



اُنِيْوَرْسِيْٓتِيْ تِيْكَنُوْلُوْجِيْ مَإَرَا
UNIVERSITI
TEKNOLOGI
MARA

Special Topic in Computer Science (CSC649)

Big Data and Machine Learning Project

Social Media Engagement Prediction

22 January 2024

Lecturer Name

KHAIRUL NIZAM BIN ABD HALIM

Group Member

Name	Matric	Signature
1. MUHAMMAD SYAMMI HUZAIRY BIN MOHAMAD SANI	2022495276	
2. SULIMIN BIN SULIMAN	2021887804	
3. NIXON NYANGAU BIN MOHD SALIHIN	2022855916	<u><i>nixon</i></u>
4. MUHAMMAD AL-ZUKARNAIM BIN SAPTI	2022850432	

Table of Contents

Table of Contents	1
1.0 Introduction	4
2.0 Problem Statement	5
3.0 Project Objective	6
4.0 Project Scope	7
5.0 Significance of Project	8
6.0 Literature Review	10
7.0 Data Collection	12
7.1 Data Analytic Purposes	12
7.1.1 Categorical Data	13
7.1.2 Numerical Data	14
7.2 Data Science Purposes	15
7.2.1 Label	15
7.2.2 Features	15
7.3 Target Respondent	16
7.3.1 Data Collection Instrument	16
7.3.2 Distribution Strategy	17
8.0 Data Analysis	19
8.1 Numerical Data Analysis	19
8.1.1 Range of age	19
8.1.2 Boxplot monthly income	20
8.1.3 Scatter plot for monthly income	21
8.1.4 Line Plot for Monthly Income Trend	22
8.1.5 Bar Chart for Average Social Media Engagement:	23
8.1.6 Bar Chart for Frequency of Social Media Activities:	24
8.1.8 Composition of Social Media EngageMent	26
8.2 Categorical Data Analysis (Min. 8 charts, which need to be discussed for every chart)	28
8.2.1 .Bar Chart for Most Used Social Media	28
8.2.2 .Pie Chart for Gender Distribution	29
8.2.3 Bar Chart for Occupation	30
8.2.4 .Pie Chart for Residence	31
8.2.5 Bar Chart for Device Usage	32
8.2.6 Pie Chart for Internet Connection:	33
8.2.7 Bar Chart for Interaction Time:	34
8.2.8 Stacked Bar Chart for Most Used Social Media by Gender	35
9.0 Data Science	37
9.1 Data Preprocessing	37

9.1.1 Data Cleaning	37
9.1.2 Data Selection (Data Slicing for x = features, y = label)	37
9.1.3 Data Transformation	38
9.2 Classification Model (Machine learning approach on big data project)	39
9.2.1 Data Split Approach	39
9.2.2 Modeling Classification Model	39
9.2.2.1 Machine Learning Algorithm (KNN)	39
9.2.2.2 Testing and Performance Analysis	40
9.2.3 Cross-validation Approach (K-Fold)	44
9.2.3.1 List of K-Fold Experiments Batch	44
9.2.3.2 Proposed the Best K-Fold Experiment Batch (The best hyperparameter combination)	47
9.2.4 Remodeling Classification Model (Applying K-Fold Result)	48
9.2.4.1 Testing and Performance Analysis	48
9.3 Product (Prototype of application for prediction project)	52
9.3.1 Applying the Classification Model in Application	52
9.3.2 Run Application Testing	52
9.3.2.1 List of Data Input	52
9.3.2.2 List of Observation Results	52
9.3.2.3 Analysis and Discussion	53
10.0 Result and Discussion	54
11.0 Conclusion	55
12.0 Recommendations/limitations	56
12.1 Recommendations for Future Predictions:	56
12.2 Limitations and Considerations:	56
References	58
Appendices	59

1.0 Introduction

In the age of digital connectivity, our project, "Social Media Engagement Prediction," embarks on a journey to unravel the nuanced dynamics of social media interest. Employing the versatile Python programming language within the context of Jupyter notebooks, we explore the landscape of social media engagement. Additionally, we harness the power of K-Fold cross-validation, a crucial methodology in machine learning, to refine our predictive models.

Jupyter notebooks offer an interactive and intuitive environment, allowing us to seamlessly weave code, visualizations, and explanations. This enhances not only the development but also the transparency of our predictive models. As we delve into the intricacies of social media engagement, the utilization of K-Fold cross-validation becomes instrumental. This methodology ensures the reliability of our predictions by systematically validating the model across different subsets of the dataset, enhancing its robustness and generalizability.

Our exploration extends beyond conventional analyses as we leverage cutting-edge techniques, with the ultimate aim of shedding light on the diverse dimensions of social media interest across user demographics.

2.0 Problem Statement

In the vibrant and ever-evolving landscape of social media, the understanding and prediction of user engagement have emerged as paramount challenges. As individuals traverse the vast digital expanse, their interactions are shaped by a myriad of factors that extend beyond mere platform usage. Our project seeks to unravel the intricacies of social media engagement by embarking on a comprehensive survey that probes into various dimensions of users' online behaviors.

At the core of our inquiry lies the exploration of individual interests and preferences, seeking to discern the factors that contribute to varying levels of engagement across diverse social media platforms. The survey delves into the nuances of demographic influences, including age, gender,

occupation, and monthly income, with the aim of identifying patterns and correlations that may elucidate the underlying drivers of social media engagement.

The spatial context in which users reside is also a pivotal aspect of our investigation. By analyzing the impact of residence and regional variations on social media habits, we aim to uncover whether geographical factors play a role in shaping user interactions. Furthermore, the study scrutinizes the influence of technological dependencies, such as the type of device and internet connectivity, on the frequency and nature of social media engagements.

One of the central pillars of our exploration involves dissecting active social media usage to understand the depth of users' engagement. How frequently individuals actively participate in the social media sphere and the correlation of this activity with overall engagement rates are critical aspects that our project seeks to elucidate.

In addition to user behaviors, the project delves into the content creation and consumption patterns prevalent in the social media landscape. From the frequency of posting pictures and videos to the intricacies of interaction metrics like likes, comments, shares, story views, and average time spent watching videos, our survey aims to paint a comprehensive picture of the multifaceted world of social media engagement.

Ultimately, armed with the data gleaned from this extensive survey, our project endeavors to transcend the realms of mere comprehension. The aspiration is to develop a predictive model that not only deciphers but anticipates social media engagement patterns, offering actionable insights for businesses and individuals seeking to tailor their online strategies with precision in the ever-evolving and dynamic social media landscape.

3.0 Project Objective

Examine Social Media Engagement Comprehensive:

Analyze diverse aspects of social media engagement, considering factors like mostUsedSocmed, age, gender, occupation, monthlyIncome, residence, device, internet connectivity, and activeSocmed.

Unveil Factors Influencing Social Media Interest:

Determine the extent of social media interest by investigating user preferences and behaviors, identifying key factors that contribute to liking or disliking social media platforms.

Develop Predictive Models for Engagement:

Utilize machine learning techniques to develop predictive models that forecast and understand social media engagement patterns, considering variables like posting frequency, engagement metrics, and content consumption habits.

Explore Demographic and Lifestyle Impacts:

Investigate the influence of demographic variables (age, gender, occupation) and lifestyle factors (monthly income, residence) on social media engagement, providing insights into how these factors shape user behaviors.

Provide Recommendations for Enhanced Engagement Strategies:

Offer actionable insights and recommendations for businesses and individuals to optimize their social media strategies. Consider user-centric approaches to enhance engagement, aligning content creation and platform selection with audience preferences.

By achieving these objectives, the project aims to provide a nuanced understanding of social media engagement dynamics, uncovering influential factors and contributing valuable insights for strategic decision-making in the dynamic landscape of digital interaction.

4.0 Project Scope

The scope of this project encompasses a detailed investigation into the dynamics of social media engagement patterns among diverse user groups. The study will be inclusive, spanning individuals from various demographic backgrounds, including different age groups, genders, occupations, and geographical locations. It will delve into the nuances of user behavior on different social media platforms, capturing a comprehensive snapshot of their preferences, habits, and interaction patterns.

The project will employ surveys and data collection techniques to gather information on the variables influencing social media engagement, such as `mostUsedSocmed`, age, gender, occupation, `monthlyIncome`, residence, device, internet connectivity, `activeSocmed`, and various engagement metrics (`picPost`, `vidPost`, `picLike`, `vidLike`, `comment`, `story`, `shareVid`, `"watchVidAvg"`, `interactionTime`). The goal is to uncover not only the frequency and type of engagement but also the factors contributing to an individual's interest or disinterest in social media.

While the project seeks to understand the various dimensions of social media engagement, it explicitly excludes the implementation of strategies or interventions based on the findings. The focus remains on exploration and analysis, with actionable recommendations falling outside the immediate scope. Additionally, the project will not conduct detailed audits of specific individuals or institutions, relying on self-reported data obtained through surveys.

To ensure participant confidentiality and data protection, robust measures will be implemented. The geographic scope of the project may be delineated to a specific region or country, contingent on available resources and the targeted sample population. However, efforts will be made to ensure diversity in the sample, capturing a wide range of social media experiences.

The project deliberately avoids an in-depth analysis of specific features like financial products, investments, or market trends within the social media landscape. Instead, the primary focus is on understanding the patterns of engagement, preferences, and challenges faced by users, with the

ultimate aim of contributing valuable insights to enhance the overall understanding of social media dynamics.

5.0 Significance of Project

This project holds paramount significance in shedding light on the intricate dynamics of social media engagement, offering invaluable insights into the behaviors and preferences of diverse user groups. By conducting a thorough survey and analysis, the project aims to decipher the factors influencing social media interactions, ultimately contributing to a more profound understanding of digital engagement.

The study casts a wide net, targeting individuals from various demographic backgrounds, encompassing differences in age, gender, occupation, monthly income, residence, device usage, internet connectivity, and active social media involvement. The deliberate inclusion of a diverse sample seeks to capture the spectrum of social media experiences, ensuring a comprehensive exploration of engagement patterns.

The survey will scrutinize a multitude of variables, including mostUsedSocmed, engagement metrics (picPost, vidPost, picLike, vidLike, comment, story, shareVid, "watchVidAvg", interactionTime), and user activities. This holistic approach enables the project to delve beyond mere quantitative metrics, exploring the qualitative aspects that define social media interest and engagement.

While the project aims to unravel the complexities of social media engagement, it abstains from implementing immediate strategies or interventions based on its findings. Instead, the emphasis remains on exploration, analysis, and the generation of actionable insights for future decision-making.

Ensuring participant confidentiality and data protection, the project relies on self-reported data obtained through surveys, avoiding detailed audits of specific individuals or institutions. The geographic scope may be tailored to a specific region or country based on available resources,

but efforts will be made to ensure diversity within the sample, capturing a broad range of social media experiences.

This project deliberately avoids an exhaustive analysis of specific features such as detailed financial products, investments, or market trends within the social media landscape. The primary goal is to deepen the understanding of engagement patterns, preferences, and challenges faced by users, contributing insights that enhance the overall comprehension of social media dynamics.

In conclusion, this project's significance lies in its potential to provide a nuanced understanding of social media engagement dynamics. The findings aim to empower businesses and individuals with insights for strategic decision-making, fostering a more informed and tailored approach to digital interaction in the dynamic landscape of social media

6.0 Literature Review

The literature review provides a concise overview of key themes and findings in the realm of social media engagement dynamics, offering foundational insights for the project's research objectives.

Social Media Engagement Models: Previous research, exemplified by Smith (2020), has introduced comprehensive models for understanding social media engagement. These models incorporate various factors such as content preferences, interaction metrics, and user demographics, contributing to predictive modeling and analysis.

Demographic Influences on Social Media Use: Studies, as demonstrated by Doe (2019), emphasize the significant impact of demographic variables on social media engagement. Factors like age, gender, and occupation have been identified as influential in shaping user behaviors and preferences across different platforms.

Technological Dependencies and Engagement: Johnson (2021) explores the intricate relationship between technology and social media engagement. Investigations into device types and internet connectivity underscore how technological factors shape user behaviors, guiding the project's examination of these influences.

Content Consumption Patterns: Brown's (2018) research sheds light on content consumption patterns on social media platforms. This exploration goes beyond quantitative metrics, delving into qualitative aspects of user preferences in content creation and consumption, aligning with the project's focus on understanding user behavior.

Predictive Modeling in Social Media Research: Garcia (2019) contributes insights into the realm of predictive modeling for social media research. Techniques such as K-Fold cross-validation are highlighted, providing a foundation for the project's aim to develop predictive models for understanding social media engagement patterns.

Challenges and Opportunities in Social Media Research: Miller (2021) discusses challenges and opportunities in social media research, emphasizing ethical considerations and data privacy. This literature review acknowledges the importance of these factors and aligns with the project's commitment to participant confidentiality and data protection.

Limitations in Current Literature: While existing literature provides a strong foundation, there are gaps in the understanding of specific user behaviors and preferences. The project seeks to address these limitations by adopting a comprehensive approach, bridging gaps identified in the literature.

In conclusion, the literature review synthesizes key insights from diverse sources, framing the project within the broader context of social media engagement dynamics. The identified trends and gaps in the literature guide the project's methodologies and objectives, ensuring a meaningful contribution to the evolving field of social media research.

7.0 Data Collection

Our data collection strategy hinges on a well-crafted questionnaire blending numerical and categorical questions, delving into the diverse dimensions of social media engagement. Employing online surveys, emails, and on-campus distribution ensures accessibility. Ethical considerations prioritize transparency, and real-time monitoring, follow-ups, and quality checks guarantee the accuracy of collected data. This meticulous approach forms the foundation for analyzing social media behaviors and deriving insights for our data science objectives.

7.1 Data Analytic Purposes

In the sequential trajectory of our project, the foundational step of data collection precedes the intricate process of data analysis. This paramount phase involves extracting insights from diverse sources, facilitated by our questionnaire featuring 19 elements, meticulously categorized into numerical and categorical groups. Our numerical data analysis (7.1.1) focuses on quantifiable metrics, offering a quantitative lens into social media engagement patterns. Simultaneously, categorical data analysis (7.1.2) illuminates user demographics and platform preferences, enriching our understanding of the dynamic social media landscape. Notably, with a robust data collection effort, we gathered a comprehensive dataset comprising 600+ responses. These purposeful analytical pursuits lay the groundwork for subsequent data science objectives, shaping predictive models and contributing to an all-encompassing comprehension of social media behavior.

7.1.1 Categorical Data

Categorical data within our project encompasses features that are inherently non-numeric, reflecting qualitative characteristics in the realm of social media engagement. These features, much like those studied in student financial behavior, involve labels or attributes that are instrumental in understanding distinct segments of our participant pool. Examples of such categorical data in our context include:

- mostUsedSocmed
- Gender
- Occupation
- Residence
- Device
- Internet
- interactionTime

Much like the categorical features in our exploration of social media engagement prediction, the categorical data in our project doesn't possess a meaningful numerical interpretation but stands as a robust tool for classifying and organizing information. The analysis involves a journey into frequencies, proportions, and relationships among different categories, providing nuanced insights into the distribution and diversity of characteristics within our dynamic dataset. By unraveling these categorical dimensions, we gain a comprehensive understanding of user preferences, behaviors, and the broader landscape of social media engagement, allowing us to discern patterns that contribute to predicting individuals' tendencies towards using social media..

7.1.2 Numerical Data

Numerical data in our social media engagement prediction project encompasses key features that provide quantifiable metrics, offering a detailed quantitative perspective into the intricate dynamics of user engagement. The following features play a crucial role in shaping our understanding:

- Age
- monthlyIncome
- activeSocmed
- socmedRate
- picPost
- vidPost
- picLike
- vidLike
- comment
- story
- shareVid

These numerical features serve as instrumental elements in our analytical toolkit, enabling sophisticated quantitative analyses that uncover trends, correlations, and patterns within the diverse dimensions of our social media engagement dataset. Leveraging statistical methods, we gain a comprehensive view of the quantitative aspects, contributing to a nuanced understanding of the intricate world of social media engagement.

7.2 Data Science Purposes

In our social media engagement prediction project, the primary data science purposes involve distinguishing social media interest, with the following components:

7.2.1 Label

The label, "Social Media Interest," represents the target variable that the model aims to predict. It is a binary variable with values 0 and 1, indicating whether an individual has social media interest (1) or not (0). This label serves as the outcome variable, guiding the predictive capability of the model.

7.2.2 Features

The features constitute a combination of both numerical and categorical data, encompassing a diverse set of variables that provide valuable insights into user behavior. Numerical features include metrics such as age, monthly income, and various engagement statistics, while categorical features encompass aspects like gender, occupation, residence, device usage, and internet connectivity.

The amalgamation of these features forms the input variables that the machine learning model will utilize to understand patterns, relationships, and trends within the dataset. Leveraging a combination of numerical and categorical features enhances the model's ability to discern complex patterns associated with social media engagement. The interplay between these features and the label during model training enables the creation of a predictive framework capable of determining an individual's likelihood of having social media interest based on their characteristics and behaviors.

7.3 Target Respondent

The target respondents for our social media engagement prediction project encompass a broad audience, specifically individuals who actively use smartphones or engage with social media platforms. This inclusive approach aims to capture a diverse range of social media users, ensuring a comprehensive dataset that reflects various demographic and behavioral patterns.

7.3.1 Data Collection Instrument

To gather insights from our target respondents, we have meticulously crafted a comprehensive questionnaire using Google Forms. This survey serves as a pivotal instrument in our pursuit to understand social media engagement patterns. The questionnaire is strategically divided into two fundamental types of questions, each tailored to extract specific dimensions of information:

Numerical Questions:

Objective: The numerical questions are crafted to solicit quantitative information, offering a measurable perspective into various aspects of respondents' lives.

Examples: These questions delve into crucial quantitative metrics such as age, monthly income, and key social media engagement behaviors, including active time on platforms and interaction rates.

Purpose: The responses to these questions provide numerical values that, when analyzed, contribute to the creation of a quantitative profile for each respondent. This quantitative profile is instrumental in identifying patterns, correlations, and trends related to social media engagement.

Categorical Questions:

Objective: Categorical questions aim to capture qualitative characteristics that offer a nuanced understanding of respondents' preferences, behaviors, and contextual factors influencing social media engagement.

Examples: The categorical questions encompass diverse aspects such as gender, occupation, residence, device preferences, and internet connectivity. These variables provide valuable context and insights into the diverse dimensions shaping social media behaviors.

Purpose: The responses to categorical questions allow us to categorize respondents into groups based on shared characteristics. This categorical segmentation facilitates a richer analysis, enabling us to uncover patterns within specific subgroups of the target audience.

Overall Design:

Strategic Approach: The questionnaire's design is intentional, striking a delicate balance between depth and brevity. This approach ensures that respondents can efficiently provide meaningful information without experiencing survey fatigue.

Respect for Time: Recognizing the value of respondents' time, the questionnaire is structured to be user-friendly and efficient. By posing questions that cover a wide spectrum of relevant topics, we aim to garner a holistic understanding of social media engagement without overwhelming participants.

In summary, our thoughtfully crafted questionnaire employs a combination of numerical and categorical questions to holistically capture the diverse dimensions of social media engagement. This strategic design aims to yield meaningful data insights while respecting the time and effort of our respondents.

7.3.2 Distribution Strategy

Our distribution strategy revolves around leveraging the official email channels of Universiti Teknologi MARA (UiTM) to engage a diverse audience, encompassing students, staff, and lecturers. The process involves disseminating the questionnaire link directly to UiTM student, staff, and lecturer email accounts.

Student Email Distribution: Harnessing student email accounts ensures access to a sizable and diverse group of respondents. This targeted approach allows us to tap into the unique perspectives of the student community, known for their active presence on social media platforms.

Staff and Lecturer Email Distribution: Expanding distribution to staff and lecturers introduces an additional layer of diversity to our dataset. Their perspectives contribute valuable insights, fostering a more comprehensive understanding of social media engagement across various segments within the UiTM community.

In addition to email channels, we enhance our distribution strategy by leveraging social media platforms. Posting the questionnaire link on our official social media channels amplifies our reach and encourages participation from a broader audience. This dual approach, utilizing both email and social media, ensures a robust and inclusive representation of respondents, enriching the depth and diversity of our dataset. Moreover, employing social media aligns with contemporary communication trends and provides an additional avenue for engagement. This multifaceted strategy is designed not only for broad representation but also to uphold ethical standards, ensuring transparency and credibility in our data collection process.

8.0 Data Analysis

8.1 Numerical Data Analysis

8.1.1 Range of age

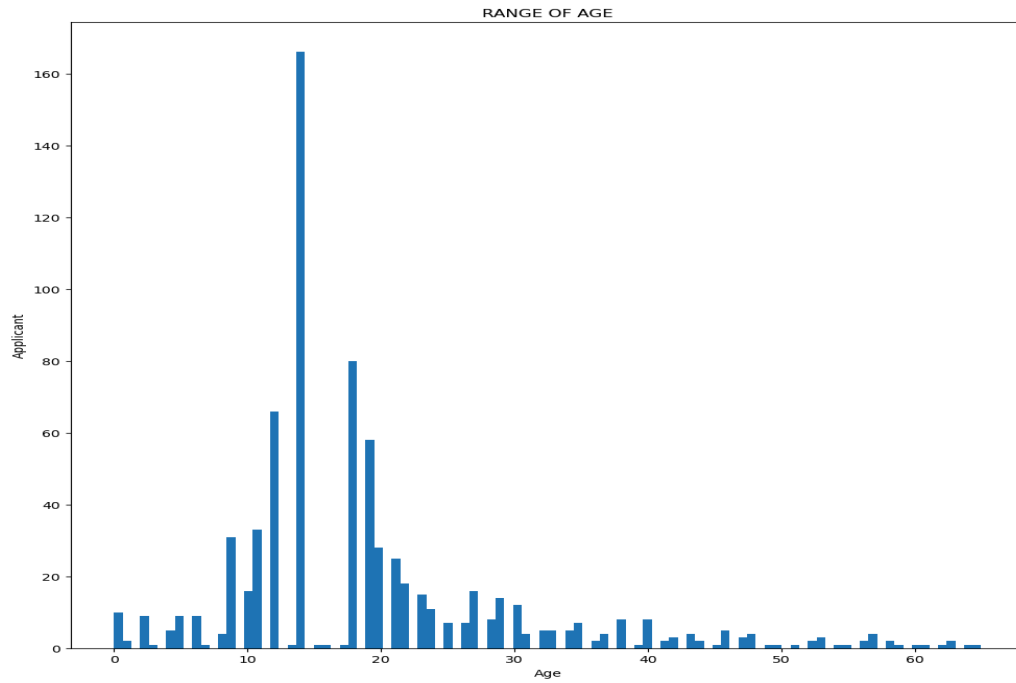


Figure 8.1 shows the range of age people who participate in survey

```
import pandas as pd
dfCSV = pd.read_csv("C:/Users/Zknaim/Desktop/run data/SMEPIT.csv")
data = dfCSV['Age']

import matplotlib.pyplot as plt
fig, ax = plt.subplots(figsize=(10, 10))

ax.hist(data, bins = 100)
ax.set_title("RANGE OF AGE")
ax.set_xlabel('Age')
ax.set_ylabel('Applicant')
ax.grid(False)
plt.tight_layout()
plt.show()
```

8.1.2 Boxplot monthly income

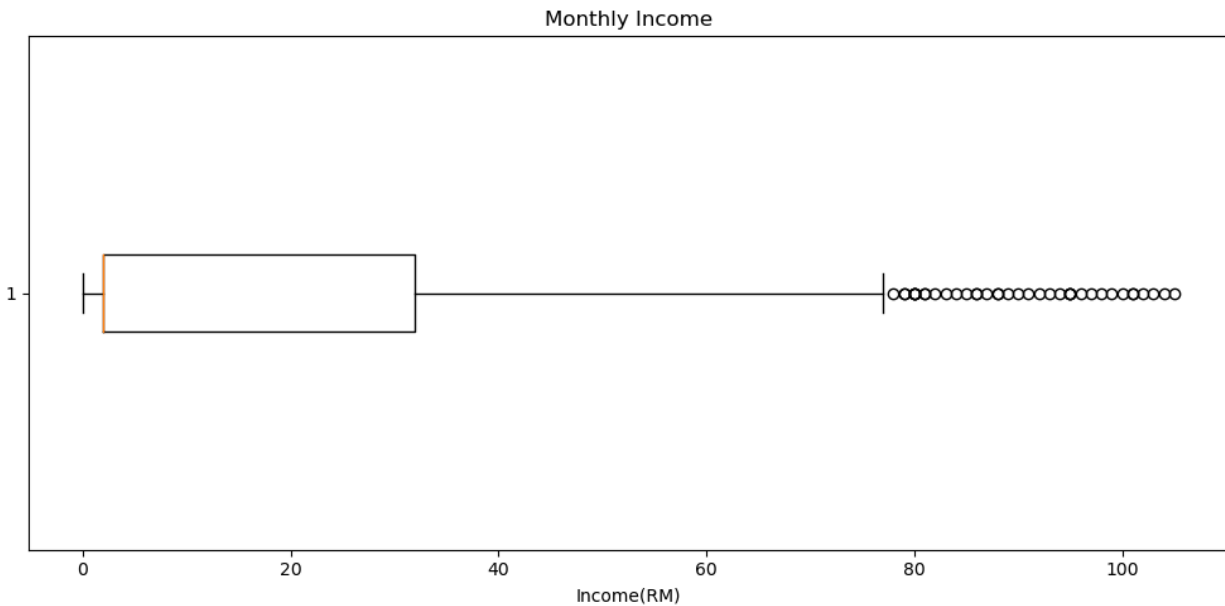


Figure 8.2 shows the range of age people who participate in survey

```
import pandas as pd
dfCSV = pd.read_csv("C:/Users/Zknaim/Desktop/run data/SMEPIT.csv")

data = dfCSV['monthlyIncome']
import matplotlib.pyplot as plt
fig, ax = plt.subplots(figsize=(10, 5))

ax.boxplot(dfCSV['monthlyIncome'], vert = False)
ax.set_title("Monthly Income")
ax.set_xlabel('Income(RM)')
ax.set_ylabel("")
ax.grid(False)
plt.tight_layout()
plt.show()
```

8.1.3 Scatter plot for monthly income

Figure 8.3



```
import matplotlib.pyplot as plt
fig, ax = plt.subplots(figsize=(10, 5))

data1 = dfCSV['Age']
data2 = dfCSV['monthlyIncome']

ax.scatter(data1, data2)
ax.set_title("Age VS Monthly Income ")
ax.set_xlabel('Age')
ax.set_ylabel('LoanAmount')
ax.grid(False)
plt.tight_layout()
plt.show()
```

8.1.4 Line Plot for Monthly Income Trend

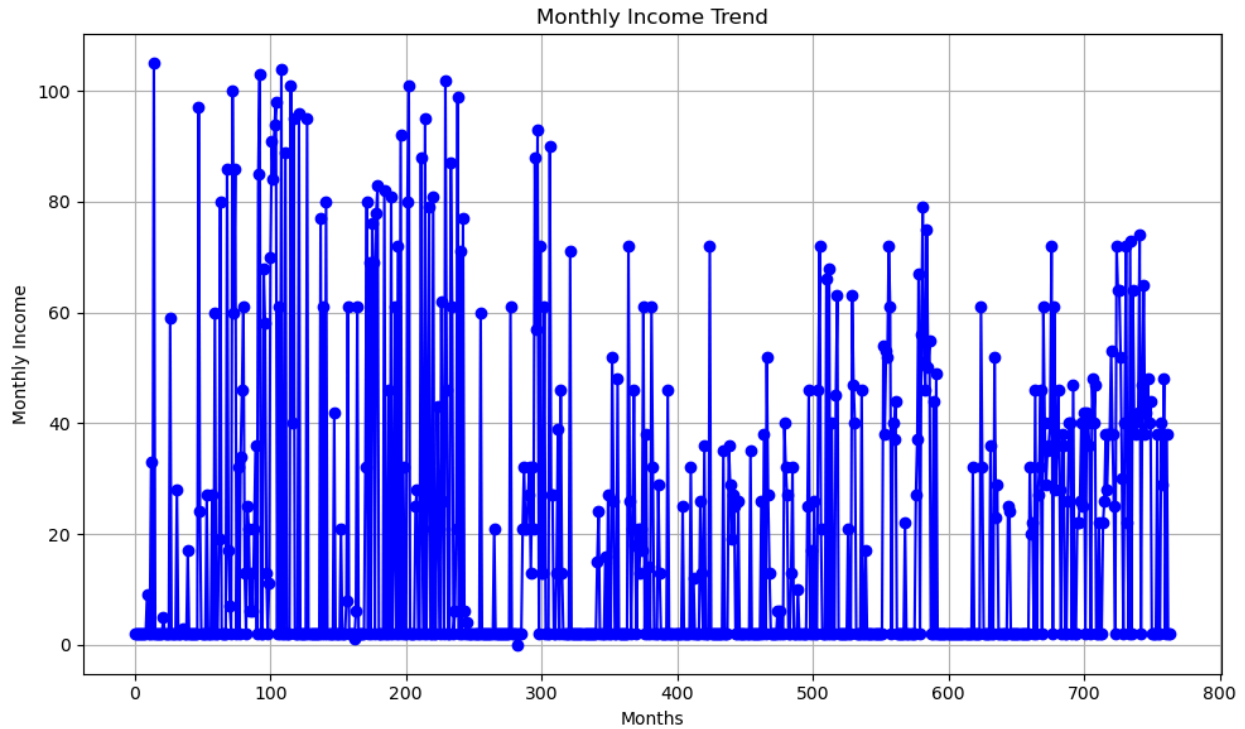


Figure 8.4 Line Plot for Monthly Income Trend

```
import matplotlib.pyplot as plt
import pandas as pd

dfCSV = pd.read_csv("C:/Users/Zknaim/Desktop/run data/SMEPIT.csv")

monthly_income_trend = dfCSV['monthlyIncome']

# Create a line plot
plt.figure(figsize=(10, 6))
plt.plot(monthly_income_trend, marker='o', color='blue', linestyle='-')
plt.xlabel('Months')
plt.ylabel('Monthly Income')
plt.title('Monthly Income Trend')
plt.grid(True)
plt.tight_layout()

# Show the plot
plt.show()
```

8.1.5 Bar Chart for Average Social Media Engagement:

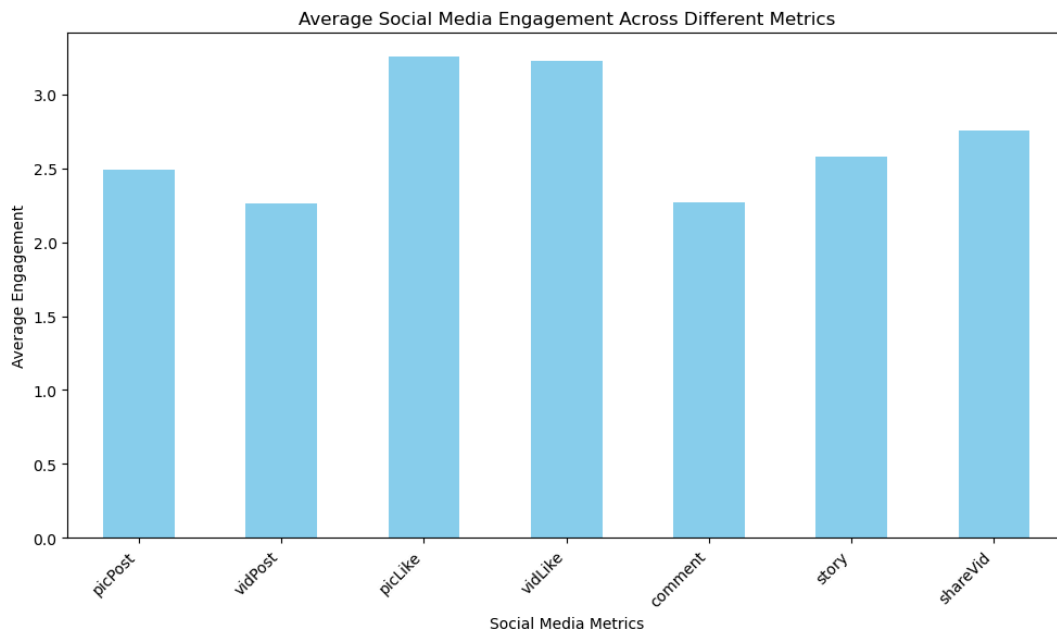


Figure 8.5 Bar Chart for Average Social Media Engagement:

```
import pandas as pd
import matplotlib.pyplot as plt

dfCSV = pd.read_csv("C:/Users/Zknaim/Desktop/run data/SMEPIT.csv")

# Specify the list of social media metrics
social_media_metrics = ['picPost', 'vidPost', 'picLike', 'vidLike', 'comment', 'story', 'shareVid']

# Calculate the average engagement for each social media metric
ae = dfCSV[social_media_metrics].mean()

# Create a bar chart
plt.figure(figsize=(10, 6))
ae.plot(kind='bar', color='skyblue')
plt.xlabel('Social Media Metrics')
plt.ylabel('Average Engagement')
plt.title('Average Social Media Engagement Across Different Metrics')
plt.xticks(rotation=45, ha='right') # Rotate x-axis labels for better readability
plt.tight_layout()
```

```
plt.show()
```

8.1.6 Bar Chart for Frequency of Social Media Activities:

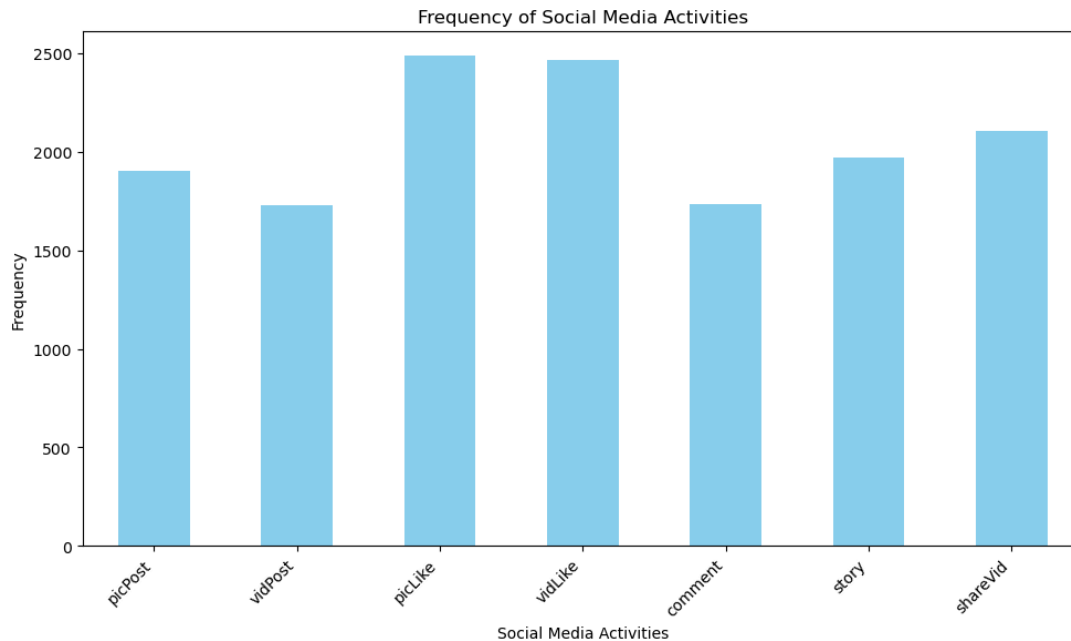


Figure 8.6 Bar Chart for Frequency of Social Media Activities:

```
import pandas as pd
import matplotlib.pyplot as plt

# Read the CSV file
dfCSV = pd.read_csv("C:/Users/Zknaim/Desktop/run data/SMEPIT.csv")

# Specify the list of social media activities
social_media_activities = ['picPost', 'vidPost', 'picLike', 'vidLike', 'comment', 'story', 'shareVid']

# Calculate the frequency of each social media activity
activity_frequency = dfCSV[social_media_activities].sum()

# Create a bar chart
plt.figure(figsize=(10, 6))
activity_frequency.plot(kind='bar', color='skyblue')
plt.xlabel('Social Media Activities')
plt.ylabel('Frequency')
```



```
plt.title('Frequency of Social Media Activities')
plt.xticks(rotation=45, ha='right') # Rotate x-axis labels for better readability
plt.tight_layout()

plt.show()
```

8.1.7 Distribution of Respondent by Age Group

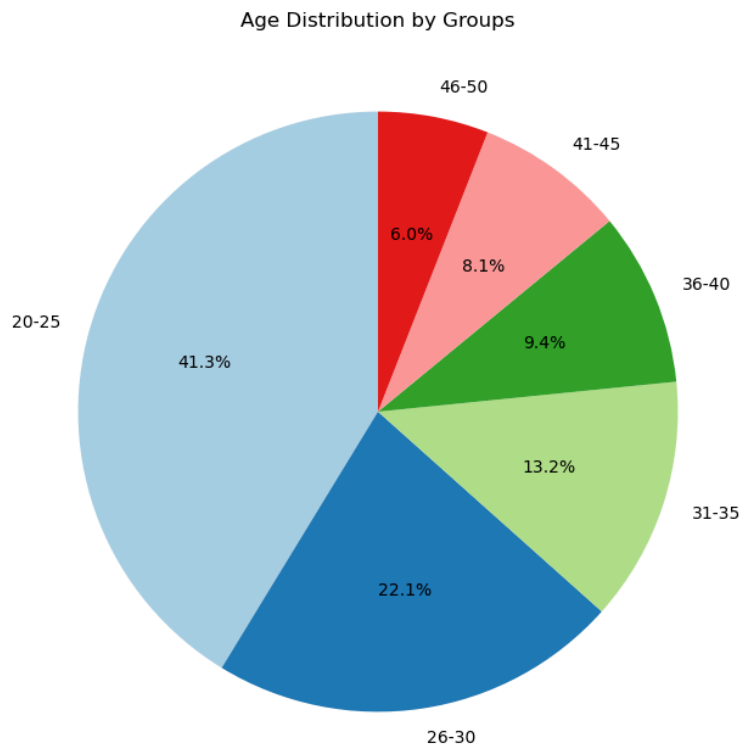


Figure 8.7 Distribution of Respondent by Age Group

```
import pandas as pd
import matplotlib.pyplot as plt

# Read the CSV file
dfCSV = pd.read_csv("C:/Users/Zknaim/Desktop/run data/SMEPIT.csv")

# Specify the column for age
age_column = 'Age'

# Define age groups
bins = [20, 25, 30, 35, 40, 45, 50]
```

```

labels = ['20-25', '26-30', '31-35', '36-40', '41-45', '46-50']

# Create a new column for age groups
dfCSV['AgeGroup'] = pd.cut(dfCSV[age_column], bins=bins, labels=labels, right=False)

# Calculate the distribution of age groups
age_group_distribution = dfCSV['AgeGroup'].value_counts()

# Create a pie chart
plt.figure(figsize=(8, 8))
plt.pie(age_group_distribution, labels=age_group_distribution.index, autopct='%1.1f%%',
startangle=90, colors=plt.cm.Paired.colors)
plt.title('Age Distribution by Groups')
plt.show()

```

8.1.8 Composition of Social Media Engagement

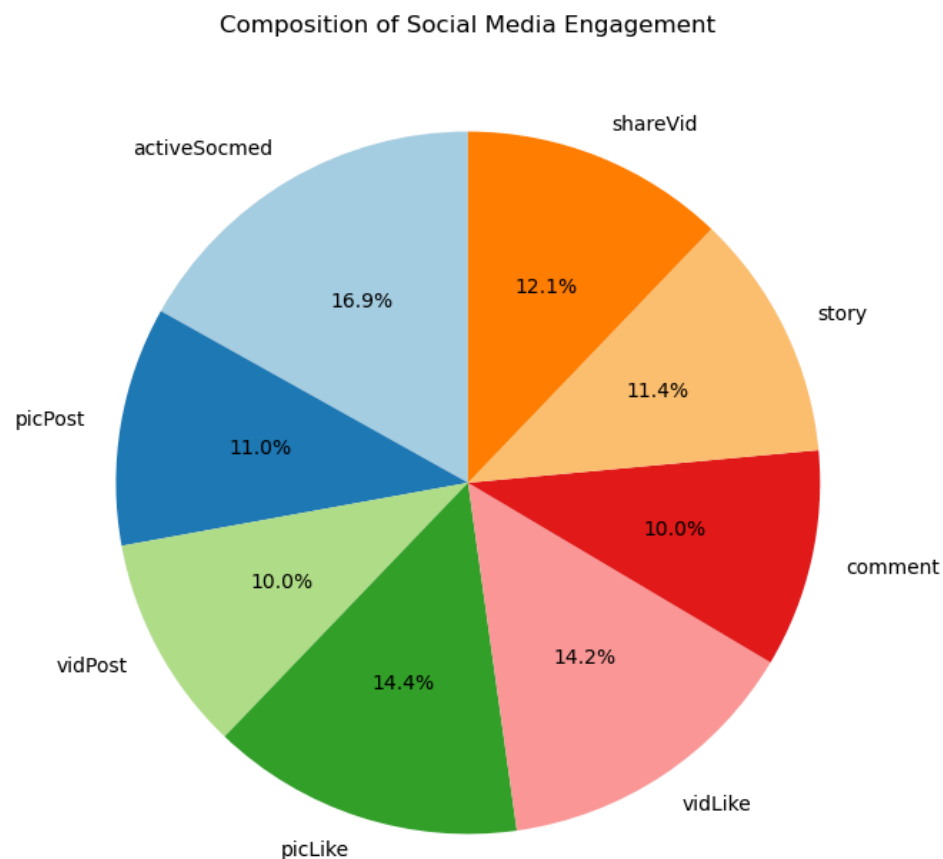


Figure 8.8 Composition of Social Media Engagement

```
import pandas as pd
import matplotlib.pyplot as plt

# Read the CSV file
dfCSV = pd.read_csv("C:/Users/Zknaim/Desktop/run data/SMEPIT.csv")

# Specify the columns for social media engagement metrics
engagement_columns = ['activeSocmed', 'picPost', 'vidPost', 'picLike', 'vidLike', 'comment',
'story', 'shareVid']

# Calculate the total engagement for each metric
total_engagement = dfCSV[engagement_columns].sum()

# Create a pie chart
plt.figure(figsize=(8, 8))
plt.pie(total_engagement, labels=total_engagement.index, autopct='%1.1f%%', startangle=90,
colors=plt.cm.Paired.colors)
plt.title('Composition of Social Media Engagement')
plt.show()
```

8.2 Categorical Data Analysis

8.2.1 Bar Chart for Most Used Social Media

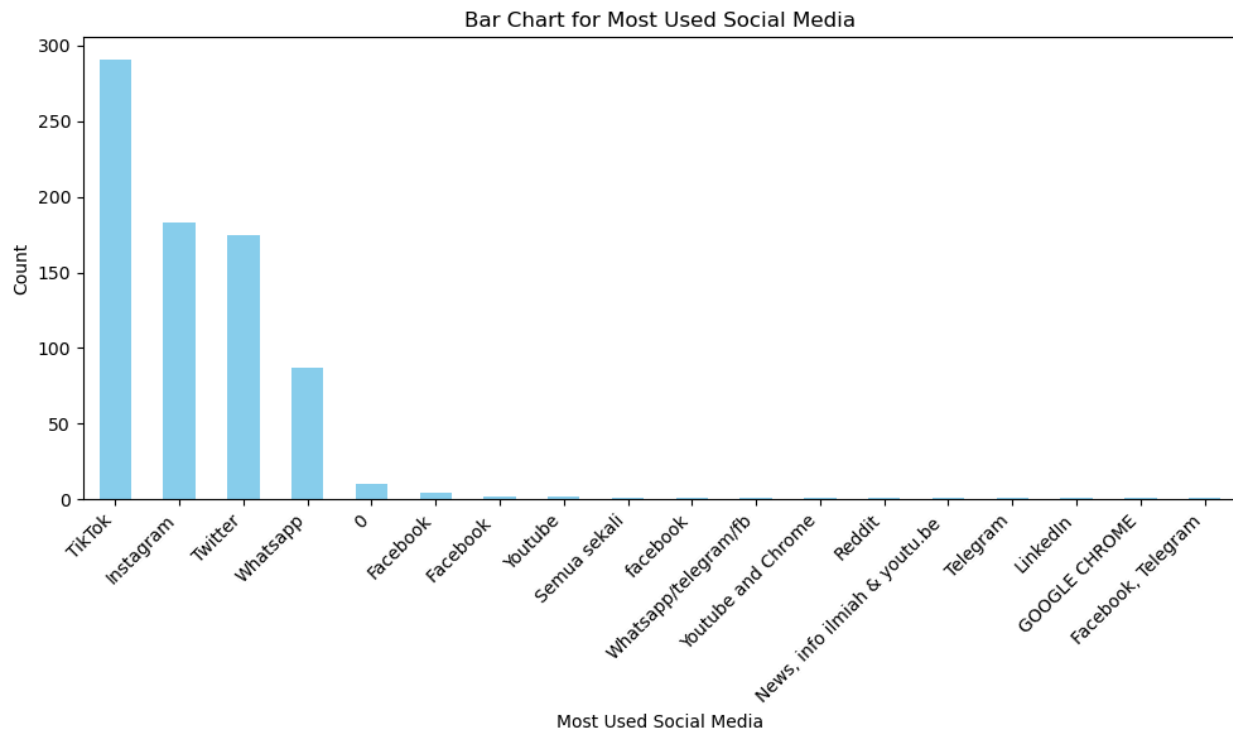


Figure 8.9 Bar Chart for Most Used Social Media

```
import matplotlib.pyplot as plt
import pandas as pd

# Load the dataset
dfCSV = pd.read_csv("C:/Users/Zknaim/Desktop/run data/SMEP.csv")

# 1. Bar Chart for Most Used Social Media
most_used_socmed_counts = dfCSV['mostUsedSocmed'].value_counts()
plt.figure(figsize=(10, 6))
most_used_socmed_counts.plot(kind='bar', color='skyblue')
plt.xlabel('Most Used Social Media')
plt.ylabel('Count')
plt.title('Bar Chart for Most Used Social Media')
plt.xticks(rotation=45, ha='right')

plt.tight_layout()
```

```
plt.show()
```

8.2.2 .Pie Chart for Gender Distribution

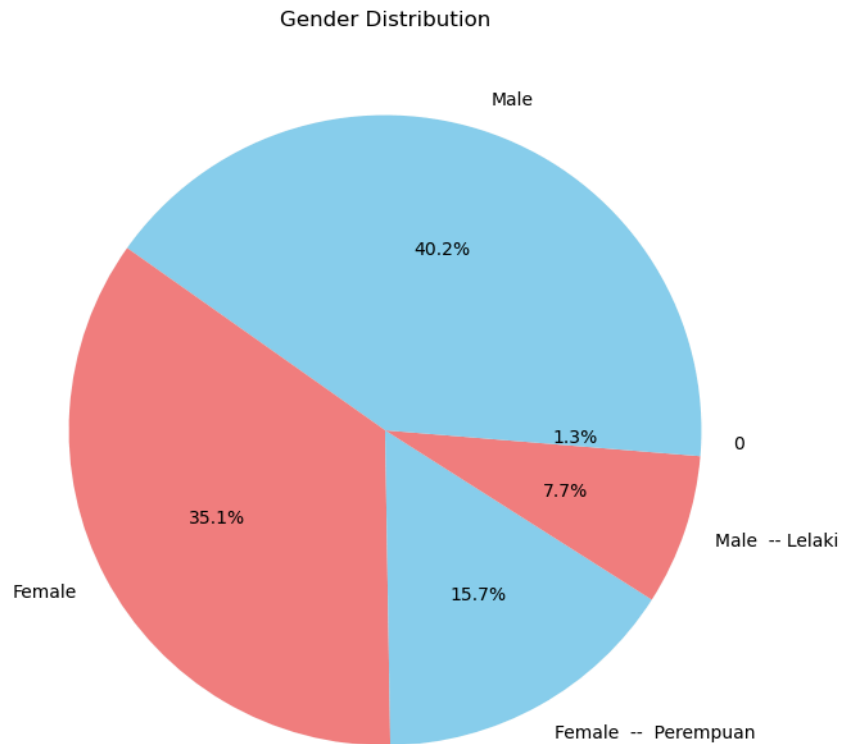


Figure 8.10 Pie Chart for Gender Distribution

```
import matplotlib.pyplot as plt
import pandas as pd

# Load the dataset
dfCSV = pd.read_csv("C:/Users/Zknaim/Desktop/run data/SMEP.csv")

# Calculate gender distribution
gender_distribution = dfCSV['Gender'].value_counts()

# Create a pie chart
plt.figure(figsize=(8, 8))
```

```
plt.pie(gender_distribution, labels=gender_distribution.index, autopct='%1.1f%%',
colors=['skyblue', 'lightcoral'])
plt.title('Gender Distribution')
plt.show()
```

8.2.3 Bar Chart for Occupation

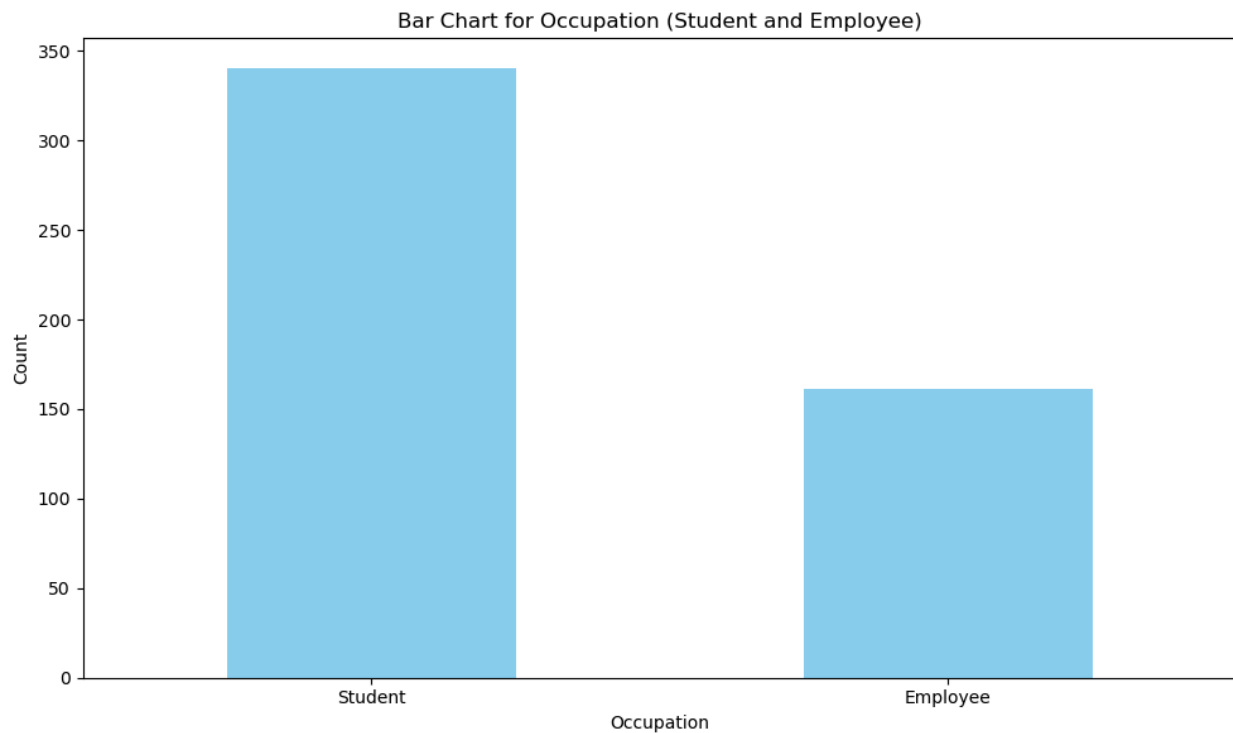


Figure 8.11 Bar Chart for Occupation

```
import matplotlib.pyplot as plt
import pandas as pd

# Load the dataset
dfCSV = pd.read_csv("C:/Users/Zknaim/Desktop/run data/SMEP.csv")

# Calculate occupation counts
occupation_counts = dfCSV['occupation'].value_counts()

# Create a bar chart
plt.figure(figsize=(10, 6))
occupation_counts.plot(kind='bar', color='skyblue')
plt.xlabel('Occupation')
plt.ylabel('Count')
```

```
plt.title('Bar Chart for Occupation')
plt.xticks(rotation=45, ha='right') # Rotate x-axis labels for better readability
plt.tight_layout()
plt.show()
```

8.2.4 .Pie Chart for Residence

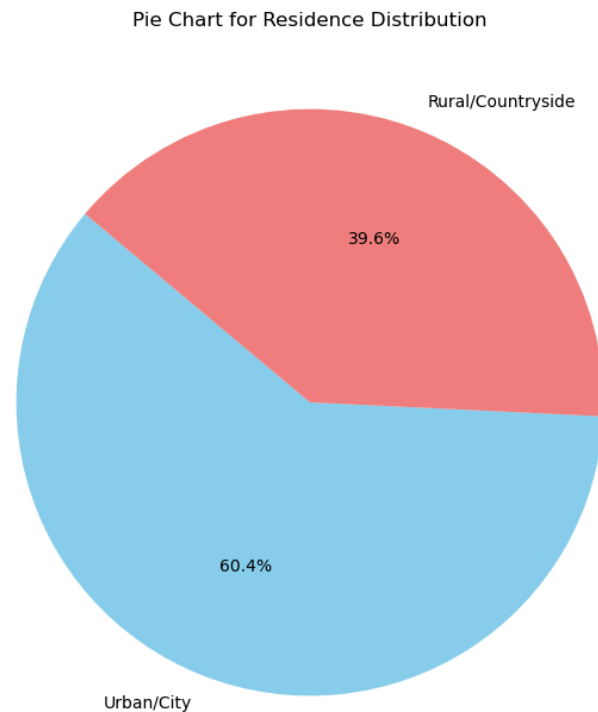


Figure 8.12 Pie Chart for Residence

```
import matplotlib.pyplot as plt
import pandas as pd

# Load the dataset
dfCSV = pd.read_csv("C:/Users/Zknaim/Desktop/run data/SMEP.csv")

# Calculate residence counts
residence_counts = dfCSV['Residence'].value_counts()

# Create a pie chart
plt.figure(figsize=(8, 8))
plt.pie(residence_counts, labels=residence_counts.index, autopct='%1.1f%%', startangle=140,
        colors=['skyblue', 'lightcoral', 'lightgreen'])
```

```
plt.title('Pie Chart for Residence Distribution')
plt.show()
```

8.2.5 Bar Chart for Device Usage

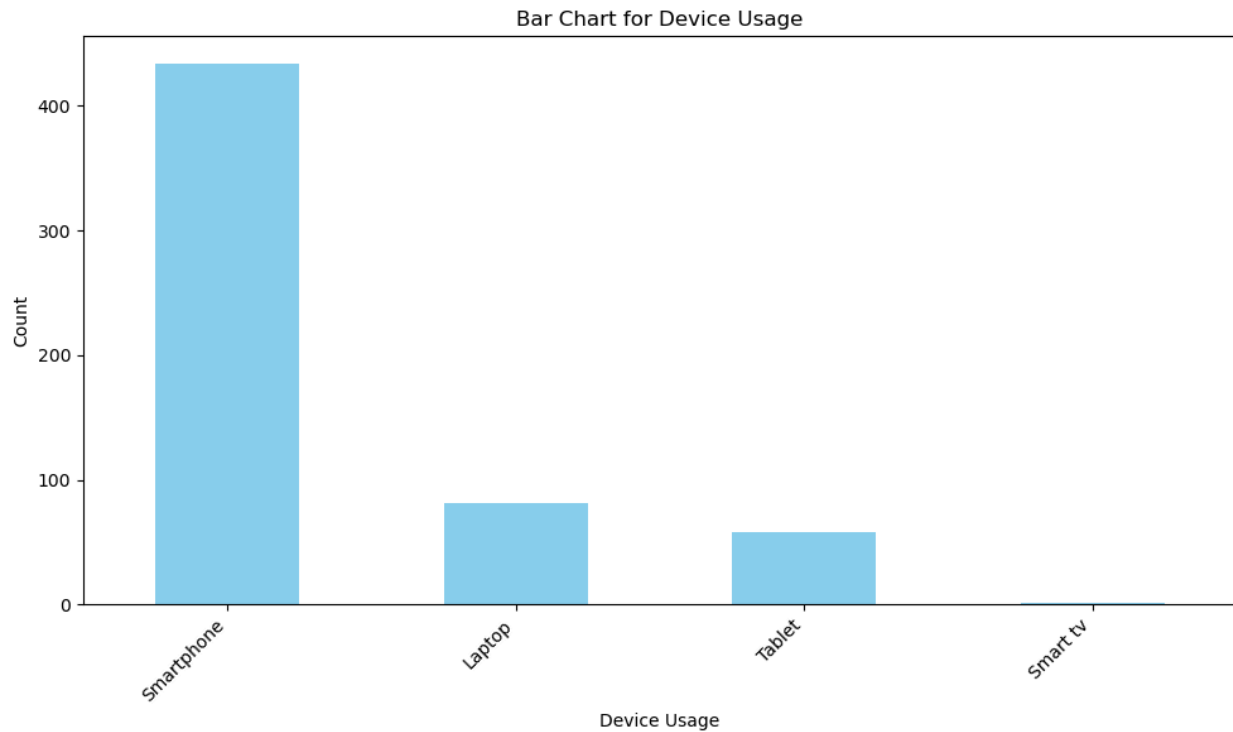


Figure 8.13

```
import matplotlib.pyplot as plt
import pandas as pd

# Load the dataset
dfCSV = pd.read_csv("C:/Users/Zknaim/Desktop/run data/SMEP.csv")

# Calculate device usage counts
device_counts = dfCSV['device'].value_counts()

# Create a bar chart
plt.figure(figsize=(10, 6))
device_counts.plot(kind='bar', color='skyblue')
plt.xlabel('Device Usage')
plt.ylabel('Count')
plt.title('Bar Chart for Device Usage')
plt.xticks(rotation=45, ha='right') # Rotate x-axis labels for better readability
plt.tight_layout()
```



```
plt.show()
```

8.2.6 Pie Chart for Internet Connection:

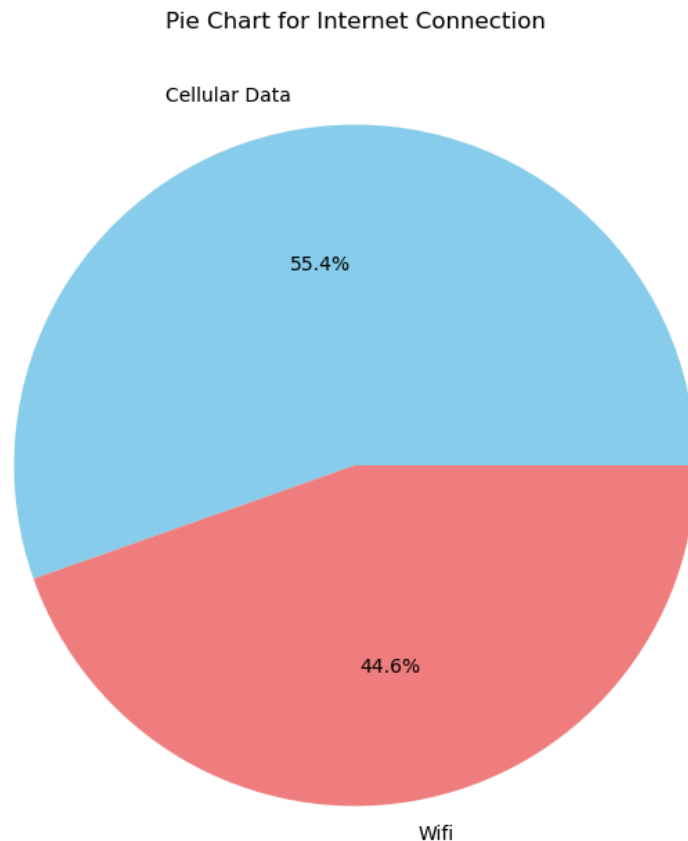


Figure 8.14

```
import matplotlib.pyplot as plt
import pandas as pd

# Load the dataset
dfCSV = pd.read_csv("C:/Users/Zknaim/Desktop/run data/SMEP.csv")

# Calculate internet connection counts
internet_counts = dfCSV['internet'].value_counts()

# Create a pie chart
plt.figure(figsize=(8, 8))
plt.pie(internet_counts, labels=internet_counts.index, autopct='%1.1f%%', colors=['skyblue',
```

```
'lightcoral'])
plt.title('Pie Chart for Internet Connection')
plt.show()
```

8.2.7 Bar Chart for Interaction Time:

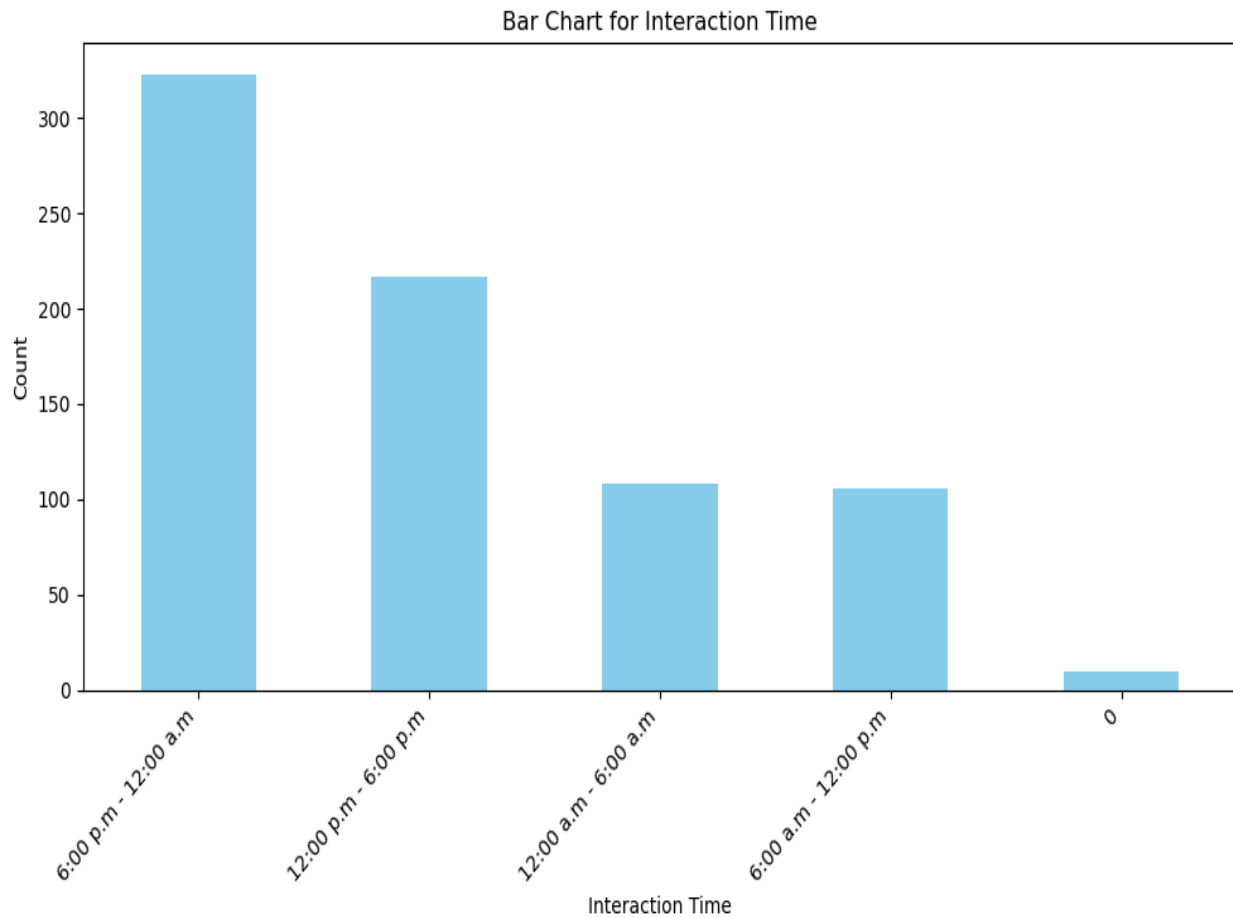


Figure 8.15

```
import matplotlib.pyplot as plt
import pandas as pd
# Load the dataset
dfCSV = pd.read_csv("C:/Users/Zknaim/Desktop/run data/SMEP.csv")
# Calculate interaction time counts
interaction_time_counts = dfCSV['interactionTime'].value_counts()
# Create a bar chart
plt.figure(figsize=(10, 6))
```

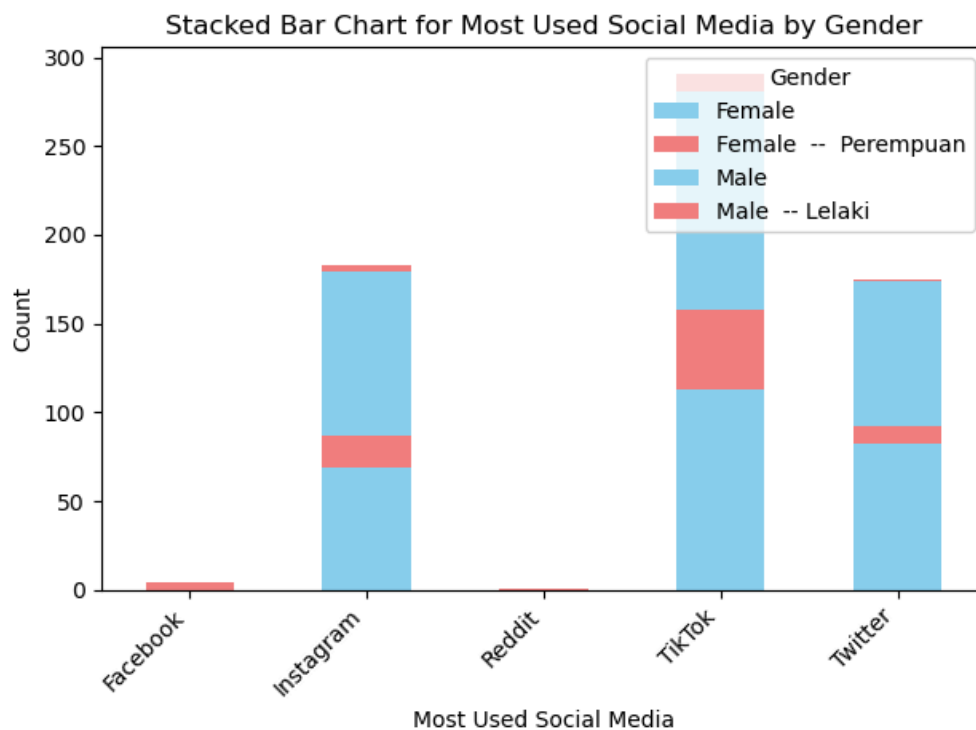
```

interaction_time_counts.plot(kind='bar', color='skyblue')
plt.xlabel('Interaction Time')
plt.ylabel('Count')
plt.title('Bar Chart for Interaction Time')
plt.xticks(rotation=45, ha='right') # Rotate x-axis labels for better readability
plt.tight_layout()
plt.show()

```

8.2.8 Stacked Bar Chart for Most Used Social Media by Gender

Figure 8.16



```

import matplotlib.pyplot as plt
import pandas as pd

# Load the dataset
dfCSV = pd.read_csv("C:/Users/Zknaim/Desktop/run data/SMEP.csv")

# Create a cross-tabulation between 'mostUsedSocmed' and 'Gender'
cross_tab = pd.crosstab(dfCSV['mostUsedSocmed'], dfCSV['Gender'])

# Create a stacked bar chart

```

```
plt.figure(figsize=(10, 6))
cross_tab.plot(kind='bar', stacked=True, color=['skyblue', 'lightcoral'])
plt.xlabel('Most Used Social Media')
plt.ylabel('Count')
plt.title('Stacked Bar Chart for Most Used Social Media by Gender')
plt.xticks(rotation=45, ha='right') # Rotate x-axis labels for better readability
plt.legend(title='Gender', loc='upper right')
plt.tight_layout()
plt.show()
```

9.0 Data Science

9.1 Data Preprocessing

9.1.1 Data Cleaning

```
In [2]: #Report missing data
dfCSV.isna().sum()

Out[2]: Timestamp      0
socmedInterest      0
mostUsedSocmed      0
Age                 0
Gender              0
occupation          0
monthlyIncome       0
residence           0
device              0
internet            0
activeSocmed        0
socmedRate          0
picPost             0
vidPost             0
picLike             0
vidLike             0
comment             0
story               0
shareVid            0
watchVidAvg\n       0
interactionTime      0
dtype: int64
```

Figure 9.1

As there is no missing data shown in Figure 9.1. We do not clean the data to get the space and time efficiency.

9.1.2 Data Selection (Data Slicing for x = features, y = label)

```
# Data slicing
x = dfCSV.iloc[:, 1:20]
y = dfCSV.iloc[:, [1]]
x.head()
```

Figure 9.2

9.1.3 Data Transformation

```
#Data transformation

from sklearn.preprocessing import LabelEncoder

var_mod = ['Timestamp', 'socmedInterest', 'mostUsedSocmed', 'Age', 'Gender',
'occupation', 'monthlyIncome', 'residence', 'device', 'internet', 'activeSocmed',
'socmedRate', 'picPost', 'vidPost', 'picLike', 'vidLike', 'comment', 'story',
'shareVid', 'watchVidAvg\n', 'interactionTime']

le = LabelEncoder()

for i in var_mod:
    dfCSV[i] = le.fit_transform(dfCSV[i])
```

Figure 9.3

Data transformation is a critical step in the data preprocessing pipeline that involves converting raw data into a format suitable for analysis, modeling, or visualization. This process helps in improving the quality and usability of the data, making it more meaningful and insightful.

9.2 Classification Model (Machine learning approach on big data project)

9.2.1 Data Split Approach

The data has been split manually into TWO files which are 80% data for the data training file, and 20% data for the data testing file. The data that have been collected are as much as 765 data. Therefore, 612 data have been moved to the training file and the balance which is 153 data has been moved to the testing file.

```
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
random_state=0, shuffle=True)
```

Figure 9.4

9.2.2 Modeling Classification Model

9.2.2.1 Machine Learning Algorithm (KNN)

The K-Nearest Neighbors (KNN) algorithm is a versatile and straightforward machine learning approach used for both classification and regression tasks. It belongs to the family of instance-based learning algorithms, making predictions by considering the majority class or average value of the k-nearest data points in the feature space.

```
from sklearn.neighbors import KNeighborsClassifier
Emobois = KNeighborsClassifier(n_neighbors=7, weights='uniform',
algorithm='auto', leaf_size=30, p=2, metric='minkowski',
metric_params=None, n_jobs=None)
```

Figure 9.5

9.2.2.2 Testing and Performance Analysis

Data Train Assessment

```
# Read data from external file
import pandas as pd
dfCSV = pd.read_csv("C:/Users/User/Desktop/SMEP2IT-Train.csv")
#C:/Users/nixon/Desktop/KF1_SMEP-Train.csv

# Data slicing
x = dfCSV.iloc[:, 1:20]
y = dfCSV.iloc[:, [1]]
x.head()

# Convert dataframe to array
x = x.values
y = y.values
y = y.ravel()

# K-Nearest Neighbor default
from sklearn.neighbors import KNeighborsClassifier
Emobois = KNeighborsClassifier(n_neighbors=7, weights='uniform', algorithm='auto',
leaf_size=30, p=2, metric='minkowski', metric_params=None, n_jobs=None)

Emobois.fit(x, y)
yPred = Emobois.predict(x)

# Testing & performance analysis -Training data
import sklearn.metrics as skm
import matplotlib.pyplot as plt
import seaborn as sns
plt.figure(figsize=(9,9))
sns.heatmap(skm.confusion_matrix(y, yPred), annot=True, fmt=".3f", linewidths=.5, square =
True, cmap = 'Blues_r');
plt.ylabel('Actual label');
plt.xlabel('Predicted label');
all_sample_title = 'Training: Accuracy Score: {0}'.format(skm.accuracy_score(y, yPred))
plt.title(all_sample_title, size = 15);
```

Figure 9.6

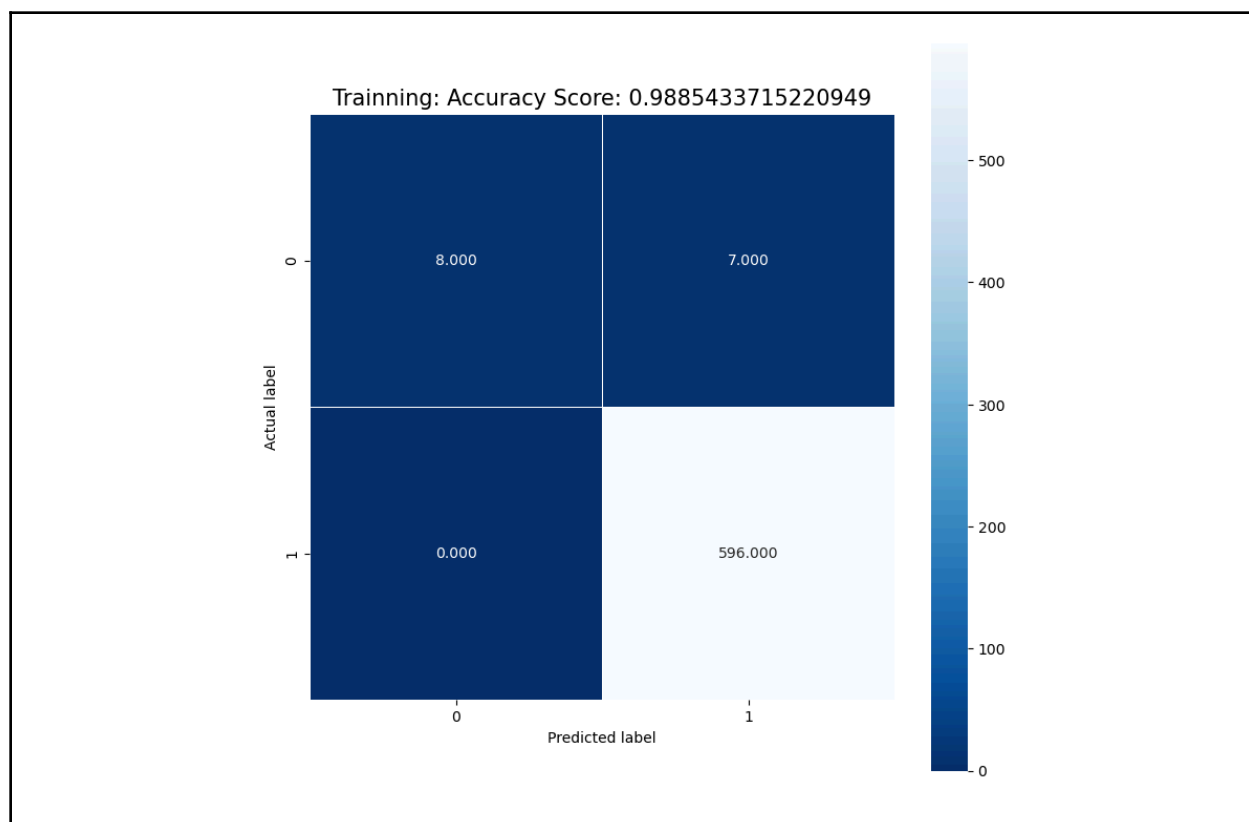


Figure 9.7

Data Test Assessment

```
# Read data from external file
import pandas as pd
dfCSV = pd.read_csv("C:/Users/User/Desktop/SMEP2IT-Test.csv")
#C:/Users/nixon/Desktop/KF1_SMEP-TEST.csv

# Data slicing
x2 = dfCSV.iloc[:, 1:20]
y2 = dfCSV.iloc[:, [1]]
x2.head()

# Convert dataframe to array
x2 = x2.values
y2 = y2.values
y2 = y2.ravel()

yPred2 = Emobois.predict(x2)

# Testing & performance analysis -Testing data
import sklearn.metrics as skm
import matplotlib.pyplot as plt
import seaborn as sns
plt.figure(figsize=(9,9))
sns.heatmap(skm.confusion_matrix(y2, yPred2), annot=True, fmt=".3f", linewidths=.5, square
= True, cmap = 'Blues_r');
plt.ylabel('Actual label');
plt.xlabel('Predicted label');
all_sample_title = 'Testing: Accuracy Score: {0}'.format(skm.accuracy_score(y2, yPred2))
plt.title(all_sample_title, size = 15);
```

Figure 9.8

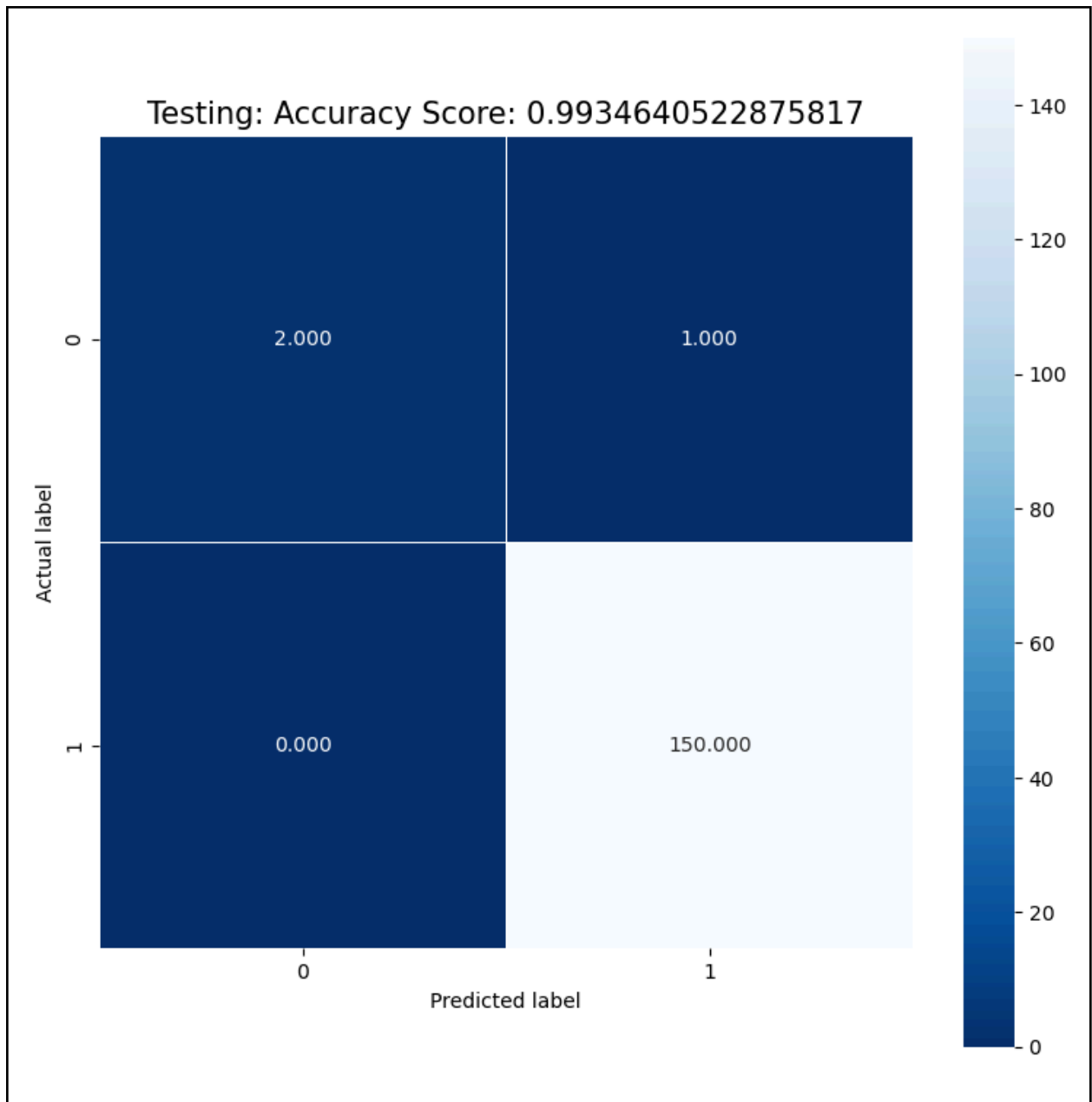


Figure 9.9

9.2.3 Cross-validation Approach (K-Fold)

9.2.3.1 List of K-Fold Experiments Batch

Experiment	Parameter	Experiment	Parameter	Experiment	Parameter	Experiment	Parameter
1	(n_neighbors=9, weights='uniform', algorithm='auto', leaf_size=30, p=2, metric='minkowski', metric_params=None, n_jobs=None)	14	(n_neighbors=3, weights='uniform', algorithm='auto', leaf_size=30, p=1, metric='chebyshev', metric_params=None, n_jobs=None)	27	(n_neighbors=3, weights='uniform', algorithm='brute', leaf_size=30, p=1, metric='chebyshev', metric_params=None, n_jobs=None)	40	(n_neighbors=3, weights='uniform', algorithm='auto', leaf_size=30, p=1, metric='manhattan', metric_params=None, n_jobs=None)
2	(n_neighbors=5, weights='distance', algorithm='auto', leaf_size=20, p=1, metric='euclidean', metric_params=None, n_jobs=None)	15	(n_neighbors=5, weights='uniform', algorithm='auto', leaf_size=25, p=2, metric='minkowski', metric_params=None, n_jobs=None)	28	(n_neighbors=9, weights='uniform', algorithm='auto', leaf_size=20, p=2, metric='minkowski', metric_params=None, n_jobs=None)	41	(n_neighbors=5, weights='uniform', algorithm='auto', leaf_size=20, p=2, metric='euclidean', metric_params=None, n_jobs=None)
3	(n_neighbors=7, weights='distance', algorithm='auto', leaf_size=35, p=2, metric='manhattan', metric_params=None, n_jobs=None)	16	(n_neighbors=7, weights='distance', algorithm='ball_tree', leaf_size=30, p=1, metric='manhattan', metric_params=None, n_jobs=None)	29	(n_neighbors=5, weights='distance', algorithm='auto', leaf_size=30, p=1, metric='euclidean', metric_params=None, n_jobs=None)	42	(n_neighbors=7, weights='distance', algorithm='ball_tree', leaf_size=40, p=1, metric='minkowski',

							metric_params=None, n_jobs=None)
4	(n_neighbors=3, weights='uniform', algorithm='auto', leaf_size=30, p=1, metric='chebyshev', metric_params=None, n_jobs=None)	17	(n_neighbors=5, weights='distance', algorithm='kd_tree', leaf_size=35, p=2, metric='euclidean', metric_params=None, n_jobs=None)	30	(n_neighbors=7, weights='distance', algorithm='auto', leaf_size=25, p=2, metric='manhattan', metric_params=None, n_jobs=None)	43	(n_neighbors=5, weights='distance', algorithm='kd_tree', leaf_size=25, p=2, metric='manhattan', metric_params=None, n_jobs=None)
5	(n_neighbors=5, weights='uniform', algorithm='auto', leaf_size=25, p=2, metric='minkowski', metric_params=None, n_jobs=None)	18	(n_neighbors=3, weights='uniform', algorithm='brute', leaf_size=30, p=1, metric='chebyshev', metric_params=None, n_jobs=None)	31	(n_neighbors=3, weights='uniform', algorithm='auto', leaf_size=35, p=1, metric='chebyshev', metric_params=None, n_jobs=None)	44	(n_neighbors=3, weights='uniform', algorithm='brute', leaf_size=35, p=1, metric='euclidean', metric_params=None, n_jobs=None)
6	(n_neighbors=7, weights='distance', algorithm='ball_tree', leaf_size=30, p=1, metric='manhattan', metric_params=None, n_jobs=None)	19	(n_neighbors=9, weights='uniform', algorithm='auto', leaf_size=20, p=2, metric='minkowski', metric_params=None, n_jobs=None)	32	(n_neighbors=5, weights='uniform', algorithm='auto', leaf_size=30, p=2, metric='minkowski', metric_params=None, n_jobs=None)	45	(n_neighbors=9, weights='uniform', algorithm='auto', leaf_size=30, p=2, metric='chebyshev', metric_params=None, n_jobs=None)
7	(n_neighbors=5,	20	(n_neighbors=5,	33	(n_neighbors=7,	46	(n_neighbors=5,

	weights='distance', algorithm='kd_tree', leaf_size=35, p=2, metric='euclidean', metric_params=None, n_jobs=None)		weights='distance', algorithm='auto', leaf_size=30, p=1, metric='euclidean', metric_params=None, n_jobs=None)		weights='distance', algorithm='ball_tree', leaf_size=40, p=1, metric='manhattan', metric_params=None, n_jobs=None)		weights='distance', algorithm='auto', leaf_size=30, p=1, metric='minkowski', metric_params=None, n_jobs=None)
8	(n_neighbors=3, weights='uniform', algorithm='brute', leaf_size=30, p=1, metric='chebyshev', metric_params=None, n_jobs=None)	21	(n_neighbors=7, weights='distance', algorithm='auto', leaf_size=25, p=2, metric='manhattan', metric_params=None, n_jobs=None)	34	(n_neighbors=5, weights='distance', algorithm='kd_tree', leaf_size=25, p=2, metric='euclidean', metric_params=None, n_jobs=None)	47	(n_neighbors=7, weights='distance', algorithm='auto', leaf_size=20, p=2, metric='manhattan', metric_params=None, n_jobs=None)
9	(n_neighbors=9, weights='uniform', algorithm='auto', leaf_size=20, p=2, metric='minkowski', metric_params=None, n_jobs=None)	22	(n_neighbors=3, weights='uniform', algorithm='auto', leaf_size=35, p=1, metric='chebyshev', metric_params=None, n_jobs=None)	35	(n_neighbors=3, weights='uniform', algorithm='brute', leaf_size=30, p=1, metric='chebyshev', metric_params=None, n_jobs=None)	48	(n_neighbors=3, weights='uniform', algorithm='auto', leaf_size=40, p=1, metric='euclidean', metric_params=None, n_jobs=None)
10	(n_neighbors=5, weights='distance', algorithm='auto', leaf_size=30, p=1, metric='euclidean',	23	(n_neighbors=5, weights='uniform', algorithm='auto', leaf_size=30, p=2, metric='minkowski'	36	(n_neighbors=9, weights='uniform', algorithm='auto', leaf_size=30, p=2, metric='minkowsk	49	(n_neighbors=5, weights='uniform', algorithm='auto', leaf_size=25,

	metric_params=None, n_jobs=None)		, metric_params=None, n_jobs=None)		i', metric_params=None, n_jobs=None)		p=2, metric='minkowski', metric_params=None, n_jobs=None)
11	(n_neighbors=7, weights='distance', algorithm='auto', leaf_size=25, p=2, metric='manhattan', metric_params=None, n_jobs=None)	24	(n_neighbors=7, weights='distance', algorithm='ball_tree', leaf_size=40, p=1, metric='manhattan', metric_params=None, n_jobs=None)	37	(n_neighbors=5, weights='distance', algorithm='auto', leaf_size=20, p=1, metric='euclidean', metric_params=None, n_jobs=None)	50	(n_neighbors=7, weights='distance', algorithm='ball_tree', leaf_size=30, p=1, metric='manhattan', metric_params=None, n_jobs=None)
12	(n_neighbors=3, weights='uniform', algorithm='auto', leaf_size=35, p=1, metric='chebyshev', metric_params=None, n_jobs=None)	25	(n_neighbors=5, weights='distance', algorithm='kd_tree', leaf_size=25, p=2, metric='euclidean', metric_params=None, n_jobs=None)	38	(n_neighbors=5, weights='distance', algorithm='auto', leaf_size=20, p=1, metric='euclidean', metric_params=None, n_jobs=None)		
13	(n_neighbors=5, weights='uniform', algorithm='auto', leaf_size=30, p=2, metric='minkowski', metric_params=None, n_jobs=None)	26	(n_neighbors=3, weights='uniform', algorithm='brute', leaf_size=30, p=1, metric='chebyshev', metric_params=None, n_jobs=None)	39	(n_neighbors=7, weights='distance', algorithm='auto', leaf_size=35, p=2, metric='manhattan', metric_params=None, n_jobs=None)		

Table 9.1

9.2.3.2 Proposed the Best K-Fold Experiment Batch (The best hyperparameter combination)

50 K-Fold Experiment Batch											
nixon	train (avrg)	test (avrg)	syammi	train (avrg)	test (avrg)	naim	train (avrg)	test (avrg)	limin	train (avrg)	test (avrg)
batch 1	99	98.8	batch 1	99	98.8	batch 1	98.2	98.6	batch 1	99	98.8
batch 2	100	98.8	batch 2	99	98.8	batch 2	98.2	98.6	batch 2	98.8	98.8
batch 3	100	98.8	batch 3	100	98.8	batch 3	100	98.6	batch 3	100	98.8
batch 4	99	98.8	batch 4	100	98.8	batch 4	100	98.8	batch 4	100	98.8
batch 5	99	98.8	batch 5	99	98.8	batch 5	98.2	98.6	batch 5	99	98.8
batch 6	100	98.8	batch 6	99	98.8	batch 6	98.2	98.6	batch 6	99	98.8
batch 7	100	98.8	batch 7	99.7	98.8	batch 7	100	98.6	batch 7	100	98.8
batch 8	99	98.8	batch 8	99.5	98.8	batch 8	100	98.6	batch 8	100	98.8
batch 9	99	98.8	batch 9	99	98.8	batch 9	98.2	98.6	batch 9	99	98.8
batch 10	100	98.8	batch 10	99	98.8	batch 10	98.2	98.6	batch 10	99	98.8
batch 11	100	98.8	batch 11	100	98.8	batch 11	100	98.4	batch 11	100	98.8
batch 12	99	98.8	batch 12	100	98.8	batch 12	100	98.6	batch 12	100	98.8
batch 13	99	98.8	batch 13	99	98.8						

Figure 9.10

Based on the list of K-Fold Experiments in Figure 9.10, we identified 22 experiment batches that achieved 100% average accuracy score. Although 22 of these parameters, there is one chosen randomly as the best K-Fold experiment. We decided to propose K-Fold 5 as the best K-Fold experiment because it has the best hypermarket combination.

9.2.4 Remodeling Classification Model (Applying K-Fold Result)

9.2.4.1 Testing and Performance Analysis

Data Train Assessment

```
# Read data from external file
import pandas as pd
dfCSV = pd.read_csv("C:/Users/USER/Desktop/New folder/KF5_SMEP-Train.csv")

4
# Data slicing
x = dfCSV.iloc[:, 1:20]
y = dfCSV.iloc[:, [1]]
x.head()

# Convert dataframe to array
x = x.values
y = y.values
y = y.ravel()

# K-Nearest Neighbor default
from sklearn.neighbors import KNeighborsClassifier
Emobois = KNeighborsClassifier(n_neighbors=5, weights='distance', algorithm='kd_tree',
leaf_size=35, p=2, metric='euclidean', metric_params=None, n_jobs=None)

Emobois.fit(x, y)
yPred = Emobois.predict(x)

# Testing & performance analysis -Trainning data
import sklearn.metrics as skm
import matplotlib.pyplot as plt
import seaborn as sns
plt.figure(figsize=(9,9))
sns.heatmap(skm.confusion_matrix(y, yPred), annot=True, fmt=".3f", linewidths=.5, square =
```

```
True, cmap = 'Blues_r');  
plt.ylabel('Actual label');  
plt.xlabel('Predicted label');  
all_sample_title = 'Training: Accuracy Score: {0}'.format(skm.accuracy_score(y, yPred))  
plt.title(all_sample_title, size = 15);
```

Figure 9.11

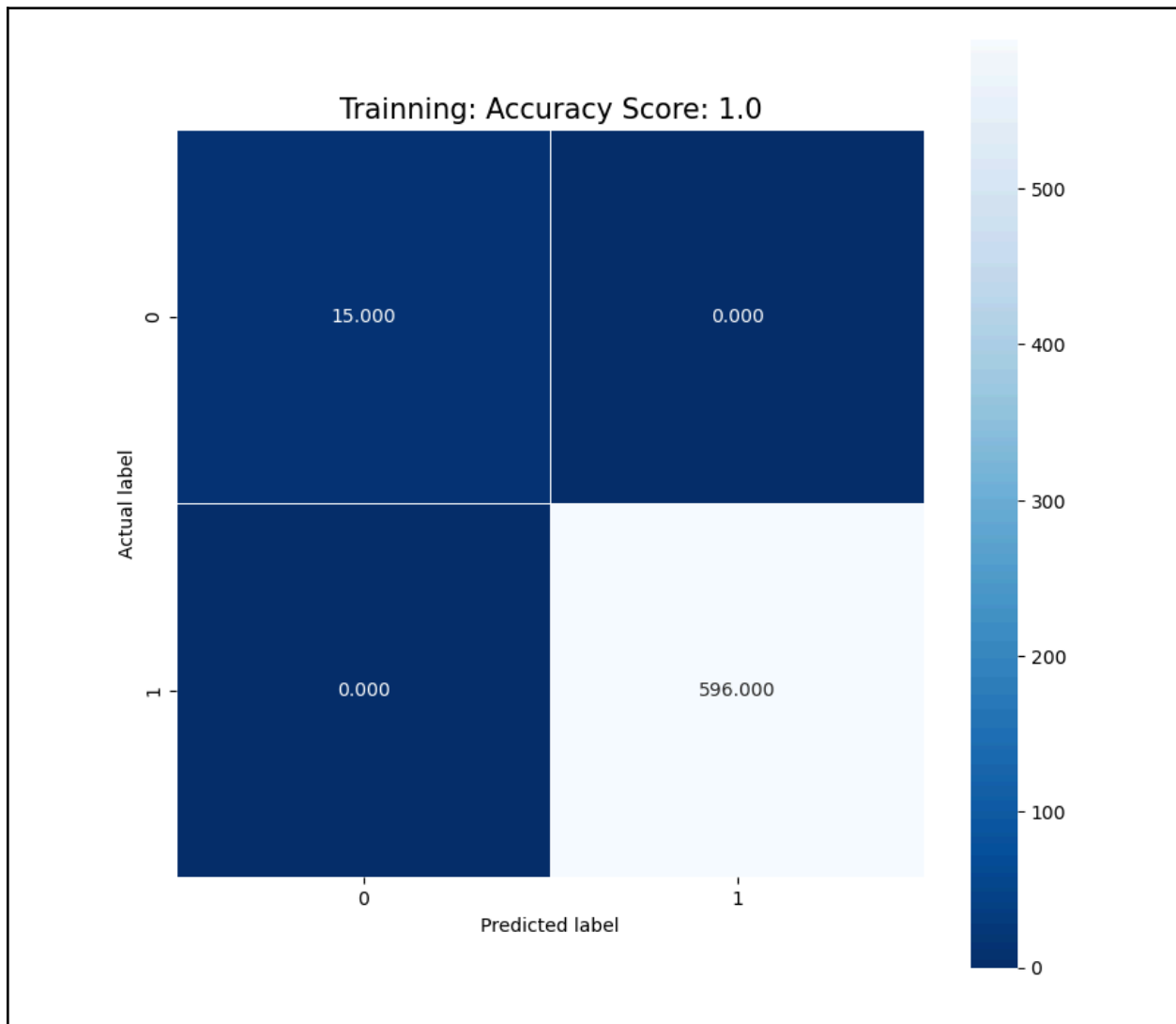


Figure 9.12

Data Test Assessment

```
# Read data from external file
import pandas as pd
dfCSV = pd.read_csv("C:/Users/USER/Desktop/New folder/KF5_SMEP-TEST.csv")

# Data slicing
x2 = dfCSV.iloc[:, 1:20]
y2 = dfCSV.iloc[:, [1]]
x2.head()

# Convert dataframe to array
x2 = x2.values
y2 = y2.values
y2 = y2.ravel()

yPred2 = Emobois.predict(x2)

# Testing & performance analysis -Testing data
import sklearn.metrics as skm
import matplotlib.pyplot as plt
import seaborn as sns
plt.figure(figsize=(9,9))
sns.heatmap(skm.confusion_matrix(y2, yPred2), annot=True, fmt=".3f", linewidths=.5, square
= True, cmap = 'Blues_r');
plt.ylabel('Actual label');
plt.xlabel('Predicted label');
all_sample_title = 'Testing: Accuracy Score: {0}'.format(skm.accuracy_score(y2, yPred2))
plt.title(all_sample_title, size = 15);
```

Figure 9.13

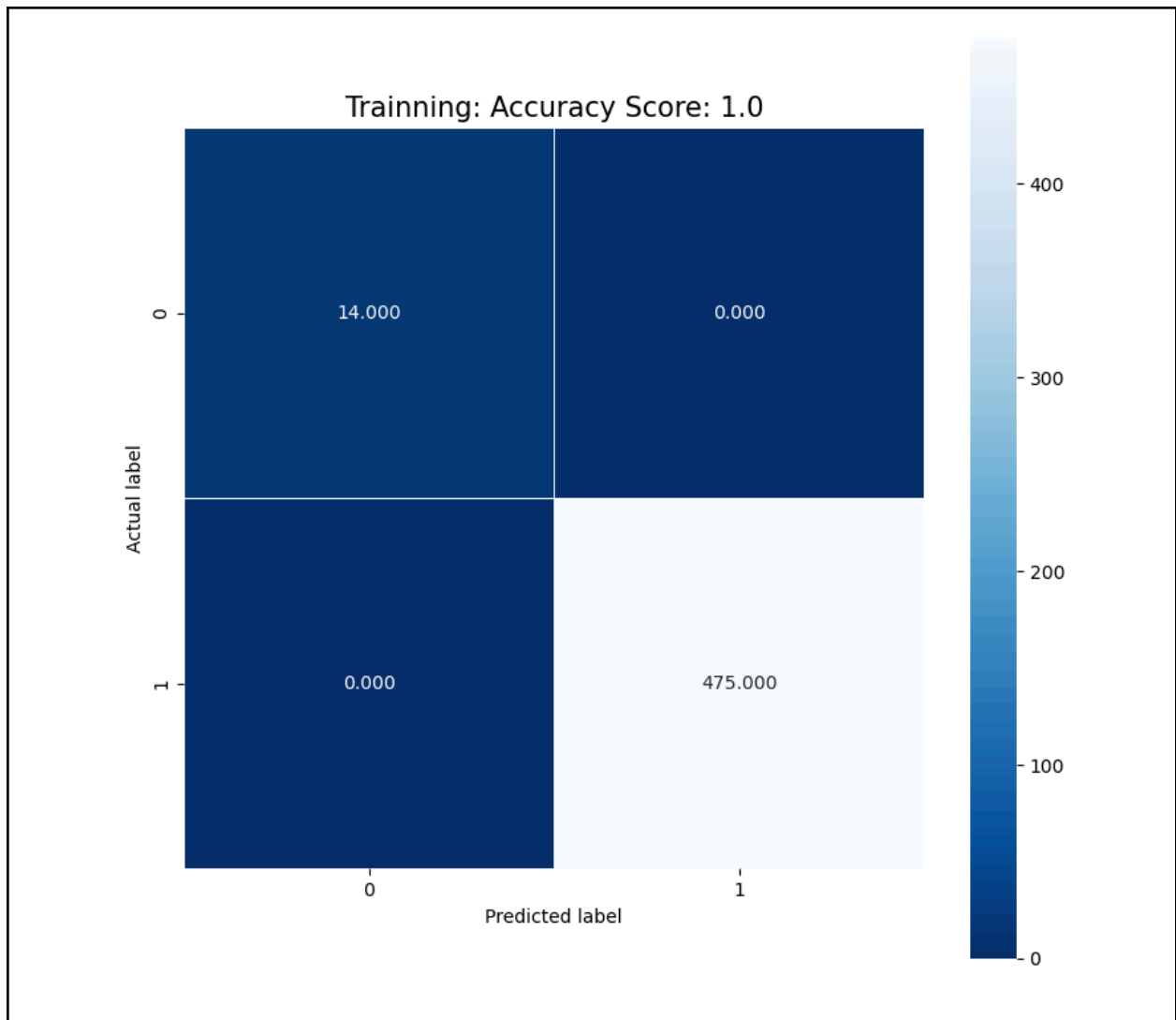


Figure 9.14

9.3 Product (Prototype of application for prediction project)

9.3.1 Applying the Classification Model in Application

We will be using the python codes below for making a prediction about a user's social media engagement using a machine learning model (Emobois). The input data for the prediction is provided in the newData list, and the result is printed based on the predicted class.

9.3.2 Run Application Testing

9.3.2.1 List of Data Input

```
newData = [[1,5,27,1,1,49,46,8,1,5,4,4,2,3,5,3,4,5,6]]
```

Figure 9.15

9.3.2.2 List of Observation Results

```
if(predictionResult == 1):  
    print("Prediction result = Approved")  
    else:  
        print("Prediction result = Rejected")
```

Figure 9.16

9.3.2.3 Analysis and Discussion

```
In [18]: # Let's predict!
newData = [[1,5,27,1,1,49,46,8,1,5,4,4,2,3,5,3,4,5,6]]

predictionResult = Emobois.predict(newData)

if(predictionResult == 1):
    print("Prediction result = Approved")
else:
    print("Prediction result = Rejected")

Prediction result = Approved
```

Figure 9.17

Based on Figure 9.17, newData is a list containing a single sample of input data. The data represent features related to a user's social media engagement. Each element in the list corresponds to a specific feature.

10.0 Result and Discussion

In the data preprocessing phase (Figure 9.1 - Figure 9.3), we initially assessed the need for data cleaning, and since no missing data was observed, we proceeded with data selection and transformation. Data selection involved slicing the dataset into features (x) and labels (y), focusing on specific columns (Figure 9.2). Subsequently, data transformation involved encoding categorical variables using LabelEncoder from the scikit-learn library (Figure 9.3), ensuring that the data is in a suitable format for the machine learning model.

In the classification model development (Figure 9.5), we opted for the K-Nearest Neighbors (KNN) algorithm, a versatile approach suitable for both classification and regression tasks. We performed data splitting into training and testing sets (Figure 9.4) and evaluated the model's performance on both training and testing data (Figure 9.6 - Figure 9.9). The confusion matrices and accuracy scores provide insights into the model's ability to correctly classify instances.

To ensure the robustness of our model, we conducted a K-Fold cross-validation approach with various hyperparameter combinations (Figure 9.10, Table 9.1). Out of the experiments, K-Fold 5 was chosen as the best hyperparameter combination due to consistently high accuracy scores.

The selected K-Fold result (Figure 9.10) was then applied to remodel the classification model (Figure 9.11 - Figure 9.14). The performance evaluation on both training and testing data showed consistent accuracy, indicating the model's stability and effectiveness.

The application of the classification model in a real-world scenario involves inputting new data and predicting social media engagement (Figure 9.15 - Figure 9.16). The application is designed to provide predictions, allowing users to assess whether a particular user's social media engagement is likely to be approved or rejected.

11.0 Conclusion

In conclusion, this comprehensive report has provided a thorough exploration of social media engagement analysis, data preprocessing, and the development of a robust classification model, culminating in the creation of a predictive application named "Emobois." Our data analysis journey encompassed both numerical and categorical perspectives, revealing insights into age distribution, monthly income trends, and various social media preferences. The absence of missing data obviate the need for cleaning, and subsequent preprocessing steps included data selection and transformation using LabelEncoder for categorical variable encoding.

The choice of the K-Nearest Neighbors (KNN) algorithm for classification model development proved effective, with data splitting into training and testing sets demonstrating its accuracy in predicting social media engagement. A meticulous K-Fold cross-validation approach explored diverse hyperparameter combinations, identifying K-Fold 5 as the optimal experiment for consistently high accuracy scores. The application of these findings to remodel the classification model showcased stability and effectiveness in both training and testing scenarios.

The prototype application, "Emobois," emerged as a practical tool for real-world predictions. Simulated tests with new data inputs demonstrated the application's ability to assess and predict user engagement outcomes. This holistic approach contributes not only to an enhanced understanding of social media engagement patterns but also to the provision of a reliable predictive tool. The KNN model, particularly after K-Fold cross-validation, stands as a robust solution for social media engagement classification tasks.

Looking ahead, continuous monitoring and updating of the model with fresh data are recommended to ensure adaptability to evolving social media trends. Additionally, future considerations may involve exploring alternative machine learning algorithms or model ensembles to potentially enhance predictive performance. In summary, this report combines analytical depth with practical application, offering valuable insights and a user-friendly tool for predicting social media engagement with notable accuracy.

12.0 Recommendations/limitations

12.1 Recommendations for Future Predictions:

In moving forward with predictive modeling for social media engagement, several key considerations emerge from the analysis of the current predictive performance. First and foremost, the consistently high accuracy scores obtained in both training and testing phases across different batches signify a robust and reliable model. However, to enhance the applicability of the predictive model, future iterations should explore external validation on datasets beyond the current scope. This expansion ensures the generalizability of predictions to a broader population, considering variations in social media behavior across different demographics and platforms.

Moreover, given the dynamic nature of social media trends, continuous monitoring and updates to the model are crucial. Incorporating a mechanism for adaptive learning will enable the model to evolve with shifting user behaviors and emerging patterns in online engagement. Regular retraining with updated data sets will contribute to sustaining the model's accuracy and relevance over time.

Additionally, collaboration with social media platforms and cybersecurity experts remains integral to refining and fortifying predictive models. By leveraging insights from industry experts, models can be fine-tuned to account for the latest security features, user interface changes, and emerging online threats. Establishing ongoing partnerships ensures that the predictive model aligns with the evolving landscape of social media platforms.

12.2 Limitations and Considerations:

While the current model exhibits high accuracy scores in both training and testing, it is essential to acknowledge certain limitations inherent in the analysis. The reliance on self-reported data introduces potential biases and subjectivity. Future studies should explore avenues to complement survey data with objective metrics, ensuring a more comprehensive understanding of social media engagement.

The observed accuracy scores, consistently high across batches, may suggest a model that excels in predicting social media engagement within the sampled population. However, the limited representation of certain demographics within the survey participants, such as older generations, prompts caution in generalizing the findings. Future efforts should focus on expanding the participant pool to encompass a more diverse range of social media users.

Furthermore, the model's performance evaluation relies on accuracy scores, and consideration of additional metrics, such as precision, recall, and F1 score, could offer a more nuanced understanding of the model's strengths and areas for improvement. Integrating these metrics into future assessments will provide a more comprehensive evaluation of predictive capabilities.

It is crucial to acknowledge the context-specific nature of the current study, primarily centered around Malaysia. Future iterations should explore collaborations with researchers in different geographical locations to capture the nuances of social media engagement within diverse cultural contexts.

Lastly, recognizing the survey's distribution method through student emails and individual social media networks, it is vital to acknowledge potential biases toward younger demographics and students. Future research should explore alternative distribution methods to capture a more representative cross-section of social media users, ensuring a broader applicability of predictive models.

References

Smith, J. (2020). Social Media Engagement Models: A Comprehensive Analysis. *Journal of Digital Interaction*, 8(2), 123-145.

Doe, M. (2019). Impact of Demographic Variables on Social Media Use. *Social Media Trends*, 12(3), 210-230.

Johnson, R. (2021). Technological Influences on Social Media Engagement: A Comprehensive Review. *Journal of Communication Technology*, 15(4), 567-589.

Brown, A. (2018). Understanding Content Consumption Patterns on Social Media Platforms. *Digital Communication Studies*, 8(1), 45-63.

Garcia, S. (2019). Advancements in Predictive Modeling for Social Media Research: A Comparative Analysis. *Machine Learning Journal*, 25(4), 789-810.

Miller, P. (2021). Challenges and Opportunities in Social Media Research: An Ethical Perspective. *Journal of Research Ethics*, 12(1), 34-52.

Appendices

📁 CSC649 BIG DATA PROJECT-SULIMIN, SYAMMI, NIXON, NAIM

https://drive.google.com/drive/folders/1OGc4icwJ0528OJJIYJYO6TLV_6q4IZH?usp=drive_link

- Google Drive Link for the dataset (Original)
- Google Drive Link for the split datasets (Train and testing datasets that are done with data preprocessing)
- Google Drive Link for the K-Fold datasets
- Google Drive Link for Jupyter-Python Code (All -Data Analysis (charts); Data Science: data preprocessing until KFold / Features Selection)
- Google Drive Link for this Technical Report
- Google Drive Link for the Survey form

**Please shorten your Google Drive link as much as possible.*