#### **Dogo Rangsang Research Journal** ISSN: 2347-7180

#### **SALES PREDICTION**

Siddheswar Sahoo 4<sup>th</sup> Year, Department of CSE, Gandhi Institute for Technology, BPUT, India siddheswar2022@gift.edu.in

Rishav Ranjan 4<sup>th</sup> Year, Department of CSE, Gandhi Institute for Technology, BPUT, India <u>rishavranjan2021@gift.edu.in</u>

<sup>3</sup> Assistant Professor, Department of CSE, Gandhi Institute for Technology, BPUT, India

#### Abstract—

This project Sales Management System uses machine learning to predict retail product sales based on historical data and product-store attributes. By comparing models like Linear Regression, Random Forest, and AdaBoost, it identifies the most accurate predictor. The system aims to improve inventory management, pricing strategies, and decision-making in real-world retail environments.

#### Keywords:

Python , PANDAS

# I. INTRODUCTION

This project leverages machine learning to predict retail product sales using historical data and features like product attributes and store characteristics. It compares models such as Linear Regression, KNN, Decision Tree, Random Forest, and AdaBoost to find the most accurate predictor. Accurate sales forecasting aids inventory management, pricing, and marketing strategies, enhancing profitability and reducing waste.

#### II. LITERATURE REVIEW

Sales prediction is crucial in retail for optimizing inventory, pricing, and marketing strategies. Traditional statistical methods often fall short in capturing complex data relationships. Recent literature highlights the effectiveness of machine learning models, such as Linear Regression, Decision Trees, Random Forests, AdaBoost, and K-Nearest Neighbors, in improving predictive accuracy. Studies emphasize the importance of incorporating both product-level features (e.g., price, visibility) and store-level attributes (e.g., outlet size, location). Tree-based and ensemble models, in particular, handle non-linearity and diverse data well. This project builds on existing research by comparing multiple models to identify the most accurate method for forecasting item outlet sales.

In addition to model selection, data preprocessing and feature engineering play a vital role in enhancing sales prediction accuracy. Literature shows that handling missing values, encoding categorical variables, and normalizing data significantly impact model performance. Feature importance analysis, especially in tree-based models like Random Forest, helps identify key drivers of sales. Studies also suggest that temporal factors such as seasonality and promotional events can improve forecast reliability when integrated into the model. Moreover, cross-validation techniques are widely used to prevent overfitting and ensure generalizability. This project adopts these best practices to develop a robust, real-world sales prediction system for retail applications.

#### III. SYSTEM DESIGN

The system design for the sales prediction model is structured to ensure efficient data processing, model training, and prediction deployment. The system begins with a data collection module that gathers historical sales data, product details, and store-specific attributes. This is followed by a data preprocessing stage, where missing values are handled, categorical variables are encoded, and numerical features are normalized to prepare the data for modeling. The feature engineering module extracts relevant insights from raw data, enhancing model performance. The core of the system includes multiple machine learning algorithms such as Linear Regression, K-Nearest Neighbors, Decision Tree, Random Forest, and AdaBoost, allowing for model comparison and selection based on performance metrics like RMSE and R<sup>2</sup> score. Once the best model is identified, it is integrated into a

# **Dogo Rangsang Research Journal** ISSN : 2347-7180

prediction module that provides sales forecasts for new data inputs. The system also includes a user interface for visualizing predictions and a backend for model updates and maintenance.

# IV. IMPLEMENTATION

The implementation of the sales prediction system involves a structured workflow that integrates data handling, model development, evaluation, and deployment. The project begins with importing and exploring the dataset, typically using Python libraries such as Pandas and NumPy for data manipulation, and Matplotlib or Seaborn for visualization. During preprocessing, missing values in features like item weight are filled using mean imputation, and inconsistent categorical data entries are standardized. Label Encoding and One-Hot Encoding are applied to transform categorical features for machine learning compatibility. The dataset is then split into training and testing subsets to evaluate model performance objectively.

Various machine learning models—including Linear Regression, K-Nearest Neighbors (KNN), Decision Tree, Random Forest, and AdaBoost—are implemented using the Scikit-learn library. Each model is trained on the processed data, and their predictions are evaluated using metrics such as Root Mean Square Error (RMSE) and R<sup>2</sup> score to determine accuracy. Hyperparameter tuning is performed for models like Random Forest and AdaBoost to optimize performance. Once the most accurate model is identified, it is saved using joblib for future use. Finally, a simple user interface or dashboard can be built using frameworks like Streamlit or Flask to allow non-technical users to input data and view predictions, completing the implementation phase.

# V. RESULTS

The results of the sales prediction system demonstrate that ensemble models, particularly Random Forest and AdaBoost, outperform simpler models like Linear Regression and K-Nearest Neighbors in terms of prediction accuracy. Random Forest achieved the lowest Root Mean Square Error (RMSE) and the highest R<sup>2</sup> score, indicating strong generalization to unseen data. The model effectively captured complex relationships between product and store-level features. Feature importance analysis revealed that Maximum Retail Price (MRP), outlet type, and product visibility were key drivers of sales. Overall, the implementation confirms that machine learning models can significantly enhance sales forecasting, enabling better retail planning and resource optimization.

# VI. CONCLUSION

In conclusion, this project successfully demonstrates the effectiveness of machine learning algorithms in predicting retail product sales using historical and contextual data. Among the tested models, Random Forest delivered the best performance, highlighting its suitability for complex datasets. Accurate sales forecasting enables businesses to make data-driven decisions, improve inventory management, and optimize pricing strategies. The system provides a scalable solution for enhancing profitability and operational efficiency in retail environments.

# ACKNOWLEDGEMENT

I would like to express my sincere gratitude to my project guide [Guide's Name] for their invaluable support, guidance, and encouragement throughout the project. I am also thankful to [Institution/Organization Name] for providing the resources and infrastructure needed to complete this work. Special thanks to my peers and faculty members for their constructive feedback and motivation, which played a crucial role in the successful completion of this project.

# REFERENCES

- <u>https://www.kaggle.com/c/walmart-recruiting-store-sales-forecasting</u>
- <u>https://www.mdpi.com/1424-8220/19/2/286</u>
- <u>https://www.jetir.org/view?paper=JETIRGH06010</u>
- <u>https://www.researchgate.net/publication/366389871\_MACHINE\_LEARNING-BASED\_SALES\_FORECASTING\_SYSTEM</u>