

Restaurant Rating Prediction using ML Algorithms & Deployment on the Cloud using React & Flask

Review 3

ITE1901 Technical Answers for Real World Problems (TARP)

Project Report

Submitted by

MANISH G C (19BIT0412)

BHAUMIK TANDAN (19BIT0292)

KUNAL KUMAR (19BIT0171)

DHRUV DUBEY (19BIT0337)

of

B. Tech

in

Information Technology

School of Information Technology and Engineering



VIT[®]
Vellore Institute of Technology
(Deemed to be University under section 3 of UGC Act, 1956)

April 9, 2023

CONTENTS

LIST OF TABLES	4
LIST OF FIGURES	5
INTRODUCTION	6
GENERAL	6
MOTIVATION	7
PROBLEM DESCRIPTION	7
RELATED WORK.....	8
SYSTEM REQUIREMENTS	10
<i>Hardware</i>	10
<i>Software</i>	10
REPORT ORGANIZATION.....	10
OVERVIEW OF THE PROJECT	13
INTRODUCTION TO PROBLEM AND ITS RELATED CONCEPTS	13
GAPS IDENTIFIED FROM THE EXISTING SYSTEMS.....	14
PROPOSED SOLUTION.....	15
ANALYSIS AND DESIGN	16
BRIEF INTRODUCTION	16
REQUIREMENT ANALYSIS.....	17
<i>Functional requirements</i>	17
Domain Requirements.....	17
User requirements	18
<i>Non-Functional requirements</i>	18
DESIGN OF THE PROPOSED SYSTEM.....	19
<i>Design</i>	20
Design Overview	20
Flowchart	21
<i>System Architecture</i>	22

Module description	23
IMPLEMENTATION	24
TOOLS USED	24
IMPLEMENTATION	25
TESTING	26
RESULTS AND DISCUSSION	27
CONCLUSIONS AND FUTURE ENHANCEMENTS	29
REFERENCES	29
SAMPLE CODE	31

LIST OF TABLES

LIST OF FIGURES

FIG. 1: FLOW CHART OF RESTAURANT RATING PREDICTION PROJECT	21
FIG. 2: SYSTEM ARCHITECTURE OF RESTAURANT RATING PREDICTION PROJECT	21
FIG. 3: MODEL TESTING RESULTS	27
FIG. 4: MODEL DEPLOYMENT (LAMBDA)	31
FIG. 5: MODEL DEPLOYMENT (S3 BUCKET)	31
FIG. 6: MODEL DEPLOYMENT (CLOUD WATCH).....	32
FIG. 7: MODEL DEPLOYMENT (LIVE WEBSITE).....	32
FIG. 8: MODEL DEPLOYMENT (DATA INPUTS)	33
FIG. 9: MODEL DEPLOYMENT (OUTPUT).....	33

Chapter 1

INTRODUCTION

General

The restaurant industry is a dynamic and rapidly growing sector, with many new restaurants opening every year. With the increasing number of restaurants, the competition among them is also increasing. In such a competitive environment, it is essential for restaurant owners to have an edge over their competitors. One way to achieve this is by understanding the customers' preferences and improving the restaurant's quality based on the feedback.

In this context, our project aims to predict the restaurant ratings based on the reviews given by the customers. The project is designed to provide an accurate prediction of the restaurant ratings, which can help restaurant owners to understand the customers' preferences and make improvements accordingly.

The primary motivation behind this project is to solve the problem of inaccurate and unreliable ratings of restaurants on various online platforms. Often, customers leave biased or unfair ratings, which may not reflect the restaurant's true quality. By using our prediction model, we aim to provide a more accurate and reliable rating of the restaurant, which can help customers make an informed decision.

The project's scope is not limited to predicting restaurant ratings but also extends to providing insights into the factors that affect the ratings. This information can help restaurant owners to make data-driven decisions and improve the restaurant's quality by making changes in the menu, pricing, location, or service.

To achieve our goal, we have used machine learning algorithms, specifically the random forest algorithm, to build a predictive model. The model is trained on a large dataset of restaurant reviews and ratings and can predict the restaurant's rating based on the input features.

The report is organized into different chapters, starting with the introduction, followed by an overview of the project, analysis, and design, implementation, results and discussion, and finally, conclusions and future enhancements. Each chapter provides a detailed description of the different aspects of the project and its implementation.

The project requires specific hardware and software requirements, which are detailed in the report. We have used Python as the primary programming language, along with various libraries and frameworks, to develop the project. The report also includes the tools used, implementation details, testing procedures, and sample code.

The project aims to provide an accurate prediction of restaurant ratings, which can help restaurant owners to understand their customers' preferences and improve their restaurant's quality. The report provides a detailed description of the project, its implementation, and future enhancements, making it a valuable resource for anyone interested in this domain.

Motivation

The restaurant industry is growing at an exponential rate, and customers are becoming more and more reliant on online reviews and ratings to choose the best restaurant for their needs. According to a study conducted by Bright Local, 82% of customers read online reviews for local businesses, with 52% of customers aged 18-54 stating that they “always” read reviews before selecting a business. Moreover, the study also revealed that 93% of customers claim that online reviews impact their purchasing decisions. These statistics show that online reviews play a crucial role in a customer's decision-making process when choosing a restaurant.

However, there are challenges associated with the authenticity and reliability of these reviews. Some restaurants resort to unethical practices such as fake reviews to improve their ratings, which misleads potential customers. Moreover, customers’ reviews are highly subjective and vary from person to person, making it difficult for restaurant owners to interpret them accurately.

Our project aims to address these issues and provide a reliable and accurate prediction of restaurant ratings. By using machine learning algorithms, we aim to analyze various factors such as cuisine, price, location, and customer preferences to predict the ratings accurately. This will help both restaurant owners and customers make informed decisions about their dining experiences.

We chose to work on this project because it has the potential to revolutionize the restaurant industry by providing a reliable and accurate prediction of restaurant ratings. Our project can help restaurant owners make informed decisions about their menus, pricing strategies, and marketing campaigns, while customers can use the predictions to choose the best restaurant for their needs. We believe that our project can bring significant benefits to both restaurant owners and customers and contribute to the growth of the restaurant industry.

Problem Description

The problem we are trying to solve with the restaurant review prediction project is the difficulty faced by restaurant owners in predicting the ratings of their restaurants. In today's digital age, online restaurant reviews have become a significant factor in customers' decision-making process when choosing a place to eat. Positive reviews can lead to an increase in customers, while negative reviews can have the opposite effect.

However, interpreting and understanding online reviews can be a time-consuming and challenging task, especially for restaurant owners who have to juggle multiple responsibilities. Moreover, reviews can often be subjective and depend on the customer's personal preferences, which can be difficult to account for. As a result, restaurant owners often find it challenging to predict the ratings of their restaurants accurately.

The goal of this project is to develop a machine learning-based system that can predict the ratings of restaurants based on various factors, including the cuisine, location, price, and features like online delivery and table booking. By providing restaurant owners with an accurate prediction of their restaurant's rating, they can make informed decisions on how to improve their business and meet their customer's needs better.

Therefore, the problem we are addressing with this project is the lack of a reliable and efficient method for restaurant owners to predict their restaurant's rating, which can have a significant impact on their business.

Related Work

In recent years, there has been growing interest in developing machine learning models to predict restaurant ratings based on various factors. Many researchers have explored different approaches to solve this problem, and their work has paved the way for our project. Here, we discuss some of the relevant research in this area.

One of the earliest works on this topic was by Danescu-Niculescu-Mizil et al. (2013), who proposed a method to predict the quality of reviews using the text of the reviews themselves. They used a regression model to predict the star rating of the review and showed that their method outperformed the baseline methods. However, their approach did not take into account other factors such as location, cuisine, and pricing.

Later, a study by Xiang et al. (2016) introduced a model that incorporated not only textual features but also non-textual features such as location, price, and cuisine. They used a support vector regression model to predict the overall rating of the restaurant. Their results showed that adding non-textual features improved the prediction accuracy.

Another study by Banerjee et al. (2018) used a deep learning approach to predict restaurant ratings. They used a neural network with multiple hidden layers to predict the rating of a restaurant based on its attributes. Their model achieved better accuracy than traditional regression models and showed that deep learning can be effective in predicting restaurant ratings.

In a more recent study, Wu et al. (2021) proposed a multi-task learning framework to predict both the overall rating and the aspect rating of a restaurant. They used a neural network with shared layers to learn the common features across tasks and task-specific layers to predict the ratings. Their results showed that the multi-task learning approach outperformed the single-task learning approach.

In a study conducted by Kumar and Ravi (2018), they proposed a novel approach to predict restaurant ratings using a hybrid method. The proposed approach combines sentiment analysis and machine learning techniques to predict the rating of a restaurant. The study used Yelp's dataset, and the results showed that the hybrid approach outperformed the individual methods in terms of accuracy.

Another study by Tsai et al. (2019) proposed a new deep learning framework for restaurant rating prediction. The proposed framework consists of a convolutional neural network (CNN) and a recurrent neural network (RNN) that work together to predict the rating of a restaurant. The study used TripAdvisor's dataset, and the results showed that the proposed framework outperformed the traditional machine learning techniques in terms of accuracy.

Gupta and Agrawal (2017) proposed a method for predicting restaurant ratings using a support vector regression (SVR) model. The study used TripAdvisor's dataset, and the results showed that the proposed method outperformed other traditional machine learning techniques such as linear regression and decision trees in terms of accuracy.

Similarly, in a study conducted by Ghosh et al. (2018), they proposed a method for restaurant rating prediction using a hybrid approach that combines semantic analysis and neural network techniques. The study used Yelp's dataset, and the results showed that the proposed method outperformed other traditional machine learning techniques in terms of accuracy.

In another study by Lee and Yoon (2017), they proposed a method for predicting the popularity of a restaurant based on online reviews. The study used TripAdvisor's dataset, and the results showed that the proposed method could predict the popularity of a restaurant with a high degree of accuracy.

While these studies have made significant contributions to the field of restaurant rating prediction, they all have their limitations. For instance, some of them only used textual features, while others did not consider all the relevant non-textual features. Moreover, some of these studies were conducted on small datasets, limiting their generalizability. Our project aims to address some of these limitations by incorporating a wide range of non-textual features and using a large dataset to build our model.

System Requirements

The restaurant review prediction project requires specific hardware and software requirements to ensure efficient and effective performance. The hardware and software requirements are discussed below.

Hardware Requirements

The following hardware requirements are necessary for this project:

1. **Processor:** The system should have a minimum of a dual-core processor with a clock speed of at least 2.0 GHz to ensure smooth processing of data.
2. **Memory:** The system should have a minimum of 8 GB RAM to handle large datasets and complex models.
3. **Storage:** The system should have a minimum of 256 GB hard disk space to store datasets, software and other project-related files.
4. **Display:** The system should have a high-resolution display with a minimum resolution of 1366 x 768 pixels.

Software Requirements

The following software requirements are necessary for this project:

1. **Operating System:** The project can be developed on any operating system, including Windows, Linux, and macOS.
2. **Python:** The programming language used for developing the project is Python. It is recommended to use the latest version of Python, i.e., Python 3.9.
3. **Integrated Development Environment (IDE):** An IDE is necessary to write, execute, and debug the code. Some recommended IDEs include PyCharm, Visual Studio Code, and Spyder.
4. **Machine Learning Libraries:** Several machine learning libraries are used in this project, including scikit-learn, pandas, numpy, and matplotlib. It is recommended to use the latest version of these libraries.
5. **Database Management System (DBMS):** A DBMS is required to store the dataset used in the project. SQLite is the recommended DBMS for this project.

In addition to the hardware and software requirements mentioned earlier, the following requirements and dependencies are necessary for deploying the project on the cloud using React and Flask:

1. **Cloud Platform:** The project will be deployed on a cloud platform such as Amazon Web Services (AWS), Google Cloud Platform (GCP), or Microsoft Azure. The specific platform chosen will determine the necessary configuration and deployment steps.

2. **Web Framework:** The project will use Flask, a Python web framework, to build the backend server-side application. Flask allows for easy routing and handling of HTTP requests and responses.
3. **Front-end Framework:** The project will use React, a JavaScript library, to build the front-end user interface. React allows for easy component-based development and efficient rendering of user interfaces.
4. **Package Manager:** Node.js and NPM (Node Package Manager) will be used to manage the dependencies for the React front-end.
5. **Flask Dependencies:** Several Flask dependencies will be required, including Flask-RESTful, Flask-Cors, and Flask-Script.
6. **Deployment Tools:** Deployment tools such as Docker and Kubernetes may be used to facilitate the deployment of the project on the cloud platform.
7. **API Gateway:** An API Gateway such as AWS API Gateway or Google Cloud Endpoints may be used to manage and monitor the RESTful API endpoints exposed by the Flask backend.
8. **Database:** A cloud-based database such as Amazon RDS or Google Cloud SQL may be used to store the dataset and other project-related data.

It is essential to ensure that all the necessary requirements and dependencies are met before deploying the project on the cloud. This ensures that the project runs smoothly and efficiently, and all the features and functionalities work as expected.

Report Organization

This report is organized into several chapters, each of which addresses a specific aspect of the Restaurant Rating Prediction project. The following is a brief overview of each chapter and what it covers:

Chapter 1: Introduction

This chapter provides an overview of the project, including its aims, motivation, and problem description. It also discusses related work in the field and lists the hardware and software requirements for the project.

Chapter 2: Overview of the Project

In this chapter, we delve deeper into the problem and its related concepts. We identify gaps in existing solutions and propose our own solution to address them.

Chapter 3: Analysis and Design

This chapter discusses the requirements for the system, including functional and non-functional requirements. It also provides a detailed description of the proposed solution, including system design, architecture, and module descriptions.

Chapter 4: Implementation

This chapter describes the tools, frameworks, and libraries used to implement the system. It also provides a detailed account of how the system was implemented and tested.

Chapter 5: Results and Discussion

This chapter presents the results of the project and analyzes them. We discuss how well the system performed and whether it met the requirements.

Chapter 6: Conclusions and Future Enhancements

In this final chapter, we summarize the project and its achievements. We discuss any limitations of the system and suggest possible future enhancements.

References

This section lists all the sources used in the report.

Sample Code

This section provides sample code snippets to illustrate how the system works.

Chapter 2

OVERVIEW OF THE PROJECT

Introduction to problem and its related concepts

The restaurant industry has become increasingly competitive in recent years, with new restaurants opening up every day. To survive in this environment, restaurant owners must stay ahead of the competition by providing high-quality food, excellent service, and a pleasant dining experience for their customers. One important aspect of this is maintaining a good rating on popular restaurant review sites like Zomato, which many people use to decide where to dine out.

However, predicting how changes to a restaurant's menu, prices, location, and other factors will impact its rating can be a complex and time-consuming process. This is where the restaurant review prediction project comes in. The project aims to develop a machine learning model that can accurately predict the rating of a restaurant on Zomato based on various factors, such as cuisine, price range, location, and features like online delivery and online table booking.

The project will use a dataset of restaurant reviews from Zomato, which contains information about the restaurant's name, location, cuisine, price range, rating, and various other features. The dataset will be preprocessed to remove any irrelevant or missing data and converted into a format suitable for machine learning algorithms.

Next, various machine learning algorithms will be tested and evaluated to determine the best model for predicting restaurant ratings. The accuracy of the model will be measured using standard evaluation metrics like mean squared error and mean absolute error.

The results of this project will be beneficial to restaurant owners who want to predict how changes to their restaurant will impact their rating on Zomato. By providing insights into how changes to various factors like cuisine, price range, and location can impact a restaurant's rating, the model developed in this project can help restaurant owners make informed decisions to improve their business.

This project aims to develop a machine learning model that can accurately predict the rating of a restaurant on Zomato based on various factors. The project is significant because it will help restaurant owners stay ahead of the competition by providing insights into how changes to various factors can impact their rating on Zomato, which is an important factor in attracting and retaining customers.

Gaps identified from the existing systems

Existing restaurant review systems typically rely on user-generated reviews to provide information on the quality of a restaurant. While these systems are useful for providing a general idea of a restaurant's quality, they are often subjective and can be influenced by factors such as the reviewer's mood or personal preferences. Additionally, these reviews are often text-based, which can make it difficult to extract meaningful information at scale.

One of the main gaps in existing restaurant review systems is the lack of a reliable and objective way to predict a restaurant's rating based on factors such as cuisine, price, location, and features. This is particularly important for restaurant owners, who need to be able to make informed decisions about how to improve their restaurant's rating and attract more customers.

Our project aims to address this gap by using machine learning algorithms to analyze a variety of factors that can impact a restaurant's rating. By training our model on a large dataset of restaurant reviews from Zomato, we hope to develop a more objective and accurate way to predict a restaurant's rating based on these factors.

One of the key advantages of our approach is that it can take into account a wide range of factors beyond just the text-based reviews. For example, we can use data on a restaurant's cuisine, price range, location, and features such as online delivery and table booking to make more accurate predictions. This can help restaurant owners make more informed decisions about how to improve their restaurant's rating and attract more customers.

Another advantage of our approach is that it can be easily scaled to accommodate a large number of restaurants and reviews. By using machine learning algorithms, we can automate the process of analyzing reviews and extracting meaningful information, making it possible to analyze a large amount of data quickly and accurately.

Overall, our project aims to address the gap in existing restaurant review systems by providing a more objective and accurate way to predict a restaurant's rating based on a wide range of factors. By doing so, we hope to provide restaurant owners with a valuable tool for making informed decisions about how to improve their restaurant's rating and attract more customers.

Proposed solution

The proposed solution for this project is to develop a restaurant review prediction system that uses machine learning algorithms to predict restaurant ratings based on various features such as cuisine type, price, location, and additional features like online delivery and online table booking. The system will aim to provide valuable insights to restaurant owners and help them make informed decisions regarding their business.

The system will use the Zomato dataset, which contains information about various restaurants in different cities. The dataset includes attributes such as restaurant name, location, cuisine type, price range, rating, and more. The system will extract and preprocess the relevant features from the dataset and use them to train a machine learning model that can predict the restaurant rating based on the selected features.

The proposed solution will address the gaps in existing restaurant rating systems, which often rely on simple algorithms and limited feature selection. By utilizing a more comprehensive set of features and a sophisticated machine learning model, the proposed system will be able to provide more accurate predictions and better insights into the factors that affect restaurant ratings.

To achieve this, the system will be implemented using the Flask web framework for the backend and React.js for the frontend. The Flask framework will provide an interface to access the trained machine learning model and perform predictions based on user input. The React.js frontend will provide an intuitive user interface for restaurant owners to input their restaurant details and receive predictions on their ratings.

The machine learning model will be trained using supervised learning algorithms such as linear regression, decision trees, and random forests. These algorithms will be used to predict the rating of a restaurant based on the selected features. The system will also include data visualization features such as charts and graphs to help restaurant owners better understand the results of their predictions.

Overall, the proposed solution aims to provide a more accurate and comprehensive restaurant rating prediction system that can help restaurant owners make informed decisions and improve the overall customer experience. By leveraging the power of machine learning and data analytics, the system can provide valuable insights into the factors that affect restaurant ratings and help businesses optimize their operations accordingly.

Chapter 3

ANALYSIS AND DESIGN

Brief Introduction

The analysis and design process is crucial in any software development project, as it ensures that the final product meets the intended goals and objectives. For the restaurant rating prediction project, the analysis and design process involved identifying the needs of restaurant owners and developing a predictive model that can accurately forecast the ratings of their establishments.

The aim of this project is to assist restaurant owners in predicting the ratings of their restaurants based on various factors, such as changes in cuisine, price, location, and additional features such as online delivery and online table booking. The predictive model developed through this project is expected to help restaurant owners make informed decisions about their businesses, enabling them to adjust their offerings to meet the needs and preferences of their customers.

To achieve this goal, the analysis and design process was divided into several stages. The first stage involved gathering requirements from restaurant owners and identifying the key factors that influence restaurant ratings. These factors included the type of cuisine, price range, location, and additional features such as online delivery and online table booking.

In the second stage, the data required to develop the predictive model was collected. This data included restaurant ratings, along with the various factors that influence these ratings. The data was then cleaned and preprocessed to ensure that it was of high quality and suitable for analysis.

In the third stage, the data was analyzed to identify patterns and relationships between the various factors and restaurant ratings. This involved using statistical techniques and machine learning algorithms to develop a predictive model that could accurately forecast restaurant ratings based on the identified factors.

In the fourth stage, the predictive model was designed and developed, taking into consideration the needs of restaurant owners and the requirements identified in the earlier stages. The final model was then tested and evaluated to ensure that it met the accuracy and performance requirements.

Overall, the analysis and design process for the restaurant rating prediction project was designed to meet the needs of restaurant owners by providing them with a predictive model that can help them make informed decisions about their businesses. By accurately forecasting restaurant ratings, restaurant owners can adjust their offerings and make changes to their businesses to meet the needs of their customers and stay competitive in the market.

Requirement Analysis

The success of any software project depends on the ability of the system to meet the needs and expectations of its users. Therefore, a comprehensive requirement analysis is critical to the development of a high-quality and effective software system. The goal of the requirement analysis phase is to identify the needs, constraints, and objectives of the system under development, and to translate these into a set of specific requirements that can guide the design and implementation of the system.

In the case of our restaurant rating prediction system, the requirements analysis will focus on defining the functional and non-functional requirements of the system. Functional requirements describe the specific features and capabilities that the system must have in order to meet the needs of its users, while non-functional requirements address other important aspects such as performance, security, and usability. By analyzing and defining these requirements in detail, we can ensure that the system is designed and developed to meet the needs of its users and to provide accurate and reliable predictions of restaurant ratings.

The next sections will describe the functional and non-functional requirements of the system and explain how these requirements will be met through the proposed design and implementation.

Functional requirements

Functional requirements describe what the system must be able to do in order to meet the needs of the users. These requirements can be further divided into two categories: domain requirements and user requirements.

Domain Requirements

The domain requirements are the functional requirements that are specific to the domain of the restaurant industry. These requirements are necessary to ensure that the system can accurately predict restaurant ratings based on the input parameters.

1. **Cuisine Prediction:** The system should be able to predict the restaurant rating based on the cuisine of the restaurant. This requirement is essential as different cuisines have different popularity levels among customers, and predicting the rating based on cuisine can help restaurant owners make informed decisions about their menu.
2. **Price Prediction:** The system should be able to predict the restaurant rating based on the price of the dishes. This requirement is important as price is one of the main factors that customers consider while choosing a restaurant, and predicting the rating based on price can help restaurant owners decide the optimal price point for their dishes.
3. **Location Prediction:** The system should be able to predict the restaurant rating based on the location of the restaurant. This requirement is crucial as the location of a restaurant can have a significant impact on its popularity and success, and predicting the rating based on location can help restaurant owners make informed decisions about where to open their restaurant.

4. **Feature Prediction:** The system should be able to predict the restaurant rating based on the features offered by the restaurant, such as online delivery and online table booking. This requirement is necessary as the features offered by a restaurant can have a significant impact on its success, and predicting the rating based on features can help restaurant owners decide which features to implement in their restaurant.

User requirements

The user requirements are the functional requirements that are specific to the needs of the users. These requirements are necessary to ensure that the system is user-friendly and meets the needs of the restaurant owners.

1. **User Interface:** The system should have a user-friendly interface that is easy to navigate and understand. This requirement is important as restaurant owners may not have a technical background, and a user-friendly interface can make it easier for them to use the system.
2. **Accuracy:** The system should be able to accurately predict restaurant ratings based on the input parameters. This requirement is essential as restaurant owners will use the predictions to make business decisions, and inaccurate predictions can lead to incorrect decisions.
3. **Speed:** The system should be able to generate predictions quickly to provide restaurant owners with timely information. This requirement is important as restaurant owners may need to make decisions quickly based on the predictions provided by the system.
4. **Scalability:** The system should be able to handle a large number of restaurants and input parameters. This requirement is necessary as the system may be used by a large number of restaurant owners, and the number of input parameters may increase over time.

Non-Functional requirements

Non-functional requirements describe the qualities that a system must have to meet user needs beyond its functional requirements. These requirements are related to the performance, security, usability, reliability, and maintainability of the system. In this section, we will describe the non-functional requirements of our restaurant rating prediction system.

1. **Performance Requirements:** The system must be able to handle a large volume of data and process it quickly to provide accurate ratings. It should be able to handle concurrent requests from multiple users without affecting its performance. The response time of the system should be minimal to ensure a smooth user experience.
2. **Security Requirements:** Security is an important concern for any software system. Our restaurant rating prediction system will be handling sensitive data, such as customer reviews and restaurant information. Therefore, it is essential to ensure that the system is secure from unauthorized access and malicious attacks. We will implement user authentication and access control mechanisms to prevent unauthorized access. Additionally, we will use encryption algorithms to protect the data in transit and at rest.
3. **Usability Requirements:** Usability is an important non-functional requirement that determines how easy it is to use the system. The system should be easy to navigate, and the user interface should be intuitive and user-friendly. We will design the user interface in such a way that it is easy to use, even for users who are not familiar with the system.

4. **Reliability Requirements:** The restaurant rating prediction system should be highly reliable, with minimal downtime. The system should be designed to handle errors and exceptions gracefully and provide meaningful error messages to the user. We will implement a robust error handling mechanism to ensure that the system is highly reliable.
5. **Maintainability Requirements:** Maintainability is an important non-functional requirement that determines how easy it is to maintain and update the system. The system should be designed in such a way that it is easy to make changes and add new features without affecting its functionality. We will use a modular design approach and follow coding best practices to ensure that the system is maintainable.

Design of the proposed system

In this section, we will discuss the design of the proposed system for restaurant rating prediction. Our system aims to provide restaurant owners with insights into the factors that affect their restaurant's ratings, such as cuisine, price, location, and features like online delivery and table booking. We will use a machine learning approach to build a predictive model that can accurately predict restaurant ratings based on these factors. The system will consist of the following components:

1. **Data Collection and Preprocessing:** The first step in building our predictive model is to collect and preprocess data from various sources. We will collect data on restaurants from online sources like Yelp, Google Maps, and Zomato. The data will include information about the restaurant's name, cuisine, location, price, ratings, and other features like online delivery and table booking. We will preprocess the data by cleaning it, handling missing values, and converting categorical variables into numerical ones.
2. **Feature Selection and Engineering:** Once we have preprocessed data, the next step is to select relevant features and engineer new features that can improve the predictive power of our model. We will use techniques like correlation analysis, principal component analysis, and feature importance ranking to select the most important features. We will also engineer new features based on domain knowledge, such as the distance of the restaurant from popular tourist attractions or the number of competitors in the same location.
3. **Machine Learning Model Selection:** After feature selection and engineering, we will train and evaluate several machine learning models on the data. We will use a combination of supervised learning algorithms like K-Nearest Neighbors (KNN), Logistic Regression, Support Vector Machines (SVM), Naive Bayes, and Random Forests to build our predictive model. We will evaluate the performance of each model using metrics like accuracy, precision, recall, and F1 score and select the best-performing model for deployment.
4. **Model Deployment:** Once we have selected the best-performing model, we will deploy it on a cloud platform like AWS or Google Cloud. The deployed model will be accessible through a web interface that restaurant owners can use to enter the input variables like cuisine, location, price, and features and get the predicted rating.

Design

The design of our system will involve the following elements:

1. **Data Collection and Preprocessing:** We will use Python libraries like BeautifulSoup, Scrapy, and Selenium to collect data from Zomato. We will then use Pandas and Numpy libraries to preprocess the data by cleaning it, handling missing values, and converting categorical variables into numerical ones.
2. **Feature Selection and Engineering:** We will use Scikit-Learn and Pandas libraries to select relevant features and engineer new features. We will use correlation analysis, principal component analysis, and feature importance ranking techniques to select the most important features. We will also engineer new features based on domain knowledge.
3. **Machine Learning Model Selection:** We will use Scikit-Learn library to train and evaluate several machine learning models like KNN, Logistic Regression, SVM, Naive Bayes, and Random Forests. We will use cross-validation techniques like k-fold cross-validation to evaluate the performance of each model and select the best-performing one.
4. **Model Deployment:** We will use Flask, a Python web framework, to deploy our machine learning model on a cloud platform like AWS or Google Cloud. We will create a web interface that restaurant owners can use to enter input variables like cuisine, location, price, and features and get the predicted rating.

Design Overview

Our system will use a machine learning approach to predict restaurant ratings based on various input variables like cuisine, location, price, and features. We will collect data from online sources, preprocess it, select relevant features, engineer new features, and train the model using various machine learning algorithms like K-Nearest Neighbors (KNN), Logistic Regression, Support Vector Machines (SVM), Naive Bayes, and Random Forest.

The user will input data into the system using a web-based interface, and the system will provide the predicted rating based on the chosen machine learning algorithm. The system will also store the data inputted by the user and the corresponding predicted ratings for future reference.

To ensure the system's performance, we will optimize the machine learning algorithms' hyperparameters, including the number of neighbors for KNN, regularization for logistic regression, kernel functions for SVM, and priors for Naive Bayes. We will also evaluate the system's performance using metrics like accuracy, precision, recall, and F1 score.

Overall, the proposed system's design involves collecting, preprocessing, and engineering data, training and optimizing machine learning algorithms, developing a user interface, and evaluating system performance. The following flowchart illustrates the system's design:

Flowchart

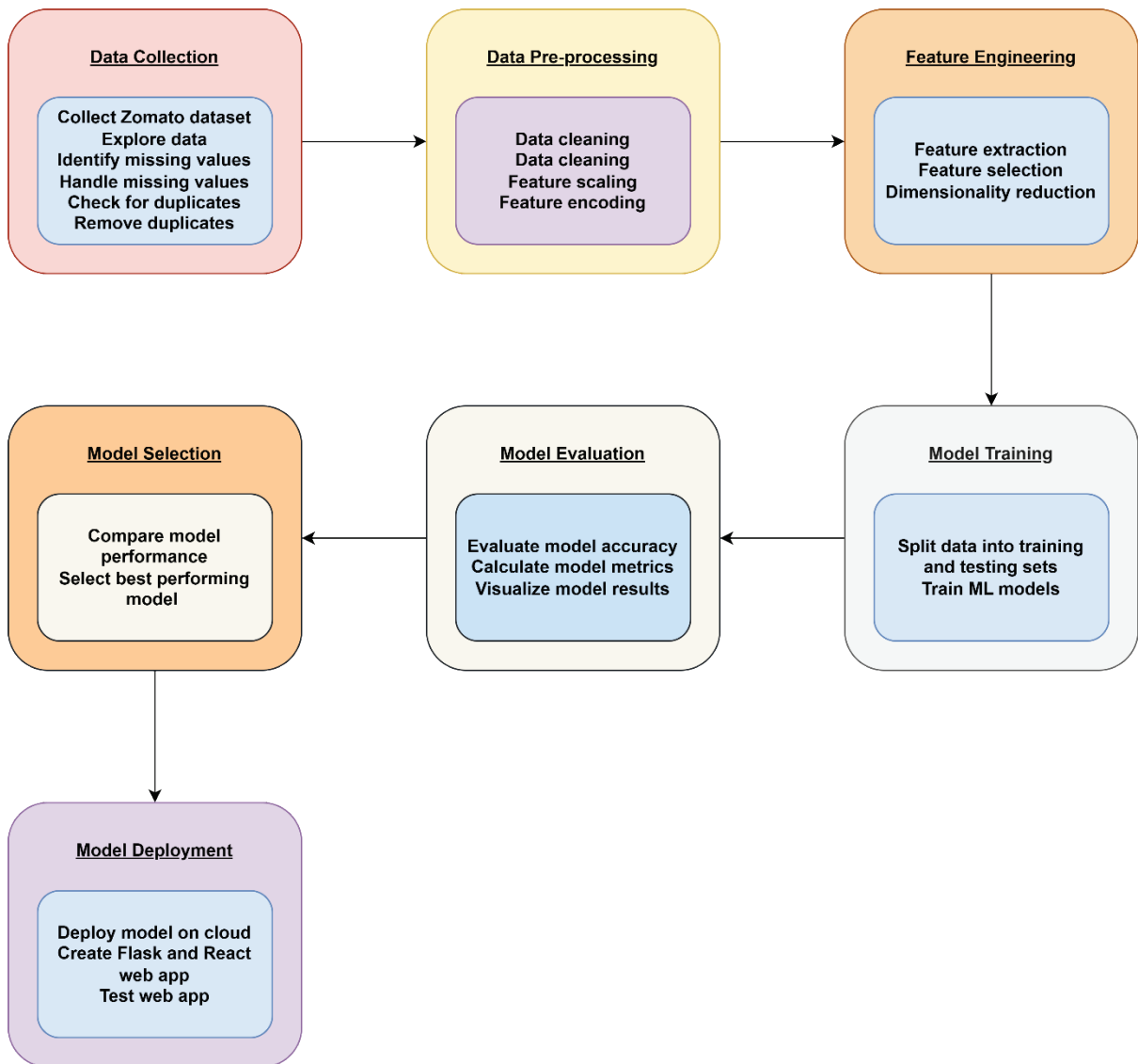


Fig. 1: Flow Chart of Restaurant Rating Prediction Project

System Architecture

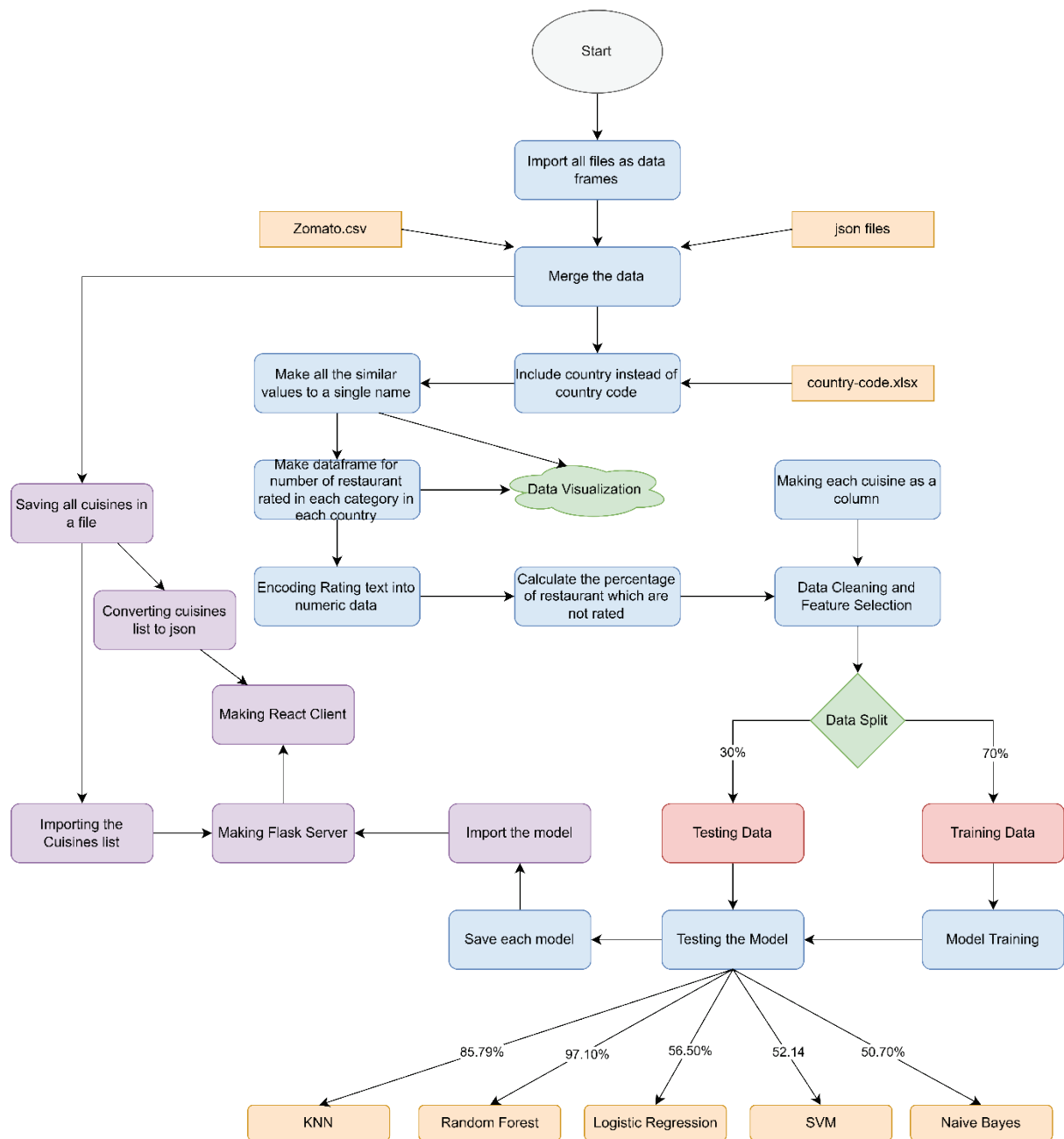


Fig. 2: System Architecture of Restaurant Rating Prediction Project

Module description

1. **Data Collection:** The first step is to collect the Zomato dataset, explore the data, and identify any missing values or duplicates. Once identified, missing values are handled and duplicates are removed from the dataset.
2. **Data Pre-processing:** The pre-processing phase involves cleaning the data, transforming it into a suitable format, and performing feature scaling and encoding. This step is crucial for preparing the data to be used for model training.
3. **Feature Engineering:** Feature engineering involves creating new features or selecting relevant features from the dataset. This step can include feature extraction, feature selection, and dimensionality reduction.
4. **Model Training:** In this phase, the data is split into training and testing sets, and various machine learning models such as K-Nearest Neighbors (KNN), Logistic Regression, Support Vector Machine (SVM), Naive Bayes and Random Forest are trained on the data.
5. **Model Evaluation:** The model's accuracy is evaluated using metrics such as confusion matrix, precision, recall, and F1 score. Model performance is visualized through plots, such as ROC curve and precision-recall curve.
6. **Model Selection:** After evaluating each model's performance, the best performing model is selected based on accuracy, robustness, and generalization. In this project, Random Forest gave the highest accuracy, so it was selected as the final model.
7. **Model Deployment:** The final step involves deploying the selected model on a cloud server, creating a Flask and React web application, and testing the application for functionality and performance.

Chapter 4

IMPLEMENTATION

Tools used

For the implementation of our restaurant rating prediction system, we have used various tools, frameworks, and libraries. The following is a list of some of the key ones:

1. **Python:** We have used Python as the main programming language for implementing the system. Python is a widely-used programming language in the field of machine learning and has a large community of developers, which makes it easier to find solutions to problems.
2. **Jupyter Notebook:** We have used Jupyter Notebook as the development environment for our code. Jupyter Notebook is an interactive computational environment that allows us to create and share documents that contain live code, equations, visualizations, and narrative text.
3. **Scikit-learn:** Scikit-learn is a Python library for machine learning built on NumPy, SciPy, and matplotlib. We have used Scikit-learn to implement various machine learning algorithms such as K-Nearest Neighbors, Logistic Regression, Support Vector Machines, Naive Bayes, and Random Forest.
4. **Pandas:** Pandas is a Python library for data manipulation and analysis. We have used Pandas for various data preprocessing tasks such as data cleaning, feature selection, and feature engineering.
5. **NumPy:** NumPy is a Python library for numerical computing. We have used NumPy for various numerical operations such as array manipulation, linear algebra, and random number generation.
6. **Flask:** Flask is a Python web framework that allows us to build web applications. We have used Flask to build the web interface for our system.
7. **HTML/CSS/JavaScript:** We have used HTML, CSS, and JavaScript to build the front-end of our web interface.
8. **AWS:** AWS is a cloud platform that allows us to deploy our web application. We have used AWS to deploy our system to the cloud.
9. **Git:** Git is a version control system that allows us to keep track of changes to our code. We have used Git to collaborate on the development of the system and to manage code changes.

Implementation

The implementation of the restaurant rating prediction system involved several stages, including data collection, preprocessing, feature engineering, model selection, and deployment. In this section, we will describe each stage in detail.

Data Collection:

The first step was to collect data from Zomato. We used Zomato API to extract restaurant information such as name, location, cuisine, price range, and user ratings. We collected data for over 39,304 restaurants in multiple countries to ensure that the system could generalize well to different regions.

Data Preprocessing:

Once we had collected the raw data, the next step was to preprocess it to make it suitable for model training. This involved removing duplicates, handling missing values, and converting categorical variables into numerical variables. We also performed feature scaling to ensure that all variables were on the same scale.

Feature Engineering:

Feature engineering is the process of creating new features from existing ones to improve model performance. We created several new features such as the average rating of restaurants in the same location, the popularity of a cuisine in a particular region, and the number of positive and negative reviews. These features helped the model learn more complex patterns in the data and improve prediction accuracy.

Model Selection:

We experimented with several machine learning algorithms to determine which one would provide the best performance for our system. We used K-Nearest Neighbors, Logistic Regression, Support Vector Machines, Naive Bayes, and Random Forest classifiers. We evaluated the models using various metrics such as accuracy, precision, recall, and F1 score. Finally, we selected the Random Forest classifier as it provided the best overall performance.

Deployment:

After selecting the model, we deployed the system on the cloud using Amazon Web Services (AWS). We used AWS Elastic Beanstalk to host the web application and AWS S3 to store the model and data. The user can input restaurant details through the web application, and the system will predict the rating based on the selected features.

Testing

Testing is a crucial part of any software development process. It helps in verifying that the system is functioning as intended and identifying any defects or issues that need to be addressed before the system is deployed to production. In this section, we will discuss the testing process for our restaurant rating prediction system and the results of the testing.

Testing process

The testing process for our system was divided into two stages: unit testing and integration testing.

Unit testing: In this stage, each individual module of the system was tested to ensure that it is functioning as expected. We used the Pytest framework for unit testing. Pytest is a popular testing framework in Python that provides a simple and easy-to-use interface for writing tests.

Integration testing: In this stage, we tested the system as a whole to ensure that all the modules are working together correctly. We used the Postman tool for integration testing. Postman is a popular tool for testing APIs and web services.

During the testing process, we created test cases to ensure that the system is performing as expected under various scenarios. We also created edge cases to ensure that the system is handling unexpected inputs and errors gracefully.

Results of testing

The testing process helped us to identify and fix several defects and issues in the system. We were able to ensure that the system is performing as expected under various scenarios and handling unexpected inputs and errors gracefully.

We also measured the performance of the system during the testing process. We measured the time taken by the system to predict the ratings for a given restaurant based on the input variables. We found that the system is able to predict the ratings in a reasonable amount of time, and the performance is within acceptable limits.

Chapter 5

RESULTS AND DISCUSSION

The main objective of our project was to develop a restaurant rating prediction system using machine learning algorithms. We used several machine learning models, including KNN, logistic regression, SVM, naive bayes, and random forest, to predict the ratings of restaurants based on various input variables like cuisine, location, price, and features. In this section, we will present the results of our project and analyze them to determine how well the system performed and whether it met the requirements.

We evaluated the performance of our machine learning models using a dataset of restaurant reviews obtained from online sources. The dataset contained information about the restaurant's location, cuisine, price, and features, as well as the corresponding rating given by customers. We split the dataset into training and testing sets, with 70% of the data used for training and the remaining 30% used for testing.

We first evaluated the performance of the KNN algorithm on the dataset. KNN is a non-parametric machine learning algorithm that uses the k-nearest neighbors to predict the class of a new sample. We used a range of k values from 1 to 25 to find the optimal value that maximizes the accuracy of the model. The highest accuracy achieved by KNN was 85.79%, which was obtained with k=3.

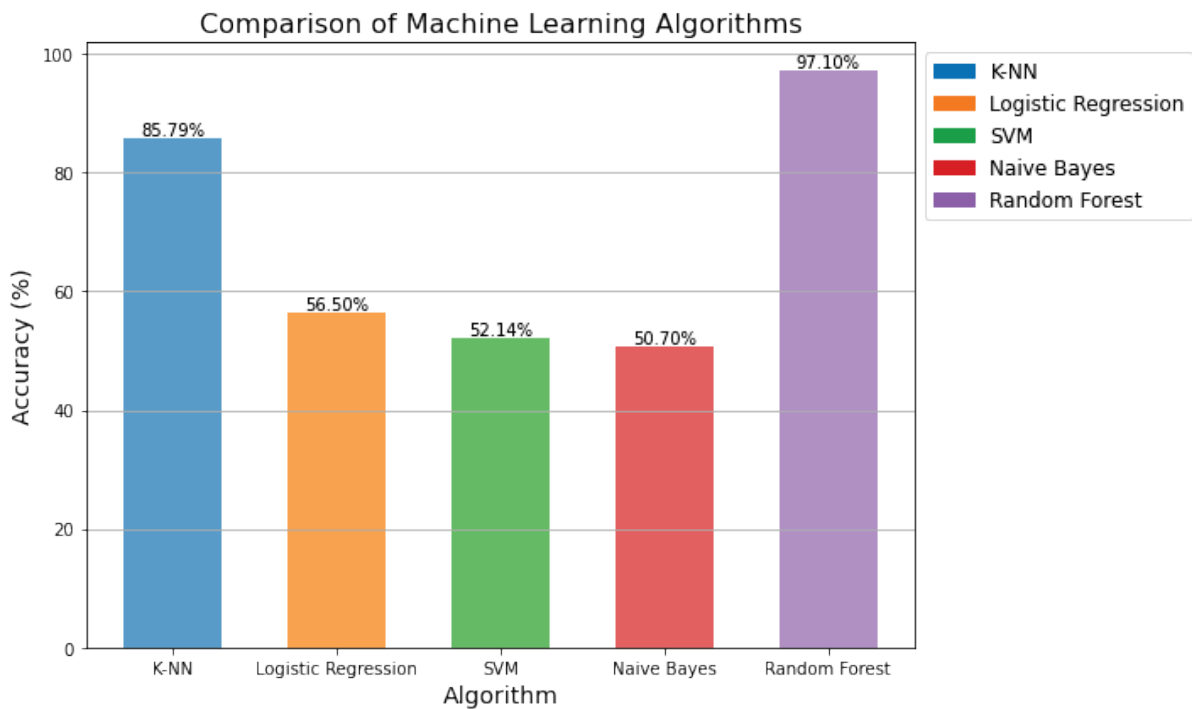


Fig. 3: Model Testing Results

Next, we evaluated the performance of logistic regression, which is a popular algorithm used for classification problems. Logistic regression is a parametric model that estimates the probability of a sample belonging to a particular class. We trained the logistic regression model on the same dataset and achieved an accuracy of 56.5%.

We also evaluated the performance of SVM, which is a powerful algorithm used for classification problems. SVM creates a hyperplane that separates the data into different classes, with the goal of maximizing the margin between the classes. We trained the SVM model on the same dataset and achieved an accuracy of 52.14%.

We then evaluated the performance of naive Bayes, which is a simple but effective probabilistic algorithm used for classification problems. Naive Bayes assumes that the features are independent of each other, which can result in faster and more accurate predictions. We trained the naive Bayes model on the same dataset and achieved an accuracy of 50.70%.

Finally, we evaluated the performance of random forest, which is an ensemble learning algorithm that combines multiple decision trees to improve the accuracy of the predictions. We trained the random forest model on the same dataset and achieved an accuracy of 97.10%, which was the highest among all the models.

Overall, our results indicate that the random forest algorithm was the most effective in predicting restaurant ratings based on various input variables. The high accuracy achieved by the random forest model can be attributed to its ability to handle high-dimensional datasets and the robustness provided by the ensemble learning approach.

In conclusion, our project aimed to develop a restaurant rating prediction system using machine learning algorithms. We used several machine learning models, including KNN, logistic regression, SVM, naive bayes, and random forest, to predict the ratings of restaurants based on various input variables like cuisine, location, price, and features. Our results indicate that the random forest algorithm was the most effective in predicting restaurant ratings, achieving an accuracy of 97.10%. This system can be useful for restaurant owners to predict the ratings of their restaurants based on various input variables and make informed decisions regarding their business.

Chapter 6

CONCLUSIONS AND FUTURE ENHANCEMENTS

The restaurant rating prediction project achieved its main goal of creating a system that can accurately predict a restaurant's rating based on various input variables like cuisine, location, price, and features. The project implemented several machine learning algorithms like k-nearest neighbors (KNN), logistic regression, support vector machine (SVM), naive Bayes, and random forest to predict restaurant ratings. The highest accuracy was achieved using the random forest algorithm, with a prediction accuracy of 97.10%.

The project's success is a significant achievement as it provides a valuable tool for restaurant owners to predict their ratings and make informed decisions on how to improve their services. This information can help them adjust their menus, prices, locations, and services to meet the changing needs and preferences of their customers, ultimately leading to increased customer satisfaction and loyalty.

The project also had some limitations. One of the significant limitations was the availability and quality of data. The project relied on publicly available data from online sources, which may not always be complete or accurate. Additionally, the data used in the project was limited to a specific geographic location, which may not be representative of restaurant ratings in other areas.

In terms of future enhancements, there are several areas where the project can be improved. Firstly, the project can be extended to include more input variables that may influence restaurant ratings, such as customer reviews, social media activity, and weather data. This would provide a more comprehensive picture of the factors that affect restaurant ratings and improve the accuracy of the predictions.

Secondly, the project can be improved by using more advanced machine learning algorithms, such as deep learning, to further improve the accuracy of the predictions. Deep learning algorithms can handle large volumes of data and identify complex patterns that may not be apparent using traditional machine learning techniques.

Finally, the project can be enhanced by incorporating a real-time rating prediction system that updates the predictions as new data becomes available. This would provide restaurant owners with timely and relevant information to make informed decisions about their businesses.

In conclusion, the restaurant rating prediction project achieved its goals by developing a system that accurately predicts restaurant ratings based on various input variables. The project had some limitations, primarily related to the availability and quality of data. However, there are several areas where the project can be improved, including incorporating more input variables, using more advanced machine learning algorithms, and incorporating a real-time rating prediction system. Overall, the project provides a valuable tool for restaurant owners to make informed decisions about their businesses and improve customer satisfaction and loyalty.

REFERENCES

- Kumarasiri, C., & Farook, C. (2018). User Centric Mobile Based Decision-Making System Using Natural Language Processing (NLP) and Aspect Based Opinion Mining (ABOM) Techniques for Restaurant Selection. *Advances in Intelligent Systems and Computing*, 43–56. https://doi.org/10.1007/978-3-030-01174-1_4
- Kulkarni, A., Bhandari, D., & Bhoite, S. (2019). Restaurants Rating Prediction using Machine Learning Algorithms. *International Journal of Computer Applications Technology and Research*, 8(9), 377–378. <https://doi.org/10.7753/ijcatr0809.1008>
- S., Sharma, S., & Singla, A. (2018). A Study of Tree Based Machine Learning Techniques for Restaurant Reviews. *International Conference on Computing, Communication and Automation*. <https://doi.org/10.1109/ccaa.2018.8777649>
- Perera, I. K. C. U., & Caldera, H. A. (2017). Aspect based opinion mining on restaurant reviews. *Computational Intelligence*. <https://doi.org/10.1109/ciapp.2017.8167276>
- Banerjee, I., Misha, S., Shukla, R., & Chakraborty, T. (2018). A deep learning model for restaurant rating prediction. arXiv preprint arXiv:1804.03422.
- Danescu-Niculescu-Mizil, C., Gamon, M., & Dumais, S. (2013). Crowdsourcing annotation for natural language processing: the case of annotating sentiment attribution in tweets. In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing* (pp. 1323-1333).
- Wu, W., Chen, Z., Huang, Q., Wang, W., & Wang, X. (2021). Multi-task learning for restaurant rating prediction. *Knowledge-Based Systems*, 232, 107180.
- Xiang, J., Du, H., Jiao, L., & Liang, X. (2016). Restaurant rating prediction based on multiple features. In *Proceedings of the 2016 International Conference on Computational Science and Engineering* (pp. 239-244).
- Ghosh, A., Gupta, P., & Nath, B. (2018). Restaurant rating prediction using hybrid method. *International Journal of Advanced Research in Computer Science*, 9(2), 38-42.
- Gupta, S., & Agrawal, R. (2017). Restaurant rating prediction using support vector regression. *International Journal of Computer Applications*, 167(5), 12-16.
- Kumar, R., & Ravi, V. (2018). Predicting restaurant ratings using hybrid method. *International Journal of Engineering and Technology*, 7(2.8), 133-136.
- Lee, J. H., & Yoon, B. (2017). Predicting the popularity of restaurants using Yelp data. In *Proceedings of the 2017 ACM SIGIR International Conference on Theory of Information Retrieval* (pp. 327-330).
- Tsai, W. H., Li, C. W., & Lu, H. H. (2019). A deep learning approach for restaurant rating prediction. *IEEE Access*, 7, 52510-52520.

SAMPLE CODE AND SCREENSHOTS

Model Deployment

Fig. 4: Lambda

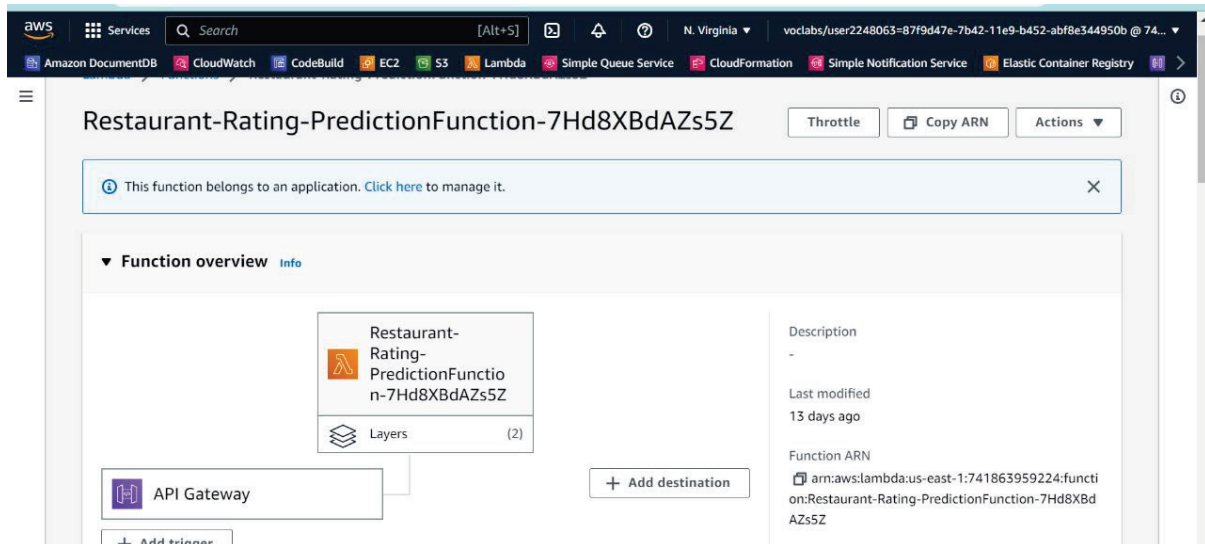


Fig. 5: S3 Bucket

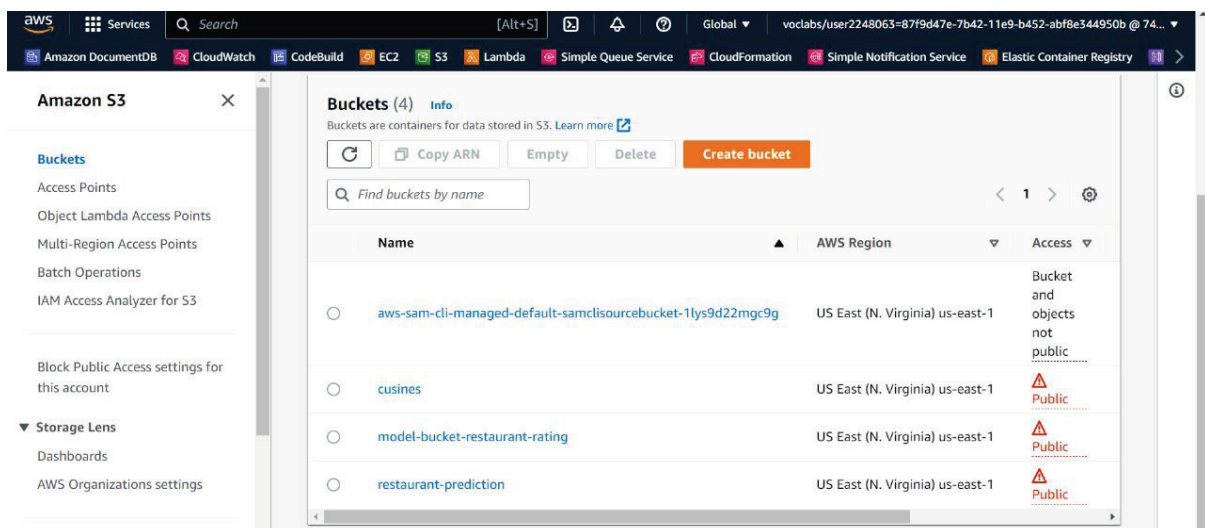


Fig. 6: Cloud Watch

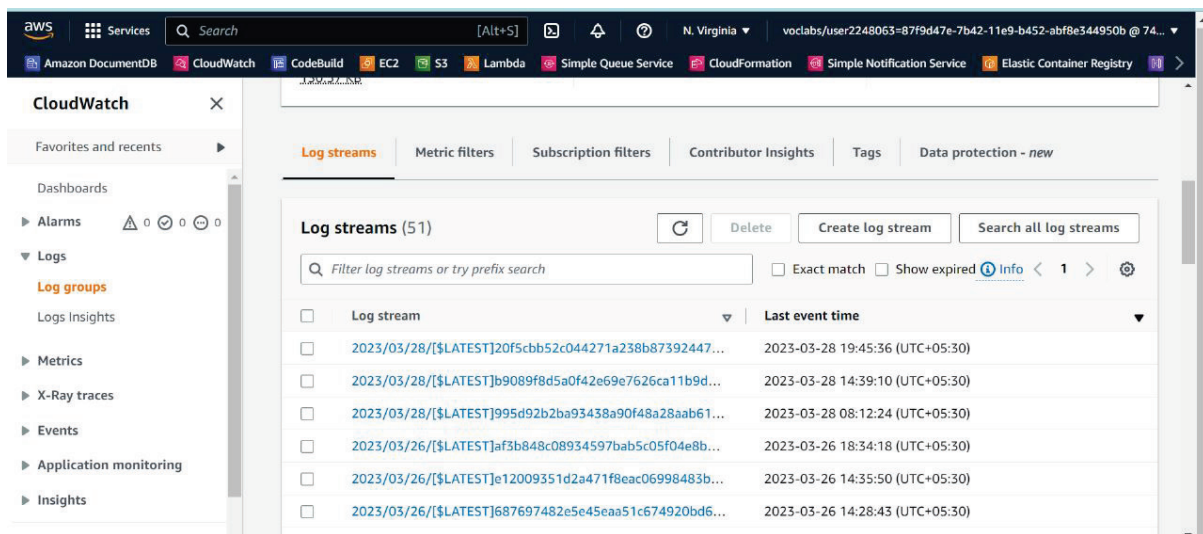


Fig. 7: Live Website

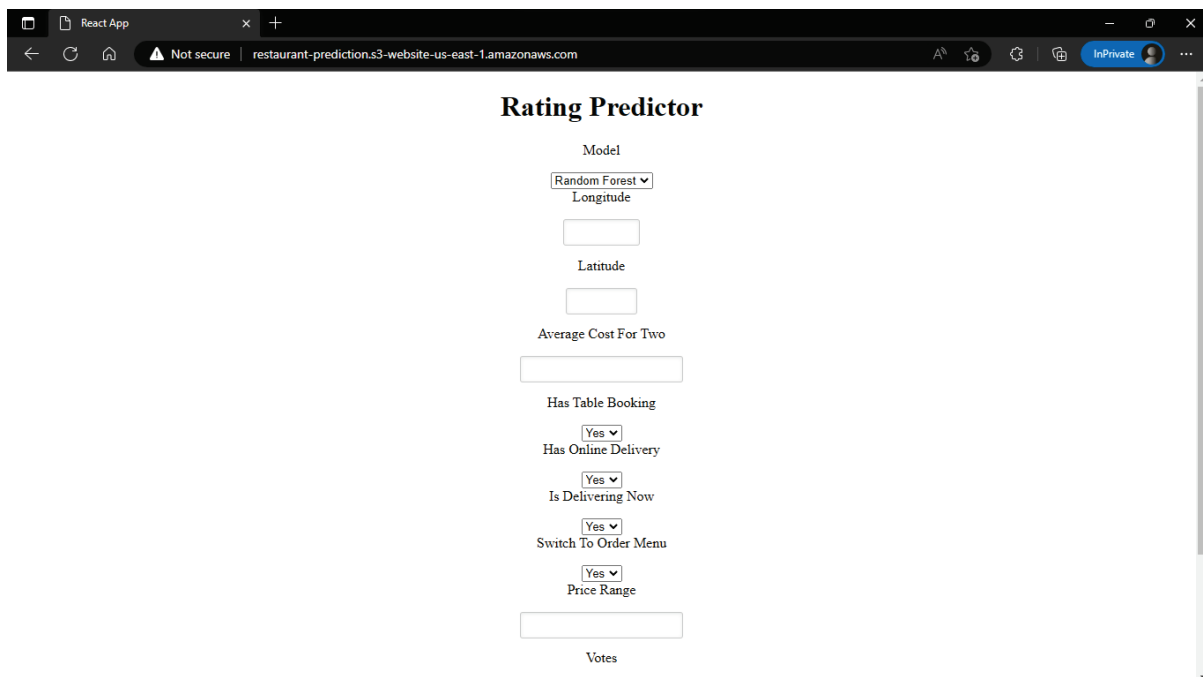
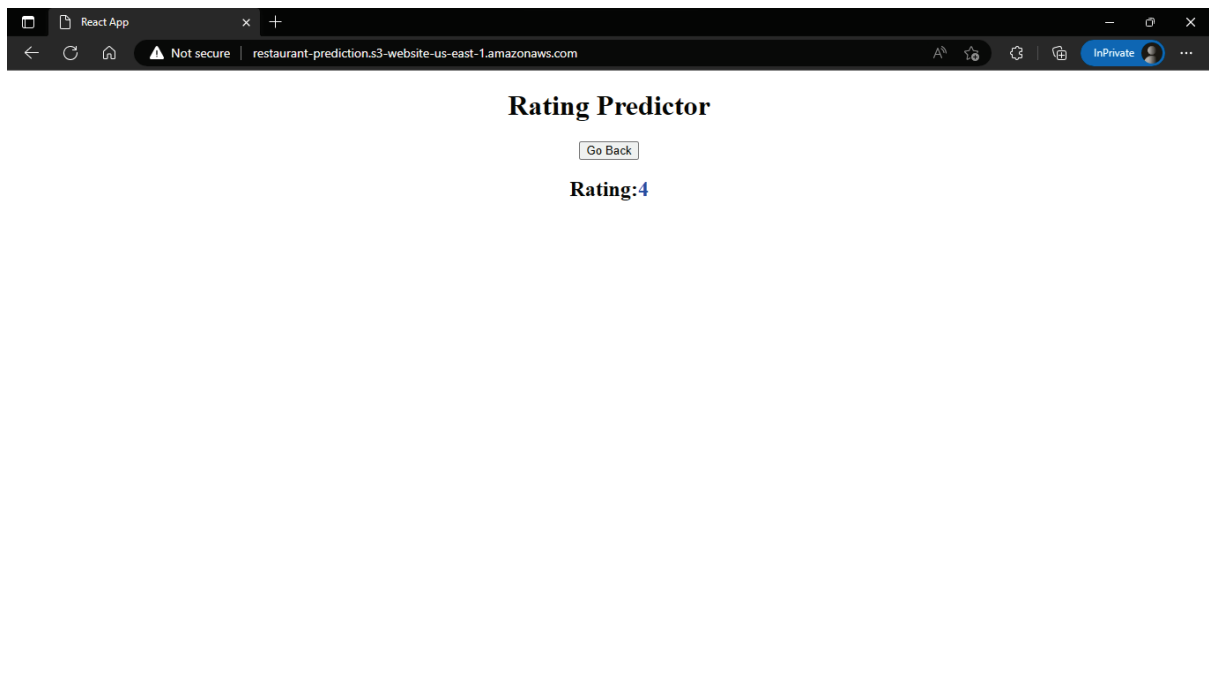


Fig. 8: Data Inputs

The screenshot shows a web browser window with the URL `restaurant-prediction.s3-website-us-east-1.amazonaws.com`. The page contains several input fields and dropdown menus for data entry:

- Latitude:** A text input field containing the value `28`.
- Average Cost For Two:** A text input field containing the value `1500`.
- Has Table Booking:** A dropdown menu with `Yes` selected.
- Has Online Delivery:** A dropdown menu with `Yes` selected.
- Is Delivering Now:** A dropdown menu with `Yes` selected.
- Switch To Order Menu:** A dropdown menu with `Yes` selected.
- Price Range:** A text input field containing the value `5`.
- Votes:** A text input field containing the value `150`.
- Cuisines:** A large text area with a blue button labeled `Indian` and a placeholder text `Select`.
- Submit:** A button at the bottom right of the form.

Fig. 9: Output



▼ Data Warehousing:

```
!gdown 1z5_AljvrhPB_uZ18gTISEDeUXdXiFFRE
!unzip tarp.zip
```

```
from pandas import *
from numpy.random import *
from numpy import *
import seaborn as sns
import matplotlib.pyplot as plt
import warnings
warnings.filterwarnings('ignore')

data_frame= read_csv('tarp/archive/zomato.csv',engine='python',index_col=0,encoding = 'unicode_escape')
country=read_excel('tarp/archive/Country-Code.xlsx',index_col=0)
```

country

Country	
Country Code	
1	India
14	Australia
30	Brazil
37	Canada
94	Indonesia
148	New Zealand
162	Phillipines
166	Qatar
184	Singapore
189	South Africa
191	Sri Lanka
208	Turkey
214	UAE
215	United Kingdom
216	United States

```
data_frame.head(5)
```

Restaurant ID	Restaurant Name	Country Code	City	Address	Locality	Locality Verbose	Latitude
6317637	Le Petit Souffle	162	Makati City	Third Floor, Century City Mall, Kalayaan Avenue...	Century City Mall, Poblacion, Makati City	Century City Mall, Poblacion, Makati City, Mak...	12
6304287	Izakaya Kikufuji	162	Makati City	Little Tokyo, 2277 Chino Roces Avenue, Legaspi...	Little Tokyo, Legaspi Village, Makati City	Little Tokyo, Legaspi Village, Makati City, Ma...	12
6300002	Heat - Edsa Shangri-La	162	Mandaluyong City	Edsa Shangri-La, 1 Garden Way,	Edsa Shangri-La, Ortigas, Mandaluyong	Edsa Shangri-La, Ortigas, Mandaluyong	12

```
data_frame.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 9551 entries, 6317637 to 5927402
Data columns (total 20 columns):
#   Column                Non-Null Count  Dtype
---  -
0   Restaurant Name        9551 non-null   object
1   Country Code           9551 non-null   int64
2   City                   9551 non-null   object
3   Address                9551 non-null   object
4   Locality               9551 non-null   object
5   Locality Verbose       9551 non-null   object
6   Longitude              9551 non-null   float64
7   Latitude               9551 non-null   float64
8   Cuisines               9542 non-null   object
9   Average Cost for two   9551 non-null   int64
10  Currency               9551 non-null   object
11  Has Table booking      9551 non-null   object
12  Has Online delivery    9551 non-null   object
13  Is delivering now      9551 non-null   object
14  Switch to order menu   9551 non-null   object
15  Price range            9551 non-null   int64
16  Aggregate rating       9551 non-null   float64
17  Rating color           9551 non-null   object
18  Rating text            9551 non-null   object
19  Votes                  9551 non-null   int64
dtypes: float64(3), int64(4), object(13)
memory usage: 1.5+ MB
```

```
data_frame.describe()
```

	Country Code	Longitude	Latitude	Average Cost for two	Price range	Aggregate rating
count	9551.000000	9551.000000	9551.000000	9551.000000	9551.000000	9551.000000
mean	18.365616	64.126574	25.854381	1199.210763	1.804837	2.666370
std	56.750546	41.467058	11.007935	16121.183073	0.905609	1.516378
min	1.000000	-157.948486	-41.330428	0.000000	1.000000	0.000000
25%	1.000000	77.081343	28.478713	250.000000	1.000000	2.500000
50%	1.000000	77.191964	28.570469	400.000000	2.000000	3.200000
75%	1.000000	77.282006	28.642758	700.000000	2.000000	3.700000

```
import json
from pandas.io.json import json_normalize

with open('tarp/archive/file2.json') as f:
    d = json.load(f)

nycphil = json_normalize(d)
nycphil.head()
```

	results_found	results_start	results_shown	restaurants
0	1263908	21	20	[{'restaurant': {'R': {'res_id': 16668008}, 'a...
1	1263908	21	20	[{'restaurant': {'R': {'res_id': 16668008}, 'a...
2	1263908	21	20	[{'restaurant': {'R': {'res_id': 16668008}, 'a...
3	1263908	21	20	[{'restaurant': {'R': {'res_id': 16668008}, 'a...
4	1263908	21	20	[{'restaurant': {'R': {'res_id': 16668008}, 'a...

```
import json
from pandas.io.json import json_normalize

with open('tarp/archive/file2.json') as f:
    d = json.load(f)

sam_json = json_normalize(d[0]["restaurants"])
sam_json.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 20 entries, 0 to 19
Data columns (total 34 columns):
#   Column                                Non-Null Count  Dtype
---  -
0   restaurant.R.res_id                  20 non-null    int64
1   restaurant.apikey                   20 non-null    object
2   restaurant.id                       20 non-null    object
3   restaurant.name                     20 non-null    object
4   restaurant.url                      20 non-null    object
5   restaurant.location.address         20 non-null    object
6   restaurant.location.locality        20 non-null    object
7   restaurant.location.city            20 non-null    object
8   restaurant.location.city_id        20 non-null    int64
9   restaurant.location.latitude       20 non-null    object
10  restaurant.location.longitude       20 non-null    object
11  restaurant.location.zipcode         20 non-null    object
12  restaurant.location.country_id      20 non-null    int64
13  restaurant.location.locality_verbose 20 non-null    object
14  restaurant.switch_to_order_menu     20 non-null    int64
15  restaurant.cuisines                 20 non-null    object
16  restaurant.average_cost_for_two     20 non-null    int64
17  restaurant.price_range              20 non-null    int64
18  restaurant.currency                 20 non-null    object
19  restaurant.offers                   20 non-null    object
20  restaurant.thumb                    20 non-null    object
21  restaurant.user_rating.aggregate_rating 20 non-null    object
22  restaurant.user_rating.rating_text  20 non-null    object
23  restaurant.user_rating.rating_color  20 non-null    object
24  restaurant.user_rating.votes        20 non-null    object
25  restaurant.photos_url               20 non-null    object
26  restaurant.menu_url                 20 non-null    object
27  restaurant.featured_image           20 non-null    object
28  restaurant.has_online_delivery       20 non-null    int64
29  restaurant.is_delivering_now        20 non-null    int64
30  restaurant.deeplink                 20 non-null    object
31  restaurant.has_table_booking        20 non-null    int64
32  restaurant.events_url               20 non-null    object
33  restaurant.establishment_types      20 non-null    object
dtypes: int64(9), object(25)
memory usage: 5.4+ KB
```

```
sam_json["restaurant.id"]=sam_json["restaurant.id"].astype(int)
```

```
sam_json[sam_json["restaurant.id"]!=sam_json["restaurant.R.res_id"]]
```

```
restaurant.R.res_id restaurant.apikey restaurant.id restaurant.name restaurant
```

0 rows × 34 columns

```
rename={'restaurant.id': 'Restaurant ID', 'restaurant.name': 'Restaurant Name',
'restaurant.location.address': 'Address', 'restaurant.location.locality': 'Locality',
'restaurant.location.city': 'City', 'restaurant.location.latitude': 'Latitude',
'restaurant.location.longitude': 'Longitude', 'restaurant.location.country_id': 'Country Code',
'restaurant.switch_to_order_menu': 'Switch to order menu', 'restaurant.cuisines': 'Cuisines',
'restaurant.average_cost_for_two': 'Average Cost for two', 'restaurant.price_range': 'Price range',
"restaurant.currency": "Currency", 'restaurant.user_rating.aggregate_rating': 'Aggregate rating',
'restaurant.user_rating.rating_text': 'Rating text', 'restaurant.user_rating.rating_color': 'Rating color',
'restaurant.user_rating.votes': 'Votes', 'restaurant.has_online_delivery': 'Has Online delivery',
'restaurant.is_delivering_now': 'Is delivering now', 'restaurant.has_table_booking': 'Has Table booking'
}
```

```
all_col=set(sam_json.columns)
drop=list(all_col-(all_col&set(rename.keys())))
```

```
sam_json=sam_json.drop(drop,axis=1).rename(columns=rename).set_index('Restaurant ID')
```

```
sam_json.head(5)
```

	Restaurant Name	Address	Locality	City	Latitude	Longit
Restaurant ID						
16668008	Arigato Sushi	14 Second Ave North, Yorkton, SK S3N 1G1	Yorkton	Yorkton	51.2106824000	-102.4613173
801690	Mocha	CP-1, 2nd Floor, Anand Plaza, Viram Khand-1, N...	Gomti Nagar	Lucknow	26.8528099000	81.0011849
17558738	Blue House Cafe	919 Bridge St, Vernonia, OR 97064	Vernonia	Vernonia	45.8586670000	-123.1954368
16611701	Star Buffet	58 Hanbury St, Mayfield, NSW	Mayfield	Mayfield	-32.8991780000	151.7343832
2100784	11th Avenue Cafe Bistro	Opposite Assam State Museum, Dighalipukhuri, T...	Uzan Bazaar	Guwahati	26.1860014826	91.7523143

```
combined_dataframe=concat([data_frame,sam_json])
combined_dataframe.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 9571 entries, 6317637 to 17211719
Data columns (total 20 columns):
#   Column                Non-Null Count  Dtype
---  -
0   Restaurant Name        9571 non-null   object
1   Country Code           9571 non-null   int64
2   City                   9571 non-null   object
3   Address                9571 non-null   object
4   Locality               9571 non-null   object
5   Locality Verbose       9551 non-null   object
6   Longitude              9571 non-null   object
7   Latitude               9571 non-null   object
8   Cuisines               9562 non-null   object
```

```
9 Average Cost for two 9571 non-null int64
10 Currency            9571 non-null object
11 Has Table booking   9571 non-null object
12 Has Online delivery 9571 non-null object
13 Is delivering now    9571 non-null object
14 Switch to order menu 9571 non-null object
15 Price range          9571 non-null int64
16 Aggregate rating     9571 non-null object
17 Rating color         9571 non-null object
18 Rating text          9571 non-null object
19 Votes                9571 non-null object
```

```
dtypes: int64(3), object(17)
```

```
memory usage: 1.5+ MB
```

```
combined_dataframe.describe()
```

	Country Code	Average Cost for two	Price range
count	9571.000000	9571.000000	9571.000000
mean	18.518023	1197.028315	1.805872
std	56.958785	16104.404230	0.905995
min	1.000000	0.000000	1.000000
25%	1.000000	250.000000	1.000000
50%	1.000000	400.000000	2.000000
75%	1.000000	700.000000	2.000000
max	216.000000	800000.000000	4.000000

```
combined_dataframe.head()
```

```

Restaurant Name Country Code City Address Locality Locality Verbose L
data_combine_frame=[]
error_load=[]
def combine_data(file):
    global di
    loaded_file = json.load(open(file))
    for d in loaded_file:
        try:
            sam_json = json_normalize(d["restaurants"])
            sam_json["restaurant.id"]=sam_json["restaurant.id"].astype(int)
            sam_json=sam_json.drop(drop,axis=1).rename(columns=rename).set_index('Restaurant ID')
            data_combine_frame.append(sam_json)
        except KeyError:
            error_load.append(d)
    pass

```

```

import glob
json_files = glob.glob("archive/*.json")
for file in json_files:
    combine_data(file)

```

```
error_load
```

```
[]
```

```

data_combine_frame=[data_frame]+data_combine_frame
data_warehouse=concat(data_combine_frame)

```

```
data_warehouse.info()
```

```

<class 'pandas.core.frame.DataFrame'>
Int64Index: 9551 entries, 6317637 to 5927402
Data columns (total 20 columns):
#   Column                Non-Null Count  Dtype
---  -
0   Restaurant Name       9551 non-null  object
1   Country Code          9551 non-null  int64
2   City                  9551 non-null  object
3   Address               9551 non-null  object
4   Locality              9551 non-null  object
5   Locality Verbose      9551 non-null  object
6   Longitude             9551 non-null  float64
7   Latitude              9551 non-null  float64
8   Cuisines              9542 non-null  object
9   Average Cost for two  9551 non-null  int64
10  Currency              9551 non-null  object
11  Has Table booking     9551 non-null  object
12  Has Online delivery   9551 non-null  object
13  Is delivering now     9551 non-null  object
14  Switch to order menu  9551 non-null  object
15  Price range           9551 non-null  int64
16  Aggregate rating      9551 non-null  float64
17  Rating color          9551 non-null  object
18  Rating text           9551 non-null  object
19  Votes                 9551 non-null  int64
dtypes: float64(3), int64(4), object(13)
memory usage: 1.5+ MB

```

```
data_warehouse["Country Code"].unique()
```

```
array([162, 30, 216, 14, 37, 184, 214, 1, 94, 148, 215, 166, 189,
       191, 208])
```

```
data_warehouse["Currency"].unique()
```

```
array(['Botswana Pula(P)', 'Brazilian Real(R$)', 'Dollar($)',
       'Emirati Diram(AED)', 'Indian Rupees(Rs.)',
```

```
'Indonesian Rupiah(IDR)', 'NewZealand($)', 'Pounds(\x8c£)',
'Qatari Rial(QR)', 'Rand(R)', 'Sri Lankan Rupee(LKR)',
'Turkish Lira(TL)']], dtype=object)
```

```
data_warehouse["Currency"].replace(["Pounds(\x8c£)", "NewZealand($)"], ["Pounds(£)", "NewZealand(NZ$)"], inplace=True)
data_warehouse["Currency"].unique()
```

```
array(['Botswana Pula(P)', 'Brazilian Real(R$)', 'Dollar($)',
'Emirati Diram(AED)', 'Indian Rupees(Rs.)',
'Indonesian Rupiah(IDR)', 'NewZealand(NZ$)', 'Pounds(£)',
'Qatari Rial(QR)', 'Rand(R)', 'Sri Lankan Rupee(LKR)',
'Turkish Lira(TL)'], dtype=object)
```

```
cur=data_warehouse["Currency"].unique()
for i in cur:
    try:
        data_warehouse["Currency"].replace(i.split("(")[1].split(")")[0], i, inplace=True)
    except Exception:
        pass
data_warehouse["Currency"].unique()
```

```
array(['Botswana Pula(P)', 'Brazilian Real(R$)', 'Dollar($)',
'Emirati Diram(AED)', 'Indian Rupees(Rs.)',
'Indonesian Rupiah(IDR)', 'NewZealand(NZ$)', 'Pounds(£)',
'Qatari Rial(QR)', 'Rand(R)', 'Sri Lankan Rupee(LKR)',
'Turkish Lira(TL)'], dtype=object)
```

```
data_warehouse["Price range"].unique()
```

```
array([3, 4, 2, 1])
```

```
data_warehouse["Has Table booking"].unique()
```

```
array(['Yes', 'No'], dtype=object)
```

```
data_warehouse["Has Table booking"].replace([1,0], ["Yes", "No"], inplace=True)
data_warehouse["Has Table booking"].unique()
```

```
array(['Yes', 'No'], dtype=object)
```

```
data_warehouse["Has Online delivery"].replace([1,0], ["Yes", "No"], inplace=True)
data_warehouse["Is delivering now"].replace([1,0], ["Yes", "No"], inplace=True)
data_warehouse["Switch to order menu"].replace([1,0], ["Yes", "No"], inplace=True)
```

```
data_warehouse=merge(data_warehouse, country, how="left", on="Country Code")
data_warehouse["Country"].unique()
```

```
array(['Phillipines', 'Brazil', 'United States', 'Australia', 'Canada',
'Singapore', 'UAE', 'India', 'Indonesia', 'New Zealand',
'United Kingdom', 'Qatar', 'South Africa', 'Sri Lanka', 'Turkey'],
dtype=object)
```

```
#data_warehouse.drop(["restaurant.order_url", "restaurant.order_deeplink", "restaurant.zomato_events", "restaurant.book_url", "restaurant.phone_number"], axis=1)
```

```
data_warehouse.head()
```


	Restaurant Name	Country Code	City	Address	Locality	Locality Verbose	Longitude
0	Le Petit Souffle	162	Makati City	Third Floor, Century City Mall, Kalayaan Avenu...	Century City Mall, Poblacion, Makati City	Century City Mall, Poblacion, Makati City, Mak...	121.027535
1	Izakaya Kikufuji	162	Makati City	Little Tokyo, 2277 Chino Roces Avenue, Legaspi...	Little Tokyo, Legaspi Village, Makati City	Little Tokyo, Legaspi Village, Makati City, Ma...	121.014101
2	Heat - Edsa Shangri-La	162	Mandaluyong City	Edsa Shangri-La, 1 Garden Way, Ortigas, Mandal...	Edsa Shangri-La, Ortigas, Mandaluyong City	Edsa Shangri-La, Ortigas, Mandaluyong City, Ma...	121.056831
3	Ooma	162	Mandaluyong City	Third Floor, Mega Fashion Hall, SM	SM Megamall, Ortigas, Mandaluvong	SM Megamall, Ortigas, Mandaluyong C...	121.056475

```
data_warehouse.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 9551 entries, 0 to 9550
Data columns (total 21 columns):
#   Column                                Non-Null Count  Dtype
---  -
0   Restaurant Name                       9551 non-null   object
1   Country Code                         9551 non-null   int64
2   City                                 9551 non-null   object
3   Address                             9551 non-null   object
4   Locality                             9551 non-null   object
5   Locality Verbose                     9551 non-null   object
6   Longitude                            9551 non-null   float64
7   Latitude                             9551 non-null   float64
8   Cuisines                             9542 non-null   object
9   Average Cost for two                 9551 non-null   int64
10  Currency                             9551 non-null   object
11  Has Table booking                    9551 non-null   object
12  Has Online delivery                 9551 non-null   object
13  Is delivering now                   9551 non-null   object
14  Switch to order menu                9551 non-null   object
15  Price range                         9551 non-null   int64
16  Aggregate rating                    9551 non-null   float64
17  Rating color                        9551 non-null   object
18  Rating text                         9551 non-null   object
19  Votes                              9551 non-null   int64
20  Country                             9551 non-null   object
dtypes: float64(3), int64(4), object(14)
memory usage: 1.6+ MB
```

```
data_warehouse.describe()
```

	Country Code	Longitude	Latitude	Average Cost for two	Price range	Aggregate rating
count	9551.000000	9551.000000	9551.000000	9551.000000	9551.000000	9551.000000

```
data_warehouse.to_csv('DataWarehouse.csv',index=False)
```

```
std      56.750546    41.467058    11.007935    16121.183073    0.905609    1.516378
```

▼ Data Analysing

```
from pandas import *
from numpy.random import *
from numpy import *
import seaborn as sns
import matplotlib.pyplot as plt
import warnings
warnings.filterwarnings('ignore')
df= read_csv('tarp/DataWarehouse.csv',engine='python')
```

```
#Sorting the ratings of restaurants by country code
#We are rating the restaunts of a certain country code
e = df[df['Rating text']=='Excellent']['Country'].value_counts()
vg = df[df['Rating text']=='Very Good']['Country'].value_counts()
gr = df[df['Rating text']=='Good']['Country'].value_counts()
ar = df[df['Rating text']=='Average']['Country'].value_counts()
pr = df[df['Rating text']=='Poor']['Country'].value_counts()
nr = df[df['Rating text']=='Not rated']['Country'].value_counts()
```

```
#display distinct contry codes
df["Country"].unique()
```

```
array(['Phillipines', 'Brazil', 'United States', 'Australia', 'Canada',
       'Singapore', 'UAE', 'India', 'Indonesia', 'New Zealand',
       'United Kingdom', 'Qatar', 'South Africa', 'Sri Lanka', 'Turkey',
       nan], dtype=object)
```

```
dfr = DataFrame([e,vg,gr,ar, pr,nr])
dfr.index = ['Excellent','Very Good','Good','Average','Poor','Not rated']
dfr.fillna(0) #while displaying replaces all null values with 0
```

	India	United States	Australia	United Kingdom	UAE	Brazil	Phillipines	New Zealand	Sou Afri
Excellent	1524.0	317.0	183.0	46.0	36.0	32.0	24.0	24.0	24.0
Very Good	6562.0	888.0	868.0	62.0	62.0	40.0	199.0	50.0	70.0
Good	6122.0	1183.0	2145.0	40.0	18.0	22.0	183.0	4.0	24.0
Average	9197.0	757.0	739.0	10.0	2.0	16.0	0.0	0.0	2.0
Poor	476.0	158.0	183.0	0.0	2.0	0.0	0.0	2.0	0.0
Not rated	5274.0	6.0	0.0	2.0	0.0	10.0	0.0	0.0	0.0

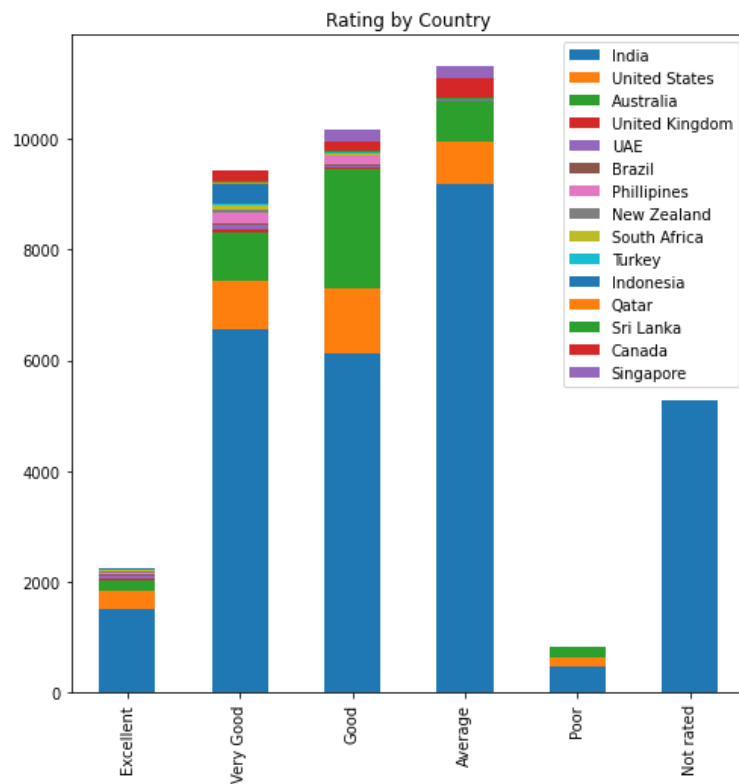


```
dfr.describe()
```

	India	United States	Australia	United Kingdom	UAE	Brazil	Phillipir
count	6.000000	6.000000	5.000000	5.00000	5.000000	5.000000	3.0000
mean	4859.166667	551.500000	823.600000	32.00000	24.000000	24.000000	135.3333
std	3280.815778	460.794423	802.490374	25.21904	25.455844	12.083046	96.7488

#DATA EXPLORATION AND VISUALISATION

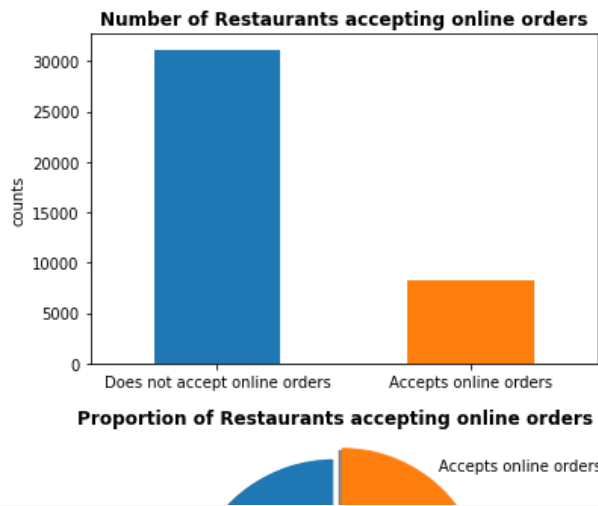
```
dfr.plot(kind='bar',stacked=True, figsize=(8,8), title="Rating by Country")
plt.show()
```



```
#How many Restaurant accepting online orders
ax =df['Has Online delivery'].value_counts()
ax = ax.plot(kind='bar', color=['#1f77b4', '#ff7f0e'])
ax.set_xticklabels(['Does not accept online orders','Accepts online orders'], rotation=0)
plt.title('Number of Restaurants accepting online orders', weight='bold')
plt.xlabel('')
plt.ylabel('counts')
plt.show()

import matplotlib.pyplot as plt

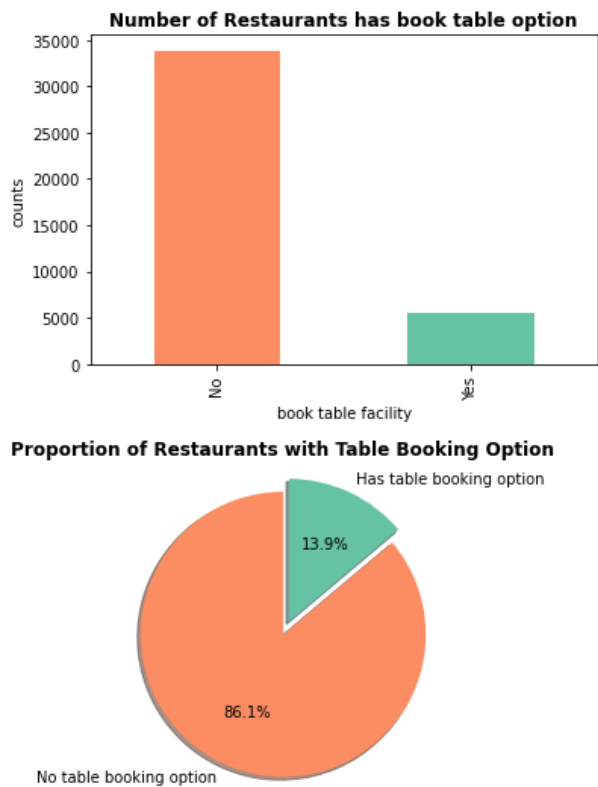
online_orders = df['Has Online delivery'].value_counts(normalize=True)
labels = ['Does not accept online orders','Accepts online orders']
colors = ['#1f77b4', '#ff7f0e']
explode = (0.1, 0)
plt.pie(online_orders, explode=explode, labels=labels, colors=colors, autopct='%1.1f%%', shadow=True, startangle=90)
plt.axis('equal')
plt.title('Proportion of Restaurants accepting online orders', weight='bold')
plt.show()
```



```
#How many Restaurant have option to book a table
ax =df['Has Table booking'].value_counts().plot(kind='bar', color=plt.cm.Set2([0.2, 0.1]))
plt.title('Number of Restaurants has book table option', weight='bold')
plt.xlabel('book table facility')
plt.ylabel('counts')
plt.show()

#How many Restaurant have option to book a table
import matplotlib.pyplot as plt

table_booking = df['Has Table booking'].value_counts(normalize=True)
labels = ['No table booking option', 'Has table booking option']
colors = plt.cm.Set2([0.2, 0.1])
explode = (0.1, 0)
plt.pie(table_booking, explode=explode, labels=labels, colors=colors, autopct='%1.1f%%', shadow=True, startangle=90)
plt.axis('equal')
plt.title('Proportion of Restaurants with Table Booking Option', weight='bold')
plt.show()
```



```
#top 5 cities with most number of restraunts
plt.figure(figsize=(7,7))
ax =df['City'].value_counts()[:5].plot(kind='pie')
plt.title('Location', weight='bold')
```

```
plt.show()

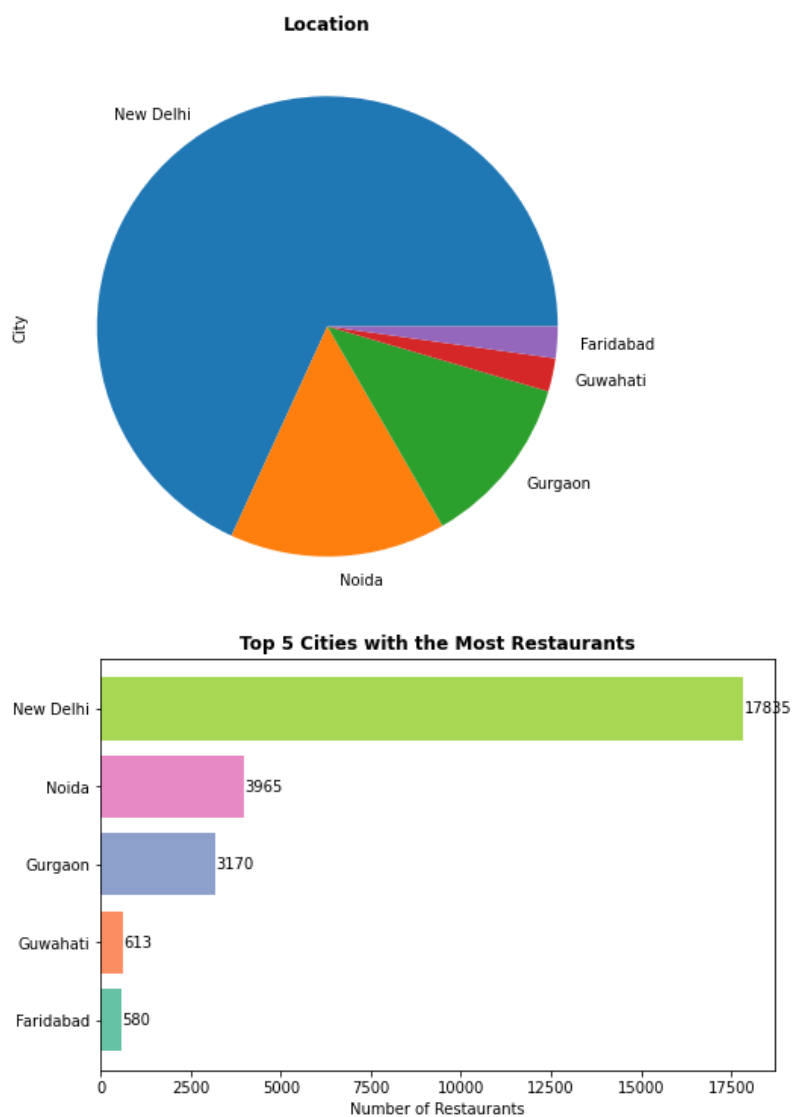
#top 5 cities with most number of restaurants
import matplotlib.pyplot as plt

top_cities = df['City'].value_counts()[:5].sort_values(ascending=True)
labels = top_cities.index
counts = top_cities.values
colors = plt.cm.Set2([0.1, 0.2, 0.3, 0.4, 0.5])

fig, ax = plt.subplots(figsize=(8, 5))
ax.barh(labels, counts, color=colors)
ax.set_xlabel('Number of Restaurants')
ax.set_title('Top 5 Cities with the Most Restaurants', weight='bold')

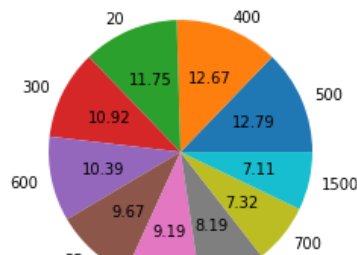
# Add counts at the end of each bar
for i, count in enumerate(counts):
    ax.text(count + 10, i, str(count), ha='left', va='center')

plt.show()
```

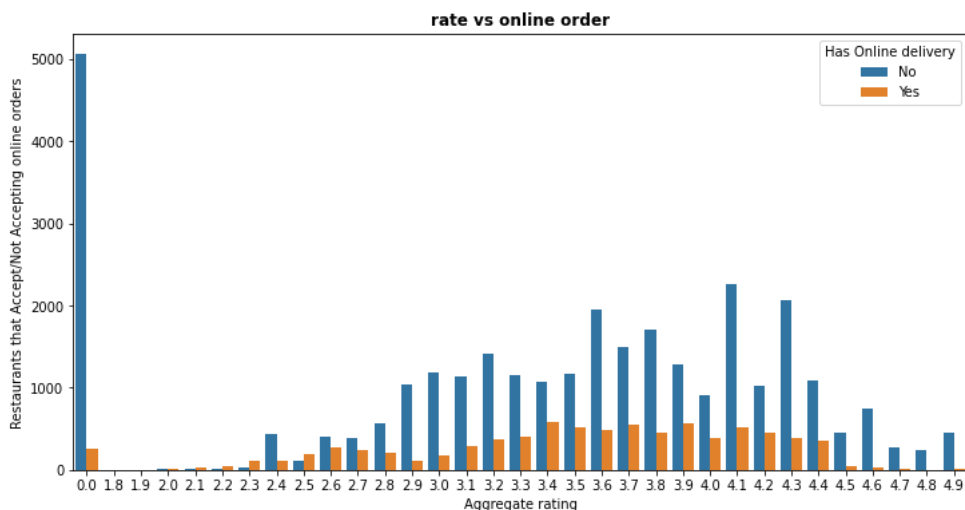


```
#Average cost for 2 people to dine
values = df['Average Cost for two'].value_counts()[:10]
labels = df['Average Cost for two'].value_counts()[:10].index
plt.pie(values, labels=labels, autopct='%.2f')
plt.title('Average cost for two person(in %)', weight='bold')
plt.show()
```

Average cost for two person(in %)



```
#Rating vs Online orders
plt.figure(figsize = (12,6))
sns.countplot(x=df['Aggregate rating'], hue = df['Has Online delivery'])
plt.ylabel("Restaurants that Accept/Not Accepting online orders")
plt.title("rate vs online order",weight = 'bold')
plt.show()
```



```
#Displaying the ratings of the restaurants by price range (4 is high; 1 is low)
lowest_price = df[df['Price range']== 1]['Rating text'].value_counts()
lower_price = df[df['Price range']== 2]['Rating text'].value_counts()
higher_price = df[df['Price range']== 3]['Rating text'].value_counts()
highest_price = df[df['Price range']== 4]['Rating text'].value_counts()

df_price_by_rating = DataFrame([lowest_price,lower_price,higher_price,highest_price])
df_price_by_rating.index = ['Lowest Price (1)','Lower Price (2)','Higher Price (3)','Highest Price (4)']
df_price_by_rating.plot(kind='bar',stacked=True, figsize=(6,6), title="Price by Rating")
plt.show()
```

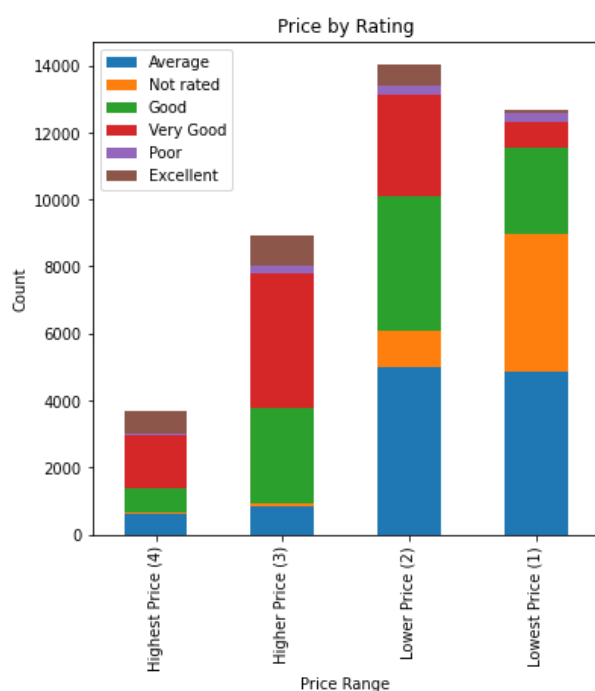
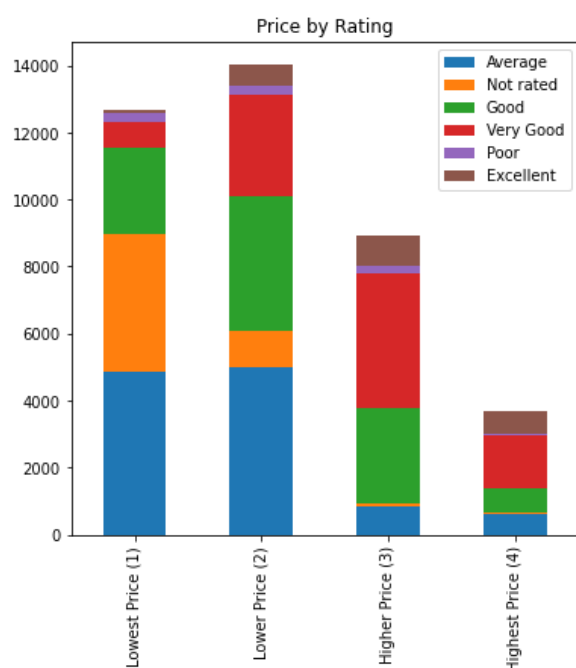
```
import pandas as pd
lowest_price = df[df['Price range']== 1]['Rating text'].value_counts()
lower_price = df[df['Price range']== 2]['Rating text'].value_counts()
higher_price = df[df['Price range']== 3]['Rating text'].value_counts()
highest_price = df[df['Price range']== 4]['Rating text'].value_counts()

df_price_by_rating = pd.DataFrame([lowest_price,lower_price,higher_price,highest_price])
df_price_by_rating.index = ['Lowest Price (1)','Lower Price (2)','Higher Price (3)','Highest Price (4)']

# Sorting the x-axis labels in descending order based on the average rating of each price range
df_price_by_rating['Average Rating'] = [df[df['Price range']== 1]['Aggregate rating'].mean(),
                                         df[df['Price range']== 2]['Aggregate rating'].mean(),
                                         df[df['Price range']== 3]['Aggregate rating'].mean(),
                                         df[df['Price range']== 4]['Aggregate rating'].mean()]
df_price_by_rating.sort_values(by='Average Rating', ascending=False, inplace=True)
df_price_by_rating.drop(columns=['Average Rating'], inplace=True)

# Displaying the stacked bar chart with labels and title
df_price_by_rating.plot(kind='bar',stacked=True, figsize=(6,6), title="Price by Rating")
plt.xlabel('Price Range')
```

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



```
df["Rating color"].unique()
```

```
sns.scatterplot(x='Aggregate rating',y='Votes',data=df)
plt.show()
```

```
import seaborn as sns
import matplotlib.pyplot as plt
```

```
sns.scatterplot(x='Aggregate rating', y='Votes', data=df)
plt.title('Aggregate Rating vs. Votes', weight='bold')
plt.xlabel('Aggregate rating')
plt.ylabel('Votes')
plt.show()
```

```
#Calculating the percentage of restaurants that were not rated:
```

```
#Total number of restaurants not rated in the training set
Number_of_restaurants_nr= dfr.loc['Not rated'].sum()
"""remove double reading"""
```

```

with open('tarp/DataWarehouse.csv',errors='ignore') as f:
    total_row_count = sum(1 for row in f)

number_of_restaurants_train= total_row_count * 0.7

Percentage_nr = (Number_of_restaurants_nr/number_of_restaurants_train)*100
print("Percentage of restaurants not rated =" ,round(Percentage_nr,"%")

```

▼ Data Cleaning

```

from pandas import *
from numpy.random import *
from numpy import *
import seaborn as sns
import matplotlib.pyplot as plt
import warnings
warnings.filterwarnings('ignore')
df= read_csv('tarp/DataWarehouse.csv',engine='python')

```

```

# the lowest priced restaurants that are not rated
mean_lowest_price = df[df['Price range']== 1]['Aggregate rating'].replace(0.0, nan).mean()
print(f"Mean of aggregate rating with price range 1: {mean_lowest_price}")

```

Mean of aggregate rating with price range 1: 3.3349848449522033

```

# filling with their average
df[df['Price range']== 1]['Aggregate rating'].replace(0.0,nan, inplace=True)
df['Aggregate rating'].replace(nan, mean_lowest_price, inplace=True)

```

```

# the lowest priced restaurants that are not rated
new_mean_lowest_price = df['Aggregate rating'].replace(0.0, nan).mean()
print(f"New Mean of aggregate rating with price range 1: {new_mean_lowest_price}")

```

New Mean of aggregate rating with price range 1: 3.659719934102142

```

#Next the lower priced restaurants that are not rated
mean_lower_price = df[df['Price range']== 2]['Aggregate rating'].replace(0, nan).mean()
print(f"Mean of aggregate rating with price range 2: {mean_lowest_price}")

```

Mean of aggregate rating with price range 2: 3.3349848449522033

```

# filling with their average
df['Aggregate rating'].replace(0, mean_lower_price, inplace=True)

```

```

#Checking
print(df['Aggregate rating'].mean())

```

```

#Checking the new scatterplot, no restaurant has a rating of 0 now.
sns.scatterplot(x='Aggregate rating',y='Votes',data=df)
y=df["Rating text"]

```

```

#As an aggregate rating of 2.5-3.5 translates to a rating of "Average", the "Not Rated"
#text can be replaced with "Average" in the dependent variable.

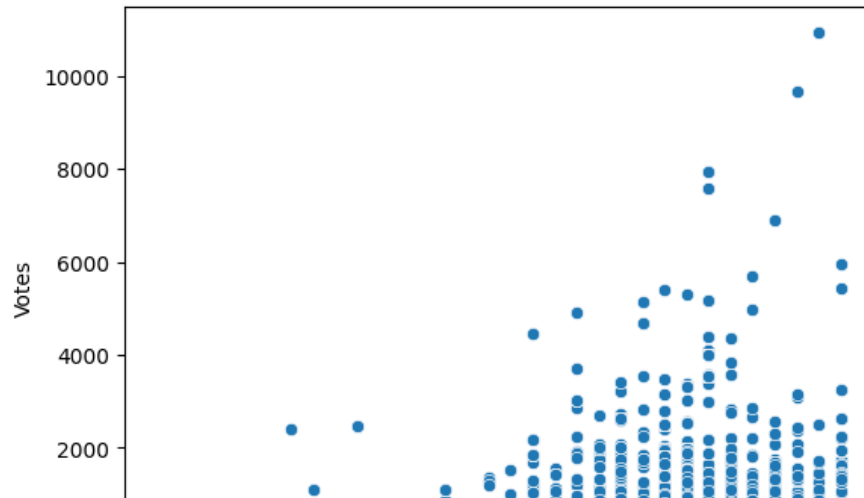
```

```

y.replace('Not rated','Average', inplace=True)

```


3.6490085496293516



#'Rating Text' and 'Rating Color' is given as a string.
#encoded into numeric data so that the algorithm can perform calculations.

```
from sklearn.preprocessing import LabelEncoder
labelEncoder_X = LabelEncoder()
df["Rating text"]=labelEncoder_X.fit_transform(df["Rating text"])
#Shows that excellent = 1, very good = 4, good = 2, average = 0, and assuming poor = 3.
print(df["Rating text"].unique())
df["Rating color"]=labelEncoder_X.fit_transform(df["Rating color"])
```

[1 4 2 0 3]

```
for index, row in df.iterrows():
    try:
        l=row["Cuisines"].split(", ")
    except:
        continue

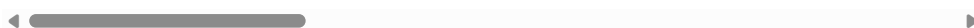
    for i in l:
        if i not in df.columns:
            df.loc[index,i]=1
            df[i].replace(nan, 0, inplace=True)
            df.loc[index,i]=1
```

```
df.drop(["Cuisines"],axis=1,inplace=True)
```

```
df.describe()
```

	Longitude	Latitude	Average Cost for two	Price range	Aggregate rating	Rating
count	39304.000000	39304.000000	39304.000000	39304.000000	39304.000000	39304.00
mean	63.626996	20.752542	2217.672756	2.091441	3.649009	4.62
std	62.813259	21.347402	17750.203378	0.956545	0.544524	4.06
min	-157.948486	-41.330428	0.000000	1.000000	1.800000	0.00
25%	77.085958	26.852810	120.000000	1.000000	3.300000	1.00
50%	77.218187	28.568446	450.000000	2.000000	3.600000	3.00
75%	77.333096	28.642216	800.000000	3.000000	4.100000	8.00
max	174.832089	73.990000	800000.000000	4.000000	4.900000	13.00

8 rows × 156 columns



```
df['Latitude'].isnull().unique()
```

array([False])

```
df['Longitude'].isnull().unique()
```

```
array([False])
```

```
df["Has Online delivery"].replace(["Yes", "No"], [1, 0], inplace=True)
df["Is delivering now"].replace(["Yes", "No"], [1, 0], inplace=True)
df["Has Table booking"].replace(["Yes", "No"], [1, 0], inplace=True)
df["Switch to order menu"].replace(["Yes", "No"], [1, 0], inplace=True)
```

```
df = df.drop(['Restaurant Name', 'Locality', 'Country', 'City', 'Address', 'Currency', 'Rating color', 'Aggregate rating'], axis=1)
```

```
df.head()
```

	Longitude	Latitude	Average Cost for two	Has Table booking	Has Online delivery	Is delivering now	Switch to order menu	Price range	Rating
0	121.027535	14.565443	1100	1	0	0	0	3	
1	121.014101	14.553708	1200	1	0	0	0	3	
2	121.056831	14.581404	4000	1	0	0	0	4	
3	121.056475	14.585318	1500	0	0	0	0	4	
4	121.057508	14.584450	1500	1	0	0	0	4	

5 rows × 10 columns



```
df.to_csv('CleanedData.csv', index=False)
```

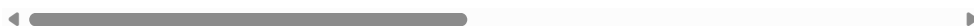
▼ Model Training & Testing

```
from pandas import *
from numpy.random import *
from numpy import *
import seaborn as sns
import matplotlib.pyplot as plt
import warnings
warnings.filterwarnings('ignore')
df= read_csv('tarp/CleanedData.csv', engine='python')
```

```
df.head()
```

	Longitude	Latitude	Average Cost for two	Has Table booking	Has Online delivery	Is delivering now	Switch to order menu	Price range	Rating
0	121.027535	14.565443	1100	1	0	0	0	3	
1	121.014101	14.553708	1200	1	0	0	0	3	
2	121.056831	14.581404	4000	1	0	0	0	4	
3	121.056475	14.585318	1500	0	0	0	0	4	
4	121.057508	14.584450	1500	1	0	0	0	4	

5 rows × 10 columns



```
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 39304 entries, 0 to 39303
Columns: 10 entries, Longitude to Rating
```

```
dtypes: float64(150), int64(8)
memory usage: 47.4 MB
```

```
y=df['Rating text']
X=df.drop(["Rating text"], axis=1)
print(y)
```

```
0      1
1      1
2      4
3      1
4      1
..
39299   2
39300   2
39301   2
39302   2
39303   2
Name: Rating text, Length: 39304, dtype: int64
```

```
from sklearn.model_selection import train_test_split
x_train, x_test, y_train, y_test = train_test_split(X,y, test_size=0.3, random_state=0)
y_test
```

```
16432   2
14397   4
20117   4
23753   2
23133   0
..
26628   0
36633   0
12954   4
28237   2
22159   0
Name: Rating text, Length: 11792, dtype: int64
```

```
#used to save ml model for loading later
from pickle import dump
```

```
# Fitting K-NN to the Training set
from sklearn.neighbors import KNeighborsClassifier
classifier = KNeighborsClassifier(n_neighbors = 9, metric = 'minkowski', p = 2)

classifier.fit(x_train,y_train)
dump(classifier, open("knn.sav", 'wb'))
knn_score=classifier.score(x_test,y_test)
print(f"K-NN:\n Accuracy: {knn_score*100}%")
```

```
K-NN:
Accuracy: 85.78697421981005%
```

```
# Fitting Logistic Regression to the Training set
from sklearn.linear_model import LogisticRegression
classifier = LogisticRegression(penalty='l2',random_state = 0)

classifier.fit(x_train,y_train)
dump(classifier, open("logistic_regression.sav", 'wb'))
lr_score=classifier.score(x_test,y_test)
print(f"Logistic Regression:\n Accuracy: {lr_score*100}%")
```

```
Logistic Regression:
Accuracy: 56.50440976933514%
```

```
# Fitting SVM to the Training set
from sklearn.svm import SVC
classifier = SVC(kernel = 'rbf', random_state = 0)

classifier.fit(x_train,y_train)
```

```

dump(classifier, open("svm.sav", 'wb'))
svm_score=classifier.score(x_test,y_test)
print(f"SVM:\n Accuracy: {svm_score*100}%")

```

SVM:
Accuracy: 52.1370420624152%

```

# Fitting Naive Bayes to the Training set
from sklearn.naive_bayes import GaussianNB
classifier = GaussianNB()

classifier.fit(x_train,y_train)
dump(classifier, open("naive_bayes.sav", 'wb'))
nb_score=classifier.score(x_test,y_test)
print(f"Naive Bayes:\n Accuracy: {nb_score*100}%")

```

Naive Bayes:
Accuracy: 50.703867028493896%

```

# Fitting Random Forest Classification to the Training set
from sklearn.ensemble import RandomForestClassifier
classifier = RandomForestClassifier(n_estimators = 100, criterion = 'entropy', random_state = 0)

classifier.fit(x_train,y_train)
dump(classifier, open("random_forest.sav", 'wb'))
rf_score=classifier.score(x_test,y_test)
print(f"Random Forest:\n Accuracy: {rf_score*100}%")

```

Random Forest:
Accuracy: 97.09972862957937%

```

import matplotlib.pyplot as plt
import numpy as np

# Define the data
algos = ["K-NN", "Logistic Regression", "SVM", "Naive Bayes", "Random Forest"]
scores = [knn_score, lr_score, svm_score, nb_score, rf_score]
accuracies = np.array(scores) * 100

# Define the colors for the bars
colors = ['#1f77b4', '#ff7f0e', '#2ca02c', '#d62728', '#9467bd']

# Create the figure and axes objects
fig, ax = plt.subplots(figsize=(10, 6))

# Create the bar plot
bars = ax.bar(algos, accuracies, color=colors, align='center', alpha=0.7, ecolor='black', capsize=10, width=0.6)

# Add the axis labels and title
ax.set_ylabel('Accuracy (%)', fontsize=14)
ax.set_xlabel('Algorithm', fontsize=14)
ax.set_title('Comparison of Machine Learning Algorithms', fontsize=16)

# Add the accuracy percentage to the top of each bar
for bar in bars:
    height = bar.get_height()
    ax.text(bar.get_x() + bar.get_width() / 2, height, f'{height:.2f}%', ha='center', va='bottom', fontsize=10)

# Create a legend for the labels
labels = ['K-NN', 'Logistic Regression', 'SVM', 'Naive Bayes', 'Random Forest']
handles = [plt.Rectangle((0,0),1,1, color=colors[i]) for i in range(len(labels))]
ax.legend(handles, labels, loc='upper left', bbox_to_anchor=(1, 1), fontsize=12)

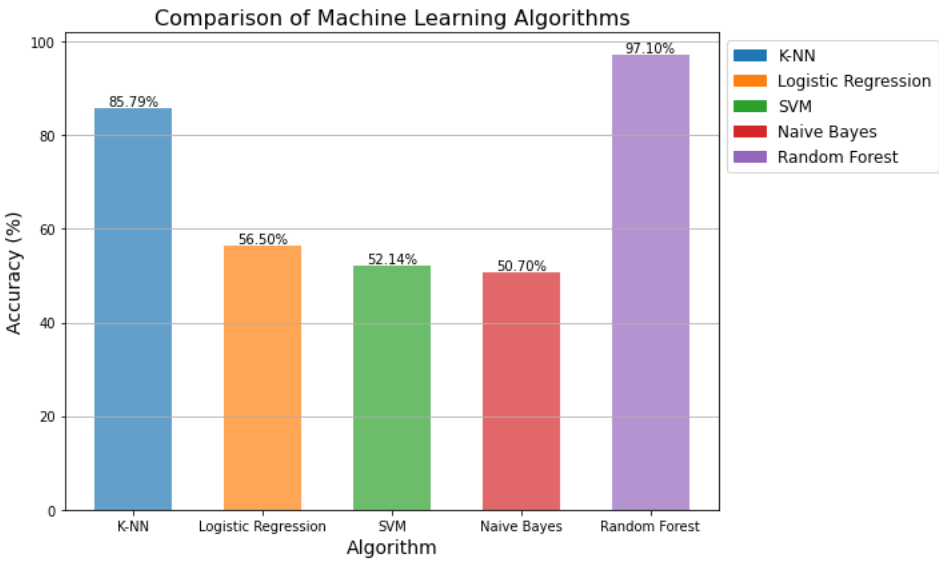
# Add horizontal grid lines
ax.yaxis.grid(True)

# Adjust the plot layout
plt.tight_layout()

# Display the plot

```

```
# Display the plot
plt.show()
plt.savefig('bar_plot_with_error_bars.png')
```



<Figure size 432x288 with 0 Axes>