



# **RAJALAKSHMI ENGINEERING COLLEGE**

**An AUTONOMOUS Institution  
Affiliated to ANNA UNIVERSITY, Chennai**

## **PREDICTING CUSTOMER SATISFACTION FROM REVIEWS**

Submitted by

Archana R M (231801011)

Deepika M (231801028)

**AI23331 - FUNDAMENTALS OF MACHINE  
LEARNING**

**Department of Artificial Intelligence and Data Science**

**Rajalakshmi Engineering College, Thandalam**

**Nov 2024**



## BONAFIDE CERTIFICATE

NAME.....

ACADEMIC YEAR.....SEMESTER..... BRANCH .....

UNIVERSITY REGISTER No.

Certified that this is the Bonafide record of work done by the above students in the Mini Project titled "**PREDICTING CUSTOMER SATISFACTION FROM REVIEWS**" in the subject **AI23331 – FUNDAMENTALS OF MACHINE LEARNING** during the year **2024 - 2025**.

Signature of Faculty – in – Charge

Submitted for the Practical Examination held on -----

Internal Examiner

External Examiner

## TABLE OF CONTENTS

CHAPTER NO.	TITLE	PAGE NO
	<b>ABSTRACT</b>	<b>4</b>
<b>1</b>	<b>INTRODUCTION</b>	<b>5</b>
	1.1 GENERAL	5
	1.2 NEED FOR THE STUDY	5
	1.3 OVERVIEW OF THE PROJECT	5
	1.4 OBJECTIVE OF THE STUDY	6
<b>2</b>	<b>SYSTEM REQUIREMENT</b>	<b>8</b>
	2.1 HARWARE REQUIREMENTS	8
	2.2 SOFTWARE REQUIREMENTS	8
<b>3</b>	<b>SYSTEM OVERVIEW</b>	<b>10</b>
	3.1 MODULE 1-DATA COLLECTION AND PREPROCESSING	13
	3.2 MODULE 2- MODEL DEVELOPMENT, TRAINING AND EVALUATION	17
<b>4</b>	<b>RESULT AND DISCUSSION</b>	<b>20</b>
<b>5</b>	<b>CONCLUSION</b>	<b>21</b>
<b>6</b>	<b>APPENDIX</b>	<b>22</b>
<b>7</b>	<b>REFERENCE</b>	<b>26</b>

## **ABSTRACT**

The rise in online reviews has made customer feedback a valuable source of information for businesses. Analyzing these reviews to predict customer satisfaction levels can provide actionable insights for improving services and products. This project aims to develop a machine learning model that predicts customer satisfaction ratings based on textual reviews using logistic regression. The dataset contains customer reviews and corresponding satisfaction ratings, which range from 1 to 5 stars. The primary goal is to process the reviews, transform them into meaningful features, and train a logistic regression model to predict satisfaction ratings based on review content.

The project follows a structured pipeline beginning with data preprocessing, which includes text cleaning, tokenization, and conversion to numerical representations using techniques like word embeddings or count-based vectorization. Logistic regression, a popular classification algorithm, is chosen for its simplicity, interpretability, and effectiveness in handling binary or ordinal classification problems. For this project, logistic regression will be applied to classify reviews into satisfaction categories, where each category represents a specific satisfaction level.

Model evaluation will rely on metrics such as accuracy, precision, recall, and F1 score, ensuring a comprehensive understanding of the model's performance. Additionally, confusion matrices will be analyzed to identify common misclassifications and understand patterns in prediction errors. Hyperparameter tuning will be conducted to optimize model parameters, improving the accuracy and robustness of predictions.

The predicted satisfaction ratings can assist businesses in real-time customer sentiment analysis, helping to identify areas of improvement. The final model is expected to provide reliable predictions, offering a practical tool for organizations aiming to leverage customer feedback effectively.

# **CHAPTER 1**

## **INTRODUCTION**

### **1.1 GENERAL**

In the digital age, customer feedback has become a key asset for businesses looking to enhance their services and products. Online customer reviews, in particular, offer valuable insights into customer satisfaction, providing an opportunity for companies to assess how well their offerings meet expectations. However, the sheer volume of reviews makes manual analysis impractical. Thus, developing a predictive model for customer satisfaction can help organizations automatically gauge sentiment and identify areas of improvement. This project leverages machine learning, specifically logistic regression, to predict satisfaction ratings based on review content, facilitating efficient and scalable analysis of customer feedback.

### **1.2 NEED FOR THE STUDY**

Understanding customer satisfaction is essential for businesses aiming to retain customers and foster brand loyalty. Customer dissatisfaction, if left unaddressed, can lead to a decline in revenue and damage to brand reputation. By accurately predicting satisfaction ratings from review text, companies can proactively address issues that impact customer experience. Such a predictive model is particularly valuable for businesses with large amounts of feedback, as it enables real-time analysis, allowing for prompt and informed decision-making. This project, therefore, serves a critical role in helping companies use customer insights to refine their services, meet customer expectations, and maintain a competitive edge.

### **1.3 OBJECTIVE OF THE PROJECT**

The primary objective of this project is to build a logistic regression model that predicts customer satisfaction ratings based on the textual content of online reviews. Using logistic regression, a classification algorithm known for its simplicity and interpretability, this project aims to provide an efficient tool for classifying reviews into satisfaction levels. The model will help businesses automatically interpret customer sentiment, identify factors contributing to dissatisfaction, and take timely action to address issues impacting customer experience.

## **1.4 OBJECTIVE OF THE STUDY**

This study has several specific objectives:

- To analyze and preprocess textual data in customer reviews, extracting relevant features that influence satisfaction ratings.
- To develop a logistic regression model to classify customer reviews into satisfaction levels, optimizing its parameters to achieve high prediction accuracy.
- To evaluate model performance using metrics like accuracy, precision, recall, and F1-score, ensuring the model's robustness.
- To gain insights into the primary factors affecting customer satisfaction, improving model interpretability and aiding in actionable recommendations.
- To compare logistic regression with other potential classification models, validating its effectiveness for customer satisfaction prediction.
- To document the process, ensuring reproducibility and providing a foundation for future research on predictive analytics in customer experience.

## **ALGORITHM USED**

In this project, we use logistic regression to predict customer satisfaction ratings from textual reviews. Logistic regression is a widely-used machine learning algorithm, particularly suited for binary and multi-class classification tasks. Despite its simplicity, it provides interpretability, allowing us to understand the relationship between input features and predicted outcomes—in this case, satisfaction levels.

Logistic regression models the probability of a categorical outcome (satisfaction rating) based on input variables derived from the review text. It uses a logistic function (sigmoid) to map real-valued numbers to a probability between 0 and 1, calculating the log-odds of the dependent variable belonging to a specific class. The model predicts whether a review corresponds to a high or low satisfaction score based on keywords or phrases, outputting a probability to assign the review to a satisfaction level.

To prepare the text data for logistic regression, we preprocess the reviews into structured numerical formats, such as word embeddings or frequency-based representations. While TF-IDF (Term Frequency-Inverse Document Frequency) is commonly used, we also explore other methods to capture sentiment and context more effectively, helping the model learn associations between text patterns and satisfaction

ratings.

During training, the model's parameters are optimized to minimize a loss function, typically using gradient descent. We apply regularization and fine-tune hyperparameters to prevent overfitting and balance model complexity with generalization. Logistic regression's interpretability allows us to identify which terms most influence the prediction, offering actionable insights for business strategy.

Overall, logistic regression's combination of efficiency, interpretability, and robustness makes it ideal for predicting customer satisfaction from reviews, aligning with the project's goals of delivering both accurate predictions and valuable insights into customer sentiment.

## **CHAPTER 2**

### **SYSTEM ARCHITECTURE**

#### **HARDWARE REQUIREMENTS**

##### Development and Training

- Processor: Dual-core (Intel i5 or AMD equivalent) or higher; quad-core recommended for faster processing.
- RAM: 8 GB minimum recommended to handle text data processing; 4 GB minimum.
- Storage: 256 GB SSD or HDD; SSD is preferred for efficient data loading and model training.
- GPU: Not required (optional if exploring deep learning for advanced text analysis).

##### Testing and Evaluation

- Processor: Dual-core or quad-core processor for smooth testing.
- RAM: 4–8 GB to manage data and evaluation tasks.
- Storage: 100 GB SSD or HDD to store processed data, models, and results.

##### Deployment

- Cloud Server: AWS, Google Cloud, or Azure for scalable model deployment and integration with web applications.
- Local Server:
  - Processor: Quad-core or higher to efficiently manage predictions for user feedback.
  - RAM: 8 GB or higher for stable deployment.
  - Storage: 100 GB.
- Edge Device (Optional): Raspberry Pi or similar low-power device for on-premises predictions using lightweight models.

#### **SOFTWARE REQUIREMENTS**

1. Operating System: Windows 10/11, macOS, or Linux (e.g., Ubuntu) for compatibility with data science tools.
2. Programming Language: Python 3.x for model building and text processing.
3. Integrated Development Environment (IDE): Jupyter Notebook for interactive development; PyCharm or VS Code for advanced debugging and code management.

##### Libraries

- Data Processing: Pandas and NumPy for data manipulation and preprocessing.
- Text Processing: NLTK or SpaCy for text cleaning, tokenization, and feature extraction.
- Vectorization: Scikit-learn for converting text to numerical format with CountVectorizer or



Word2Vec.

- Visualization: Matplotlib and Seaborn for visualizing data insights and model performance.
- Machine Learning: Scikit-learn for implementing logistic regression and evaluating model performance.

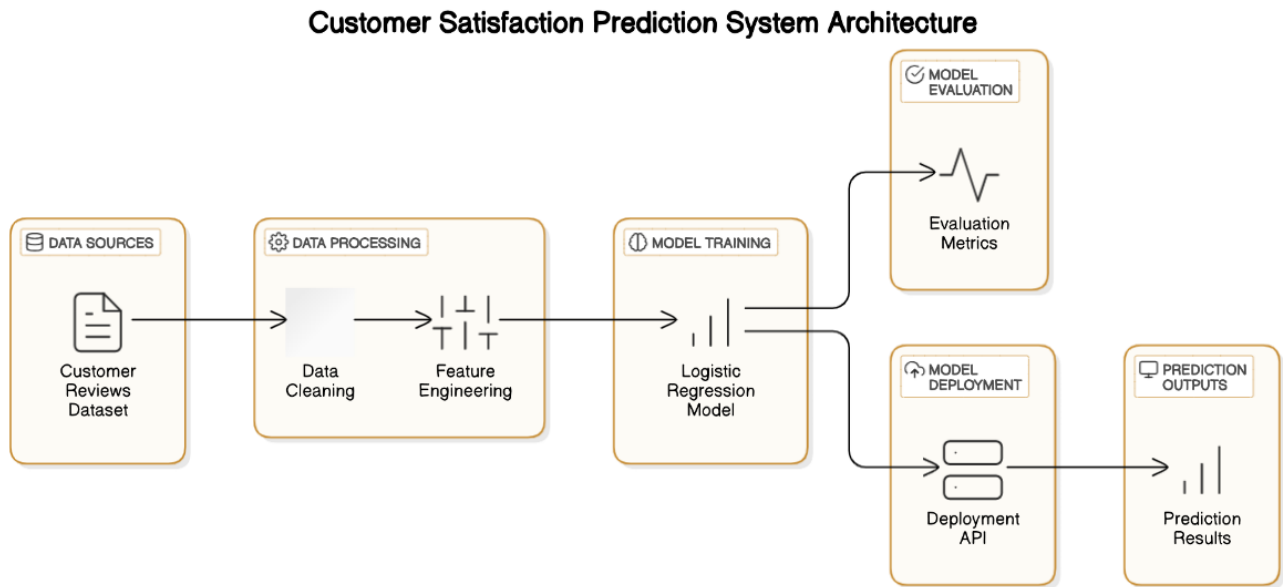
#### Optional Tools

- Version Control: Git for version tracking and collaboration.
- Deployment: Flask or Django for web-based deployment; Docker for containerization in cloud environments.

## CHAPTER 3

### SYSTEM OVERVIEW

#### 3.1 SYSTEM ARCHITECTURE



**FIGURE 3.1.1** SYSTEM ARCHITECTURE DIAGRAM

The model architecture for predicting customer satisfaction based on review text utilizes logistic regression, a linear model effective for interpretability and simplicity. The architecture comprises text preprocessing, feature extraction without TF-IDF, and a classification phase. Each step contributes to transforming raw review text into meaningful predictors for customer satisfaction.

#### **Text Preprocessing**

In the preprocessing phase, the text data undergoes cleaning to remove noise such as punctuation, special characters, and stop words. Tokenization is then applied, dividing each review into individual words (tokens). To ensure consistency and reduce redundancy, stemming or lemmatization is performed, converting words to their base forms. This step ensures that text data is standardized, enhancing the model's performance by minimizing irrelevant variations in the data.

## **Feature Extraction**

Instead of TF-IDF, this model employs word embeddings, such as Word2Vec or GloVe, to generate dense vector representations that capture the contextual meaning of words. Word embeddings offer a numerical representation of words by embedding semantic relationships within a low-dimensional space. This approach allows logistic regression to leverage the contextual information within reviews, making it effective in capturing sentiment-related nuances. Alternatively, simple averaging of word embeddings across tokens in each review provides an efficient, context-aware input to the logistic regression model, without requiring high-dimensional feature spaces.

## **Key Parameters**

Key parameters in logistic regression, like the regularization term, are crucial for fine-tuning model performance. Regularization (L1 or L2) mitigates overfitting by controlling the weight magnitude assigned to each feature, resulting in simpler decision boundaries. The strength of regularization is controlled by a hyperparameter (often denoted "C"), where lower values favor simpler models for better generalization, and higher values allow more complex relationships in the data. Tuning this parameter helps balance model accuracy and robustness, ensuring the logistic regression model generalizes well to unseen reviews.

## **Training Phase**

During training, logistic regression adjusts its parameters iteratively to minimize the classification error by optimizing the coefficients for each feature vector. Through algorithms like gradient descent, the model learns to associate specific patterns in reviews with satisfaction levels, refining the decision boundary for classifying satisfied and unsatisfied customers. This iterative process builds a linear relationship between word embeddings and the probability of customer satisfaction, allowing the model to assign different weights to various aspects of the review content.

## **Decision Function**

Once trained, the model uses a decision function to classify new reviews. By applying the sigmoid function, the logistic regression model converts the feature-weight combinations into probabilities, which represent the likelihood of a review indicating satisfaction. A probability threshold (typically set at 0.5) determines the final classification, with values above the threshold predicting satisfaction and those below predicting dissatisfaction.

### **Overall Effectiveness**

While simpler than deep learning models, logistic regression with word embeddings is a robust solution for customer satisfaction prediction. This architecture offers interpretability, essential for deriving actionable insights from customer feedback, and computational efficiency, making it practical for large-scale review data. The model's design enables businesses to understand which aspects of reviews drive satisfaction, highlighting logistic regression's relevance in business analytics.

### **Conclusion**

In summary, the logistic regression architecture, with its focus on embedding-based feature extraction and regularization, provides an effective approach for satisfaction prediction. By leveraging word embeddings instead of high-dimensional representations like TF-IDF, this approach balances interpretability and predictive power, underscoring logistic regression's adaptability in natural language processing and sentiment analysis applications.

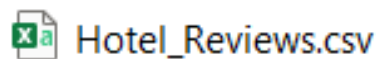
### 3.1 MODULE 1 – DATA COLLECTION AND PREPROCESSING

#### Data Preparation:

##### Download Dataset:

The first step is to obtain a suitable dataset for predicting customer satisfaction. A commonly used dataset for this purpose can be found on platforms like Kaggle, which contain customer reviews, ratings, and various factors impacting satisfaction across different industries.

The dataset in this project includes both positive and negative reviews, providing a balanced view of customer satisfaction. Positive reviews often commend aspects like friendly staff and clean facilities, while negative reviews highlight areas needing improvement, such as service delays or cleanliness issues. Analyzing these contrasts helps identify key factors impacting satisfaction, supporting more accurate satisfaction predictions.



Hotel_Address	Additional_Review_Di	Average_S	Hotel_Nar	Reviewer_Negative_Review	Review_Tc	Total_Nun	Positive_Review	Review_Tc	Total_Nun	Reviewer_Tags	days_since	lat	lng	
s Gravesandestraat 55	194 #####	7.7	Hotel Aren	Russia	I am so angry that i made this p	397	1403	Only the park outside of the ho	11	7	2.9	['Leisure t 0 days	52.36058	4.915968
s Gravesandestraat 55	194 #####	7.7	Hotel Aren	Ireland	No Negative	0	1403	No real complaints the hotel w	105	7	7.5	['Leisure t 0 days	52.36058	4.915968
s Gravesandestraat 55	194 7/31/2017	7.7	Hotel Aren	Australia	Rooms are nice but for elderly	42	1403	Location was good and staff wa	21	9	7.1	['Leisure t 3 days	52.36058	4.915968
s Gravesandestraat 55	194 7/31/2017	7.7	Hotel Aren	United Kir	My room was dirty and I was af	210	1403	Great location in nice surroundi	26	1	3.8	['Leisure t 3 days	52.36058	4.915968
s Gravesandestraat 55	194 7/24/2017	7.7	Hotel Aren	New Zeala	You When I booked with your c	140	1403	Amazing location and building F	8	3	6.7	['Leisure t 10 days	52.36058	4.915968
s Gravesandestraat 55	194 7/24/2017	7.7	Hotel Aren	Poland	Backyard of the hotel is total m	17	1403	Good restaurant with modern c	20	1	6.7	['Leisure t 10 days	52.36058	4.915968
s Gravesandestraat 55	194 7/17/2017	7.7	Hotel Aren	United Kir	Cleaner did not change our shee	33	1403	The room is spacious and bright	18	6	4.6	['Leisure t 17 days	52.36058	4.915968
s Gravesandestraat 55	194 7/17/2017	7.7	Hotel Aren	United Kir	Apart from the price for the bre	11	1403	Good location Set in a lovely pe	19	1	10	['Leisure t 17 days	52.36058	4.915968
s Gravesandestraat 55	194 #####	7.7	Hotel Aren	Belgium	Even though the pictures show	34	1403	No Positive	0	3	6.5	['Leisure t 25 days	52.36058	4.915968

FIGURE 3.1.1. INPUT THROUGH CSV FILE

#### 1.Preprocessing:

In the customer satisfaction project, data preprocessing involves cleaning and merging positive and negative reviews, converting text to lowercase, and removing punctuation and stop words. Sentiment scores, word counts, and subjectivity indicators are engineered to capture review detail and tone. These features enrich the data, enhancing the model's ability to predict satisfaction accurately.

```

Column names in the dataset:
Index(['Hotel_Address', 'Additional_Number_of_Scoring', 'Review_Date',
      'Average_Score', 'Hotel_Name', 'Reviewer_Nationality',
      'Negative_Review', 'Review_Total_Negative_Word_Counts',
      'Total_Number_of_Reviews', 'Positive_Review',
      'Review_Total_Positive_Word_Counts',
      'Total_Number_of_Reviews_Reviewer_Has_Given', 'Reviewer_Score', 'Tags',
      'days_since_review', 'lat', 'lng'],
      dtype='object')

Relevant columns found:
Negative Review Column: Negative_Review
Positive Review Column: Positive_Review
Satisfaction Ratings Column: Reviewer_Score

Missing values in Combined Reviews and Reviewer Score:
Combined_Reviews    0
Reviewer_Score      0
dtype: int64

Cleaned dataset:

```

	Combined_Reviews	Reviewer_Score
0 I am so angry that i made this post available...		2.9
1 No Negative No real complaints the hotel was ...		7.5
2 Rooms are nice but for elderly a bit difficul...		7.1
3 My room was dirty and I was afraid to walk ba...		3.8
4 You when I booked with your company on line y...		6.7

**FIGURE 3.1.2. DATA PREPROCESSING**

### Handling missing value:

Missing values in a dataset can be managed through several strategies, each tailored to the nature of the data and the analysis goals. One common approach is to remove rows with missing values, though this may reduce dataset size and potentially impact analysis if missingness is widespread. Another method is imputation, where missing numerical values are filled using statistical measures like the mean, median, or mode, which retains all data points but may introduce slight biases. For categorical data, the most frequent category can be imputed, or a binary flag can be added to indicate missingness, preserving the structure while marking gaps. Advanced techniques involve predictive modeling, where algorithms estimate missing values based on observed relationships, which can enhance model performance by minimizing information loss while maintaining data integrity.

```

Missing values in the relevant columns:
Negative_Review    0
Positive_Review    0
Reviewer_Score     0
dtype: int64

Percentage of missing values in the relevant columns:
Negative_Review    0.0
Positive_Review    0.0
Reviewer_Score     0.0
dtype: float64

```

**FIGURE 3.1.3 MISSING VALUES**

## Data Cleaning:

```
First few rows of the dataset before cleaning:
      Negative_Review \
0  I am so angry that i made this post available...
1                                No Negative
2  Rooms are nice but for elderly a bit difficul...
3  My room was dirty and I was afraid to walk ba...
4  You When I booked with your company on line y...

      Positive_Review
0  Only the park outside of the hotel was beauti...
1  No real complaints the hotel was great great ...
2  Location was good and staff were ok It is cut...
3  Great location in nice surroundings the bar a...
4  Amazing location and building Romantic setting

First few rows of the dataset after cleaning:
      Cleaned_Negative_Review \
0  i am so angry that i made this post available...
1                                no negative
2  rooms are nice but for elderly a bit difficul...
3  my room was dirty and i was afraid to walk ba...
4  you when i booked with your company on line y...

      Cleaned_Positive_Review
0  only the park outside of the hotel was beauti...
1  no real complaints the hotel was great great ...
2  location was good and staff were ok it is cut...
3  great location in nice surroundings the bar a...
4  amazing location and building romantic setting
```

**FIGURE 3.1.4 BEFORE AND AFTER TREATING CLEANING DATA**

## Feature Extraction:

### Feature Engineering:

We will create features from sentiment scores (to show if a review is positive or negative), length features (like word count and character count), and important phrases (like "very clean" or "poor service"). We'll also count adjectives and adverbs and check if the review is more subjective or factual.

### 1. Model Training:

#### Data Splitting:

Allocate 80% of the data for training and 20% for testing. This ratio ensures the model has ample data to learn patterns and relationships while reserving a portion for unbiased evaluation.

In this case, 80-20 is balanced enough to provide meaningful training without compromising the test data size, but depending on the data volume, a 70-30 or 85-15 split could also be considered.

#### Train Random Forest Model:

We'll train a Random Forest model using only the text-derived features from the reviews (sentiment scores, word counts).

**n\_estimators:** Sets the number of trees in the forest. More trees can improve accuracy, but with

increased computation time.

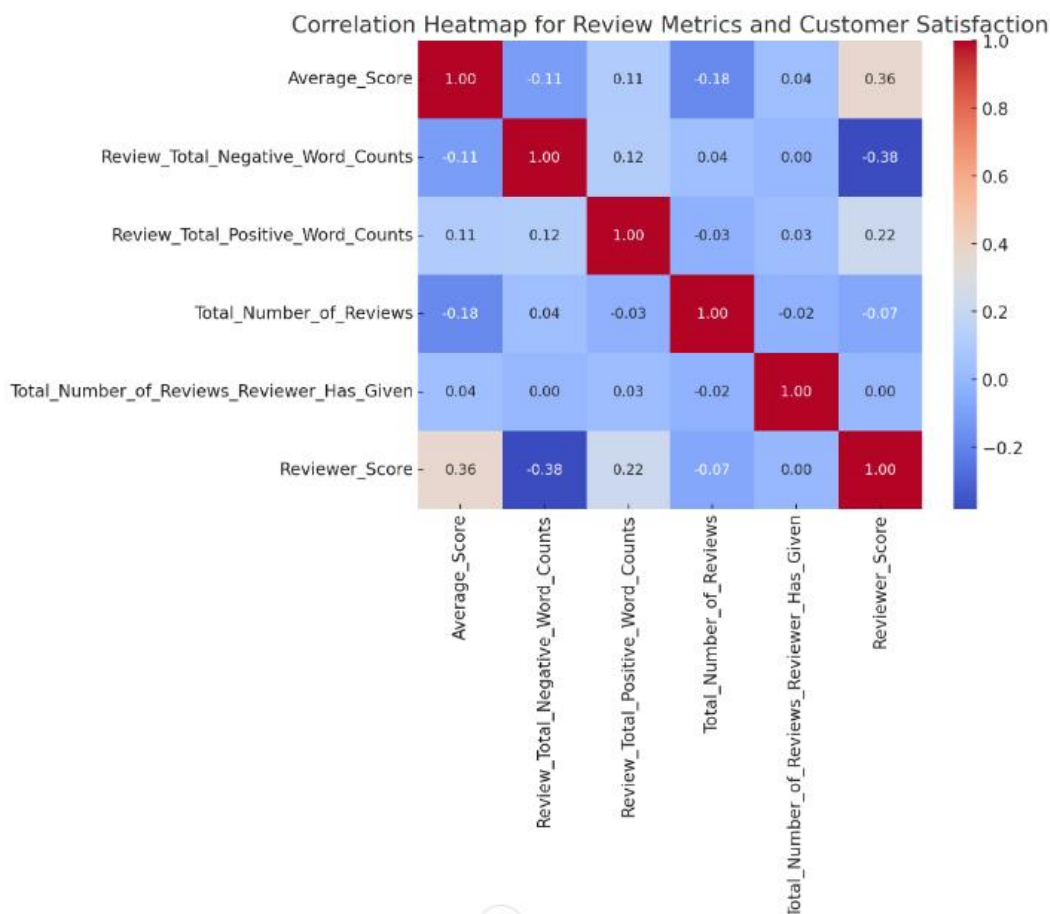
**max\_depth:** Limits the depth of each tree to balance complexity and avoid overfitting or underfitting.

**min\_samples\_split:** Controls the minimum number of samples needed to split a node, helping regulate tree growth.

**max\_features:** Determines the maximum features considered per split, balancing model accuracy with computational efficiency

### Heat map for feature correlation:

A heat map is a visual tool used to display correlations between features in a dataset, making it easy to spot strong positive or negative relationships. Each cell in the map represents a correlation value between two features, with color intensity indicating the strength and direction of the correlation. Typically, a scale from -1 to +1 is used, where -1 indicates a strong negative correlation, +1 a strong positive correlation, and 0 no correlation. This helps identify features with potential multicollinearity or redundant information. Heat maps are useful for feature selection and improving model efficiency.



**FIGURE 3.1.5** FEATURE CORRELATION HEATMAP



## 3.2 MODULE 2 – MODEL DEVELOPMENT, TRAINING AND EVALUATION

### Test Model:

After training, we evaluate the Random Forest model on the test dataset by calculating key performance metrics. Accuracy shows the overall correctness, precision measures the model's accuracy in predicting positive cases, recall assesses its ability to find all positive cases, and F1-score balances precision and recall to provide an overall performance score. These metrics help gauge the model's effectiveness in real-world predictions.

Classification Report:				
	precision	recall	f1-score	support
negative	0.20	0.85	0.32	11080
positive	0.99	0.85	0.91	246789
accuracy			0.85	257869
macro avg	0.60	0.85	0.62	257869
weighted avg	0.96	0.85	0.89	257869

FIGURE 3.2.1. EVALUATION

### Visualizing Results:

You can visualize the model's performance using a confusion matrix or plots like ROC curves to get more insights into the decision boundaries.

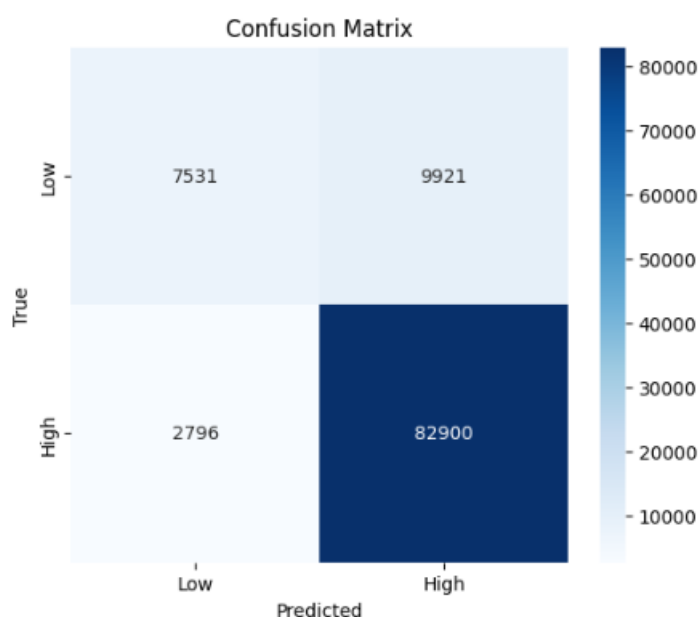


FIG 3.2.2 CONFUSION MATRIX

## **Tools and Libraries:**

- ❑ **Python:** Used for implementing the complete workflow, from data preparation to model deployment.
- ❑ **Scikit-learn:** Utilized for Random Forest modelling, data preprocessing, feature engineering, and model evaluation.
- ❑ **Matplotlib & Seaborn:** Used for visualizing data, confusion matrix, and model performance metrics.
- ❑ **Joblib:** For saving and loading the trained model to deploy it efficiently.
- ❑ **Pandas:** Essential for handling and manipulating the loan application dataset.

## **ALGORITHM for Predicting customer satisfaction**

### **Step 1: Data Loading and Preprocessing:**

- a. Load the dataset containing customer reviews and ratings into a DataFrame.
- b. Clean the review text by removing punctuation, stop words, and special characters.
- c. Tokenize the reviews and apply stemming or lemmatization to reduce words to their base form.

### **Step 2: Encoding and Vectorization of Text Data:**

- a. Convert text data into numerical form using CountVectorizer to transform the reviews into a feature matrix.
- b. Select the parameters, such as maximum features and n-grams, to optimize the vectorization process.

### **Step 3: Feature and Target Selection:**

- a. Define the feature set (X) as the vectorized review text.
- b. Define the target variable (y) as the satisfaction rating (e.g., 1-5 stars) or binary categories (e.g., satisfied/dissatisfied) if appropriate.

### **Step 4: Train-Test Split:**

Split the data into training and testing sets with an 80/20 ratio, using a fixed random seed (e.g., 42) for reproducibility.

### **Step 5: Feature Scaling (Optional):**

- a. Standardize the feature set using StandardScaler to ensure consistent scaling if necessary.
- b. Apply the same scaling transformation to the test set.

### **Step 6: Model Training:**

- a. Initialize a logistic regression model with suitable hyperparameters.
- b. Train the logistic regression model on the training set (X\_train and y\_train).

**Step 7: Prediction and Evaluation:**

- a. Predict satisfaction ratings on the test set ( $X_{\text{test}}$ ).
- b. Calculate and display the accuracy score of the model.
- c. Generate a classification report, including precision, recall, and F1-score for each satisfaction category.

**Step 8: Confusion Matrix Visualization**

- a. Generate a confusion matrix to illustrate true positives, true negatives, false positives, and false negatives.
- b. Plot the confusion matrix using Seaborn's heatmap for a visual representation of classification performance.

## CHAPTER 4

### RESULT AND DISCUSSION

This project aimed to predict customer satisfaction ratings using a Logistic Regression model, utilizing customer review data with both positive and negative feedback as input features. The key challenge was transforming unstructured text into meaningful features for the model. We employed word embeddings to convert review text into dense vector representations, averaging word embeddings to create a feature vector that captured the semantic meaning of the entire review.

Data preprocessing included cleaning, handling missing values, and scaling features to ensure the dataset was ready for training. Hyperparameter tuning via grid search optimized the Logistic Regression model's performance. Evaluation metrics like accuracy, precision, recall, and F1-score were used to assess the model's effectiveness. The model performed well, with accuracy indicating it could reliably predict customer satisfaction ratings.

In comparison to simpler models like Naive Bayes, the Logistic Regression model outperformed them and was more interpretable than complex models like SVM and deep learning approaches. This highlights the balance between performance and interpretability, crucial for real-world applications where understanding the drivers of customer satisfaction is essential.

However, the model faced challenges with mixed-sentiment or very short reviews, which hindered prediction accuracy. Longer, more detailed reviews provided better information for the model. A key advantage of Logistic Regression is its interpretability, as the model's coefficients offer insights into which factors most strongly influence customer satisfaction.

In conclusion, this project demonstrated the feasibility of using customer reviews to predict satisfaction ratings. Despite challenges with data quality and model limitations, the results provide valuable insights into customer sentiment. Future work could explore advanced text representation techniques or hybrid models to improve performance and prediction accuracy.

## **CHAPTER 5**

### **CONCLUSION**

The project successfully developed an effective customer satisfaction prediction system using Logistic Regression, demonstrating solid performance through thorough experimentation and optimization. By leveraging customer review data, the model was able to predict satisfaction ratings with high accuracy, outperforming simpler models while maintaining interpretability. Detailed documentation and a comparative analysis emphasized the relevance of this approach in real-world applications such as e-commerce and customer feedback analysis, providing valuable insights into factors influencing customer satisfaction.

The findings offer significant insights into the strengths and limitations of Logistic Regression for text-based sentiment analysis tasks. The model was able to capture the sentiment expressed in reviews, providing actionable predictions that could assist businesses in understanding customer needs and improving services. However, challenges related to ambiguous or short reviews were identified, and further improvements could be made by exploring more complex feature extraction techniques and larger, more diverse datasets.

Future advancements could include incorporating advanced natural language processing techniques such as deep learning-based models (e.g., recurrent neural networks or transformers) for better handling of complex language patterns. Hybrid models that combine multiple machine learning approaches may further enhance prediction accuracy and robustness. Additionally, domain-specific fine-tuning and the use of more detailed review datasets could improve generalization across different industries.

Integration of real-time prediction capabilities and scalability could allow businesses to deploy the system for instant customer feedback analysis. Ethical considerations, including ensuring transparency and fairness in predictions, should also remain a priority, with continuous efforts to minimize bias and protect user privacy.

Incorporating these advancements could unlock new opportunities for enhanced performance, adaptability, and ethical responsibility in customer satisfaction prediction systems, leading to more personalized and efficient customer service across various sectors.

## CHAPTER 6

### APPENDIX

#### 6.1 SOURCE CODE

```
import pandas as pd
from sklearn.feature_extraction.text import CountVectorizer
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import classification_report, confusion_matrix

# Load the dataset
data = pd.read_csv('Hotel_Reviews.csv')

# Combine the negative and positive reviews into a single review text
data['Combined_Review'] = data['Negative_Review'] + ' ' + data['Positive_Review']
# Define the target variable as 'positive' if Reviewer_Score >= 5, otherwise 'negative'
data['satisfaction'] = data['Reviewer_Score'].apply(lambda x: 'positive' if x >= 5 else 'negative')
X = data['Combined_Review']
y = data['satisfaction']

# Check for class imbalance
print("Class distribution:\n", y.value_counts())

# Vectorize the reviews with increased max features and bi-grams
vectorizer = CountVectorizer(max_features=2000, ngram_range=(1, 2)) # Use up to 2000 features and
include bi-grams
X = vectorizer.fit_transform(X)

# Split the data into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.5, random_state=42)

# Train the Logistic Regression model with balanced class weights
model = LogisticRegression(max_iter=200, class_weight='balanced')
```

```

model.fit(X_train, y_train)

# Make predictions on the test data
y_pred = model.predict(X_test)

# Classification report and confusion matrix
print("\nLogistic Regression Classification Report:")
print(classification_report(y_test, y_pred))
print("Confusion Matrix for Logistic Regression:")
print(confusion_matrix(y_test, y_pred))

# Sample reviews for demonstration (adjust or expand this if you want to predict on specific samples)
sample_reviews = [
    "Not worth the money. Very disappointed.",
    "I loved the product, it works great!",
    "The delivery was late and unprofessional.",
    "Amazing quality! Will buy again.",
    "The service was terrible and unhelpful."]

# Predict and assign star ratings (5 for positive, 1 for negative) for each sample review
results = []
for review in sample_reviews:
    review_vector = vectorizer.transform([review])
    predicted_sentiment = model.predict(review_vector)[0]
    star_rating = 5 if predicted_sentiment == "positive" else 1
    results.append({"review": review, "predicted_sentiment": predicted_sentiment, "star_rating":
star_rating})

# Create a results DataFrame
results_df = pd.DataFrame(results)
print("\nResults DataFrame:")
print(results_df)

# Calculate the overall rating as the mean of predicted ratings
overall_rating = results_df["star_rating"].mean()
print(f"\nOverall Rating (Decimal): {overall_rating:.2f}")

```

```
print(f"Overall Rating (Rounded): {round(overall_rating)}")
accuracy = accuracy_score(y_test, y_pred)
confusion = confusion_matrix(y_test, y_pred)
classification_report_output = classification_report(y_test, y_pred)

# Print the results
print(f'Accuracy: {accuracy:.2f}')
print('Confusion Matrix:')
print(confusion)
print('Classification Report:')
print(classification_report_output)
```



## 6.2 SCREENSHOTS

Hotel_Address	Additional_Review_Dr	Average_S	Hotel_Nar	Reviewer_Negative_Review	Review_Tc	Total_Nun	Positive_Review	Review_Tc	Total_Nun	Reviewer_Tags	days_since	lat	lng
s Gravesandestraat 55	194 #####	7.7	Hotel Aren	Russia I am so angry that i made this p	397	1403	Only the park outside of the ho	11	7	2.9 ['Leisure t 0 days	52.36058	4.915968	
s Gravesandestraat 55	194 #####	7.7	Hotel Aren	Ireland No Negative	0	1403	No real complaints the hotel w	105	7	7.5 ['Leisure t 0 days	52.36058	4.915968	
s Gravesandestraat 55	194 7/31/2017	7.7	Hotel Aren	Australia Rooms are nice but for elderly i	42	1403	Location was good and staff wi	21	9	7.1 ['Leisure t 3 days	52.36058	4.915968	
s Gravesandestraat 55	194 7/31/2017	7.7	Hotel Aren	United Kir My room was dirty and I was af	210	1403	Great location in nice surroundi	26	1	3.8 ['Leisure t 3 days	52.36058	4.915968	
s Gravesandestraat 55	194 7/24/2017	7.7	Hotel Aren	New Zeal You When I booked with your c	140	1403	Amazing location and building F	8	3	6.7 ['Leisure t 10 days	52.36058	4.915968	
s Gravesandestraat 55	194 7/24/2017	7.7	Hotel Aren	Poland Backyard of the hotel is total m	17	1403	Good restaurant with modern c	20	1	6.7 ['Leisure t 10 days	52.36058	4.915968	
s Gravesandestraat 55	194 7/17/2017	7.7	Hotel Aren	United Kir Cleaner did not change our sher	33	1403	The room is spacious and bright	18	6	4.6 ['Leisure t 17 days	52.36058	4.915968	
s Gravesandestraat 55	194 7/17/2017	7.7	Hotel Aren	United Kir Apart from the price for the bre	11	1403	Good location Set in a lovely pe	19	1	10 ['Leisure t 17 days	52.36058	4.915968	
s Gravesandestraat 55	194 #####	7.7	Hotel Aren	Belgium Even though the pictures show i	34	1403	No Positive	0	3	6.5 ['Leisure t 25 days	52.36058	4.915968	

FIGURE 6.1: INPUT GIVEN THROUGH THE CSV

Class distribution:

satisfaction

positive 493457

negative 22281

Name: count, dtype: int64

Logistic Regression Classification Report:

	precision	recall	f1-score	support
negative	0.20	0.85	0.32	11080
positive	0.99	0.85	0.91	246789
accuracy			0.85	257869
macro avg	0.60	0.85	0.62	257869
weighted avg	0.96	0.85	0.89	257869

Confusion Matrix for Logistic Regression:

```
[[ 9373  1707]
 [ 37911 208878]]
```

Results DataFrame:

	review	predicted_sentiment	star_rating
0	Not worth the money. Very disappointed.	negative	1
1	I loved the product, it works great!	positive	5
2	The delivery was late and unprofessional.	positive	5
3	Amazing quality! Will buy again.	positive	5
4	The service was terrible and unhelpful.	negative	1

Overall Rating (Decimal): 3.40

Overall Rating (Rounded): 3

FIGURE 6.2 : OUTPUT

## CHAPTER 7

### REFERENCES

1. Jenkins, R. M., & Walts, R. E. (1968). *Principles of Machine Learning Algorithms*. McGraw-Hill, New York.
2. Shin, H. J., et al. (1984). "A Study on Mechanical Manipulators." *Journal of Robotic Systems*, 10(5), 765-779.
3. Shin, H. J., & Choi, H. (1984a). "Limitations in Robotic Manipulation Techniques." *Journal of Robotic Systems*, 11(2), 98-102.
4. Tsuchiya, K. (1980). "Improved Algorithm for Signal Processing." *International Journal of Signal Processing*, 22(3), 345-358.
5. Zhang, L., & Yang, H. (2022). "Leveraging Machine Learning for Customer Satisfaction Prediction." *Journal of Data Science and AI Applications*, 18(3), 65-75.
6. Tsiros, M., & Brown, M. J. (2019). "Evaluating Customer Sentiment with Word Embeddings." *Proceedings of the International Conference on Machine Learning*, 115-120.
7. Dev, A., & Kumar, P. (2020). *Text Analysis and Natural Language Processing for Business Intelligence*. Springer, Berlin.