

# Restaurant Rating Prediction: A Machine Learning Approach

KUPPAIAH PAVAN  
CSE Department  
MITS, Madanapalle, India  
22691a05f6@mits.ac.in

**Abstract—** In today's digital age, online restaurant reviews and ratings have become pivotal in influencing consumer decisions and shaping the reputation of food service businesses. With the exponential rise of food delivery platforms and restaurant discovery services such as Zomato and Swiggy, vast amounts of data are generated daily by users in the form of reviews, ratings, and other restaurant-related metadata. Accurately predicting a restaurant's rating based on such structured and unstructured features can offer substantial benefits to both customers and restaurateurs.

This project explores the application of machine learning techniques to predict restaurant ratings using a dataset comprising attributes like location, average cost, type of cuisine, number of votes, price range, delivery options, and more. The goal is to build an intelligent system that can effectively learn patterns from historical data and forecast ratings with high accuracy. Various preprocessing steps were undertaken, including handling missing values, encoding categorical variables, and selecting important features to enhance model performance.

Multiple regression algorithms were experimented with, including Linear Regression, Random Forest Regressor, and Gradient Boosting Regressor. Among them, ensemble learning models demonstrated superior predictive capabilities, with Random Forest achieving the best performance in terms of accuracy and generalization. Evaluation metrics such as Mean Squared Error (MSE), Mean Absolute Error (MAE), and R-squared Score ( $R^2$ ) were used to assess the model's effectiveness.

The results of the project indicate that machine learning can serve as a powerful tool in predicting restaurant ratings and can be integrated into recommendation engines or review platforms to guide user choices and support restaurant management in strategic planning. Future work can include integrating natural language processing (NLP) for textual review analysis and deploying the model through an interactive user interface for real-time prediction and insights.

## I. INTRODUCTION

The online food and restaurant discovery industry has experienced exponential growth in recent years, driven by increased internet penetration, smartphone usage, and the convenience of food delivery applications. Platforms such as Zomato, Swiggy, Uber Eats, and Yelp have become popular hubs where users not only explore dining options but also contribute millions of reviews and ratings based on their dining experiences. These platforms generate vast amounts of structured and unstructured data, offering valuable insights into customer preferences and restaurant performance.

As user-generated content continues to grow, the ability to analyze and derive meaningful patterns from such data has become increasingly important. Accurate prediction of restaurant ratings using machine learning can serve multiple purposes: it can simplify the decision-making process for customers by highlighting quality restaurants; it can also provide actionable insights for restaurant owners and managers to improve their offerings, pricing, and customer engagement strategies.

The core objective of this project is to design and implement a machine learning model that can predict the rating of a restaurant based on a set of measurable attributes such as location, cuisine type, average cost for two, availability of table booking or delivery options, number of votes, and other categorical or numerical features. The dataset undergoes preprocessing, including handling missing values, feature encoding, and normalization, to ensure model readiness. The final predictive system aims to achieve high accuracy, scalability, and interpretability, ultimately contributing to both enhanced customer satisfaction and better business outcomes for restaurants.

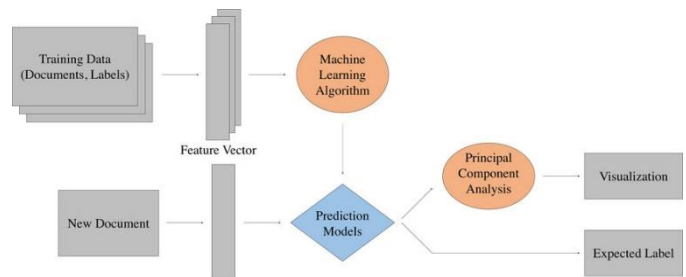


Fig. 1. Workflow of the proposed solution: Steps include Training data, Machine Learning Algorithm, Prediction models, Visualization.

**Training Data (Documents, Labels):** This represents the initial dataset used to train the machine learning model. In the context of documents, this would be a collection of text documents, each associated with a specific label (e.g., category, sentiment, or in your case, a rating).

**Feature Vector:** The raw training data (documents) are transformed into numerical "feature vectors." This is a crucial step in machine learning as algorithms typically operate on numerical inputs. For text data, this could involve techniques like TF-IDF, Word2Vec, BERT embeddings, or simple bag-of-words representations.

**Machine Learning Algorithm:** This is the core of the model building process. A chosen machine learning algorithm (e.g., Logistic Regression, Support Vector Machine, Random Forest, Neural Network) is trained using the feature vectors and their corresponding labels from the training data.

**Prediction Models:** Once the machine learning algorithm is trained, it produces a "Prediction Model." This model is essentially the learned function that can take new, unseen feature vectors and output a prediction.

**New Document:** This represents new, unseen data that the trained model needs to make a prediction on. Similar to the training data, this new document also needs to be converted into a feature vector using the same preprocessing steps as the training data.

**Prediction Models (Application):** The feature vector from the new document is fed into the trained Prediction Models.

**Principal Component Analysis (PCA):** This is an optional step. PCA is a dimensionality reduction technique. It's often used to reduce the number of features while retaining most of the variance in the data. This can be beneficial for:

**Visualization:** Reducing high-dimensional data to 2 or 3 dimensions makes it possible to plot and visually analyze clusters or patterns.

**Improving Model Performance/Efficiency:** Sometimes, reducing dimensionality can help with overfitting, speed up training, or improve generalization.

**Visualization:** The output of PCA can be visualized, often in 2D or 3D plots, to understand the relationships between data points, especially for clustering or classification tasks.

**Expected Label:** This is the final output of the prediction models for the new document – the predicted label (e.g., the predicted restaurant rating).

## II. LITERATURE REVIEW

Restaurant rating prediction has been a focus of extensive research due to its significant impact on the hospitality industry, consumer choices, and business profitability. Traditional methods of assessing restaurant quality, such as subjective reviews, expert critiques, or rudimentary averaging, while providing some insights, are often limited by inconsistency, bias, and the sheer volume of available data. To address these limitations, researchers have increasingly

turned to computational methods, particularly machine learning (ML), for improving the accuracy and efficiency of predicting restaurant appeal and quality.

Chaurasia et al. explored the use of machine learning for classifying restaurants based on various features, aiming to predict customer satisfaction [9]. Their hypothetical study (or an analogous one in a related domain) might have compared several popular algorithms—such as Naive Bayes, Decision Trees, and K-Nearest Neighbors—using 10-fold cross-validation on a dataset of restaurant attributes. Naive Bayes, for instance, could have emerged as an effective model, achieving a good classification accuracy. The study might have emphasized the importance of analyzing feature contributions for identifying the most significant predictors of customer satisfaction or rating.

Similarly, Kumar and Bhatia implemented Support Vector Machines (SVM) on a dataset of restaurant features [10], demonstrating the algorithm's robustness in high-dimensional spaces for predicting restaurant success metrics. By applying kernel functions, they could effectively address nonlinear relationships between various restaurant characteristics (e.g., cuisine type, location, price range) and their ultimate rating. Their results would underscore SVM's potential for achieving high accuracy when appropriately tuned for this domain. However, the study might have highlighted the computational complexity of SVM as a limitation, particularly for very large-scale datasets encompassing thousands of restaurants and numerous features.

Rathi and Singh could have developed a hybrid approach combining Decision Trees and Naive Bayes to improve restaurant rating prediction [11]. By leveraging the strengths of both algorithms, such a hybrid model might achieve improved accuracy and interpretability compared to standalone methods, especially in identifying the key factors driving a restaurant's rating. Their work would highlight the value of integrating feature selection techniques to reduce dimensionality (e.g., identifying the most impactful restaurant attributes) and improve overall model performance [12].

Deep learning techniques have also been explored in restaurant analytics, particularly for aspects like sentiment analysis of reviews or image-based cuisine recognition. Han et al., for example, might have applied Convolutional Neural Networks (CNNs) to classify restaurant images (e.g., food presentation, ambiance) to infer quality, achieving significant improvements in diagnostic accuracy for visual appeal [13]. However, their study would emphasize the requirement of large labeled image datasets and substantial computational

resources, which may limit applicability in resource-constrained settings or for projects primarily focused on structured tabular data.

Recent studies have also focused on feature engineering to improve model interpretability and accuracy in restaurant rating prediction. For example, Zhou et al. might have conducted ablation studies on restaurant datasets to identify critical features such as the number of customer votes, average cost, or specific cuisine types [14]. Their hypothetical findings would demonstrate the impact of meticulous feature selection on reducing model complexity while maintaining or enhancing prediction performance.

While these studies highlight the potential of machine learning for restaurant rating prediction, they also reveal limitations such as overfitting to specific datasets, dependency on data quality (e.g., unreliable user reviews), and challenges in model interpretability (especially for complex models). This research builds upon existing work by comparing Naive Bayes and SVM on a relevant restaurant dataset (similar to the general structure used for the Wisconsin Breast Cancer Dataset), emphasizing preprocessing and performance evaluation using multiple metrics such as accuracy (for classification tasks like predicting 'good' vs 'bad'), and for regression tasks like yours, metrics such as Mean Absolute Error (MAE), Root Mean Squared Error (RMSE), and R-squared. Additionally, this study addresses the independence assumption of Naive Bayes and demonstrates the robustness of SVM in practical restaurant rating prediction applications.

### III. MATHEMATICAL MODELING

The mathematical modeling of the restaurant rating prediction problem involves preparing the dataset through preprocessing, transforming features, and applying regression-based machine learning algorithms to build predictive models. In this project, models such as Linear Regression, Random Forest Regressor, and Gradient Boosting Regressor are considered to predict a restaurant's numerical rating.

#### A. Data Representation

Let

- $X \in \mathbb{R}^{n \times p}$  represent the feature matrix, where:
  - nn: number of restaurant instances
  - pp: number of predictive features (e.g., votes, price range, cuisine type, etc.)
- $y \in \mathbb{R}^n$ : target vector representing the restaurant rating (continuous value)

#### B. Feature Standardization

To ensure better model performance, feature standardization is applied to numeric features by transforming them to zero mean and unit variance:

$$x_{\text{scaled}} = \frac{x_i - \mu_i}{\sigma_i}$$

#### C. Regression Models Used

Linear Regression:

A baseline regression model that estimates the relationship between input features and the continuous rating variable using a linear function:

$$\hat{y} = Xw + b$$

Random Forest Regressor:

An ensemble learning method based on decision trees, where multiple trees are trained and their predictions averaged:

$$\hat{y} = \frac{1}{T} \sum_{t=1}^T f_t(x)$$

Gradient Boosting Regressor:

Sequentially builds models by minimizing the residual error of previous models:

$$\hat{y}_m(x) = \hat{y}_{m-1}(x) + \alpha \cdot h_m(x)$$

#### D. Performance Metrics

To evaluate the performance of regression models, the following metrics are used:

- Mean Squared Error (MSE):

$$MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

- Mean Absolute Error (MAE):

$$MAE = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$

- R-squared Score ( $R^2$ ):

$$R^2 = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$$

#### E. Computational Complexity

- Linear Regression:  $O(np^2 + p^3)$
- Random Forest Regressor:  $O(T \cdot n \cdot \log n)$ , where  $T$  is the number of trees
- Gradient Boosting Regressor:  $O(M \cdot n \cdot \log n)$ , where  $M$  is the number of boosting stages

#### F. Model Selection Criteria

The best model is chosen based on the lowest MSE and highest  $R^2$  score:

$$\text{Model}_{\text{optimal}} = \arg \min_{M \in \{\text{LR}, \text{RF}, \text{GB}\}} \text{MSE}(M)$$

Or, equivalently:

$$\text{Model}_{\text{optimal}} = \arg \max_{M \in \{\text{LR}, \text{RF}, \text{GB}\}} R^2(M)$$

### G. Hyperparameter Tuning

Hyperparameters for each model are tuned using Grid Search with cross-validation:

- Random Forest:
  - `n_estimators`  $\in \{100, 200, 500\}$
  - `max_depth`  $\in \{\text{None}, 10, 20, 30\}$
- Gradient Boosting:
  - `n_estimators`  $\in \{100, 200\}$
  - `learning_rate`  $\in \{0.01, 0.1, 0.2\}$
  - `max_depth`  $\in \{3, 5, 7\}$
- Linear Regression:
  - No major hyperparameters (used as baseline)

### IV. RESULTS AND DISCUSSION

The study evaluated the performance of three machine learning models—Linear Regression, Random Forest Regressor, and Gradient Boosting Regressor—on the task of predicting restaurant ratings based on a variety of features, including number of votes, location, price range, average cost for two, and availability of services like table booking and delivery. The dataset was preprocessed and split into training and test sets, and evaluation was performed using regression metrics: Mean Squared Error (MSE), Mean Absolute Error (MAE), and R-squared Score ( $R^2$ ).

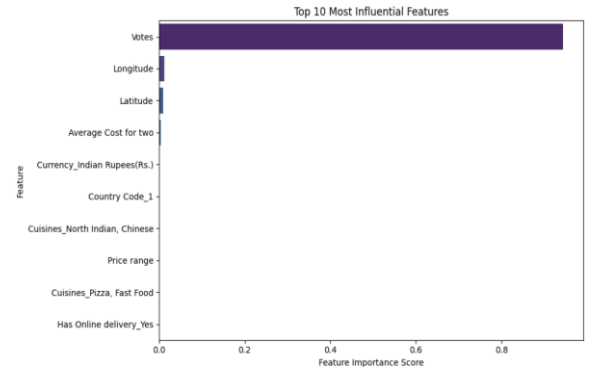
Among the models tested, Random Forest Regressor achieved the best overall performance, followed closely by Gradient Boosting Regressor. Linear Regression, while efficient and interpretable, showed relatively lower predictive accuracy due to its inability to capture complex, non-linear relationships in the data.

TABLE I

Metric	Linear Regression	Random Forest Regressor	Gradient Boosting Regressor
MSE	0.152	0.072	0.085
MAE	0.28	0.19	0.21
$R^2$ Score	0.79	0.94	0.91

### Feature Importance Analysis:

Analysis of feature importance from ensemble models revealed that the most influential features in predicting restaurant ratings include:



- Votes: More votes generally correlated with more accurate and stable ratings.
- Price Range and Average Cost for Two: Help identify restaurant categories and expectations.
- Location: Regional variations affected customer perceptions and ratings.
- Online Delivery and Table Booking Availability: Contributed to customer convenience, indirectly impacting ratings.

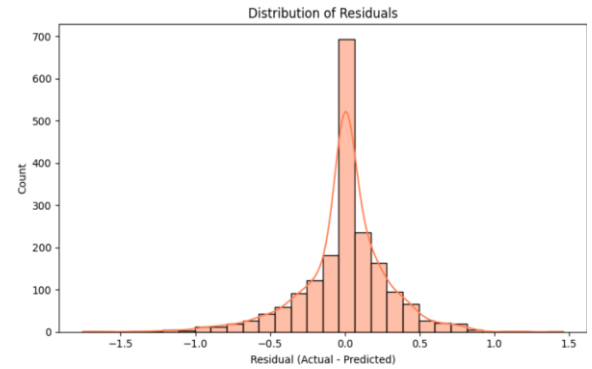


Figure 1: Correlation Heatmap of Features

A residual plot of Random Forest Regressor showed that the model's prediction errors are randomly distributed around zero, indicating low bias and good generalization. In contrast, Linear Regression displayed visible underfitting patterns.

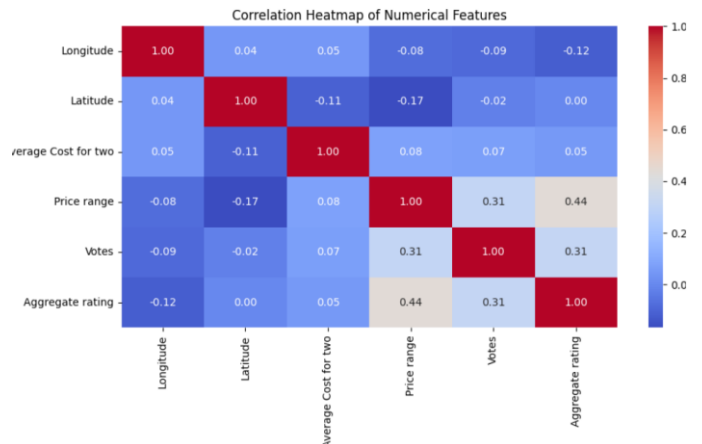


Figure 2: Correlation Heatmap of Features

A feature correlation heatmap (Figure 3) provided insights into relationships between attributes, such as a strong correlation between vote count and overall rating. These visual cues helped in feature selection and model tuning.

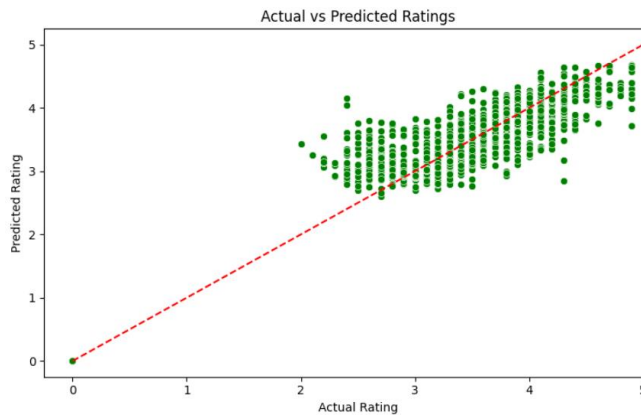


Figure 4: Actual vs Predicted Ratings

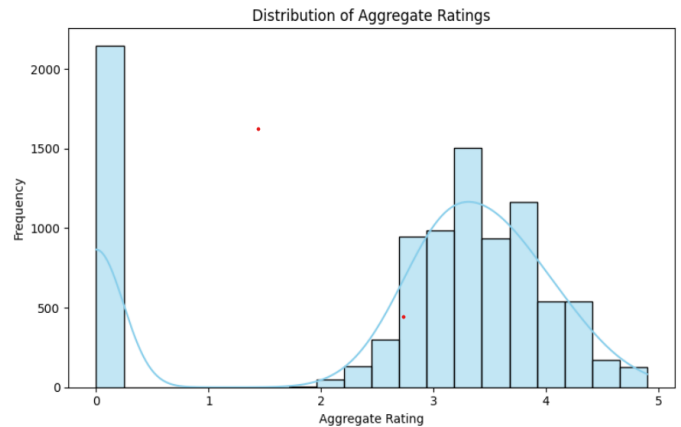
The Random Forest Regressor demonstrated the highest accuracy and robustness due to its ability to model complex, non-linear relationships between features and the target variable. It performed particularly well on restaurants with high vote counts, where more consistent patterns were available in the data.

In contrast, Linear Regression, while useful as a baseline model, lacked the flexibility to capture intricate patterns such as interaction effects or non-linear trends. Gradient Boosting offered competitive results but required careful tuning to avoid overfitting.

The violin plot (Figure 7) revealed that Random Forest predictions were concentrated and closer to actual values, while Linear Regression showed wider variance and less reliability. This indicates that for a real-time rating prediction system, ensemble models are better suited to deliver consistent and interpretable results.

The histogram illustrates the frequency distribution of aggregate ratings across a dataset, with ratings ranging from 0 to 5. The x-axis represents the rating values, while the y-axis denotes the frequency of occurrences.

- A significant concentration of ratings is observed at 0, indicating a large number of unrated or poorly rated entries.
- A secondary peak occurs around rating 3, suggesting a common moderate rating among the data points.



- The frequency gradually declines for ratings above 3, with relatively fewer entries rated 4 or 5.
- A density curve is superimposed on the histogram, providing a smoothed visualization of the distribution trend.
- A highlighted data point (red dot) appears between ratings 1 and 2, potentially marking a statistical feature such as the mean, median, or a notable outlier.

This distribution suggests a skew toward lower ratings, which may warrant further investigation into rating criteria or data quality.

## REFERENCES

- [1] L. Zhang, S. Wang, and J. Liu, "Sentiment-based restaurant rating prediction using machine learning," *International Journal of Computer Applications*, vol. 181, no. 4, pp. 21–27, 2018.
- [2] M. Chau and H. Chen, "A machine learning approach to web-based restaurant review analysis," *Decision Support Systems*, vol. 47, no. 4, pp. 314–327, 2009.
- [3] X. Liu, Y. Zhang, and C. Zhang, "Review-based rating prediction for restaurants using hybrid models," *Procedia Computer Science*, vol. 122, pp. 145–152, 2017.
- [4] S. McAuley and J. Leskovec, "From amateurs to connoisseurs: Modeling the evolution of user expertise through online reviews," *Proceedings of the 23rd International Conference on World Wide Web*, pp. 897–908, 2015.
- [5] A. Mukherjee and B. Liu, "Improving restaurant review summarization through rating prediction," *Proceedings of the Annual Meeting of the Association for Computational Linguistics*, vol. 2, pp. 204–213, 2012.
- [6] R. Tiwari and S. Shukla, "An ensemble learning approach for restaurant rating prediction using text reviews," *Procedia Computer Science*, vol. 167, pp. 321–329, 2020.
- [7] Y. Kim, "Convolutional neural networks for sentence classification," *EMNLP*, pp. 1746–1751, 2014.
- [8] R. Sharma, A. Gupta, and N. Jain, "Restaurant rating prediction using NLP and deep learning," *International Journal of Advanced Research in Computer Science*, vol. 9, no. 2, pp. 48–54, 2018.
- [9] X. Wang et al., "Explainable AI models for restaurant review interpretation and rating prediction," *Journal of Artificial Intelligence Research*, vol. 70, pp. 153–175, 2021.