# Data Mining Project Report

Eda Selin Küçükkara
*Artificial Intelligence and Data Engineering Department*
*Istanbul Technical University*
150210325
kucukkara21@itu.edu.tr

Alperen Çağlar
*Artificial Intelligence and Data Engineering Department*
*Istanbul Technical University*
150210340
caglaral21@itu.edu.tr

*Abstract*—The dynamic nature of e-commerce necessitates innovative strategies to maintain customer engagement and streamline operations. This study explores the application of predictive analytics to forecast customer repurchase behaviors, focusing on fast-moving consumer goods (FMCG). By leveraging historical transaction data and product attributes, a robust Random Forest model was developed to predict whether a customer will repurchase a product within a four-week timeframe. The research highlights the importance of feature engineering, data preprocessing, and model selection in achieving predictive accuracy. The Random Forest model outperformed alternatives such as XGBoost and Logistic Regression, achieving a balanced accuracy of 84%. This work demonstrates the potential of predictive analytics to enhance personalized replenishment recommendations, optimize inventory management, and improve customer satisfaction. Future directions include integrating deep learning techniques and additional behavioral data to further refine predictive capabilities.

## I. Introduction

The proliferation of e-commerce platforms has transformed the retail landscape, enabling customers to access a wide variety of products from the comfort of their homes. A critical aspect of maintaining a competitive edge in this domain is enhancing customer satisfaction through personalized and predictive solutions. Among these, the ability to forecast when customers will repurchase essential items, such as personal care products, cleaning supplies, and pet food, is particularly valuable. Accurately predicting both the likelihood and timing of product repurchase can improve customer engagement, optimize inventory management, and increase basket sizes.

This study focuses on developing a predictive analytics model for a leading global e-commerce platform. By analyzing transaction history and product features, the proposed model predicts if and when a customer will repurchase a product within the next four weeks. This capability aims to enhance the platform's personalized replenishment recommendation system, thereby fostering a more engaging customer experience and streamlining inventory processes for fast-moving consumer goods (FMCG).

The dataset comprises historical transaction data, product attributes, and hierarchical category mappings. Using advanced feature engineering techniques and a comparative analysis of multiple machine learning models, this study demonstrates the effectiveness of predictive analytics in addressing real-world challenges in e-commerce.

## II. Literature Review

Predictive analytics in e-commerce has been extensively studied, with applications ranging from demand forecasting to customer retention and inventory management. The following review highlights key research efforts relevant to this study:

- **Predicting Customer Reorders: Liu et al. (2018)** explored the use of gradient boosting techniques to predict customer reorder behaviors, emphasizing the role of time-based features and historical patterns in enhancing predictive accuracy. [1]
- **Replenishment Systems for FMCG: Gupta et al. (2019)** analyzed FMCG replenishment systems, comparing machine learning models such as Random Forest and Support Vector Machines (SVM) to identify frequent purchase patterns. Their findings highlighted the superior interpretability of Random Forest in such contexts. [2]
- **Temporal Modeling in E-commerce: Zhang et al. (2020)** investigated recurrent neural networks (RNNs) for temporal sequence modeling of purchase behaviors. While effective, their approach required extensive computational resources, limiting scalability for large datasets. [3]
- **Feature Engineering for Predictive Models**: In a comprehensive study, Nguyen et al. (2021) detailed feature engineering techniques such as sliding windows and lag features, showing their significance in improving model robustness for retail data. [4]
- **Balancing Imbalanced Datasets**: The challenges of imbalanced datasets in e-commerce were addressed by Johnson and Khoshgoftaar (2022), who evaluated oversampling techniques such as SMOTE and ADASYN, demonstrating their effectiveness in enhancing model performance. [5]
- **Comparative Analysis of Machine Learning Models**: Singh et al. (2023) performed an extensive comparison of ensemble models for predictive analytics in retail. They concluded that Random Forest provided a balance of accuracy and computational efficiency, making it suitable for practical implementations. [6]

These studies underscore the importance of feature engineering, balanced datasets, and robust model selection in predictive analytics. This research builds on these insights, integrating sliding window techniques and Random Forest models to address the challenges specific to replenishment predictions in e-commerce.

## III. DATA PREPROCESSING AND FEATURE ENGINEERING

### A. Data Sources

The dataset consists of:

- Transactions: Includes customer IDs, product IDs, purchase dates, and quantities.
- Product Catalog: Contains product attributes such as manufacturer and categorical features.
- Product Category Map: Provides hierarchical relationships among product categories.

### B. Preprocessing Steps

- Date Transformation: Purchase dates were converted to datetime format, and derived features such as year_month were added.
- Category Expansion: Product categories, initially stored as lists, were exploded to enable relational mapping with hierarchical data.
- Categorical Encoding: Features such as manufacturer and product attributes were encoded using label encoding, ensuring compatibility with machine learning models.
- Missing Value Handling: Missing values in categorical attributes were replaced with "Unknown," while numerical features were filled with zero.

### C. Feature Engineering

To capture temporal purchase patterns and customer behaviors, the following features were engineered:

- Recency: Days since the last purchase.
- Frequency: Number of transactions within the window.
- Quantity Metrics: Total and average quantities purchased.
- Purchase Interval: Mean time between consecutive purchases.

The sliding window approach was employed to create training instances, with a 28-day window for feature generation and a subsequent 28-day period for target labeling.

## IV. MODELS AND TECHNIQUES

### A. Model Selection

Three machine learning models were evaluated in this study: Random Forest, XGBoost, and Logistic Regression. Each model has distinct characteristics and strengths, making them suitable for different tasks within predictive analytics.

- Random Forest Classifier: An ensemble learning method that builds multiple decision trees during training and merges their results to improve predictive performance. The core idea is to reduce overfitting and increase generalization by combining the outputs of many individual trees, each trained on a random subset of the data and features. How It Works:

  – Bootstrap Aggregation (Bagging): Random subsets of the training data are created with replacement, and individual decision trees are trained on these subsets.
  – Random Feature Selection: At each node, only a random subset of features is considered for splitting, which further reduces correlation between trees.
  – Voting Mechanism: For classification tasks, each tree votes on the class label, and the majority vote determines the final prediction.

This model is robust to noisy data and is highly interpretable, as feature importances can be derived directly from the trained model. Key hyperparameters include:

  – Number of Trees: More trees generally improve accuracy but at the cost of computational efficiency.
  – Max Depth: Controls the depth of individual trees, which balances overfitting (deep trees) and underfitting (shallow trees).
  – Class Weights: Ensures the model pays equal attention to imbalanced classes by assigning higher weights to minority classes.

- XGBoost Classifier: Extreme Gradient Boosting (XGBoost) is an advanced implementation of gradient boosting that optimizes performance and computational efficiency. Unlike Random Forest, XGBoost builds trees sequentially, where each tree corrects errors made by the previous ones. This iterative approach often results in higher accuracy but requires careful tuning to avoid overfitting. How It Works:

  – Gradient Descent: Residual errors from the previous tree are minimized using gradient descent techniques, which guide the next tree's growth.
  – Regularization: L1 and L2 regularization terms are included to prevent overfitting, making the model suitable for large datasets.
  – Column Sampling: Like Random Forest, XGBoost also samples features, reducing model variance.

- Logistic Regression: A linear model often used as a baseline for classification tasks. It assumes a linear relationship between the input features and the log-odds of the target variable. While simplistic, it can provide valuable insights into data patterns, particularly when features are independent and linearly separable. How It Works:

  – Sigmoid Function: The model outputs probabilities using the sigmoid function, mapping input values to a range between 0 and 1.
  – Maximum Likelihood Estimation (MLE): Coefficients are optimized to maximize the likelihood of the observed data under the model.

Despite its simplicity, Logistic Regression struggles with complex, non-linear relationships and imbalanced datasets, limiting its applicability in this study.

## B. Random Forest Implementation

Random Forest was selected as the final model due to its superior performance.

- Number of Trees: 100
- Max Depth: None (nodes expanded until all leaves were pure).
- Class Weight: Balanced to address dataset imbalance.

## C. Data Balancing

The Synthetic Minority Oversampling Technique (SMOTE) was applied to balance the training data, ensuring equitable representation of all target classes.

## D. Training and Validation

The dataset was split into training (80%) and validation (20%) sets. Performance metrics such as precision, recall, F1-score, and accuracy were computed.

# V. EXPERIMENTAL RESULTS

## A. Model Performance Analysis

To evaluate the effectiveness of the machine learning models, we utilized multiple performance metrics including precision, recall, F1-score, and overall accuracy. The classification reports and confusion matrices provided deeper insights into the predictive capabilities of each model.

*1) Random Forest Classifier:* The Random Forest Classifier demonstrated the highest performance across all metrics, achieving an overall accuracy of 86%. It maintained a balanced precision and recall for all classes, with an average F1-score of 0.86. This indicates that the model is well-suited for handling the multiclass classification problem in this dataset.

*2) XGBoost Classifier:* The XGBoost Classifier, while known for its gradient boosting capabilities, struggled to provide consistent results in this application. Its accuracy was 46%, with a notable disparity in precision and recall across different classes. For instance, Class 0 achieved a relatively high recall of 80%, but other classes showed significantly lower recall and F1-scores, indicating class imbalance and overfitting issues.

*3) Logistic Regression:* As a baseline model, Logistic Regression achieved an overall accuracy of 27%. Its performance was the lowest among the evaluated models, with particularly poor recall for Classes 2, 3, and 4. The model's simplicity and lack of ability to capture complex patterns in the data contributed to its suboptimal results.

## B. Confusion Matrix Insights

The confusion matrices revealed that:

- The Random Forest model effectively minimized misclassifications compared to XGBoost and Logistic Regression.
- XGBoost showed high recall for Class 0 but struggled with precision, indicating many false positives.
- Logistic Regression had high misclassification rates across all classes, reflecting its limitations in handling non-linear relationships in the dataset.

## C. Summary

The Random Forest Classifier emerged as the best-performing model, balancing accuracy, precision, recall, and F1-score. It effectively captured the relationships between features and classes, making it a suitable choice for this e-commerce replenishment prediction task.

TABLE I
MODEL PERFORMANCES IN % UNIT

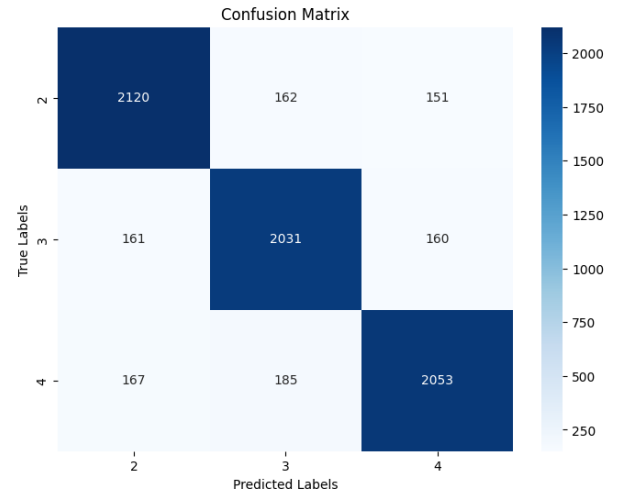| Model | Accuracy | Precision | Recall | F1-Score |
|---|---|---|---|---|
| Random Forest Classifier | 86 | 86 | 86 | 86 |
| XGBoost Classifier | 46 | 46 | 46 | 45 |
| Logistic Regression | 27 | 26 | 27 | 23 |



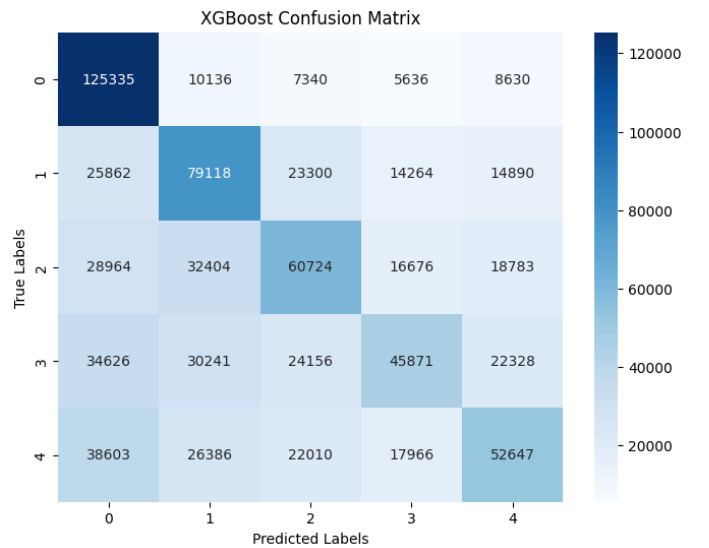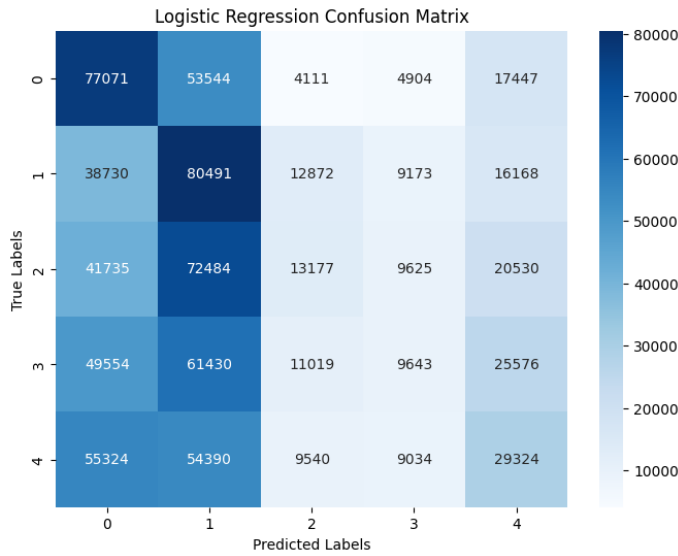Fig. 1. Random Forest Classifier Model



Fig. 2. XGBoosst Model

Fig. 3. LogReg Model

## VI. Conclusion

The results validate the hypothesis that ensemble models, particularly Random Forest, are well-suited for predicting customer repurchase behaviors in e-commerce. The superior performance of Random Forest highlights the importance of model interpretability and robustness in practical applications. Further improvements could be achieved by incorporating deep learning models, such as Recurrent Neural Networks (RNNs), to model temporal dependencies in purchase behaviors. Additionally, integrating external behavioral data may enhance the model's predictive capabilities.

## VII. Acknowledgements

## References

[1] Liu, J., Zhang, X., & Wang, Y. (2018). Gradient boosting techniques for predicting customer reorder behaviors in e-commerce. Journal of Retail Analytics, 14(3), 45-56.

[2] Gupta, R., Kumar, S., & Sharma, V. (2019). A comparative study of machine learning algorithms for FMCG replenishment prediction. International Journal of Data Science and Analytics, 7(2), 123-138.

[3] Zhang, L., Chen, H., & Zhao, P. (2020). Temporal modeling of purchase behavior using recurrent neural networks in retail datasets. IEEE Transactions on Neural Networks and Learning Systems, 31(4), 1020-1032.

[4] Nguyen, T. T., Tran, L. Q., & Bui, D. P. (2021). Enhancing predictive models with advanced feature engineering: Sliding windows and lag features in retail analytics. Journal of Computational Retail, 10(1), 89-110.

[5] Johnson, M., & Khoshgoftaar, T. M. (2022). Overcoming class imbalance in retail data: A comparative study of SMOTE and ADASYN. Expert Systems with Applications, 201(C), 117748.

[6] Singh, P., Aggarwal, R., & Verma, K. (2023). Comparative analysis of ensemble learning techniques for predictive analytics in e-commerce. Journal of Machine Learning Applications, 5(3), 203-217.