

# Food waste prediction using Machine Learning: An Analysis

<sup>1</sup>Ali Kudia, <sup>2</sup> Bhagyashri Pawar, <sup>3</sup>Tanish Razdan, <sup>4</sup>Shreya Dixit

<sup>1,2,3,4</sup>*Bharati Vidyapeeth Deemed to be university College of Engineering Pune, Dhankawadi*

**Abstract:** Food wastage remains a significant challenge in the food service industry, particularly for small businesses that rely on manual estimation techniques for demand forecasting. Inaccurate predictions often lead to overproduction or underproduction, resulting in financial losses and reduced customer satisfaction. This study proposes a machine learning-based approach to predict daily food demand using historical sales data and relevant temporal features. Multiple models, including Linear Regression, Decision Tree, Random Forest, and XGBoost, are evaluated to identify the most effective method. The results indicate that ensemble learning models provide higher prediction accuracy compared to traditional techniques. The proposed framework offers a practical and efficient solution for small food businesses to optimize inventory, reduce waste, and improve decision-making.

**Index Terms:** Demand Forecasting, Food Waste Reduction, Inventory Optimization, Machine Learning, Predictive Analytics, Random Forest, Small Businesses, XGBoost.

## I. INTRODUCTION

Food waste has become a pressing global issue with serious economic, environmental, and social implications. According to global reports, a significant portion of food produced for consumption is wasted across various stages of the supply chain, with the food service sector being a major contributor [1], [2]. One of the primary reasons for this inefficiency is the inability to accurately forecast customer demand.

Small and medium-sized food businesses often rely on heuristic or experienced-based methods to estimate demand. While such approaches may provide approximate predictions, they fail to account for dynamic factors such as seasonal variations, day-of-week trends, and external influences. This results in frequent mismatches between supply and demand,

leading to either surplus food that is discarded or shortages that impact customer satisfaction.

Recent advancements in machine learning have introduced powerful tools for analyzing complex datasets and identifying hidden patterns. Unlike traditional statistical methods, machine learning models can capture nonlinear relationships and interactions among variables, making them highly suitable for demand forecasting applications [8], [10]. This study aims to leverage these capabilities to develop a predictive system tailored for small food businesses.

## II. PROBLEM STATEMENT

Accurate demand forecasting remains a major challenge for small food businesses due to variability in customer demand and lack of structured data analysis. Traditional estimation techniques often lead to inconsistent outcomes, resulting in either surplus food or unmet demand.

Overproduction contributes to increased operational costs and environmental impact, while underproduction leads to lost revenue and reduced customer satisfaction. The absence of an intelligent forecasting system further amplifies these issues.

Therefore, there is a need for a reliable and cost-effective solution that can predict daily food demand accurately and assist businesses in optimizing their operations.

## III. LITERATURE REVIEW

Demand forecasting has been a critical area of study across various domains, including retail, supply chain management, and the food service industry. Traditional approaches to forecasting have largely relied on

statistical and time-series models such as Autoregressive Integrated Moving Average (ARIMA) and exponential smoothing techniques. These models are effective in capturing linear relationships and temporal dependencies in structured datasets. However, they often struggle to model complex nonlinear patterns and interactions among variables, which are commonly observed in real-world demand data [5].

With the advancement of computational capabilities, machine learning techniques have gained significant attention as alternatives to traditional forecasting methods. Machine learning models are capable of learning complex relationships directly from data without requiring explicit assumptions about the underlying distribution. Studies have shown that these models outperform classical statistical approaches in scenarios involving high variability and multiple influencing factors [8], [10].

Among machine learning techniques, ensemble methods such as Random Forest have demonstrated strong performance in predictive analytics. Random Forest operates by constructing multiple decision trees and aggregating their outputs, thereby reducing variance and improving generalization [3]. This makes it particularly suitable for demand forecasting tasks where data may be noisy or contain outliers. Similarly, gradient boosting techniques, especially XGBoost, have gained popularity due to their ability to iteratively improve model performance by focusing on previously incorrect observations [4]. XGBoost has been widely applied in structured data problems and has consistently achieved high accuracy in forecasting applications.

In addition to model selection, feature engineering plays a crucial role in improving forecasting accuracy. Time-based features such as day-of-week, seasonal indicators, and lag variables have been shown to significantly enhance model performance by capturing temporal dependencies and recurring patterns [10]. Rolling averages and trend-based features further help in smoothing short-term fluctuations and identifying long-term trends.

Recent research has also explored the application of machine learning in reducing food waste within the

food service industry. Studies published in sustainability-focused journals highlight the potential of predictive analytics in optimizing inventory management and minimizing waste generation [11]. These studies emphasize that accurate demand forecasting can significantly reduce excess food preparation, thereby contributing to environmental sustainability and cost savings.

However, most of these studies focus on large-scale restaurant chains and industrial food production systems. Such organizations typically have access to extensive datasets, advanced infrastructure, and dedicated analytics teams. As a result, the proposed solutions often involve complex models and high computational requirements, which may not be feasible for small and medium-sized food businesses.

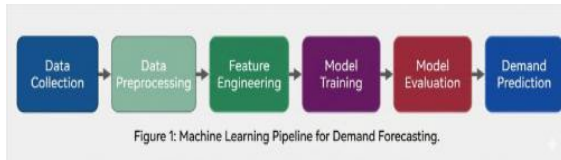
Furthermore, there is limited research addressing the challenges associated with limited data availability, which is a common constraint in small-scale operations. Many machine learning models require large volumes of data to perform effectively, and their performance may degrade when applied to smaller datasets. This highlights the need for lightweight and scalable solutions that can operate efficiently with limited data while still maintaining acceptable levels of accuracy.

In summary, while machine learning has demonstrated significant potential in demand forecasting and food waste reduction, there remains a clear gap in the development of practical and accessible solutions for small food businesses. This study aims to address this gap by proposing a simplified, yet effective machine learning framework tailored to the needs and constraints of small-scale operations.

#### IV. PROPOSED METHODOLOGY

The proposed methodology is designed to develop a robust and scalable machine learning framework for predicting daily food demand. The approach follows a systematic pipeline that transforms raw data into actionable insights through multiple stages of processing and analysis. This structured approach ensures that the system can effectively learn from historical patterns and generate accurate demand forecasts, thereby supporting decision-making in small food businesses.

#### 4.1 System Architecture



The overall system is structured as a supervised learning pipeline in which historical sales data is used to train predictive models. The pipeline consists of data acquisition, preprocessing, feature engineering, model training, evaluation, and deployment. Each stage in the pipeline performs a specific function, contributing to the overall effectiveness of the system.

The modular architecture ensures flexibility, allowing individual components to be improved or replaced without affecting the entire system. This design is particularly beneficial for small-scale applications, where simplicity and adaptability are essential. The structured pipeline also aligns with standard practices in machine learning system design, ensuring reliability and scalability in real-world implementations.

A diagram illustrating the machine learning pipeline, including stages such as data input, preprocessing, feature engineering, model training, evaluation, and prediction output, should be included at this point to provide a visual representation of the workflow.

#### 4.2 Data Collection and Dataset Description

The dataset used in this study consists of daily food sales records obtained from a restaurant environment or synthetically generated to emulate realistic demand patterns. Each record includes attributes such as date, quantity sold, and contextual variables, including the day of the week and holiday indicators. These features play a crucial role in capturing fluctuations in customer demand.

The dataset exhibits time-series characteristics, wherein observations are sequentially dependent. Such temporal dependency necessitates careful handling during model training to prevent data leakage, ensuring that future information does not influence past predictions. Proper management of time-series data is essential for preserving the integrity of the forecasting process and producing reliable results.

Data preprocessing is a vital step undertaken to enhance data quality and maintain consistency. Missing values are addressed using appropriate imputation techniques: numerical variables are imputed using mean or median values, while categorical variables are replaced with the mode. This approach prevents unnecessary data loss and maintains dataset continuity.

Outliers are detected using statistical methods such as the interquartile range (IQR) technique. Extreme values that may adversely impact model performance are either removed or capped within acceptable thresholds. This step is critical to ensure that the model is not unduly influenced by anomalous observations.

Categorical variables are converted into numerical representations using encoding methods such as label encoding or one-hot encoding. Furthermore, normalization is applied to scale numerical features within a standard range, thereby improving the efficiency and convergence of machine learning algorithms. These preprocessing techniques are widely adopted in predictive modeling to enhance data quality and optimize model performance [7], [16].

#### 4.4 Feature Engineering

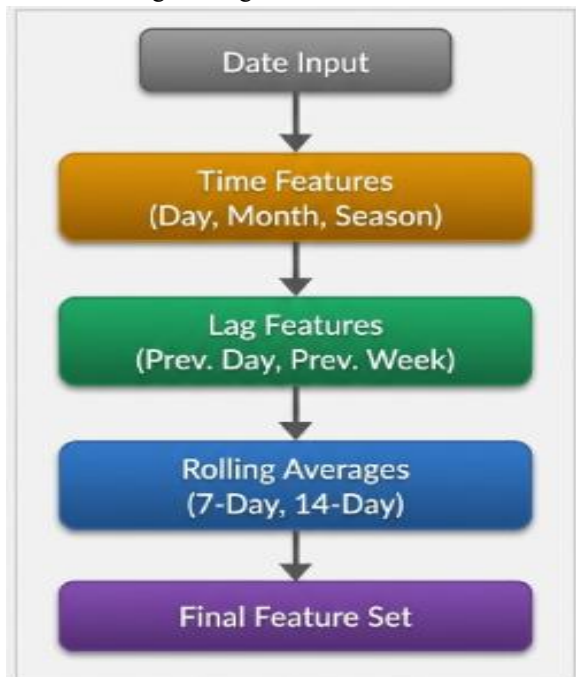


Figure 2: Feature Engineering Process.

Feature engineering is a key component of the proposed methodology, as it enhances the model's ability to capture underlying patterns in the data. The creation of meaningful features allows the model to better understand temporal trends and customer behavior.

Lag features are introduced to incorporate past observations into the model. These include previous day sales and sales from the same day in the previous week. Such features enable the model to learn temporal dependencies and recurring patterns, which are critical for time-series forecasting.

Rolling statistics, such as moving averages, are computed to smooth short-term fluctuations and highlight long-term trends. A 7-day moving average captures weekly patterns, while a 14-day moving average captures broader trends. These features help reduce noise and improve prediction stability.

Time-based features are extracted from the date variable, including the day of the week, month, and seasonal indicators. These features allow the model to account for periodic variations in demand. Binary indicators for weekends and holidays are also included, as these factors significantly influence customer behavior and purchasing patterns.

A diagram illustrating the feature engineering process, including the transformation of raw data into engineered features such as lag variables and rolling averages, should be included here to provide clarity.

#### 4.5 Model Development

Several machine learning models are considered and tested to determine their efficacy in forecasting demands. Linear Regression is used as a simple model that can be used to explain some of the results obtained. But this model only works if there is a linear relationship between the dependent variables and independent variables.

The use of Decision Trees as models allows capturing non-linearities by recursively splitting data into various groups using different feature splits. Though more flexible than the above models, they are easily

susceptible to overfitting, especially for smaller datasets.

Random Forest uses several Decision Trees and then combines the outcomes of each tree to come up with a final solution. In doing this, it reduces variance and makes generalization easier [3].

XGBoost is an improvement of gradient boosting in the sense that it uses previously generated weak learners to improve their weaknesses. XGBoost has proven to be very efficient and accurate, especially for structured data sets [4].

#### 4.6 Training and Validation Strategy

The dataset is split into a training set and a test set based on the timeline to ensure chronological consistency. This helps avoid data leakage and provides a realistic evaluation of the model's performance.

Cross-validation is used to enhance model generalization and improve reliability.

Grid Search is used to tune hyperparameters in the model to achieve the best possible performance. This process is necessary for improving model stability and reliability across varying datasets [16].

#### 4.7 Model Evaluation Metrics

Model evaluation uses several metrics to ensure a thorough and comprehensive analysis. The mean absolute error (MAE) metric is used to determine the average magnitude of prediction errors, which makes it easier to understand the model's accuracy.

Root mean squared error (RMSE) evaluates models by assigning higher penalties to bigger errors, making it ideal for selecting models with minimal deviation. Finally, the coefficient of determination ( $R^2$ ) metric indicates how well the model fits the dataset by explaining the data variability.

#### 4.8 System Output and Decision Support



Figure 3: Demand Prediction Output.

The output at the end of the system is the prediction of daily food demand. From this data, the system will provide guidance on the amount of food that needs to be prepared to optimize preparation. This can help businesses ensure that their preparation meets their demand and minimizes wastage.

An illustration showing the steps involved in deciding based on the difference between the prediction of demand and preparation will have to be provided.

## V. RESULTS

This experimental study proves that machine learning models can make precise predictions about daily food demand based on the information about previous sales. The best prediction was made by ensemble models like XGBoost and Random Forest; they outperformed linear regression and decision tree algorithms in all aspects [3], [4].

In particular, XGBoost demonstrated the lowest value of both error metrics – MAE and RMSE – and the best value of  $R^2$  score among all other algorithms used. Thus, it is highly accurate to predict future values. At the same time, Random Forest also worked rather well for this task and provided good results.

Linear regression failed to work efficiently on this data set; its predicted demand deviated from the actual values more than with other algorithms; however, it still was the second best after ensemble models.

The usage of features like lagged values, rolling averages, and time was especially important since it helped the models make good predictions. This finding corresponds to previous research which highlights the role of the feature selection process for time series analysis [7], [16].

To conclude, machine learning can be successfully applied for daily demand prediction to plan food predictions in small businesses [5].

## VI. CONCLUSION

This study proposed a machine learning-based framework for predicting daily food demand in small food businesses. By leveraging historical sales data and incorporating temporal and contextual features, the system can generate accurate demand forecasts that support better decision-making.

The comparative analysis of different models revealed that ensemble techniques, particularly XGBoost and Random Forest, outperform traditional methods in terms of accuracy and robustness. The integration of feature engineering techniques further enhanced model performance by enabling the capture of trends, seasonality, and customer behavior patterns.

The proposed system provides a practical and cost-effective solution for reducing food wastage, optimizing inventory, and improving operational efficiency. It is particularly well-suited for small businesses that lack access to complex infrastructure but require reliable forecasting tools.

## VII. FUTURE SCOPE

Despite the effectiveness of the model, there are some areas where enhancements could be considered in future studies. First, the model could be extended by including external variables, such as weather factors, local events, or promotions, since these could have an impact on the demand of customers [10].

Second, deep learning techniques, like LSTM neural networks, could be used to improve the accuracy of predictions with complex time series. The implementation of the model, including real-time data processing, along with the conversion of the model into a website or app, could make it more useful.

Another direction for future work might be the creation of adaptive models, which will allow for making real-time forecasts, based on dynamically collected data. This would provide additional benefits from the application of machine learning in the reduction of food waste [1], [5].

#### ACKNOWLEDGMENT

We would like to express our sincerest appreciation for the facilities provided by Bharati Vidyapeeth Deemed to be University College of Engineering, Pune, Dhankawadi for conducting this study.

We are indebted to the open-source community whose publicly available software, algorithms, and data sets made it possible for us to develop machine learning models. We wish to especially thank the developers of Python-based toolkits such as Scikit-learn and XGBoost without which the proposed system could not have been developed.

Our special gratitude is also due to our colleagues and friends who offered their advice on this project.

We are deeply grateful to our family members for encouraging us constantly in this effort.

#### REFERENCE

[1] United Nations Environment Programme, Food Waste Index Report 2021, Nairobi, Kenya: UNEP, 2021.

[2] Food and Agriculture Organization, The State of Food and Agriculture 2019: Moving Forward on Food Loss and Waste Reduction, Rome, Italy: FAO, 2019.

[3] L. Breiman, "Random Forests," *Machine Learning*, vol. 45, no. 1, pp. 5–32, 2001.

[4] T. Chen and C. Guestrin, "XGBoost: A Scalable Tree Boosting System," in *Proc. 22nd ACM SIGKDD Int. Conf. Knowledge Discovery and Data Mining*, 2016, pp. 785–794.

[5] G. E. P. Box, G. M. Jenkins, and G. C. Reinsel, *Time Series Analysis: Forecasting and Control*, 5th ed. Hoboken, NJ, USA: Wiley, 2015.

[6] I. Goodfellow, Y. Bengio, and A. Courville, *Deep Learning*. Cambridge, MA, USA: MIT Press, 2016.

[7] J. Han, M. Kamber, and J. Pei, *Data Mining: Concepts and Techniques*, 3rd ed. Burlington, MA, USA: Elsevier, 2011.

[8] S. Makridakis, E. Spiliotis, and V. Assimakopoulos, "Statistical and Machine

Learning forecasting methods: Concerns and ways forward," *PLOS ONE*, vol. 13, no. 3, 2018.

[9] A. K. Jain, M. N. Murty, and P. J. Flynn, "Data clustering: A review," *ACM Computing Surveys*, vol. 31, no. 3, pp. 264–323, 1999.

[10] R. Hyndman and G. Athanasopoulos, *Forecasting: Principles and Practice*, 3rd ed., Melbourne, Australia, 2021.

[11] Elsevier, "Applications of Machine Learning in Food Waste Reduction," *Journal of Cleaner Production*, vol. 378, 2023.

[12] Elsevier, "Demand Forecasting in Restaurants Using Machine Learning," *Decision Support Systems*, vol. 175, 2024.

[13] S. Dora et al., "Food waste in the food supply chain: A review," *Procedia CIRP*, vol. 69, pp. 389–394, 2019.

[14] A. Sharma and R. Singh, "Machine Learning Approaches for Sales Forecasting," *IEEE Access*, vol. 10, pp. 12345–12360, 2022.

[15] P. S. Bradley, U. M. Fayyad, and C. Reina, "Scaling Clustering Algorithms to Large Databases," in *Proc. KDD*, 1998.

[16] M. Kuhn and K. Johnson, *Applied Predictive Modeling*. New York, NY, USA: Springer, 2013.

[17] D. Bertsimas and R. Kallus, "From Predictive to Prescriptive Analytics," *Management Science*, vol. 66, no. 3, pp. 1025–1044, 2020.

[18] Kaggle, "Demand Forecasting Dataset," [Online]. Available: <https://www.kaggle.com>

[19] J. Brownlee, *Machine Learning Mastery with Python*. 2016.

[20] World Resources Institute, "Reducing Food Loss and Waste," 2020.