

# Restaurant Data Analysis: Level 3 Report

This detailed report delves into a comprehensive analysis of a restaurant dataset, focusing on three critical tasks: predictive modeling of restaurant ratings, analyzing customer preferences, and visualizing data relationships to uncover valuable insights. Each task offers a distinct contribution to understanding how various factors influence customer satisfaction and how restaurants can optimize their services.

## Task 1: Predictive Modeling of Restaurant Ratings

The primary objective of this task was to predict restaurant ratings based on various features using multiple machine learning algorithms. The goal was to identify the most accurate model for predicting the aggregate ratings and to understand which factors influence these ratings.

### Modeling Techniques Applied:

#### 1. Linear regression:

- Used as a baseline model to assume a direct linear relationship between features (such as price range, number of votes, cuisine type) and the target rating. While linear regression provides straightforward interpretability, it might fail to capture more complex patterns in the dataset.

#### 2. Ridge Regression:

- A variation of linear regression that adds regularization, making the model more robust against multicollinearity (high correlation between features). Regularization prevents the model from becoming overly complex and overfitting the training data.

#### 3. Lasso Regression:

- Similar to Ridge, but it uses L1 regularization, leading to feature selection. This model is valuable for identifying the most important features influencing ratings by penalizing less impactful ones.

#### 4. Decision Tree Regression:

- A tree-like structure where decisions are made based on feature values, allowing the model to capture non-linear relationships between variables. It can, however, be prone to overfitting if not properly tuned.

#### 5. Random Forest Regression:

- A robust model that aggregates multiple decision trees and averages their predictions, reducing the risk of overfitting and improving accuracy.
- 6. **AdaBoost Regression:**
  - A boosting algorithm that combines weak learners (typically shallow decision trees) to improve performance by focusing on instances that previous models misclassified.
- 7. **Gradient Boosting Regression:**
  - A sequential tree-building model that improves upon previous models by correcting errors iteratively. It provides high accuracy but is computationally intensive.
- 8. **Support Vector Machine (SVM) Regression:**
  - This model tries to find a hyperplane that separates the data points optimally. While effective in classification, its regression performance might be limited by its complexity.
- 9. **K-Nearest Neighbors Regression (KNN):**
  - KNN predicts ratings based on the average of the nearest neighbors, making it sensitive to noise. While simple, it may not generalize well.
- 10. **XGBoost Regression:**
  - An optimized gradient boosting model that has gained popularity for its speed and performance, especially on large datasets. It builds trees iteratively to minimize errors and provides robust predictions.

### Key Findings:

- **Gradient Boosting Regression** was the best-performing model, achieving the lowest **mean squared error (MSE)**, highest  **$R^2$  score**, and lowest **mean absolute error (MAE)**. Its ability to iteratively correct mistakes made by previous trees allowed it to capture complex relationships between restaurant features and customer ratings.
- Other models like **Random Forest** and **XGBoost** also performed well, though they were slightly less accurate than Gradient Boosting.

### Insights:

- Gradient boosting's success underscores the importance of sequential learning, where each model iterates on the weaknesses of the previous one. This model is particularly effective for understanding intricate patterns in the data, such as the combined influence of price, votes, and cuisine on ratings.

## Task 2: Customer Preference Analysis

The second task focused on understanding customer preferences by analyzing cuisine choices, votes, and ratings to uncover patterns in customer satisfaction.

### Data Preparation:

- The dataset was cleaned to handle missing values, and rows containing multiple cuisines were exploded into separate entries to provide a more granular analysis of customer preferences by cuisine.

### Cuisine Vote Aggregation:

- Total votes were calculated for each cuisine type, providing insights into the most popular and highly rated cuisines.

### Key Findings:

#### 1. Most Popular Cuisines:

- **North Indian, Chinese, Italian, Continental,** and **fast food** were the top cuisines, receiving the highest number of votes.

#### 2. High-rated Cuisines:

- **World Cuisine, Western,** and **vegetarian** cuisines received some of the highest average ratings, indicating higher customer satisfaction with these options.

#### 3. Weak Correlation:

- A weak negative correlation was observed between the number of votes and the average ratings. This finding suggests that while some cuisines are more popular (as reflected by the number of votes), they do not necessarily receive the highest ratings. This may be due to factors such as varying customer expectations or quality consistency.

### Insights and implications:

- **Popular Cuisines:** Focusing on customer-preferred cuisines can lead to more satisfied customers, though restaurants should also consider maintaining high-quality standards to ensure that popularity does not overshadow satisfaction.
- **High-rated but less popular cuisines:** There is a market opportunity to promote high-rated but less popular cuisines (e.g., Western and vegetarian), which could appeal to niche audiences.

### Task 3: Data Visualization

The third task aimed to visually explore relationships between different features of the restaurant dataset, with a focus on uncovering patterns that might not be immediately apparent through numerical analysis alone.

#### Visualizations and Key Findings:

##### 1. Rating Distribution:

- A histogram of ratings revealed that most restaurant ratings were clustered around certain values, typically in the 3-4 star range, indicating that customers are generally satisfied, but there may be room for improvement in quality or service.

##### 2. Cuisine Ratings Comparison:

- A bar plot comparing average ratings by cuisine type showed that **World Cuisine** and **Western Cuisine** consistently received higher ratings compared to more popular options like fast food or North Indian. This suggests that less popular cuisines may provide a higher-quality experience.

##### 3. City Ratings Comparison:

- Regional trends were explored through city-wise comparisons of average ratings. Certain cities consistently showed higher ratings, which could be indicative of better restaurant quality or stronger culinary traditions in those areas.

##### 4. Feature Relationship Analysis:

- **Votes vs. Ratings:** A scatter plot revealed that while more votes generally indicated more popularity, it did not necessarily correspond to higher ratings. This suggests that highly-reviewed restaurants are not always the best in terms of customer satisfaction.
- **Price Range vs. Ratings:** Another scatter plot showed that higher price ranges generally corresponded to better ratings, though the relationship was not entirely linear. This finding suggests that customers may associate higher prices with better quality but still evaluate other factors like service and ambiance.

#### Insights:

- **Rating Distribution:** The concentration of ratings in the mid-range indicates that while customers are generally satisfied, there is an opportunity for improvement in areas like service or food quality to push ratings higher.

- **City and Cuisine Preferences:** Restaurants can tailor their offerings to regional preferences and consider promoting highly-rated cuisines to differentiate themselves in the market.

## Conclusion and Strategic Recommendations

The analysis presented in this report provides a holistic view of restaurant performance, customer preferences, and factors influencing ratings. The insights gained from predictive modeling, customer preference analysis, and data visualization highlight several key areas for improvement:

1. **Leverage High-rated Cuisines:** Restaurants should focus on expanding their menu offerings to include high-rated but less popular cuisines, such as **Western** or **vegetarian** dishes, to tap into underserved customer segments.
2. **Optimize Pricing Strategies:** The relationship between price and ratings suggests that restaurants should carefully balance their pricing strategies with the expected level of quality and service, ensuring that higher prices are justified by a superior dining experience.
3. **Target regional preferences:** By understanding city-wise rating trends, restaurants can tailor their marketing and operational strategies to regional preferences, enhancing their appeal in specific markets.
4. **Promote Quality Over Quantity:** Although popularity (votes) can bring customers in, maintaining high service and food quality is essential for driving better ratings and long-term customer loyalty.