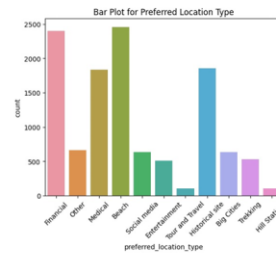
count      11720.000000
mean      13.61081
std       9.072250
min       0.000000
25%       8.000000
50%       12.000000
75%       18.000000
max       270.000000
```



```
In [161]: # Frequency table
print(Data['preferred_location_type'].value_counts())

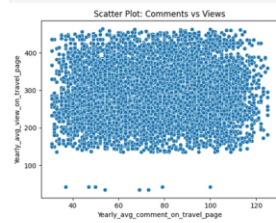
# Bar plot
sns.countplot(x=Data['preferred_location_type'])
plt.title('Bar Plot for Preferred Location Type')
plt.xticks(rotation=45)
plt.show()

Beach      2405
Financial  2296
Historical site 1856
Medical     1811
Other       864
Big Cities  636
Social media 632
Trekking    528
Entertainment 559
Hill Stations 388
Tour and Travel 384
Name: preferred_location_type, dtype: int64
```



```
In [162]: # Scatter plot
sns.scatterplot(x=Data['yearly_avg_comment_on_travel_page'], y=Data['yearly_avg_view_on_travel_page'])
plt.title('Scatter Plot: Comments vs Views')
plt.xlabel('Yearly_avg_comment_on_travel_page')
plt.ylabel('Yearly_avg_view_on_travel_page')
plt.show()

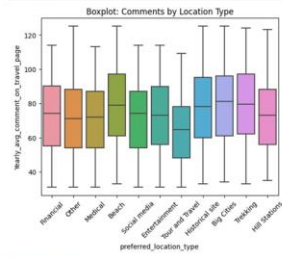
# Correlation
correlation = Data[['yearly_avg_comment_on_travel_page', 'yearly_avg_view_on_travel_page']].corr()
print('Correlation Matrix:\n', correlation)
```



```
Correlation Matrix:
Yearly_avg_comment_on_travel_page    Yearly_avg_comment_on_travel_page    \
Yearly_avg_view_on_travel_page        1.000000
Yearly_avg_view_on_travel_page        0.844778

Yearly_avg_comment_on_travel_page    Yearly_avg_view_on_travel_page
Yearly_avg_comment_on_travel_page    0.844778
Yearly_avg_view_on_travel_page        1.000000

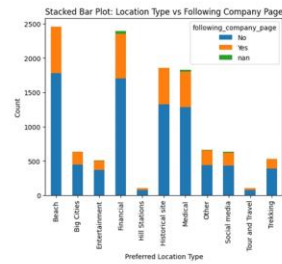
In [105]: # Boxplot
sns.boxplot(data[preferred_location_type], y=data['yearly_avg_comment_on_travel_page'])
plt.title('Boxplot: Comments by Location Type')
plt.xticks(rotation=45)
plt.show()
```



```
In [106]: # Cross-tabulation
crosstab = pd.crosstab(data[preferred_location_type], data[following_company_page])
print(crosstab)

# Stacked Bar plot
crosstab.plot(kind='bar', stacked=True)
plt.title('Stacked Bar Plot: Location Type vs Following Company Page')
plt.xlabel('Preferred Location Type')
plt.ylabel('Count')
plt.show()

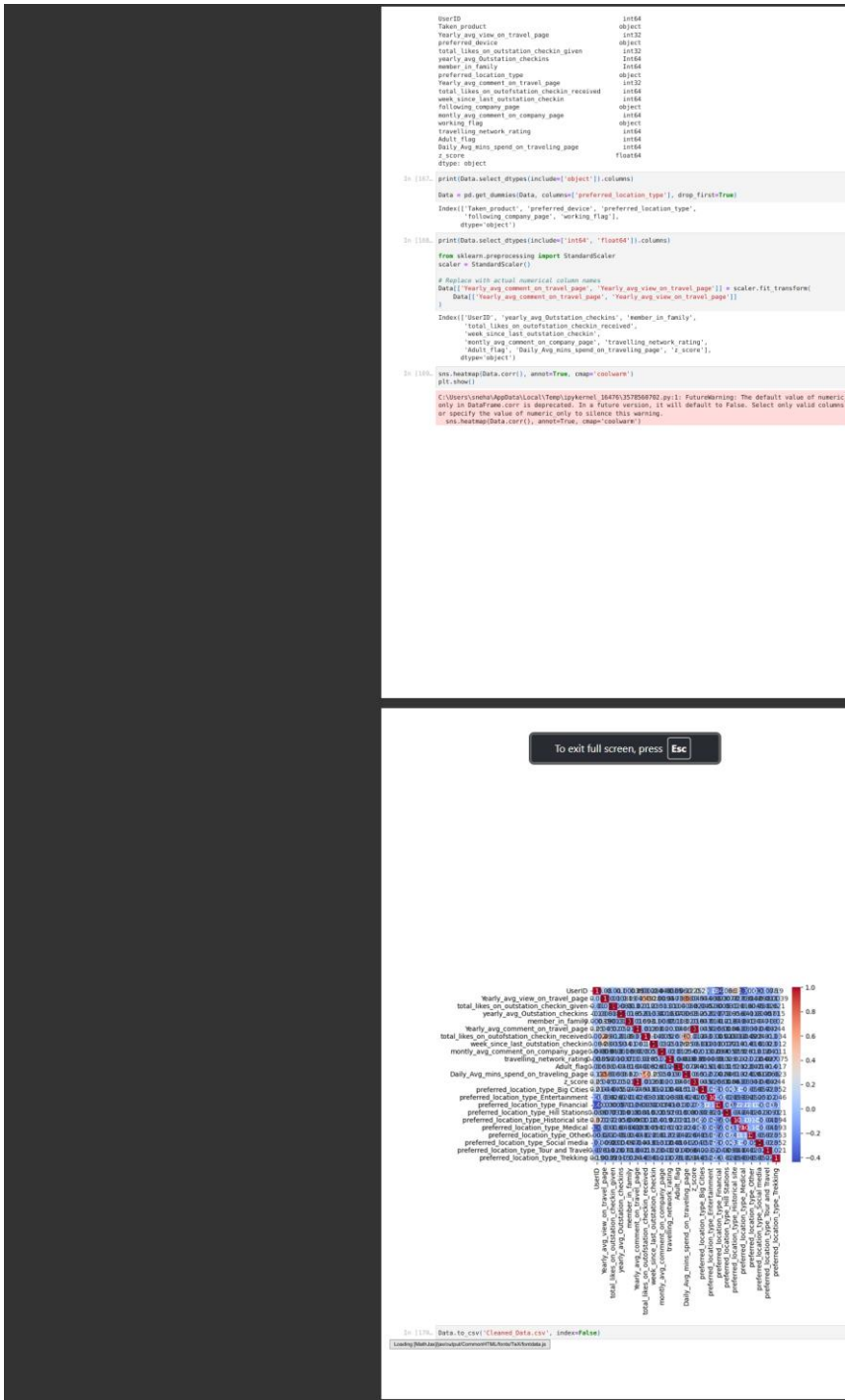
following_company_page    No    Yes    nan
preferred_location_type
Beach                    1789    675     0
Big Cities                464    192     0
Entertainment            372    128    19
Financial                1764    402    47
Hill Stations             76     32     0
Historical site          1139    320     0
Medical                  1284    515    32
Other                     441    218     5
Social media             434    198     0
Tour and Travel           74     20     0
Treking                  392    136     0
```



```
In [107]: print(data.isnull().sum())

UserID                    0
Taken_product            0
Yearly_avg_view_on_travel_page    0
Preferred_Device         53
Total_likes_on_outstation_checkin_given    0
Yearly_avg_outstation_checkin    75
member_in_family         0
preferred_location_type    0
Yearly_avg_comment_on_travel_page    0
Total_likes_on_outstation_checkin_received    0
week_since_last_outstation_checkin    0
following_company_page     0
monthly_avg_comment_on_company_page    0
working_flag              0
travelling_network_rating    0
Adult_Flag                0
Daily_Avg_mins_spend_on_travelling_page    0
dtype: int64

In [108]: print(data.dtypes)
```



```
In [7]: import numpy as np
import matplotlib.pyplot as plt
import pandas as pd
import seaborn as sns

In [8]: from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import classification_report, confusion_matrix

# 1. Load data
data = pd.read_csv('D:\VBA\Semester 4\Project\Data\TM\Python Implementation')

# 2. Feature selection (replace with your actual feature columns)
feature_cols = [
    'Yearly_avg_comment_on_travel_page',
    'Yearly_avg_view_on_travel_page',
    'total_likes_on_outstation_checkin_given',
    'Yearly_avg_outstation_checkins',
    'number_in_family',
    # add more relevant features as needed
]
X = data[feature_cols]
y = data['Taken_product'] # updated target variable

# 3. Train-test split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# 4. Model training
model = LogisticRegression(max_iter=1000)
model.fit(X_train, y_train)

# 5. Evaluation
y_pred = model.predict(X_test)
print("Classification Report:\n", classification_report(y_test, y_pred))
print("Confusion Matrix:\n", confusion_matrix(y_test, y_pred))

Classification Report:
              precision    recall  f1-score   support

     No         0.84         1.00         0.91         1968
     Yes         0.80         0.80         0.80           376

 accuracy         0.84         0.84         0.84         2344
 macro avg         0.82         0.90         0.86         2344
 weighted avg         0.70         0.84         0.77         2344

Confusion Matrix:
[[1968   0]
 [ 376   0]]
```

```
c:\Users\sneha\AppData\Local\Programs\Python\Python310\lib\site-packages\sklearn\
n\metrics\_classification.py:1531: UndefinedMetricWarning: Precision is ill-def
ined and being set to 0.0 in labels with no predicted samples. Use 'zero_divisi
on' parameter to control this behavior.
  _warn_prf(average, modifier, f"(metric.capitalize()) is", len(result))
c:\Users\sneha\AppData\Local\Programs\Python\Python310\lib\site-packages\sklearn\
n\metrics\_classification.py:1531: UndefinedMetricWarning: Precision is ill-def
ined and being set to 0.0 in labels with no predicted samples. Use 'zero_divisi
on' parameter to control this behavior.
  _warn_prf(average, modifier, f"(metric.capitalize()) is", len(result))
c:\Users\sneha\AppData\Local\Programs\Python\Python310\lib\site-packages\sklearn\
n\metrics\_classification.py:1531: UndefinedMetricWarning: Precision is ill-def
ined and being set to 0.0 in labels with no predicted samples. Use 'zero_divisi
on' parameter to control this behavior.
  _warn_prf(average, modifier, f"(metric.capitalize()) is", len(result))

As of now the Data is imbalanced

In [9]: print(data['Taken_product'].value_counts())

No         1968
Yes         376
Name: Taken_product, dtype: int64

In [10]: print(np.unique(y_pred, return_counts=True))

(array(['No'], dtype=object), array([2344], dtype=int64))

In [11]: model = LogisticRegression(max_iter=1000, class_weight='balanced')
model.fit(X_train, y_train)

Out[11]:
LogisticRegression
LogisticRegression(class_weight='balanced', max_iter=1000)

In [12]: from imblearn.over_sampling import SMOTE
smote = SMOTE(random_state=42)
X_resampled, y_resampled = smote.fit_resample(X_train, y_train)
model.fit(X_resampled, y_resampled)

Out[12]:
LogisticRegression
LogisticRegression(class_weight='balanced', max_iter=1000)

In [13]: pip install imbalanced-learn
```



```
WARNING: Failed to remove contents in a temporary directory 'C:\Users\sneha\AppData\Local\Programs\Python\Python310\Lib\site-packages\umpy'.
You can safely remove it manually.

[notice] A new release of pip is available: 25.0.1 -> 25.1.1
[notice] To update, run: python.exe -m pip install --upgrade pip

In [17]: from sklearn.metrics import classification_report, confusion_matrix

print("Classification Report:\n", classification_report(y_test, y_pred))
print("Confusion Matrix:\n", confusion_matrix(y_test, y_pred))

Classification Report:
              precision    recall  f1-score   support

     No       0.84        1.00        0.91       1968
     Yes       0.90        0.60        0.75        376

 accuracy          0.84        0.84       0.84       2344
 macro avg         0.42        0.50        0.46       2344
 weighted avg         0.70        0.84        0.77       2344

Confusion Matrix:
[[1968   0]
 [ 376   0]]

c:\Users\sneha\AppData\Local\Programs\Python\Python310\Lib\site-packages\sklearn\metrics\_classification.py:1531: UndefinedMetricWarning: Precision is ill-defined and being set to 0.0 in labels with no predicted samples. Use 'zero_division' parameter to control this behavior.
  _warn_prf(average, modifier, f"{metric.capitalize()} is", len(result))
c:\Users\sneha\AppData\Local\Programs\Python\Python310\Lib\site-packages\sklearn\metrics\_classification.py:1531: UndefinedMetricWarning: Precision is ill-defined and being set to 0.0 in labels with no predicted samples. Use 'zero_division' parameter to control this behavior.
  _warn_prf(average, modifier, f"{metric.capitalize()} is", len(result))
c:\Users\sneha\AppData\Local\Programs\Python\Python310\Lib\site-packages\sklearn\metrics\_classification.py:1531: UndefinedMetricWarning: Precision is ill-defined and being set to 0.0 in labels with no predicted samples. Use 'zero_division' parameter to control this behavior.
  _warn_prf(average, modifier, f"{metric.capitalize()} is", len(result))

Your model is still predicting only the majority class ("No") for all test samples, even after using SMOTE. Hence trying to check with different model

In [18]: from sklearn.ensemble import RandomForestClassifier

rf = RandomForestClassifier(class_weight='balanced', random_state=42)
rf.fit(X_resampled, y_resampled)
y_pred_rf = rf.predict(X_test)
print(classification_report(y_test, y_pred_rf))
print(confusion_matrix(y_test, y_pred_rf))
```

```
              precision    recall  f1-score   support

     No       0.98        0.98        0.98       1968
     Yes       0.92        0.89        0.90        376

 accuracy          0.97        0.97       0.97       2344
 macro avg         0.95        0.94        0.94       2344
 weighted avg         0.97        0.97        0.97       2344

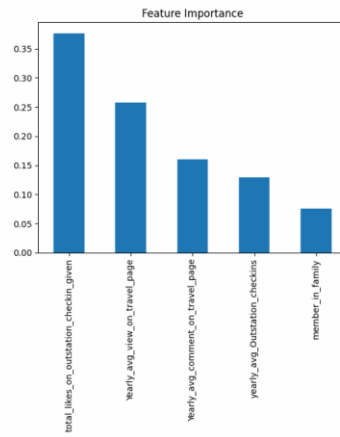
[[1937   31]
 [   41  335]]
```

This output shows that your **Random Forest model** is performing much better than Logistic Regression for the imbalanced data:

- Random Forest model is able to **detect both classes well**, including the minority class ("Yes").
- The model is **not biased toward the majority class** anymore.
- **Precision and recall for both classes are high**, indicating a balanced and effective model.

```
In [19]: importances = rf.feature_importances_
feature_importance = pd.Series(importances, index=feature_cols).sort_values(ascending=False)
print(feature_importance)
feature_importance.plot(kind='bar')
plt.title('Feature Importance')
plt.show()

total_likes_on_outstation_checkin_given    0.376483
Yearly_avg_view_on_travel_page             0.257311
Yearly_avg_comment_on_travel_page          0.160333
yearly_avg_outstation_checkins             0.129739
member_in_family                          0.076335
dtype: float64
```



Each bar represents how much a particular feature contributed to the model's decision-making process for predicting whether a user took the product. Features with higher importance scores had a greater impact on the model's predictions, helping to identify which aspects of user engagement and behavior are most influential in driving ticket purchase decisions. This insight can guide airlines to focus their marketing efforts on the most impactful engagement metrics.

```
In [20]: from sklearn.tree import DecisionTreeClassifier
from sklearn.metrics import classification_report, confusion_matrix

# Train the Decision Tree model on the resampled data
dt = DecisionTreeClassifier(class_weight='balanced', random_state=42)
dt.fit(X_resampled, y_resampled)

# Predict on the test set
y_pred_dt = dt.predict(X_test)
```

```
# Evaluate the model
print("Classification Report:\n", classification_report(y_test, y_pred_dt))
print("Confusion Matrix:\n", confusion_matrix(y_test, y_pred_dt))

Classification Report:
              precision    recall  f1-score   support

     No       0.98       0.97       0.97       1968
     Yes       0.84       0.88       0.86        376

 accuracy          0.95       2344
 macro avg       0.91       0.92       0.92       2344
 weighted avg     0.95       0.95       0.95       2344

Confusion Matrix:
[[1965  43]
 [ 46  330]]
```

- The Decision Tree model is **performing very well** for both classes.
- It is able to **detect the minority class ("Yes")** with high recall (0.88) and good precision (0.84).
- The model is **not biased toward the majority class** and provides a balanced performance.
- **Overall accuracy is high (95%)**, indicating strong predictive power.
- **Overall Performance:**  
Both models perform very well, but the **Random Forest** has slightly higher accuracy, precision, recall, and F1-score for both classes, especially for the minority class ("Yes").
- **Minority Class Detection:**
  - **Random Forest** detects more true positives for "Yes" (335 vs. 330) and has fewer false negatives (41 vs. 46) compared to the Decision Tree.
  - **Precision and recall for "Yes"** are higher in Random Forest, meaning it is better at correctly identifying users who took the product.
- **Generalization:**
  - **Random Forest** is an ensemble of many decision trees, which helps it generalize better and avoid overfitting.
  - **Decision Tree** is more prone to overfitting and may not perform as well on unseen data.
- **Interpretability:**
  - **Decision Tree** is easier to interpret and visualize.
  - **Random Forest** is less interpretable but provides more robust and reliable predictions.

## Conclusion

- **Random Forest** is the better model for your data, offering higher accuracy and better detection of both classes, especially the minority class.
- **Decision Tree** is still a strong model and useful for understanding feature splits, but for deployment or business decisions, Random Forest is preferred due to its superior performance.

#### Model Validation

```
In [21]: from sklearn.model_selection import cross_val_score
scores = cross_val_score(rf, X, y, cv=5, scoring='f1_weighted')
print('Cross-validated F1 scores:', scores)
print('Mean F1 score:', scores.mean())

Cross-validated F1 scores: [0.93584994 0.96312288 0.98216236 0.98836489 0.94831
859]
Mean F1 score: 0.9617917300718826

cross-validated F1 scores show that your Random Forest model is performing
consistently well across different splits of your data.
```

#### Analysis

- **High and Stable Scores:**  
All F1 scores are above 0.93, with a mean of ~0.96. This indicates your model is robust and not overfitting to any particular fold.
- **Low Variance:**  
The scores do not fluctuate much, showing your model generalizes well to unseen data.
- **F1 Score:**  
The F1 score balances precision and recall, which is especially important for imbalanced datasets like yours.

```
In [22]: from sklearn.model_selection import GridSearchCV

# Define the parameter grid
param_grid = {
    'n_estimators': [100, 200],
    'max_depth': [None, 10, 20],
    'min_samples_split': [2, 5, 10],
    'min_samples_leaf': [1, 2, 4],
    'class_weight': ['balanced']
}

# Initialize the model
rf = RandomForestClassifier(random_state=42)

# Set up GridSearchCV
grid_search = GridSearchCV(
    estimator=rf,
    param_grid=param_grid,
    cv=5,
    scoring='f1_weighted',
    n_jobs=-1,
    verbose=2
)
```

```

)
# Fit on the resampled training data
grid_search.fit(X_resampled, y_resampled)

# Best parameters and score
print('Best Parameters:', grid_search.best_params_)
print('Best F1 Score:', grid_search.best_score_)

# Use the best estimator to predict on the test set
best_rf = grid_search.best_estimator_
y_pred_best = best_rf.predict(X_test)
print('Classification Report:\n', classification_report(y_test, y_pred_best))
print('Confusion Matrix:\n', confusion_matrix(y_test, y_pred_best))

Fitting 5 folds for each of 54 candidates, totalling 270 fits
Best Parameters: {'class_weight': 'balanced', 'max_depth': None, 'min_samples_
leaf': 1, 'min_samples_split': 2, 'n_estimators': 200}
Best F1 Score: 0.9748222973821372
Classification Report:

```

	precision	recall	f1-score	support
No	0.98	0.98	0.98	1968
Yes	0.92	0.89	0.91	376
accuracy			0.97	2344
macro avg	0.95	0.94	0.94	2344
weighted avg	0.97	0.97	0.97	2344

```

Confusion Matrix:
[[1938  30]
 [ 40 336]]

```

The tuned Random Forest model maintains excellent performance on both classes, especially the minority class ("Yes"). High precision and recall for both classes indicate the model is both accurate and reliable. The model is well-balanced and generalizes well, as shown by the high cross-validated F1 score and strong test results.

```
In [23]: import joblib

# Save the model
joblib.dump(best_rf, 'random_forest_model.pkl')

# To load the model later:
loaded_model = joblib.load('random_forest_model.pkl')

Out[23]: ['random_forest_model.pkl']

In [24]: joblib.dump(smote, 'smote_object.pkl')

Out[24]: ['smote_object.pkl']

In [28]: import matplotlib.pyplot as plt

class_counts = data['taken_product'].value_counts()
plt.figure(figsize=(6,4))
class_counts.plot(kind='bar', color=['skyblue', 'salmon'])
```

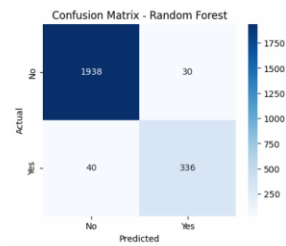
```
plt.title('Class Distribution in the Dataset')
plt.xlabel('Class')
plt.ylabel('Count')
plt.show()
```



The dataset is highly imbalanced, with a much larger number of "No" responses compared to "Yes". This imbalance can negatively impact model performance, making it important to apply techniques like SMOTE to ensure both classes are learned effectively.

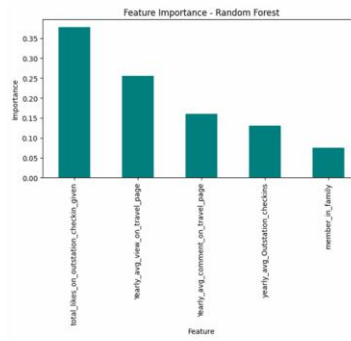
```
In [27]: import seaborn as sns
from sklearn.metrics import confusion_matrix

cm = confusion_matrix(y_test, y_pred_best)
plt.figure(figsize=(8,4))
sns.heatmap(cm, annot=True, fmt='d', cmap='Blues', xticklabels=['No', 'Yes'], yticklabels=['Actual', 'Predicted'])
plt.title('Confusion Matrix - Random Forest')
plt.show()
```



The Random Forest model performs very well, with high precision and recall for both classes. The confusion matrix shows that the model correctly identifies most "Yes" and "No" cases, with very few misclassifications. This indicates the model is reliable for predicting product uptake.

```
In [28]: feature_importance = pd.Series(best_rf.feature_importances_, index=feature_cols)
plt.figure(figsize=(8,4))
feature_importance.plot(kind='bar', color='teal')
plt.title('Feature Importance - Random Forest')
plt.xlabel('Feature')
plt.ylabel('Importance')
plt.show()
```



The "total likes on outstation check-ins given" and "Yearly average view on travel page" are the most influential features in predicting product uptake. This insight can help the airline focus its marketing efforts on customers who are highly engaged with these activities.