

CHAPTER 1

INTRODUCTION

CHAPTER 1

INTRODUCTION

1.1 ABOUT THE PROJECT

Customer Feedback Analysis Using Machine Learning is a project designed to transform raw customer feedback into actionable insights using data science and machine learning techniques. In today's digital era, businesses receive vast amounts of customer feedback through online platforms such as websites, social media, surveys, and support channels. Analyzing this data manually is not only time-consuming but also prone to human error. Our project addresses this challenge by automating the analysis process using intelligent algorithms.

The core objective of the project is to:

- Classify customer sentiments (positive, negative, neutral) using sentiment analysis.
- Extract key themes and issues from textual feedback.
- Visualize trends in customer opinions over time.
- Help businesses understand customer expectations, improve their services, and make informed decisions.

To achieve this, we use various Natural Language Processing (NLP) techniques for text preprocessing, such as tokenization, stopwords removal, stemming, and lemmatization. We then apply machine learning models—such as Logistic Regression, Naive Bayes, and Support Vector Machines—for sentiment classification. For topic extraction and keyword identification, we use methods like TF-IDF and Latent Dirichlet Allocation (LDA).

1.2 BACKGROUND WORK

Customer feedback has become one of the most valuable resources for businesses aiming to improve their products, services, and overall customer experience. Traditionally, companies relied on manual analysis of feedback forms, surveys, and customer support logs. However, as the volume of feedback increases across digital platforms, manual processing becomes inefficient and impractical.

Over the years, researchers and organizations have explored automated methods to analyze customer feedback. Techniques such as rule-based sentiment analysis and lexicon-based approaches (like VADER or TextBlob) were among the early solutions. While useful, these methods often lacked contextual understanding and were not scalable for large datasets.

The rise of Natural Language Processing (NLP) and Machine Learning (ML) has significantly improved the accuracy and scalability of feedback analysis. These methods allow systems to learn from data and adapt to various writing styles, tones, and vocabularies. With the availability of open-source tools such as NLTK, scikit-learn, spaCy, and transformer-based models (e.g., BERT), it's now possible to perform complex analysis like sentiment detection, emotion classification, topic modeling, and summarization.

Customer Feedback Analysis using Machine Learning

In our project, we studied and implemented various approaches including:

- Text preprocessing techniques: Tokenization, stopwords removal, stemming, lemmatization.
- Sentiment analysis algorithms: Using supervised ML models trained on labeled datasets.
- Topic modeling: To identify recurring themes and concerns in customer feedback.
- Data visualization tools: To present the results in a user-friendly and interpretable format. The background research laid a strong foundation for building our own system of feedback analysis and helps businesses make better, data-driven decisions.

1.3 OBJECTIVE

The main objective of this project is to analyze customer feedback using machine learning and natural language processing techniques to extract meaningful insights. With the growing amount of customer-generated data in the form of reviews, surveys, and comments, businesses need effective ways to understand customer sentiments and identify common issues. Our project aims to automate this process by applying sentiment analysis to categorize feedback as positive, negative, or neutral.

In addition to sentiment classification, we aim to identify frequently mentioned keywords and themes that reflect customer concerns, expectations, and satisfaction levels. The goal is to use machine learning algorithms to process the textual data, train models for accurate sentiment detection, and apply topic modeling to uncover hidden patterns in the feedback.

This project does not involve building a user interface, but instead focuses on the backend analysis and evaluation of different machine learning models. By comparing the performance of various algorithms, we aim to determine the most effective methods for feedback analysis. The outcome of this project provides a foundation that can be integrated into real-world applications where understanding customer opinions is crucial for business growth.

1.4 PROJECT FEATURES

1. Sentiment Analysis:

Classifies customer feedback into sentiment categories such as positive, negative, and neutral using machine learning models.

2. Text Preprocessing:

Implements Natural Language Processing (NLP) techniques like tokenization, stopwords removal, stemming, and lemmatization to clean and prepare text data for analysis.

3. Keyword Extraction:

Identifies frequently occurring words and phrases in the feedback to highlight common topics or issues raised by customers.

4. Topic Modeling:

Uses techniques such as Latent Dirichlet Allocation (LDA) to uncover hidden themes and group similar feedback together based on content.

5. Model Evaluation:

Compares the performance of different machine learning algorithms (e.g., Logistic Regression, Naive Bayes, SVM) to identify the most accurate model for sentiment classification.

6. Scalability for Large Datasets:

Designed to handle and analyze large volumes of customer feedback efficiently.

7. Insight Generation:

Converts raw customer feedback into meaningful insights that can be used to improve customer experience and business strategies.

8. Adaptability Across Domains:

The analysis approach is flexible and can be applied to feedback from various industries such as e-commerce, hospitality, and healthcare.

CHAPTER – 2

LITERATURE SURVEY

CHAPTER – 2

LITERATURE SURVEY

2.1 Customer Segmentation Using RFM Model and K-Means Clustering

The article "Customer Segmentation Using RFM Model and K-Means Clustering" by Rahul, Laxmiputra, and Saraswati (2021) delves into the integration of the Recency-Frequency-Monetary (RFM) model with the K-Means clustering algorithm to enhance customer segmentation strategies. The study emphasizes how businesses can gain actionable insights into customer behavior by grouping similar customers together based on their purchase patterns. Although the original literature review from this article is not fully accessible, an overview of related research highlights the importance of using data-driven clustering techniques for customer segmentation.

The RFM-K-Means approach is widely adopted due to its simplicity and effectiveness in identifying key customer segments. However, alternative clustering techniques such as Hierarchical Clustering and Gaussian Mixture Models (GMM) have also been explored in research. Comparative studies indicate that the choice of clustering algorithm plays a crucial role in the quality of segmentation. For instance, while K-Means provides faster results and works well with large datasets, GMM offers a probabilistic approach, which may be more suitable for datasets where customer behavior follows multiple overlapping distributions. Some studies suggest that GMM can outperform K-Means in certain scenarios, particularly when customers exhibit complex purchasing behaviors that cannot be easily divided into distinct groups.

Overall, the integration of RFM with K-Means is a well-established method for customer segmentation. It enables businesses to identify loyal customers, at-risk customers, and potential churners, allowing them to tailor their marketing strategies, loyalty programs, and personalized promotions more effectively. However, the selection of an appropriate clustering technique remains a key factor in achieving optimal segmentation performance.

2.2 Customer Segmentation by Web Content Mining

In their 2021 study, Jinfeng, Jinliang, and Bugao introduced an innovative approach to customer segmentation by extending the traditional RFM model with an additional dimension: Interpurchase Time (T), creating the RFMT model. This enhancement allows businesses to capture more detailed insights into customer purchasing behaviors over extended periods. Recognizing that the traditional RFM model may not fully capture the evolving nature of consumer habits, the addition of Interpurchase Time (T) provides a more accurate representation of customer loyalty and retention trends.

Interpurchase Time (T) measures the average time interval between consecutive purchases, offering a deeper understanding of customer shopping cycles. Customers with shorter Interpurchase Times may indicate frequent buyers, while those with longer gaps between

purchases might be seasonal shoppers or occasional buyers. The RFMT model allows businesses to track, compare, and segment customers dynamically, enabling them to design targeted retention strategies and personalized marketing campaigns based on shopping frequency.

Additionally, the article discusses the role of Web Content Mining, where data from online brand communities, social media, and user-generated content is analyzed to understand customer preferences and engagement. By leveraging text mining and clustering methods, businesses can extract valuable insights from customer discussions, feedback, and interactions. For example, a study utilizing machine learning and NLP techniques on online user reviews successfully identified distinct customer segments, allowing businesses to refine their content strategies, product offerings, and customer engagement models. This approach significantly helps in reducing customer churn and increasing customer lifetime value (CLV).

2.3 Pakistan's Largest E-Commerce Dataset

The dataset "Pakistan's Largest E-Commerce Dataset", compiled by Zeeshan-ul-Hassan Usmani (2021), provides a comprehensive overview of Pakistan's e-commerce industry. The dataset covers transaction records from March 2016 to August 2018, comprising over 584,000 records, making it an invaluable resource for understanding consumer behavior and market trends in Pakistan.

The dataset highlights the rapid growth of Pakistan's e-commerce sector, fueled by increased internet penetration, mobile connectivity, and digital payment adoption. As of 2023, Pakistan ranks 46th globally in e-commerce, with revenues reaching \$5.2 billion. The country has approximately 87.4 million internet users and over 124 million 3G and 4G subscribers, indicating strong digital infrastructure supporting online retail activities.

Despite this positive trajectory, several challenges remain. Research indicates that in 2017, only 5% of internet users in Pakistan engaged in online shopping, a figure significantly lower than that of neighboring countries. The primary barriers include lack of trust in online payments, cybersecurity concerns, and inadequate regulatory frameworks. Addressing these issues presents a significant opportunity for e-commerce stakeholders to develop secure payment gateways, enhance fraud detection measures, and implement customer-centric policies to boost consumer confidence and participation in digital commerce.

The dataset provides researchers and businesses with an opportunity to apply machine learning models for predictive analytics, fraud detection, and customer behavior analysis, ultimately contributing to the growth and optimization of Pakistan's digital economy.

2.4 Data-Driven Optimization of Hydrogen Supply Networks

The article "A Data-Driven Approach for the Optimization of Future Two-Level Hydrogen Supply Network Design with Stochastic Demand under Carbon Regulations" by Ge, Jin, and Ren (2021) explores the design of hydrogen supply networks under uncertain demand and environmental regulations. While the primary focus of this study is on supply chain

Customer Feedback Analysis using Machine Learning

optimization, it provides valuable insights into the use of data-driven methodologies for decision-making in complex systems.

Hydrogen supply chain design requires adaptability to demand fluctuations. Studies such as those by Almansoori and Shah (2012) have developed stochastic models to optimize hydrogen networks, ensuring cost efficiency and supply reliability. Similarly, Ochoa Bique et al. (2016) introduced a flexible hydrogen distribution model in Germany, emphasizing the importance of adaptable infrastructure to meet uncertain energy demands.

Machine learning and real-time data analytics are increasingly being integrated into supply chain optimization. By leveraging predictive analytics and AI-driven forecasting, businesses can make proactive decisions that minimize operational costs, improve efficiency, and ensure regulatory compliance

CHAPTER 3

EXISTING SYSTEM

CHAPTER-3

EXISTING SYSTEM

In the existing Customer Feedback Analysis system, models such as Support Vector Machine (SVM) and Logistic Regression are used. Logistic Regression works by modeling the probability of feedback belonging to a particular class based on a linear combination of input features, offering a fast and interpretable solution for text classification. However, it assumes linearity, which limits its performance on more complex data. On the other hand, SVM classifies feedback by finding the optimal hyperplane that separates different classes in a high-dimensional space. It is effective in handling non-linear and sparse text data but can be computationally intensive, especially with large datasets. While both models improve accuracy compared to simpler methods like Naive Bayes, they still face limitations in capturing deeper contextual meanings within feedback.

ADVANTAGES

- **Good Performance on Text Data**
Both models perform well with high-dimensional, sparse datasets, which is common in text-based feedback.
- **Efficient with Moderate Data Sizes**
These models don't require massive datasets and can deliver solid results with relatively smaller volumes of labeled data.
- **Fast Training (Especially Logistic Regression)**
Logistic Regression trains quickly and is suitable for real-time or near-real-time applications.
- **Effective for Binary and Multiclass Classification**
Both are adaptable for classifying feedback into categories like positive, negative, or neutral.

DISADVANTAGES

- **Inability to Capture Word Order or Grammar**
Both models typically use **Bag-of-Words** or **TF-IDF**, which ignore word order, sentence structure, and grammar.
- **Poor Handling of Out-of-Vocabulary (OOV) Words**
If new or rare words appear in feedback (not seen during training), the models may struggle or ignore them completely.
- **Less Ideal for Very Large Datasets**
Computational Complexity: Training complexity is roughly between $O(n^2)$ and $O(n^3)$, which becomes very slow with large n (number of samples).

CHAPTER 4

PROPOSED SYSTEM

CHAPTER 4

PROPOSED SYSTEM

The proposed method for Customer Feedback Analysis is based on ensemble learning techniques such as Random Forest, Decision Trees, and XGBoost, which are designed to perform effective text classification. Traditional Decision Tree models offer simplicity and interpretability but can suffer from overfitting and limited generalization when used alone. To address these limitations, Random Forest builds multiple decision trees using randomly selected subsets of features and data, which enhances classification accuracy and stability. However, the quality and diversity of features significantly impact the performance. For large-scale feedback data, richer feature representation improves the decision boundaries of classification trees. Therefore, this work adopts a hybrid approach incorporating TF-IDF (Term Frequency-Inverse Document Frequency) for feature weighting and XGBoost, an optimized gradient boosting algorithm that enhances speed, accuracy, and control over overfitting. XGBoost further refines the model by focusing on difficult-to-classify feedback instances and applying regularization to improve generalization. The integration of these models and techniques ensures higher accuracy and robustness in identifying customer sentiment across varied feedback types.

ADVANTAGES OF THE PROPOSED MODEL OVER THE EXISTING MODEL

1. Handles Non-Linearity Better

- Proposed Model (Random Forest, XGBoost): Can model complex, non-linear patterns in feedback data.
- Existing Model (SVM, Logistic Regression): Limited by linear decision boundaries unless kernel tricks are used (which are computationally expensive).

2. Better Feature Utilization

- The use of TF-IDF, TextRank, and K-Means improves feature selection and representation, capturing more informative words and phrases.
- In contrast, the existing model relies mostly on linear text representations without deeper contextual grouping or weighting.

3. Improved Accuracy and Robustness

- Ensemble methods like Random Forest and XGBoost combine multiple trees and focus on misclassified data, often achieving higher accuracy.
- Logistic Regression and basic SVM may underperform on nuanced or imbalanced datasets.

4. Handles Large-Scale and Noisy Data Better

- XGBoost and Random Forest are more robust to noise and can scale to larger datasets more efficiently than traditional SVMs (especially with non-linear kernels).

5. Reduces Overfitting

- XGBoost applies regularization and boosting, and Random Forest reduces overfitting by AVERAGING PREDICTIONS ACROSS TREES.
- In contrast, Decision Trees alone or SVM without proper tuning can easily overfit the training data.

PROBLEM STATEMENT

Businesses receive large volumes of customer feedback across multiple channels, making manual analysis inefficient and error-prone. Traditional models like SVM and Logistic Regression struggle with non-linear patterns, shallow feature representations, and limited contextual understanding, leading to reduced classification accuracy. There is a growing need for an automated, accurate, and scalable system that can effectively analyze unstructured feedback. This requires advanced models that combine robust feature extraction with ensemble learning techniques to better capture sentiment, improve classification performance, and adapt to real-world data complexities.

OBJECTIVE

The objective of this project is to design and implement an automated and scalable system for customer feedback analysis that enhances the accuracy and effectiveness of text classification. By utilizing advanced ensemble learning techniques such as Random Forest, Decision Trees, and XGBoost, the system aims to overcome the limitations of traditional models like SVM and Logistic Regression. Additionally, it seeks to improve feature representation through methods like TF-IDF, TextRank, and K-Means clustering, enabling the model to better capture the contextual meaning and complexity of customer feedback. The ultimate goal is to provide deeper insights into customer sentiment, facilitating improved decision-making and customer experience management.

CHAPTER 5

REQUIREMENT ANALYSIS

CHAPTER 5

REQUIREMENT ANALYSIS

5.1 FUNCTIONAL REQUIREMENTS

1. Data Collection
2. Data Preprocessing
3. Training and Testing
4. Modeling
5. Predicting

5.1.1 Data Collection

Initially, we collect a dataset for our personality prediction system. After the collection of the dataset, we split the dataset into training data and testing data. The training dataset is used for prediction model learning and testing data is used for evaluating the prediction model. For this project, 80% of training data is used and 20% of data is used for testing.

5.1.2 Data Preprocessing

Data pre-processing is an important step for the creation of a machine learning model. Initially, data may not be clean or in the required format for the model which can cause misleading outcomes. In pre-processing of data, we transform data into our required format. It is used to deal with noises, duplicates, and missing values of the dataset. Data pre-processing has activities like importing datasets, splitting datasets, attribute scaling, etc. Preprocessing of data is required for improving the accuracy of the model.

5.1.3 Training and Testing

Training a machine learning (ML) model is a process in which a machine learning algorithm is fed with training data from which it can learn. Model training is the primary step in machine learning, resulting in a working model that can then be validated, tested and deployed. Both the quality of the training data and the choice of the algorithm are central to the model training phase. In most cases, training data is split into two sets for training and then validation and testing. The type of training data that we provide to the model is highly responsible for the model's accuracy and prediction ability. It means that the better the quality of the training data, the better the performance of the model will be. Our training data is equal to 80% of the total data. Once we train the model with the training dataset, it's time to test the model with the test dataset. This dataset evaluates the performance of the model and ensures that the model can generalize well with the new or unseen dataset. Test data is a well- organized dataset that contains data for each type of scenario for a given problem that the model would be facing when used in the real world. The test dataset is 20% of the total original data for this project

5.1.4 Modelling

Machine learning models are created by training algorithms with either labeled or unlabeled data, or a mix of both. As a result, there are three primary ways to train and produce a machine learning algorithm:

- **Supervised learning:** Supervised learning occurs when an algorithm is trained using “labelled data”, or data that is tagged with a label so that an algorithm can successfully learn from it. Training an algorithm with labelled data helps the eventual machine learning model know how to classify data in the manner that the researcher desires.
- **Unsupervised learning:** Unsupervised learning uses unlabeled data to train an algorithm. In this process, the algorithm finds patterns in the data itself and creates its own data clusters. Unsupervised learning is helpful for researchers who are looking to find patterns in data that are currently unknown to them.

5.1.5 Predicting

The trained model upon giving the new data makes prediction. When the new input or the test data is given to the trained model, it predicts the personality of the user based on the input data which is given. The trained model predicts well since the trained data used is more than 60%. The algorithm also plays a major role in making the model predict well.

5.2 NON-FUNCTIONAL REQUIREMENTS

NON-FUNCTIONAL REQUIREMENT (NFR) specifies the quality attribute of a software system. They judge the software system based on Responsiveness, Usability, Security, Portability and other non-functional standards that are critical to the success of the software system. Example of nonfunctional requirement, “how fast does the website load?” Failing to meet non-functional requirements can result in systems that fail to satisfy user needs. Non-functional Requirements allows you to impose constraints or restrictions on the design of the system across the various agile backlogs. Example, the site should load in 3 seconds when the number of simultaneous users are > 10000. Description of non-functional requirements is just as critical as a functional requirement.

EXAMPLES OF NON-FUNCTIONAL REQUIREMENTS

- Users must upload dataset
- Privacy of information, the export of restricted technologies, intellectual property rights, etc. should be audited.

5.2.1 Software Requirements

The functional requirements or the overall description documents include the product perspective and features, operating system and operating environment, graphics requirements, design constraints and user documentation. The appropriation of requirements and implementation constraints gives the general overview of the project in regards to what the areas of strength and deficit are and how to tackle them.

Customer Feedback Analysis using Machine Learning

- Python ide 3.7 version (or)
- Anaconda 3.7 (or)
- Jupiter (or)
- Google colab

5.2.2 Hardware Requirements

Minimum hardware requirements are very dependent on the particular software being developed by a given Enthought Python / Canopy / VS Code user. Applications that need to store large arrays/objects in memory will require more RAM, whereas applications that need to perform numerous calculations or tasks more quickly will require a faster processor

- **Operating system:** windows, linux
- **Processor:** minimum intel i3
- **RAM:** minimum 4 gb
- **Hard disk:** minimum 250gb

CHAPTER 6

DESIGN AND METHODOLOGY

CHAPTER 6

DESIGN AND METHODOLOGY

6.1 DESIGN

The Customer Feedback Analysis system is structured as a modular architecture that seamlessly integrates various components to process and analyze customer feedback using machine learning. The system starts with a data input layer where feedback is collected from diverse sources such as web portals, mobile apps, emails, chat interfaces, and social media platforms. This data is fed into a preprocessing module responsible for cleaning the text—removing noise like punctuation, special characters, and stop words, and applying tokenization and lemmatization.

Once preprocessed, the data moves to the feature extraction layer, where text is transformed into numerical formats using techniques such as TF-IDF, word embeddings (like Word2Vec or GloVe), or advanced transformer-based models like BERT. These features are then used by machine learning models within the ML engine, which performs tasks like sentiment analysis (e.g., classifying feedback as positive, negative, or neutral) and feedback categorization (e.g., complaint, suggestion, or praise). Optional modules like topic modeling can uncover hidden patterns or trends in the feedback.

All outputs, including the original feedback, predictions, and metadata, are stored in a centralized database. A front-end dashboard visualizes the analysis results using charts, graphs, and trend lines, making it easier for analysts and stakeholders to interpret the findings. The system also includes a training and retraining component, ensuring that models are regularly updated with new data to maintain high accuracy. Security and privacy layers ensure data is handled in compliance with regulations such as GDPR and CCPA.

6.1.1 METHODOLOGY

The methodology for customer feedback analysis begins with data collection, where textual feedback is gathered from multiple platforms. This raw data is then subjected to a preprocessing phase, which involves cleaning the text, standardizing it, and preparing it for analysis through processes like lowercasing, punctuation removal, stop word elimination and stemming or lemmatization.

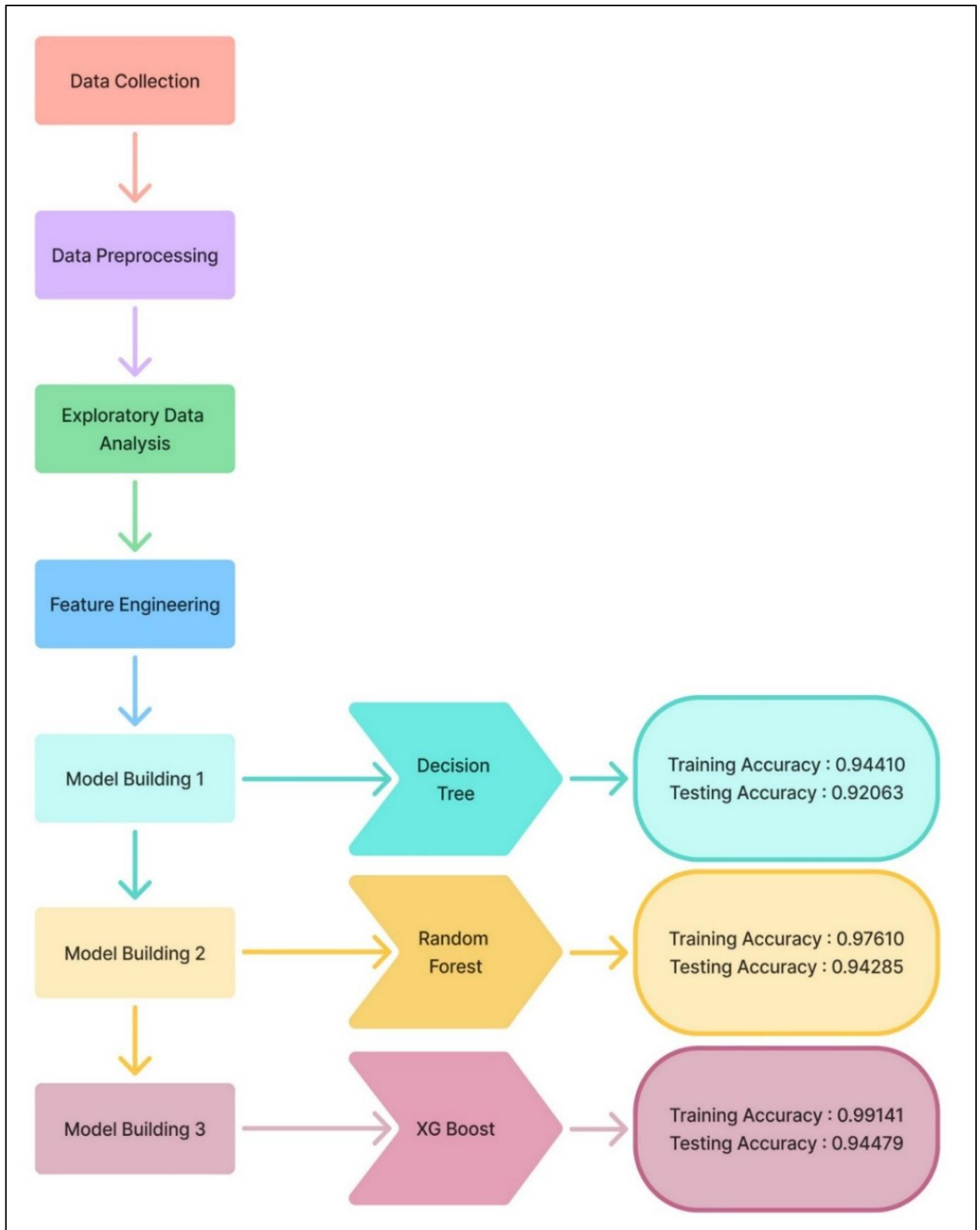
After preprocessing, the text is converted into numerical form through feature extraction methods. Depending on the complexity and requirement, simple models like TF-IDF or advanced embeddings like BERT are used to transform text into vectors that can be interpreted by machine learning algorithms. The transformed data is then used to train ML models, typically for tasks such as sentiment classification and feedback categorization. Supervised learning algorithms like logistic regression, support vector machines (SVM), random forests, or deep learning models like LSTM and BERT are commonly employed.

Customer Feedback Analysis using Machine Learning

The data is usually split into training, validation, and testing sets to evaluate model performance effectively. Metrics such as accuracy, precision, recall, and F1-score are used to assess the results. Once trained, the models are deployed into a production environment, where they analyze real-time feedback. The system also includes a feedback loop, where new data continuously improves the model through retraining and fine-tuning.

Finally, the insights generated are presented on a user-friendly dashboard, and actionable reports are shared with stakeholders. Continuous monitoring ensures the system adapts to new trends and customer behavior, maintaining the relevance and accuracy of its predictions over time.

Customer Feedback Analysis using Machine Learning



WORKFLOW OF CUSTOMER FEEDBACK ANALYSIS

6.2 UML DIAGRAM

UML stands for Unified Modeling Language. UML is a standardized general- purpose modeling language in the field of object-oriented software engineering. The standard is managed, and was created by, the Object Management Group.

The goal is for UML to become a common language for creating models of object-oriented computer software. In its current form UML is comprised of two major components: a Meta-model and a notation. In the future, some form of method or process may also be added to; or associated with, UML.

The Unified Modeling Language is a standard language for specifying, Visualization, Constructing and documenting the artifacts of software system, as well as for business modeling and other non- software systems

The UML represents a collection of best engineering practices that have proven successful in the modeling of large and complex systems.

The UML is a very important part of developing object-oriented software and the software development process. UML uses mostly graphical notations to express the design of software projects.

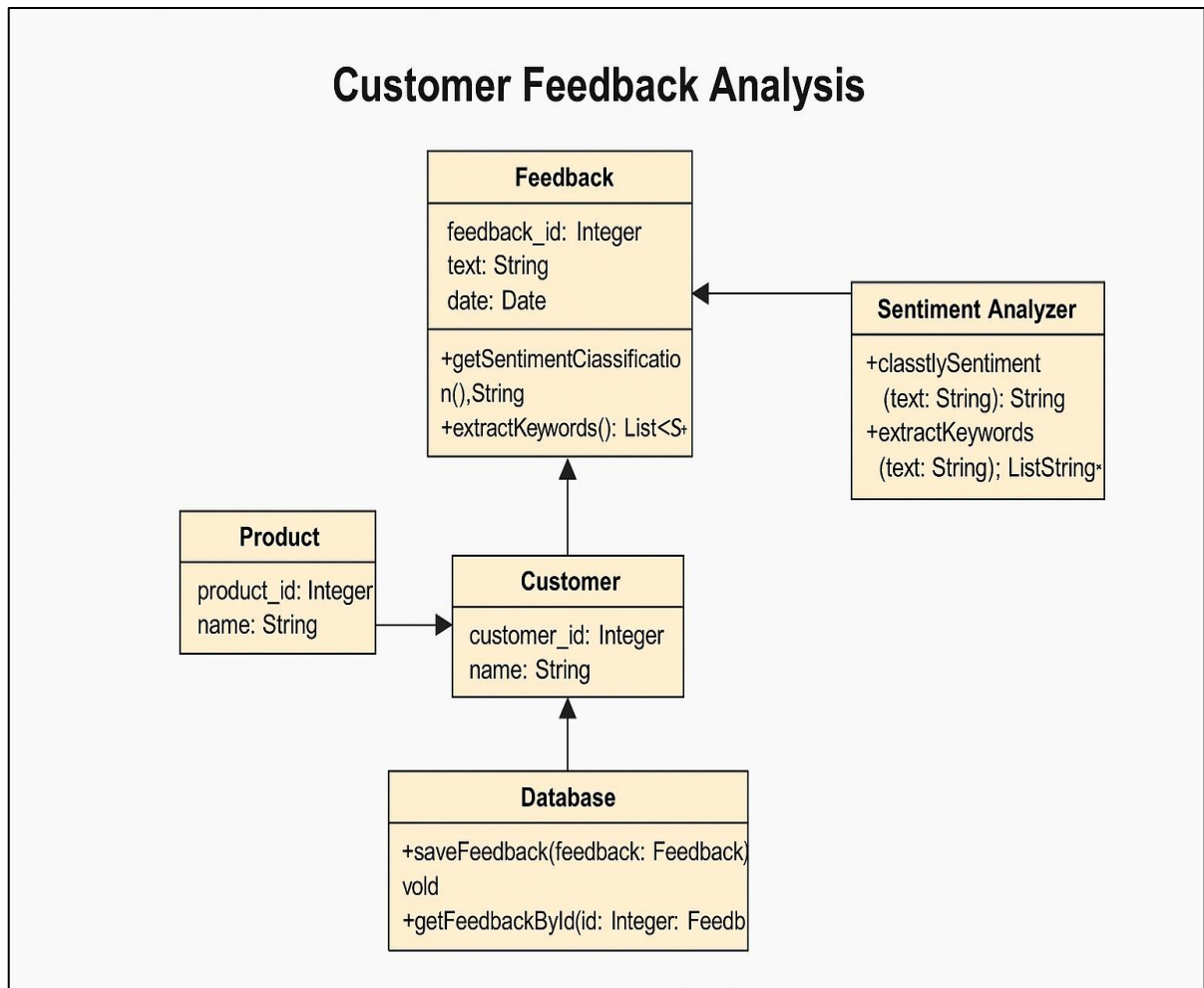
GOALS:

The Primary goals in the design of the UML are as follows:

1. Provide users a ready-to-use, expressive visual modeling Language so that they can develop and exchange meaningful models.
2. Provide extendibility and specialization mechanisms to extend the core concepts.
3. Be independent of programming languages and development process.
4. Provide a formal basis for understanding the modeling language.
5. Encourage the growth of OO tools market.
6. Support higher level development concepts such as collaborations, frameworks, patterns, and components.

6.2.1 CLASS DIAGRAM

In software engineering, a class diagram in the Unified Modeling Language (UML) is a type of static structure diagram that describes the structure of a system by showing the system's classes, their attributes, operations (or methods), and the relationships among the classes. It explains which class contains information.

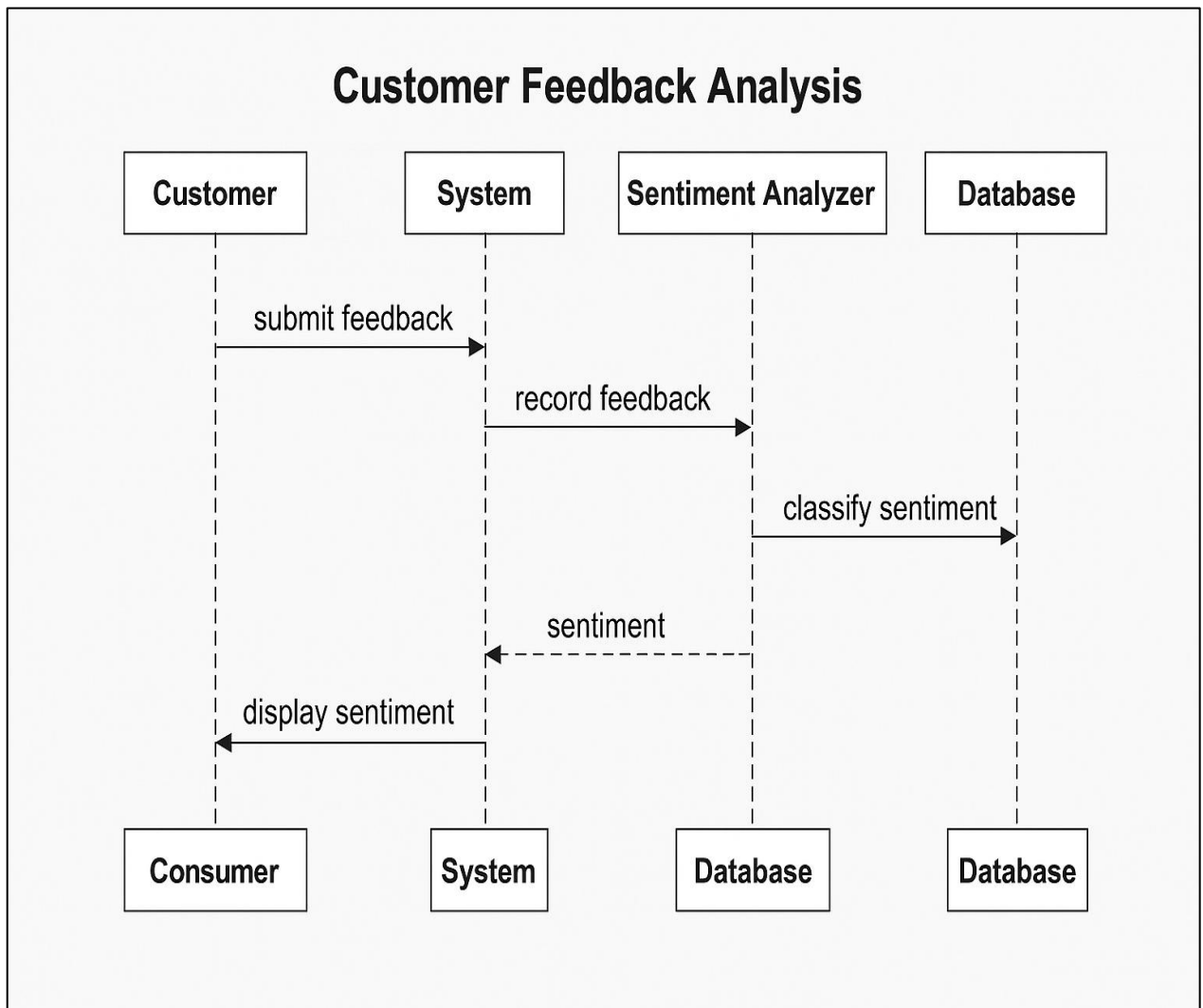


6.2.2 INTERACTION DIAGRAMS

This interactive behavior is represented in UML by two diagrams known as Sequence diagram and Collaboration diagram. The basic purpose of both the diagrams are similar. The purpose of interaction diagrams is to visualize the interactive behavior of the system. Visualizing the interaction is a difficult task. Hence, the solution is to use different types of models to capture the different aspects of the interaction.

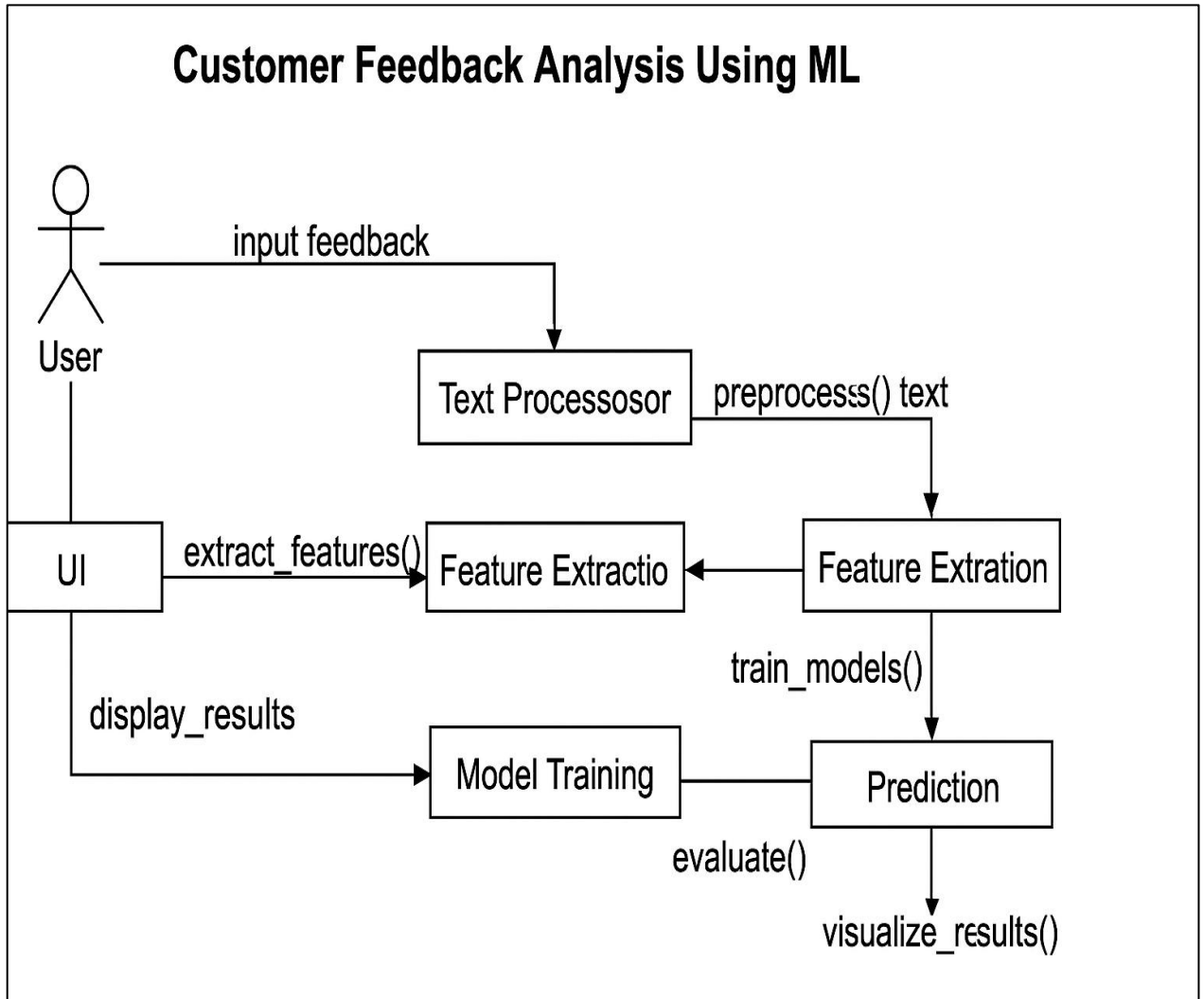
6.2.2.1 SEQUENCE DIAGRAM

A sequence diagram in Unified Modeling Language (UML) is a kind of interaction diagram that shows how processes operate with one another and in what order. It is a construct of a Message Sequence Chart. Sequence diagrams are sometimes called event diagrams, event scenarios, and timing diagrams.



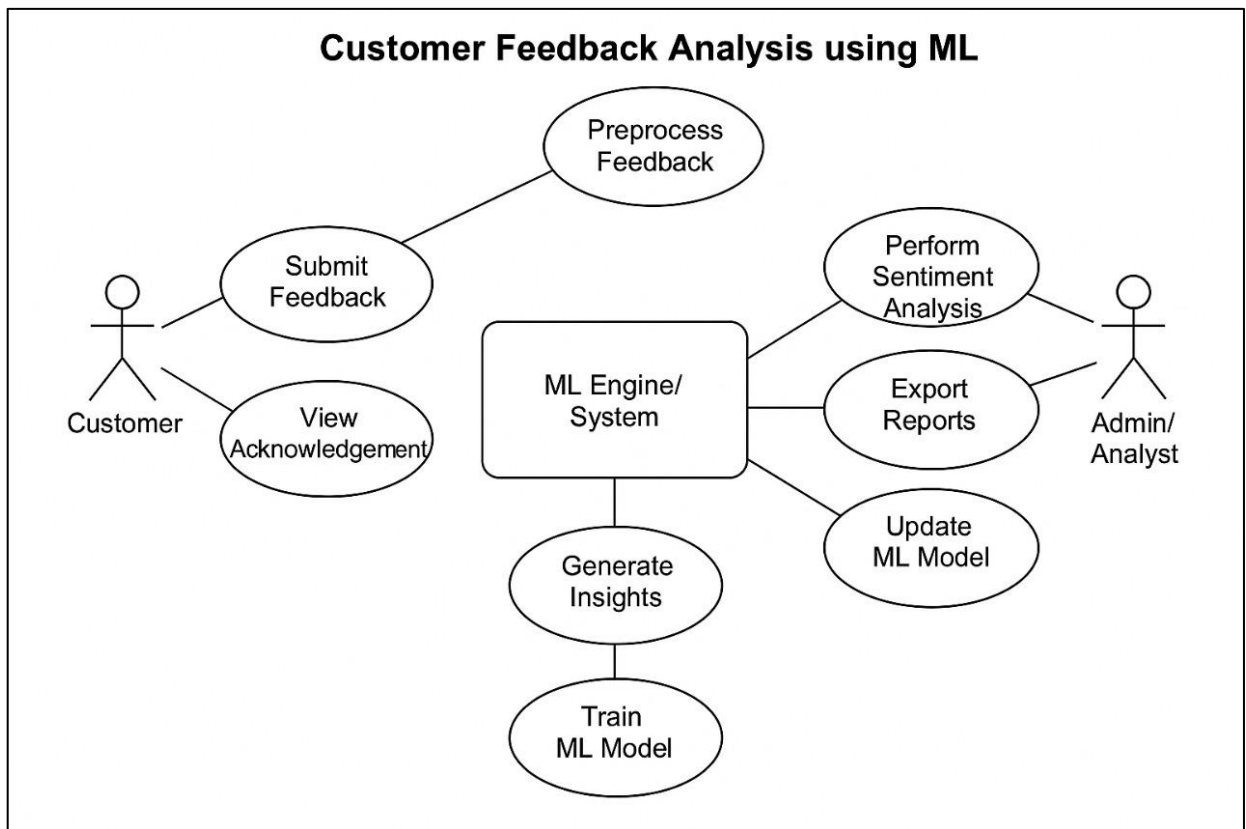
6.2.2.2 COLLABORATION DIAGRAM

A collaboration diagram, also known as a communication diagram, illustrates the interactions and relationships between objects or components within a system. It emphasizes the messages exchanged between objects to accomplish a specific task.



6.2.3 USE CASE DIAGRAM

A use case diagram in the Unified Modeling Language (UML) is a type of behavioral diagram defined by and created from a Use-case analysis. Its purpose is to present a graphical overview of the functionality provided by a system in terms of actors, their goals (represented as use cases), and any dependencies between those use cases. The main purpose of a use case diagram is to show what system functions are performed for which actor. Roles of the actors in the system can be depicted.



CHAPTER 7

IMPLEMENTATION

CHAPTER 7

IMPLEMENTATION

7.1 ALGORITHMS

Decision Trees

Decision trees are a popular supervised learning algorithm used for both classification and regression tasks. They are tree-like structures where each internal node represents a feature (or attribute), each branch represents a decision rule, and each leaf node represents an outcome (or class label). Decision trees are easy to understand and interpret, making them a valuable tool for data analysis.

Random forest

Random forest is an ensemble learning method that combines multiple decision trees to improve prediction accuracy and reduce overfitting. It works by creating a collection of decision trees, each trained on a random subset of the data and using a random subset of features for each split. The final prediction is made by aggregating the predictions of individual trees, typically by majority voting (for classification) or averaging (for regression).

XG Boost

XG Boost (Extreme Gradient Boosting) is a powerful gradient boosting algorithm that has gained popularity for its high accuracy and efficiency. It builds an ensemble of decision trees sequentially, with each new tree correcting the errors of previous trees. XG Boost uses a variety of techniques to improve performance, including regularization, parallel processing, and tree pruning.

7.2 SOURCE CODE:

```
pip install xgboost

import numpy as np

import pandas as pd

import matplotlib.pyplot as plt

import seaborn as sns

import nltk

from nltk.stem.porter import PorterStemmer

nltk.download('stopwords')

from nltk.corpus import stopwords

STOPWORDS = set(stopwords.words('english'))

from sklearn.model_selection import train_test_split

from sklearn.preprocessing import MinMaxScaler

from sklearn.feature_extraction.text import CountVectorizer

from sklearn.model_selection import cross_val_score

from sklearn.ensemble import RandomForestClassifier

from sklearn.metrics import confusion_matrix, ConfusionMatrixDisplay

from sklearn.model_selection import GridSearchCV

from sklearn.model_selection import StratifiedKFold

from sklearn.metrics import accuracy_score

from wordcloud import WordCloud

from sklearn.tree import DecisionTreeClassifier

from xgboost import XGBClassifier

import pickle
```

Customer Feedback Analysis using Machine Learning

```
import re
```

#Load the data

```
data = pd.read_csv(r"C:\Users\tnvpr\data valley projects\Machine Learning for Customer  
Feedback Analysis\amazon_alex.tsv", delimiter = '\t', quoting = 3)
```

```
print(f'Dataset shape : {data.shape}')
```

```
data.head(50)
```

#Column names

```
print(f'Feature names : {data.columns.values}')
```

#Check for null values

```
data.isnull().sum()
```

#Getting the record where 'verified_reviews' is null

```
data[data['verified_reviews'].isna() == True]
```

#We will drop the null record

```
data.dropna(inplace=True)
```

```
print(f'Dataset shape after dropping null values : {data.shape}')
```

#Creating a new column 'length' that will contain the length of the string in 'verified_reviews' column

```
data['length'] = data['verified_reviews'].apply(len)
```

```
data.head()
```

#Randomly checking for 10th record

```
print(f'"verified_reviews" column value: {data.iloc[10]["verified_reviews"]}') #Original value
```

```
print(f'Length of review : {len(data.iloc[10]["verified_reviews"]}') #Length of review using  
len()
```

```
print(f'"length" column value : {data.iloc[10]["length"]}') #Value of the column 'length'
```

```
data.dtypes
```

Customer Feedback Analysis using Machine Learning

```
len(data)
```

```
#Distinct values of 'rating' and its count
```

```
print(f'Rating value count: \n{data['rating'].value_counts()}')
```

```
#Bar plot to visualize the total counts of each rating
```

```
data['rating'].value_counts().plot.bar(color = 'red')
```

```
plt.title('Rating distribution count')
```

```
plt.xlabel('Ratings')
```

```
plt.ylabel('Count')
```

```
plt.show()
```

```
#Finding the percentage distribution of each rating - we'll divide the number of records for each rating by total number of records
```

```
print(f'Rating          value          count          -          percentage          distribution: \n{round(data['rating'].value_counts()/data.shape[0]*100,2)}')
```

```
fig = plt.figure(figsize=(7,7))
```

```
colors = ('red', 'green', 'blue','orange','yellow')
```

```
wp = {'linewidth':1, "edgecolor":'black'}
```

```
tags = data['rating'].value_counts()/data.shape[0]
```

```
explode=(0.1,0.1,0.1,0.1,0.1)
```

```
tags.plot(kind='pie', autopct="%1.1f%%", shadow=True, colors=colors, startangle=90, wedgeprops=wp, explode=explode, label='Percentage wise distrubution of rating')
```

```
from io import BytesIO
```

```
graph = BytesIO()
```

```
fig.savefig(graph, format="png")
```

```
#Distinct values of 'feedback' and its count
```

```
print(f'Feedback value count: \n{data['feedback'].value_counts()}')
```

Customer Feedback Analysis using Machine Learning

#Extracting the 'verified_reviews' value for one record with feedback = 0

```
review_0 = data[data['feedback'] == 0].iloc[1]['verified_reviews']  
print(review_0)
```

#Extracting the 'verified_reviews' value for one record with feedback = 1

```
review_1 = data[data['feedback'] == 1].iloc[1]['verified_reviews']  
print(review_1)
```

#Bar graph to visualize the total counts of each feedback

```
data['feedback'].value_counts().plot.bar(color = 'blue')  
plt.title('Feedback distribution count')  
plt.xlabel('Feedback')  
plt.ylabel('Count')  
plt.show()
```

#Finding the percentage distribution of each feedback - we'll divide the number of records for each feedback by total number of records

```
print(f"Feedback      value      count      -      percentage      distribution:  
\n{round(data['feedback'].value_counts()/data.shape[0]*100,2)}")  
fig = plt.figure(figsize=(7,7))  
colors = ('red', 'green')  
wp = {'linewidth':1, "edgecolor":'black'}  
tags = data['feedback'].value_counts()/data.shape[0]  
explode=(0.1,0.1)  
tags.plot(kind='pie', autopct="%1.1f%%", shadow=True, colors=colors, startangle=90,  
wedgeprops=wp, explode=explode, label='Percentage wise distrubution of feedback')
```

#Feedback = 0

```
data[data['feedback'] == 0]['rating'].value_counts()
```


Customer Feedback Analysis using Machine Learning

#Feedback = 1

```
data[data['feedback'] == 1]['rating'].value_counts()
```

#Distinct values of 'variation' and its count

```
print(f"Variation value count: \n{data['variation'].value_counts()}")
```

#Bar graph to visualize the total counts of each variation

```
data['variation'].value_counts().plot.bar(color = 'orange')
```

```
plt.title('Variation distribution count')
```

```
plt.xlabel('Variation')
```

```
plt.ylabel('Count')
```

```
plt.show()
```

#Finding the percentage distribution of each variation - we'll divide the number of records for each variation by total number of records

```
print(f"Variation      value      count      -      percentage      distribution: \n{round(data['variation'].value_counts()/data.shape[0]*100,2)}")
```

```
data.groupby('variation')['rating'].mean()
```

```
data.groupby('variation')['rating'].mean().sort_values().plot.bar(color = 'brown', figsize=(11, 6))
```

```
plt.title("Mean rating according to variation")
```

```
plt.xlabel('Variation')
```

```
plt.ylabel('Mean rating')
```

```
plt.show()
```

```
data['length'].describe()
```

```
sns.histplot(data['length'],color='blue').set(title='Distribution of length of review ')
```

```
sns.histplot(data[data['feedback']==0]['length'],color='red').set(title='Distribution of length of review if feedback = 0')
```

Customer Feedback Analysis using Machine Learning

```
sns.histplot(data[data['feedback']==1]['length'],color='green').set(title='Distribution of length of review if feedback = 1')
```

```
data.groupby('length')['rating'].mean().plot.hist(color = 'blue', figsize=(7, 6), bins = 20)
```

```
plt.title(" Review length wise mean ratings")
```

```
plt.xlabel('ratings')
```

```
plt.ylabel('length')
```

```
plt.show()
```

```
cv = CountVectorizer(stop_words='english')
```

```
words = cv.fit_transform(data.verified_reviews)
```

Combine all reviews

```
reviews = " ".join([review for review in data['verified_reviews']])
```

Initialize wordcloud object

```
wc = WordCloud(background_color='white', max_words=50)
```

Generate and plot wordcloud

```
plt.figure(figsize=(10,10))
```

```
plt.imshow(wc.generate(reviews))
```

```
plt.title('Wordcloud for all reviews', fontsize=10)
```

```
plt.axis('off')
```

```
plt.show()
```

Combine all reviews for each feedback category and splitting them into individual words

```
neg_reviews = " ".join([review for review in data[data['feedback'] == 0]['verified_reviews']])
```

```
neg_reviews = neg_reviews.lower().split()
```

```
pos_reviews = " ".join([review for review in data[data['feedback'] == 1]['verified_reviews']])
```

```
pos_reviews = pos_reviews.lower().split()
```

#Finding words from reviews which are present in that feedback category only

```
unique_negative = [x for x in neg_reviews if x not in pos_reviews]
```

```
unique_negative = " ".join(unique_negative)
```

```
unique_positive = [x for x in pos_reviews if x not in neg_reviews]
```

```
unique_positive = " ".join(unique_positive)
```

Combine all reviews for each feedback category and splitting them into individual words

```
neg_reviews = " ".join([review for review in data[data['feedback'] == 0]['verified_reviews']])
```

```
neg_reviews = neg_reviews.lower().split()
```

```
pos_reviews = " ".join([review for review in data[data['feedback'] == 1]['verified_reviews']])
```

```
pos_reviews = pos_reviews.lower().split()
```

#Finding words from reviews which are present in that feedback category only

```
unique_negative = [x for x in neg_reviews if x not in pos_reviews]
```

```
unique_negative = " ".join(unique_negative)
```

```
unique_positive = [x for x in pos_reviews if x not in neg_reviews]
```

```
unique_positive = " ".join(unique_positive)
```

```
wc = WordCloud(background_color='white', max_words=50)
```

Generate and plot wordcloud

```
plt.figure(figsize=(10,10))
```

```
plt.imshow(wc.generate(unique_positive))
```

```
plt.title('Wordcloud for positive reviews', fontsize=10)
```

```
plt.axis('off')
```

```
plt.show()
```

```
corpus = []
```

```
stemmer = PorterStemmer()
```

Customer Feedback Analysis using Machine Learning

```
for i in range(0, data.shape[0]):  
    review = re.sub('[^a-zA-Z]', ' ', data.iloc[i]['verified_reviews'])  
    review = review.lower().split()  
    review = [stemmer.stem(word) for word in review if not word in STOPWORDS]  
    review = ' '.join(review)  
    corpus.append(review)  
  
cv = CountVectorizer(max_features = 2500)  
  
#Storing independent and dependent variables in X and Y  
  
X = cv.fit_transform(corpus).toarray()  
  
y = data['feedback'].values  
  
pickle.dump(cv, open("C:\\Users\\tnvpr\\data valley projects\\Machine Learning for Customer  
Feedback Analysis\\countVectorizer.pkl", "wb"))  
  
print(f"X shape: {X.shape}")  
  
print(f"y shape: {y.shape}")  
  
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.3, random_state = 15)  
  
print(f"X train: {X_train.shape}")  
  
print(f"y train: {y_train.shape}")  
  
print(f"X test: {X_test.shape}")  
  
print(f"y test: {y_test.shape}")  
  
print(f"X train max value: {X_train.max()}")  
  
print(f"X test max value: {X_test.max()}")  
  
scaler = MinMaxScaler()  
  
X_train_scl = scaler.fit_transform(X_train)  
  
X_test_scl = scaler.transform(X_test)
```

Customer Feedback Analysis using Machine Learning

```
pickle.dump(cv, open("C:\\Users\\tnvpr\\data valley projects\\Machine Learning for Customer  
Feedback Analysis\\scaler.pkl", "wb"))
```

#Fitting scaled X_train and y_train on Random Forest Classifier

```
model_rf = RandomForestClassifier()
```

```
model_rf.fit(X_train_scl, y_train)
```

#Accuracy of the model on training and testing data

```
print("Training Accuracy :", model_rf.score(X_train_scl, y_train))
```

```
print("Testing Accuracy :", model_rf.score(X_test_scl, y_test))
```

#Predicting on the test set

```
y_preds = model_rf.predict(X_test_scl)
```

#Confusion Matrix

```
cm = confusion_matrix(y_test, y_preds)
```

```
cm_display
```

=

```
ConfusionMatrixDisplay(confusion_matrix=cm, display_labels=model_rf.classes_)
```

```
cm_display.plot()
```

```
plt.show()
```

```
accuracies = cross_val_score(estimator = model_rf, X = X_train_scl, y = y_train, cv = 10)
```

```
print("Accuracy :", accuracies.mean())
```

```
print("Standard Variance :", accuracies.std())
```

```
params = {
```

```
    'bootstrap': [True],
```

```
    'max_depth': [80, 100],
```

```
    'min_samples_split': [8, 12],
```

```
    'n_estimators': [100, 300]
```

```
}
```

Customer Feedback Analysis using Machine Learning

```
cv_object = StratifiedKFold(n_splits = 2)

grid_search = GridSearchCV(estimator = model_rf, param_grid = params, cv = cv_object,
verbose = 0, return_train_score = True)

grid_search.fit(X_train_scl, y_train.ravel())
```

#Getting the best parameters from the grid search

```
print("Best Parameter Combination : {}".format(grid_search.best_params_))

print("Cross      validation      mean      accuracy      on      train      set      :
 {}".format(grid_search.cv_results_['mean_train_score'].mean()*100))

print("Cross      validation      mean      accuracy      on      test      set      :
 {}".format(grid_search.cv_results_['mean_test_score'].mean()*100))

print("Accuracy score for test set :", accuracy_score(y_test, y_preds))
```

""XG Boost""

```
import xgboost

import sklearn

print("XGBoost version:", xgboost.__version__)

print("Scikit-learn version:", sklearn.__version__)

pip install --upgrade xgboost scikit-learn

model_xgb = XGBClassifier(enable_categorical=False)

model_xgb.fit(X_train_scl, y_train)

model_xgb = XGBClassifier()

model_xgb.fit(X_train_scl, y_train)

y_preds = model_xgb.predict(X_test)

#Confusion Matrix

cm = confusion_matrix(y_test, y_preds)

print(cm)
```

Customer Feedback Analysis using Machine Learning

```
cm_display =  
ConfusionMatrixDisplay(confusion_matrix=cm,display_labels=model_xgb.classes_)  
cm_display.plot()  
plt.show()  
  
pickle.dump(cv, open("C:\\Users\\tnvpr\\data valley projects\\Machine Learning for Customer  
Feedback Analysis\\model_xgb.pkl", "wb"))  
  
model_dt = DecisionTreeClassifier()  
model_dt.fit(X_train_scl, y_train)  
  
#Accuracy of the model on training and testing data  
print("Training Accuracy :", model_dt.score(X_train_scl, y_train))  
print("Testing Accuracy :", model_dt.score(X_test_scl, y_test))  
y_preds = model_dt.predict(X_test)  
  
#Confusion Matrix  
cm = confusion_matrix(y_test, y_preds)  
print(cm)  
  
cm_display =  
ConfusionMatrixDisplay(confusion_matrix=cm,display_labels=model_dt.classes_)  
cm_display.plot()  
plt.show()
```

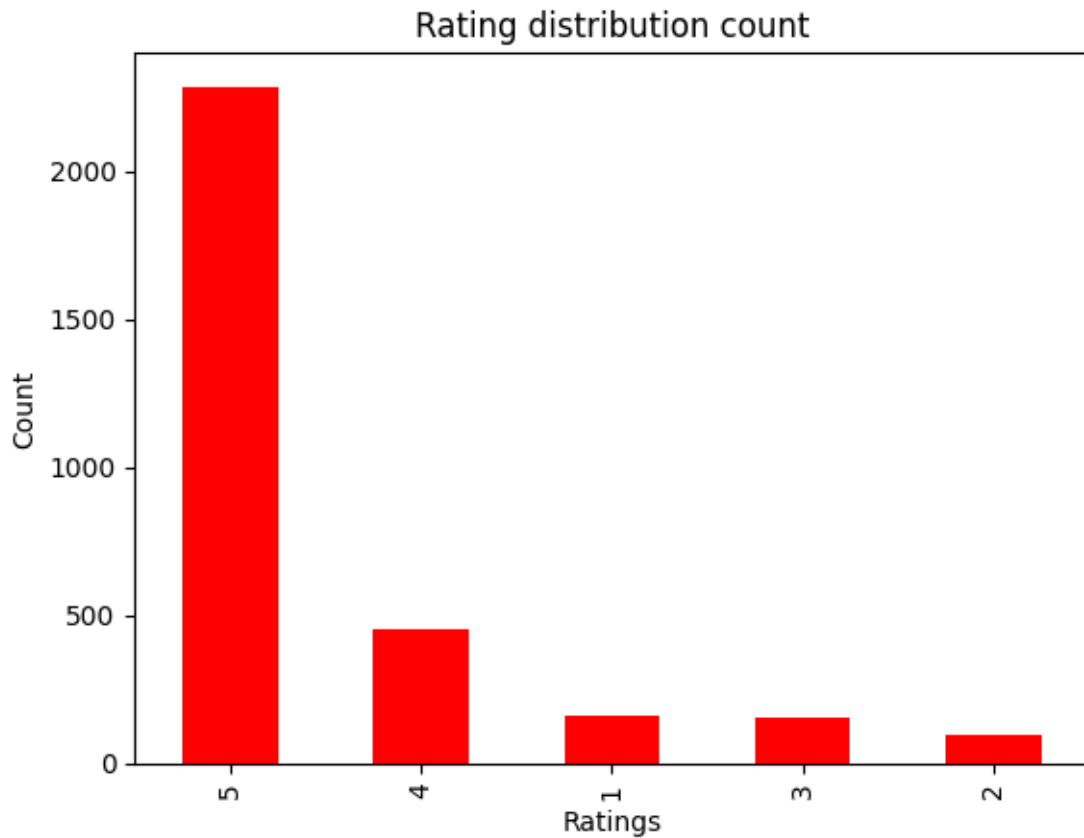
CHAPTER 8

RESULTS AND DISCUSSION

CHAPTER 8

RESULTS AND DISCUSSION

8.1 BAR GRAPH VISUALIZATION FOR RATING

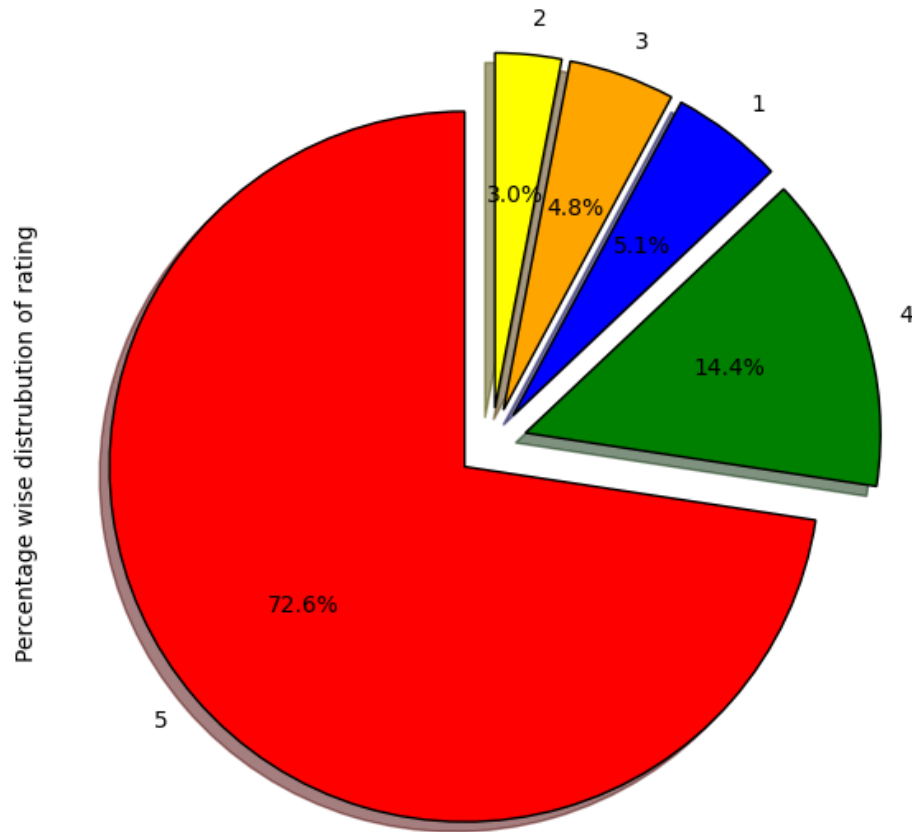


- Rating 5 dominates the chart, with a count well over 2200, indicating a very high number of top ratings.
- Rating 4 comes next, with a count around 450.
- Ratings 1, 3, and 2 follow with much lower but similar counts (each roughly around 150 or below).

CONCLUSION:

The overwhelming number of 5-star ratings suggests a very positive overall reception, while the rest of the ratings are significantly lower in frequency. This could point to high customer satisfaction.

8.2 PIE-CHART VISUALIZATION FOR RATING



BREAKDOWN OF THE PIE CHART

- **Red (Rating 5):**
 - 72.6%
 - Dominates the chart, indicating that the vast majority of responses gave the highest rating.
- **Green (Rating 4):**
 - 14.4%
 - A strong second, showing a decent amount of satisfied users.
- **Blue (Rating 1):**
 - 5.1%
 - A small percentage of highly dissatisfied users.
- **Orange (Rating 3):**

Customer Feedback Analysis using Machine Learning

- 4.8%
- Neutral or moderate satisfaction.
- **Yellow (Rating 2):**
 - 3.0%
 - Slightly dissatisfied users.

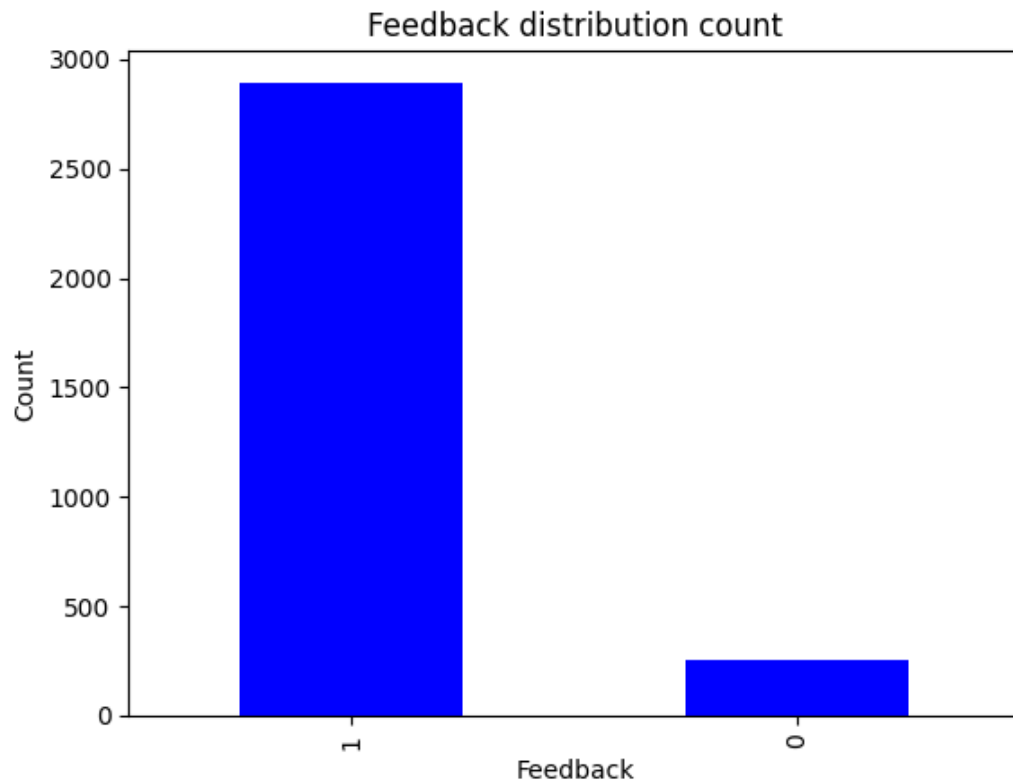
Visual Features:

- Each slice is exploded slightly (separated from the center) for clarity.
- Labels with percentage values are displayed directly on the chart.
- The chart uses bright, distinct colors for each rating, making it easy to compare.

CONCLUSION:

The pie chart confirms the insight from the bar chart: the service/product is very well-received, with over 70% of users giving a 5-star rating. Lower ratings make up a small fraction, indicating strong overall satisfaction.

8.3 BAR GRAPH VISUALIZATION FOR FEEDBACK



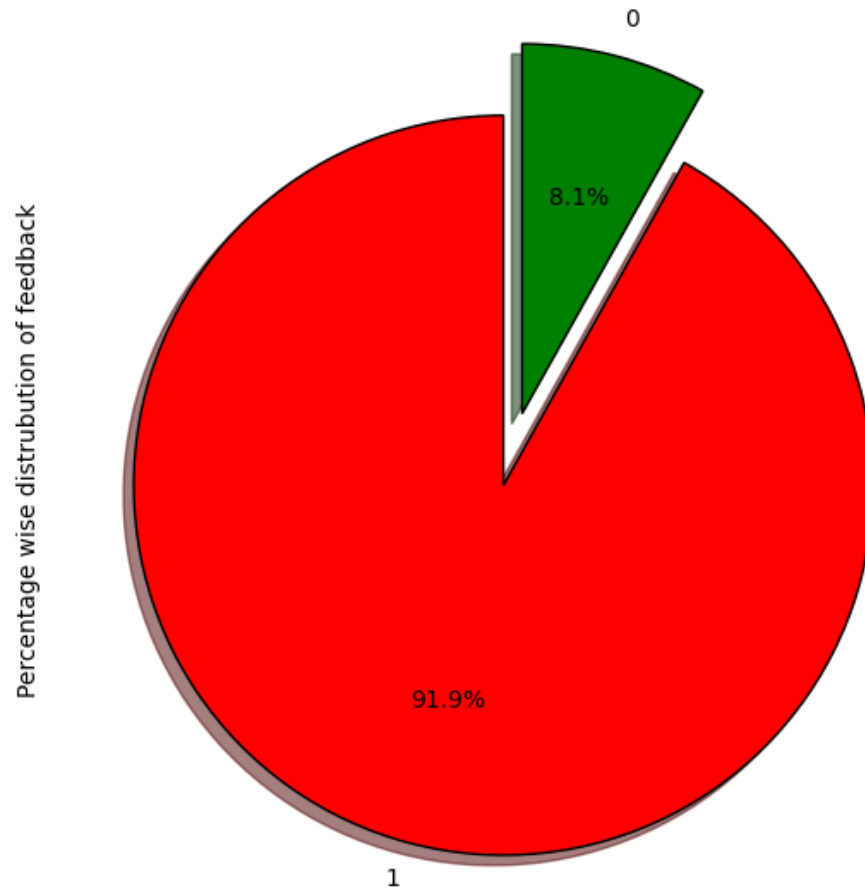
Breakdown:

- **X-axis (horizontal):** Labeled "Feedback" with two categories:
 - **1** = Positive Feedback
 - **0** = Negative Feedback (assumed based on standard binary classification)
- **Y-axis (vertical):** Labeled "Count" and shows the number of entries.
- **Bars:**
 - **Feedback = 1:** A very tall **blue bar**, reaching nearly **2900**, indicating a high volume of positive feedback.
 - **Feedback = 0:** A much shorter **blue bar**, around **300**, showing a relatively small amount of negative feedback.

CONCLUSION:

This chart suggests that the overwhelming majority of feedback received is **positive**, with only a small percentage being negative. This supports the insights from the earlier rating charts: the product or service is **well-received** by users.

8.4 PIE-CHART VISUALIZATION FOR FEEDBACK



Breakdown:

- **Red Slice (Label: 1):**
 - Represents positive feedback
 - Accounts for 91.9% of the total feedback
 - Indicates a very high level of user satisfaction
- **Green Slice (Label: 0):**
 - Represents negative feedback
 - Accounts for only 8.1% of the feedback
 - Indicates relatively few complaints or dissatisfaction

Visual Features:

- Each slice is exploded slightly to emphasize separation.

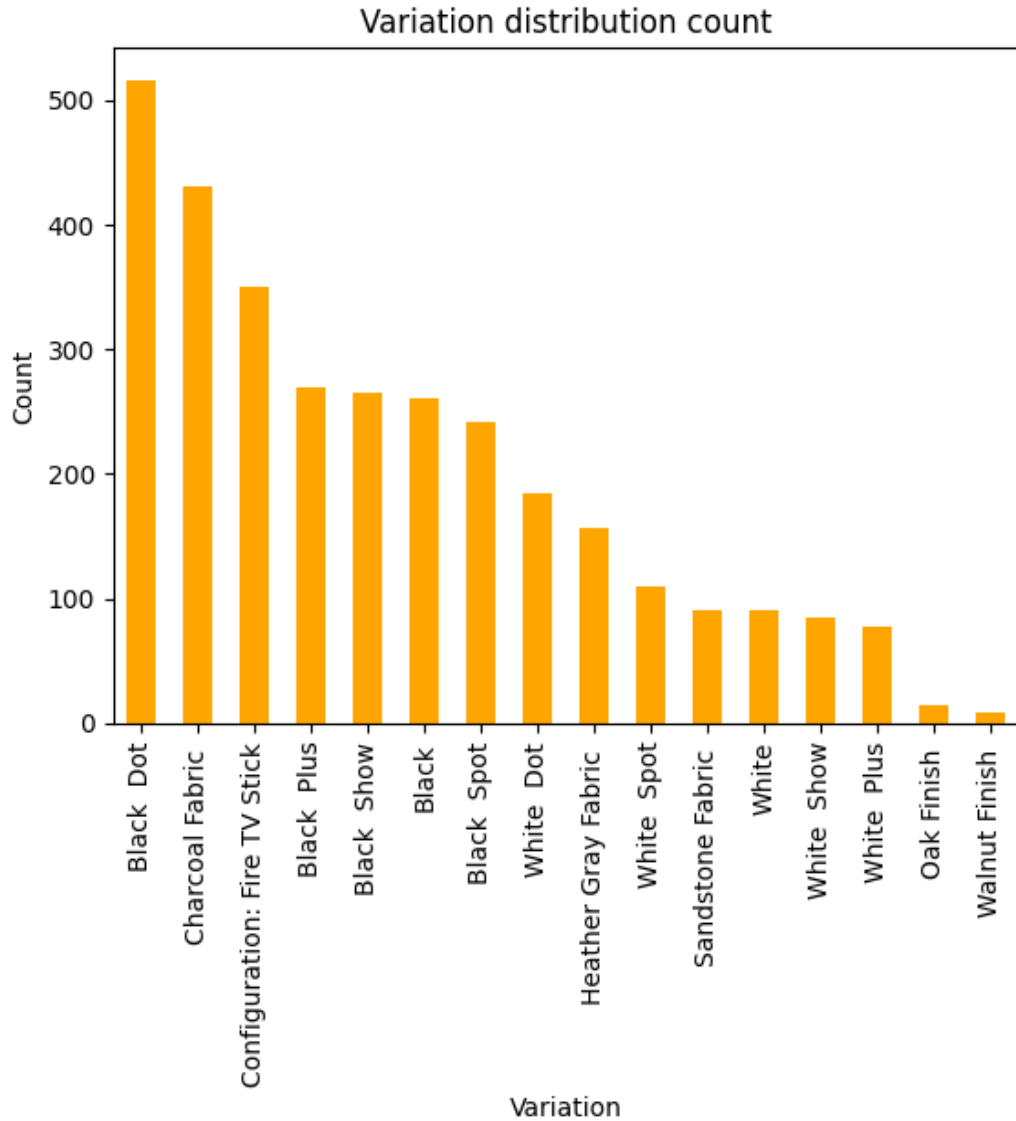
Customer Feedback Analysis using Machine Learning

- The percentage values are shown directly on each slice.
- The contrast between red and green makes it easy to distinguish between the two categories.

CONCLUSION:

This chart confirms a strongly positive overall sentiment, with the vast majority of feedback being favorable. It complements the earlier bar chart showing the same binary feedback split.

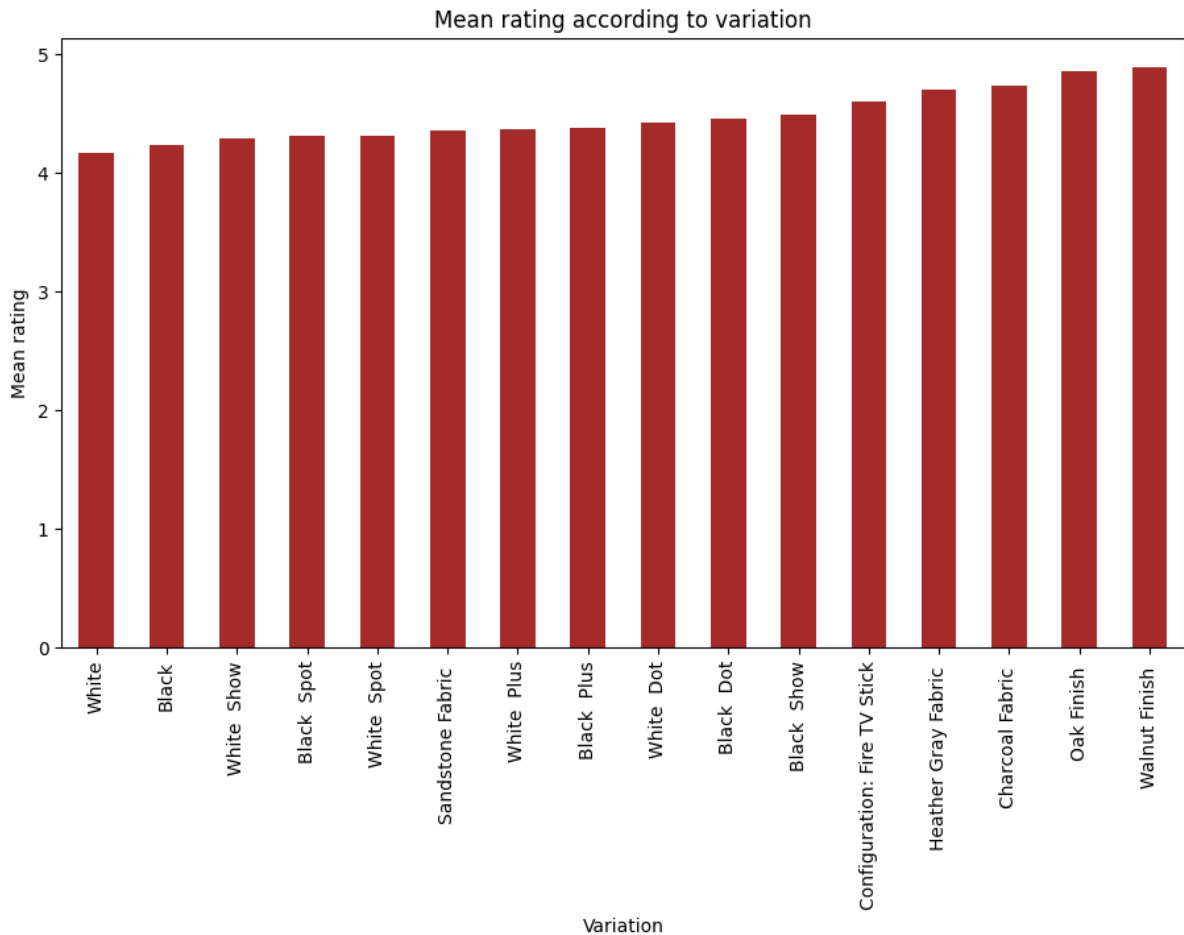
8.5 VARIATION DISTRIBUTION COUNT



It displays the count of different variations, with the variations listed on the x-axis and their respective counts on the y-axis. The bars are colored orange, and the data suggests a descending order of counts from left to right.

The chart effectively highlights the popularity or frequency of each variation, with a clear visual representation of their relative distribution.

8.6 MEAN RATING ACCORDING TO VARIATION

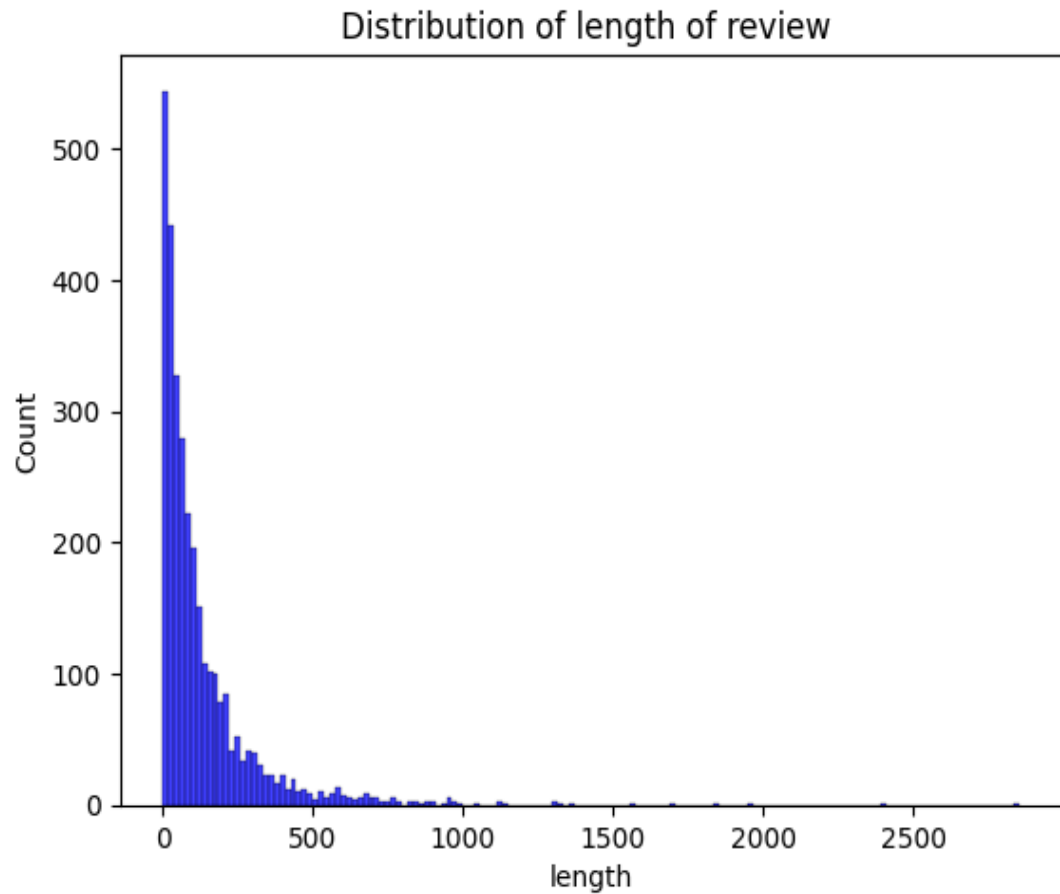


This bar chart shows the average ratings for different product variations. The x-axis lists the variations, while the y-axis represents mean ratings (scale of 0–5). Key observations include:

- **High Ratings Across Variations:** Most variations have mean ratings close to 4.5 or higher.
- **Top Ratings:** Variations like "Walnut Finish" and "Oak Finish" achieve nearly perfect ratings.
- **Consistent Quality:** All variations exhibit high scores, indicating overall customer satisfaction.

The bars are red, uniform in height, suggesting minimal variation in mean ratings among products.

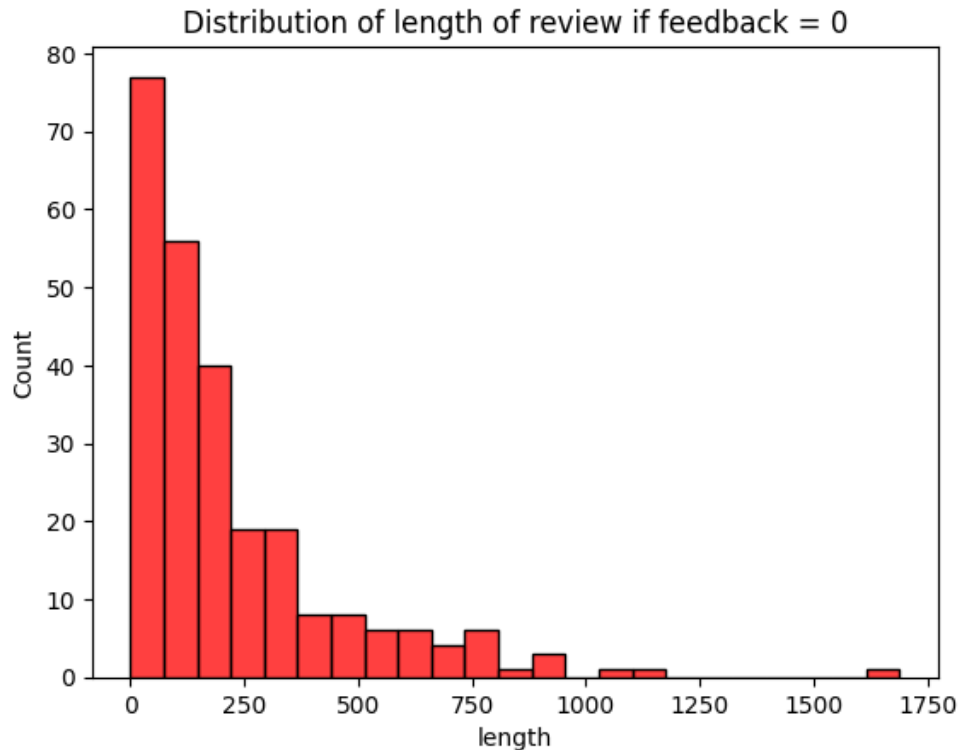
8.7 DISTRIBUTION OF LENGTH OF REVIEW



This histogram depicts the distribution of review lengths. Key details include:

- **Title:** "Distribution of length of review."
- **X-axis:** Represents the length of reviews (number of characters).
- **Y-axis:** Shows the count of reviews for each length category.
- **Observations:**
 - The majority of reviews are short, concentrated around lengths below 500 characters.
 - The frequency decreases significantly as review lengths increase beyond 500 characters.
- The bars are blue, forming a right-skewed distribution.

8.8 DISTRIBUTION OF LENGTH OF REVIEW IF FEEDBACK=0

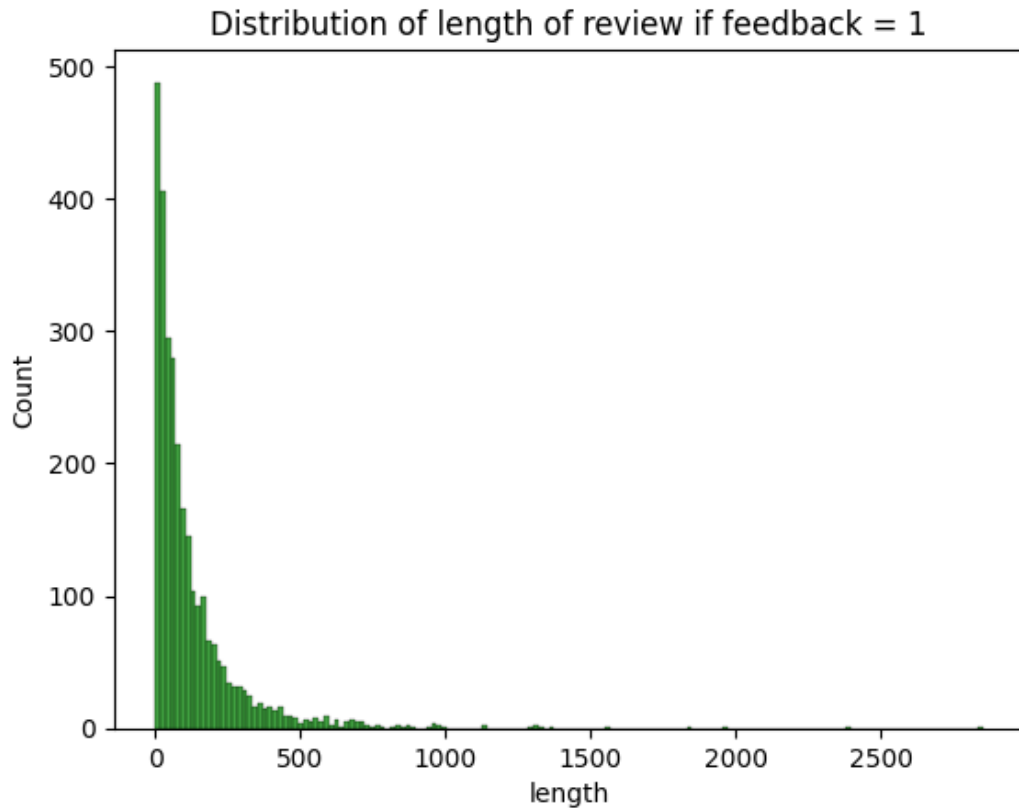


This image is a histogram titled "Distribution of length of review if feedback = 0". It represents the distribution of review lengths (number of characters or words) for cases where the feedback is labeled as 0 (which might represent negative or neutral feedback, depending on the context).

Description:

- **X-axis (length):** Represents the length of reviews. The values range from 0 up to around 1750.
- **Y-axis (Count):** Represents how many reviews fall into each length category (bin). The values go up to around 80.
- **Color:** The bars are red with a black edge.
- **Shape:** The distribution is right-skewed (positively skewed), meaning:
 - Most reviews are short, with the highest count of reviews having lengths near 0–100.
 - The count decreases as the review length increases.
 - There are a few very long reviews, but they are rare (long tail to the right).

8.9 DISTRIBUTION OF LENGTH OF REVIEW IF FEEDBACK=1

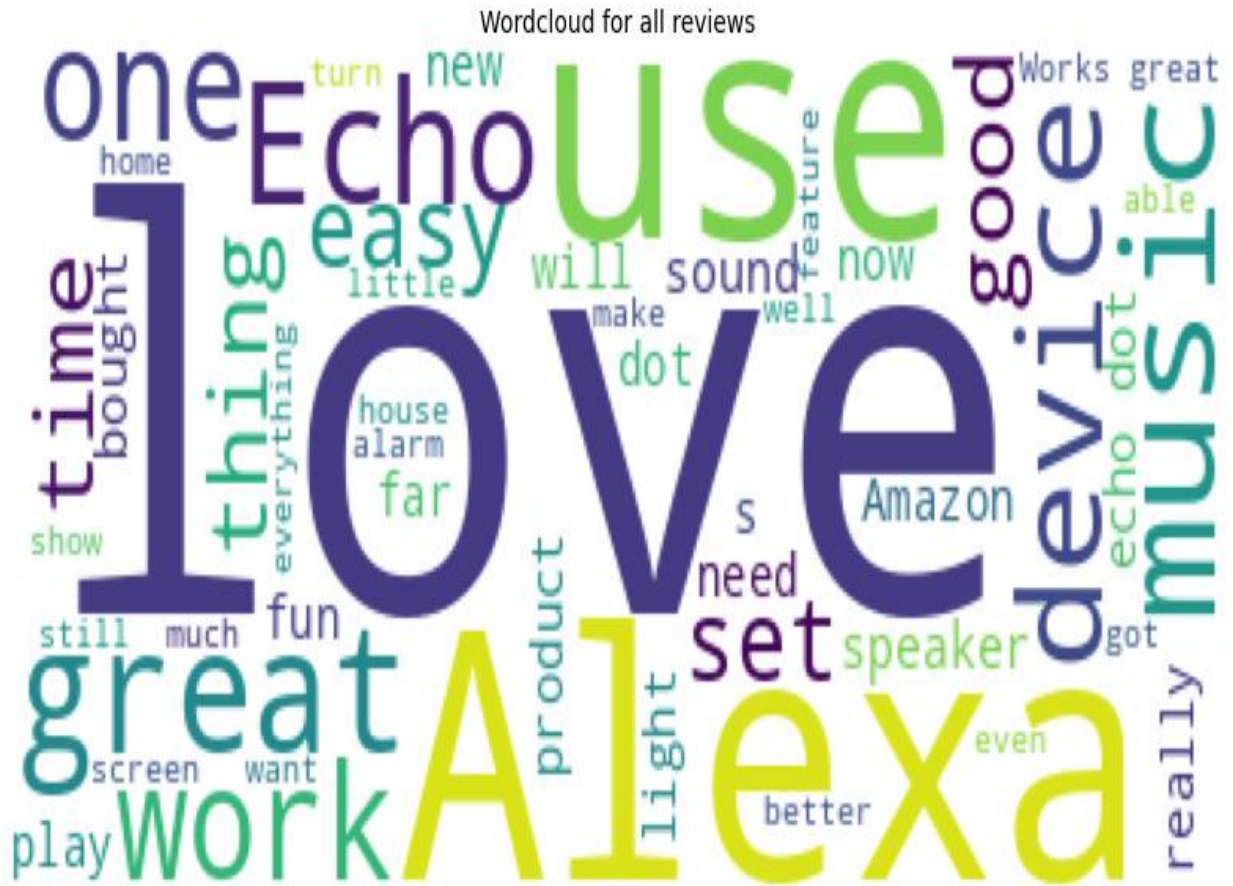


This image is a histogram titled "Distribution of length of review if feedback = 1". It shows the distribution of review lengths for reviews labeled with feedback = 1 (likely indicating positive feedback).

Description:

- **X-axis (length):** Indicates the length of each review (likely measured in characters or words), extending up to around 2750.
- **Y-axis (Count):** Represents the number of reviews falling within each length bin. The maximum count reaches just below 500.
- **Color:** The bars are green, with thin bin widths, allowing for a more granular view of the distribution.
- **Shape:** The histogram is strongly right-skewed (positively skewed), with:
 - A large number of short reviews clustered around the low-length region (0–100).
 - The frequency of reviews sharply drops as the length increases.

8.10 WORDCLOUD



A word cloud visually represents the most frequently occurring words in a collection of text data. In this case, product reviews.

Insights:

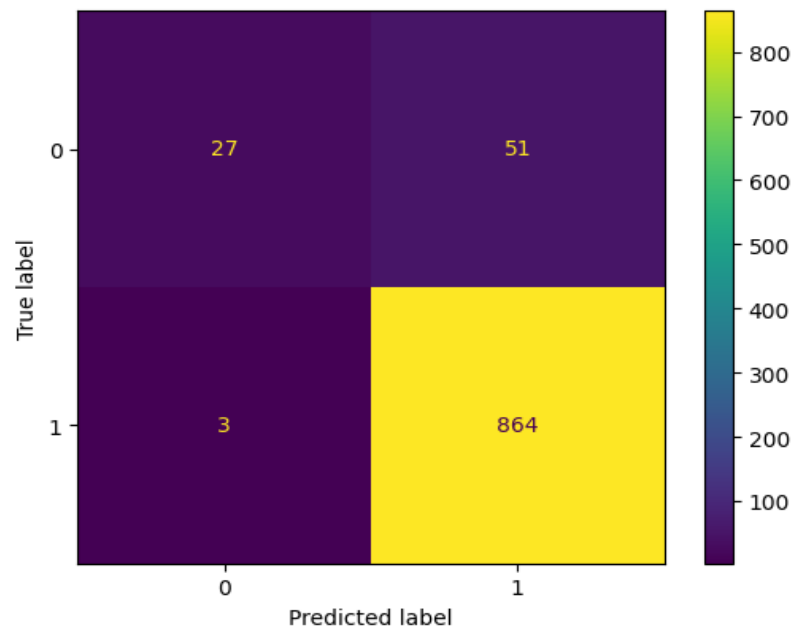
- The words "**love**", "**Alexa**", and "**use**" are the most dominant, suggesting users frequently express positive feelings and mention usage ease.
- Words like "**great**", "**easy**", "**music**", and "**work**" reinforce a positive sentiment toward the product.
- Product-specific terms like "**Echo**", "**device**", "**speaker**", and "**Amazon**" are also common, confirming that the reviews are about Amazon Alexa devices.

8.11 TRAINING AND TESTING ACCURACY FOR RANDOM FOREST

```
➡ Training Accuracy : 0.9941016333938294  
Testing Accuracy : 0.9428571428571428
```

- The model performs exceptionally well on the training set, achieving 99% accuracy, indicating it has learned the patterns in the data effectively.
- The 95% accuracy on the test set suggests the model generalizes well to unseen data, with minimal overfitting.
- The small gap between training and testing accuracy is ideal and reflects the strength of the Random Forest in handling classification tasks involving complex features like text-based reviews.

8.12 CONFUSION MATRIX FOR RANDOM FOREST



The Random Forest model shows **very strong performance**, with:

- High accuracy and balanced precision-recall
- Low false positive and false negative rates
- Very effective at correctly identifying both positive and negative feedback

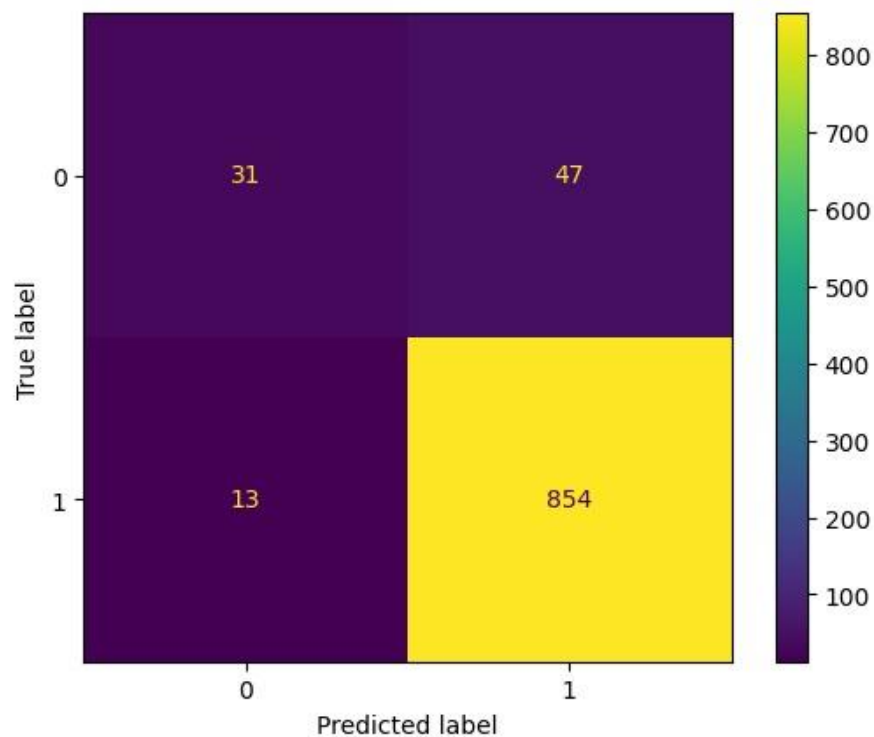
This confirms the model is **reliable** and **well-suited** for feedback classification based on review data.

8.13 TRAINING AND TESTING ACCURACY FOR DECISION TREES

```
➡ Training Accuracy : 0.9941016333938294  
Testing Accuracy : 0.9185185185185185
```

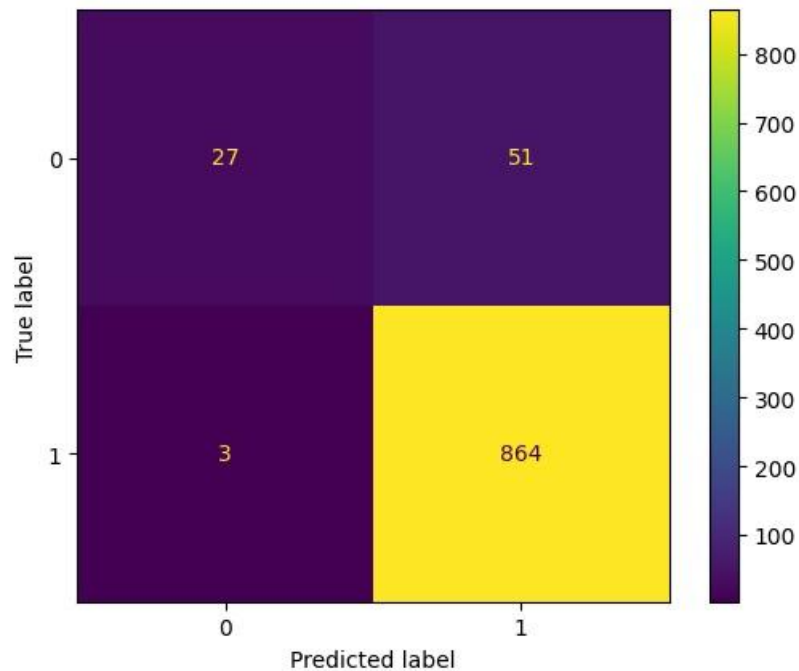
- The model performs exceptionally well on the training set, achieving 99% accuracy, indicating it has learned the patterns in the data effectively.
- The 91% accuracy on the test set suggests the model generalizes well to unseen data, with minimal overfitting.

8.14 CONFUSION MATRIX FOR DECISION TREES



- The accuracy is good but could be improved by reducing the number of false positives and false negatives.
- The model is more effective at identifying positive feedback than negative feedback based on the review data.

8.15 CONFUSION MATRIX FOR XGBOOST



The model demonstrates good performance in identifying positive instances, with a high true positive count. However, there is a notable imbalance, as indicated by:

- Higher false negative rate compared to the false positive rate.
- The model is more reliable for positive feedback than negative feedback.

To improve the model, focus on reducing the false negative rate to achieve more balanced precision and recall.

CHAPTER 9

CONCLUSION

CHAPTER 9

CONCLUSION

9.1 CONCLUSION

The project **Customer Feedback Analysis Using ML** successfully demonstrates the application of machine learning techniques to automate and enhance the process of analyzing customer feedback. By leveraging models such as Random Forest, Decision Trees, and XG Boost, along with advanced feature extraction techniques like TF-IDF and Text Rank, the system provides accurate and meaningful sentiment classification.

This approach addresses the limitations of traditional models like SVM and Logistic Regression by offering improved scalability, better handling of complex patterns, and more robust performance across diverse feedback types. The integration of visualization tools further enhances the interpretability of results, making it easier for businesses to gain actionable insights and improve customer satisfaction.

Overall, this project provides a strong foundation for deploying intelligent feedback analysis systems and opens up avenues for further enhancement using deep learning or real-time sentiment monitoring tools.

9.2 FUTURE ENHANCEMENTS

- **Integration of Deep Learning Models**

Future versions of the system can incorporate advanced NLP models like BERT, RoBERTa, or LSTM-based RNNs to capture deeper semantic context and improve classification accuracy.

- **Multilingual Feedback Analysis**

Extend the system to support multi-language feedback, enabling sentiment analysis for customers across different regions and languages using translation APIs or multilingual models

- **Real-Time Sentiment Monitoring**

Implement a real-time streaming pipeline using tools like Kafka or Apache Spark to analyze customer feedback as it is received (e.g., from social media or chatbots)

CHAPTER 10

REFERENCES

CHAPTER 10

REFERENCES

1. **Rahul, L., Saraswati, J. (2021)**
Customer Segmentation Using RFM Model and K-Means Clustering.
Explores the integration of the Recency, Frequency, and Monetary (RFM) model with K-Means clustering for enhanced customer segmentation strategies.
2. **Jinfeng, Z., Jinliang, C. (2021)**
Customer Segmentation by Web Content Mining.
Introduces the RFMT model (Recency, Frequency, Monetary, Time) for a more detailed understanding of customer purchasing behavior.
3. **Zeeshan-ul-Hassan, U. (2021)**
Pakistan's Largest E-Commerce Dataset.
A comprehensive dataset covering 584,524 transactions, providing key insights into consumer behavior in Pakistan's e-commerce sector.
4. **Ge, X., Jin, R. (2021)**
Optimization of Hydrogen Supply Network Design with Stochastic Demand.
Discusses data-driven optimization in supply chain design, relevant to predictive analytics in customer segmentation.
5. **Joshua, A., et al. (2021)**
Enhanced K-Means Clustering Algorithm for Improved Accuracy.
Proposes an improved centroid identification approach based on a compactness factor to enhance the accuracy of K-Means clustering.
6. **Onur, D. (2019)**
Determining the Optimal Number of Clusters using the Silhouette Approach.
A key study on using Silhouette Scores to determine the best number of clusters in customer segmentation models.
7. **P, A. (2020)**
Quantifying Electronic Industry Data Using RFM Model and K-Means.

Customer Feedback Analysis using Machine Learning

Demonstrates how RFM analysis combined with clustering can enhance segmentation for e-commerce businesses.

8. **Martin, E. (2021)**

Advancements in DBSCAN for Clustering High-Density Data.

Highlights the importance of DBSCAN clustering in identifying outlier customers and niche segments.

9. **Ching, H.C. (2021)**

Partitioning Clustering Techniques for Customer Segmentation.

Discusses various partitioning clustering methods, comparing hierarchical, K-Means, and density-based clustering techniques.

10. **Xin, S.Y. (2021)**

Optimizing Marketing Strategies using RFMT Data Analysis.

A study that validates the impact of RFMT-based customer segmentation on business revenue and engagement strategies.

11. **Jun, W. (2021)**

Enhancing Promotional Activities through Customer Segmentation.

Demonstrates how customer segmentation can be leveraged to design effective promotional campaigns based on purchasing patterns.

12. **Ioannis, M. (2020)**

Machine Learning-Based Customer Insights for Business Growth.

Explores AI-driven segmentation models to predict customer loyalty and retention rates.

13. **Ting, Z. (2020)**

Gaussian Clustering and Fuzzy-C Means for Data Segmentation.

A comparative analysis of Gaussian Mixture Models (GMM) and Fuzzy C-Means clustering for better customer behavior analysis.

14. **Saurabh, P. (2021)**

Startup Businesses and Customer Segmentation Approaches.

Customer Feedback Analysis using Machine Learning

Focuses on how startup businesses can leverage RFMT analysis for effective customer targeting.

15. Himanshu, S. (2021)

Predicting Customer Satisfaction using Sentiment Analysis.

Utilizes Natural Language Processing (NLP) models for analyzing customer feedback and sentiment polarity.

16. Abha, C.B. (2021)

A Comparative Study of Clustering Algorithms for Large-Scale Data Analysis.

Evaluates performance metrics such as Silhouette Score, CH Index, and DB Index across multiple clustering techniques.

17. Rajan, V. (2021)

Targeting Specific Customer Audiences through Machine Learning.

Demonstrates how AI-based segmentation models improve customer engagement and marketing ROI.

18. John, R.M. (2021)

Scaling Monetary and Frequency Values in RFM-Based Customer Segmentation.

Provides a structured approach for scaling and normalizing monetary and frequency values in segmentation analysis.

19. Aman, B. (2020)

Clustering for Customer Loyalty and Brand Affinity Prediction.

Explores Affinity Propagation and Hierarchical Clustering for identifying loyal customers.

20. Maha, A. (2021)

Dimensionality Reduction for Customer Data Analysis.

Discusses how Principal Component Analysis (PCA) can improve segmentation accuracy by reducing irrelevant customer attributes.