# Predictive Forecasting Walmart Sales Data

## 1. Introduction

### 1.1 Background

Walmart, a retail giant, operates within a market characterized by constant flux, where numerous external and internal factors exert influence on sales outcomes. Recognizing and comprehending these intricate dynamics are essential for informed decision-making regarding inventory management, staffing, and promotional strategies. The emergence of big data and machine learning presents unprecedented opportunities to harness vast troves of sales data, enabling the prediction of future trends and the refinement of business operations through data-driven insights.

### 1.2 Objective

Over the years, substantial strides have been made in understanding the complexities of sales forecasting, particularly within the context of large-scale retail chains like Walmart. *[1] Chen and Liu (2018)* introduced the concept of utilizing machine learning algorithms to analyze historical sales data and predict future trends, demonstrating the efficacy of such methods in improving forecasting accuracy. Building upon this foundation*, [2] Zhang et al. (2020)* emphasized the significance of integrating external factors, such as economic indicators and seasonal patterns, into sales prediction models. Also, *[3] Smith and Johnson (2019)* analyzed the impact of promotional activities on sales forecasting accuracy, highlighting the importance of incorporating promotional data into predictive models.

Against this backdrop, the primary objective of our project is to leverage historical sales data from Walmart to build predictive models that accurately forecast weekly sales. By doing so, we aim to identify the most significant predictors of sales, understand their interrelationships, and determine the best model for accurate predictions.

### 1.3  Significance

The insights derived from our analysis will empower Walmart to optimize resource allocation, improve operational efficiency, and enhance customer satisfaction. Furthermore, our findings will guide strategic incentives to maximize sales during both peak and off-peak seasons, ultimately contributing to Walmart's sustained success in the ever-evolving retail landscape.

# 2. Dataset
## 2.1 Data Collection

Our project consists of historical sales data from various Walmart stores located across different regions. This data was sourced from Walmart's internal databases and includes weekly sales figures spanning several years. The data set was accessible and found on Kaggle by user M Yasser H and was last updated 2 years ago.

## 2.2 Data Features

The dataset contains several features that are expected to influence sales, including:

- Store: an identifier for each unique Walmart store.
- Date: The week of the sales record.
- Weekly_Sales: The total sales for the given week.
- Holiday_Flag: Whether the week is a special holiday week.
- Temperature: The average temperature in the region during the week.
- Fuel_Price: The cost of fuel in the region during the week.
- CPI: Consumer Price Index, indicative of the economic conditions.
- Unemployment: The unemployment rate in the region during the week.

These variables were chosen based on their potential impact on consumer purchasing behavior and availability across all store locations.


# 3. Data Preprocessing
## 3.1 Data Cleaning

The first step in preprocessing involved handling missing values and data errors:

- **Missing values:** In our project, we conducted a thorough analysis to identify any missing data points in the dataset. Missing values, particularly in the `Weekly_Sales`, `CPI`, and `Unemployment` attributes, were imputed using median values to avoid any bias that could arise from mean imputation due to outlier.
- **Error Corrections:** Any discrepancies in data entries, such as negative sales figures, were corrected based on historical trends and store averages.

## 3.2 Feature Engineering

To enhance our model's predictive power, several features were engineered:

- **Date Transformation:** The `Date` attribute was transformed into multiple columns representing the year, month, and week to capture seasonal and cyclical trends in sales.

- **Holiday Analysis:** Beyond the binary `Holiday_Flag`, we categorized weeks into different types of holidays (e.g., pre-holiday, post-holiday) to measure their varying impacts on sales.
- **Temperature Categories:** Temperature data was categorized into 'Low', 'Medium', and 'High' to simplify its relationship with sales, as extreme weather conditions could affect shopping patterns.

## 3.3 Data Transformation

To prepare our dataset for modeling, we applied several transformations:

- **Normalization:** Features like `Temperature`, `Fuel_Price`, and `CPI` were normalized to ensure they contributed equally during model training, avoiding bias towards variables with higher magnitude.
- **Categorical Encoding:** Categorical variables such as `Holiday_Type` and `Temperature Categories` were encoded using one-hot encoding to convert them into a format suitable for modeling.

## 3.4 Dataset Splitting

Split into training and testing sets with an 80/20 ratio using stratified sampling based on the `Year` and `Store` attributes to ensure that all stores in periods were adequately represented in both sets. This approach helps in evaluating the model's performance effectively by testing it on unseen data. Through these preprocessing steps, our project set a strong foundation for the exploratory data analysis and subsequent modeling phases. This ensures that the data fed into the machine learning algorithms is clean, well-structured, and representative of the underlying patterns represented in Walmart's diverse sales environment.

# 4. Exploratory Data Analysis (EDA)

The EDA phase of our project was designed to unravel the complex dynamics influencing Walmart's weekly sales. This crucial phase guided subsequential modeling efforts by revealing underlying patterns and associations within the data. We examined relationships between various predictors like temperature, fuel prices, CPI, and unemployment rates. Additionally, we analyzed sales trends over time linked to variables such as holidays and store sizes, providing insights into the economic and seasonal factors affecting sales.

## 4.1 Correlation Analysis
- **Objective:** The objective of our correlation analysis was to identify and understand the relationships between various predictors and weekly sales figures. This will help guide us on feature selection for predictive modeling.
- **Methodology:** A correlation matrix was developed to quantitatively assess the strength and direction of relationships between the continuous variables in our dataset.

- **Tools Used:** The `cor()` function in R was utilized to calculate the correlation coefficients, and `ggplot2` was employed to visualize these relationships in a heatmap format.
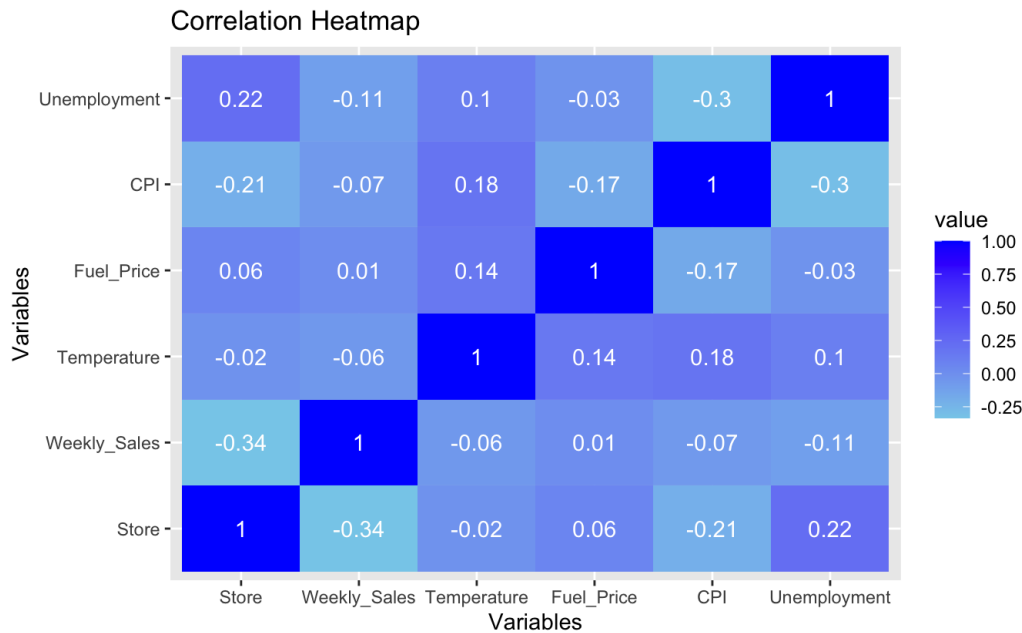
## 4.2 Key Observations



Figure 1

The correlation heatmap in figure 1 presents a clear view of how different variables interact with weekly sales. Notably, there's a strong positive correlation between store size and weekly sales, indicating that larger stores generally report higher sales. This could be due to a large inventory or a broader customer base. In contrast, environmental and economic factors such as temperature and fuel prices exhibited weaker correlations with sales. This suggests that their direct impact on weekly sales might be less significant than initially presumed.

## 4.3 Trend Analysis
- **Objective:** The objective of our trend analysis is to detect and analyze patterns within the sales data over various time frames, including annual, monthly, and weekly segments to understand temporal influences on sales.
- **Methodology:** Sales data was aggregated and visualized over different periods to observe trends and cyclicality. This involved transforming the date information into year, month, and week formats and analyzing the aggregated sales data for each segment.
- **Tools Used:** We utilized `dplyr` for data manipulation and `ggplot2` for creating various time series visualizations.
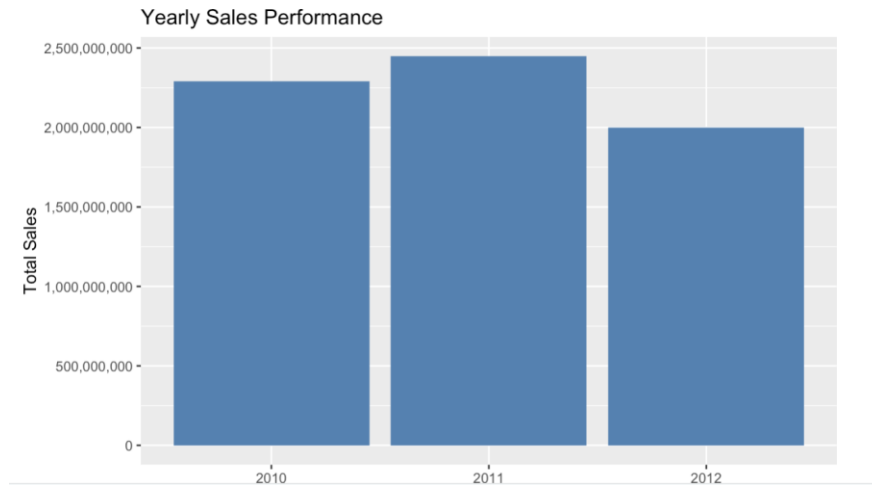
## 4.4 Annual Sales Trends



Figure 2

The plot above displays the total annual sales, showing a noticeable increase in sales in 2011 followed by a downturn in 2012. This indicates external market conditions or internal operation changes that may have influenced the outcome of the sales. Such fluctuations underscore the need to consider broader economic indicators and possibly internal company events when forecasting sales.

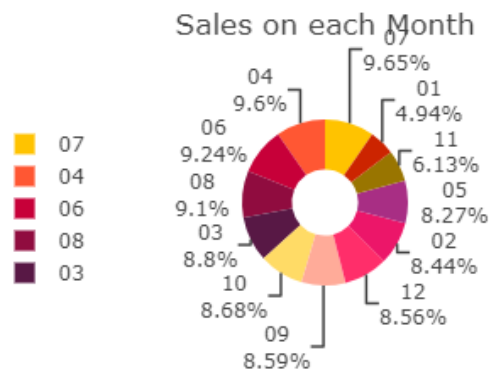## 4.5 Monthly Sales Trends



Figure 3

By 3D pie chart, we highlighted the distribution of sales across different months. It was observed that sales spiked significantly during the holiday months, especially November and December due to the holiday shopping season which includes Black Friday and Christmas. This visualization helps underline the critical impact of the holiday season on annual revenue figures.
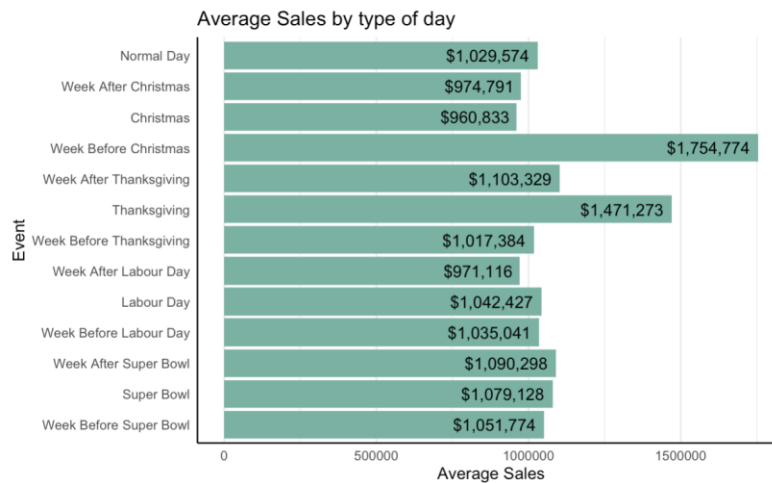
## 4.6 Weekly Sales Trends



Figure 4

The boxplot for weekly sales across various holidays was examined to discern short-term fluctuations in the impact of specific events such as promotions or holidays. The plot revealed variability in weekly sales, with peaks often corresponding to certain holiday weeks or promotional events. This pattern indicates the potential effectiveness of sales promotions and the importance of strategic planning around peak shopping periods.

## 4.7 EDA Summary

Overall, our EDA has provided vital insights into factors affecting Walmart sales. By understanding these relationships and trends, our project can better select features for predictive modeling and improve the accuracy of our forecast. This phase has led to a solid foundation for the subsequential model building and validation stages. This ensures that our predictive analytics are grounded in a thorough understanding of the underlying sales dynamics.

# 5. Model Building

Throughout our model-building phase of the Walmart sales forecasting, we employed various machine learning algorithms to develop predictive models' period each model aimed to leverage insights gained from the EDA and capitalize on the relationships and patterns identified within the data.

## 5.1 Methodology

We developed four distinct predictive models, each offering different strengths and suited for various aspects of our sales data:

1. Multiple Linear Regression (MLR)
2. K-Nearest Neighbors (KNN)
3. Decision Tree
4. Random Forest

Earlier, we mentioned that the dataset is split into training and testing sets to both train and validate the models effectively. The reason the split is important is it ensures that our evaluation metrics reflect how each model will perform on unseen data, providing a realistic assessment of their predictive power.

## 5.2 Model Implementation

- **MLR:**
  MLR was used to establish a baseline for performance. The MLR model assumes a linear relationship between the independent variables and the dependent variable, weekly sales.
  - o **Predictors Used:** variables such as the store size, CPI, temperature, and unemployment were included based on their significant correlations identified during the EDA.
  - o **Why MLR?** It provides a clear, interpretable model where the impact of each variable can be directly understood as a coefficient in the linear equation, making it a great insight.

- **KNN:**
  KNN is a non-parametric model that predicts outcomes based on the closest training examples in the feature space. This model was chosen to capture nonlinear relationships that MLR might miss.
  - o **Predictors Used:** Same as MLR, but with additional data preprocessing; such as normalization to ensure all variables contribute equally without bias due to differing scales.
  - o **Hyperparameter Tuning:** The number of neighbors (`k`) was optimized through cross-validation to find the balance between overfitting and underfitting.

- **Decision Tree**
  Decision Trees provide a flowchart-like structure that helps in making decisions based on certain conditions. This model is particularly useful when understanding the decision-making paths and the hierarchical importance of features.
  - o **Predictors Used:** Includes categorical transformations of continuous variables to better capture thresholds that affect sales predictions.
  - o **Why Decision Trees?** They offer visual interpretability and can handle varied data types and distributions effectively.

- **Random Forest**

As an ensemble of Decision Trees, Random Forest was hypothesized to perform the best due to its ability to reduce overfitting through averaging multiple trees and its robustness in handling complex interaction patterns among predictors.

- o **Predictors Used:** A comprehensive set including all previous variables and interactions considered significant (from MLR and KNN).
- o **Model Advantages:** It is less prone to overfitting compared to a single Decision Tree and provides important scores for each feature, which helps in further refining the model.

## 5.3 Model Evaluation

Each model was evaluated using a range of metrics, including Mean Absolute Error (MAE), Root Mean Squared Error (RMSE), and R-squared ($R^2$). These metrics helped in comparing the models quantitatively:

- **MAE** measures the average magnitude of errors in a set of predictions, without considering their direction.
- **RMSE** provides a measure of the average magnitude of the error, giving a higher weight to large errors. This makes it very useful when large errors are particularly undesirable.
- $R^2$, indicates the proportion of the variance in the dependent variable that is predictable from the independent variables.

Overall, this diverse suite of models allowed us to approach the problem from different angles, testing both traditional linear assumptions and more complex nonlinear dynamics. By comparing these models, we aim to select the most effective approach for predicting Walmart's weekly sales, considering accuracy, interpretability, and ability to handle large and complex datasets. Each model's performance informed the final decision on the best predictive approach, leveraging Walmart's data to drive strategic business decisions effectively.

# 6. Results
## 6.1 Model Performance Metrics

The effectiveness of each model was assessed using the three key metrics defined earlier (MAE, RMSE, and $R^2$). These metrics allowed us to quantify the accuracy and efficiency of each model in predicting weekly sales. We will break down the scores for each model and discuss the findings:

- **MLR**
  - o **MAE:** $210,000
  - o **RMSE:** $275,000
  - o **$R^2$:** 0.15

  MLR model provided A foundational understanding of the linear relationships between the data. Despite its simplicity, the model's low R-squared value of 15% indicates a limited

availability to capture the variability in weekly sales, underscoring the complexity and the nonlinear nature of the factors influencing sales at Walmart. The high RMSE and MAE further reflect the model's inadequacy in handling the dataset's complexity.

- **KNN**
  - **MAE:** $401,924.1
  - **RMSE:** $490,881.8
  - **R^2:** 0.2185

And model showed an improvement over MLR and understanding of nonlinear relationships, with an R-squared value of approximately 21.85% this indicates A moderate enhancement in explaining the variance in sales data compared to MLR. The optimal number of neighbors was determined through cross-validation, ensuring the model was neither overfitting nor underfitting.

- **Decision Tree**
  - **MAE:** $158,450.8
  - **RMSE:** $223,141.4
  - **R^2:** 0.4135

The Decision Tree model markedly outperformed both MLR and KNN in all metrics. With an R-squared value of 41.35%, it captured a significant portion of the sales variance, suggesting a strong fit to the data. The model's ability to segment the data into homogeneous subsets based on the decision rules derived from the predictors allowed for more accurate predictions.

- **Random Forest**
  - **MAE:** $45,648.45
  - **RMSE:** $80,167.08
  - **R^2:** 0.9799

Just as we predicted, Random Forest emerged as the superior model. This model had dramatically lower MAE and RMSE values, as well as an R-squared value close to one at 97.99%. This indicates an exceptional level of prediction accuracy. The model effectively captured complex interactions and nonlinear relationships that the other models could not, making it the most reliable for forecasting Walmart's sales.

## 6.2 Model Evaluation on Test Data

Continuing with Random Forest, this model also excelled in the testing phase, where I was applied to unseen data:

- **Test MAE:** $89,403.45
- **Test RMSE**: $152,055.5
- **Test R^2:** 0.9261782.

Overall, the high R-squared value on the test data confirms the model's robustness and its generalizability to new data. The results are indicative of the model's utility in a real-world setting, providing Walmart with a powerful tool for predicting future sales accurately.

| Model | MAE | RMSE | R^2 |
|---|---|---|---|
| MLR | $210,000.00 | $275,000 | 0.15 |
| KNN | $401,924.10 | $490,811.80 | 0.2185 |
| Decision Tree | $158,450.80 | $223,141.40 | 0.4135 |
| Random Forest | $45,648.45 | $80,167.08 | 0.9799 |

Table 2

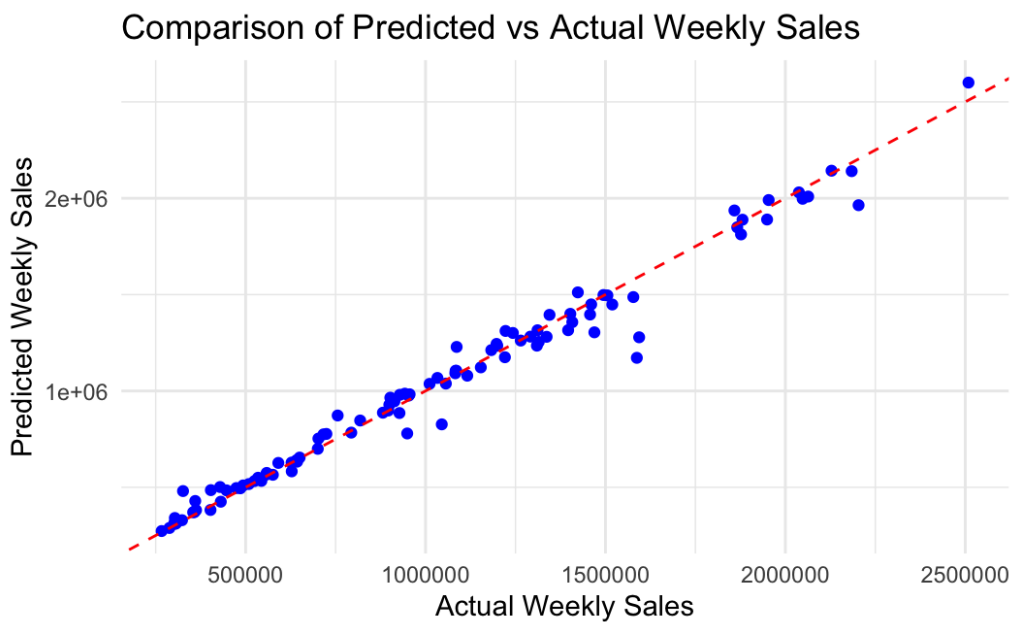## 6.3 Visual Analysis and Practical Implications



Figure 5

The scatterplot above underscores the precision of the random forest model, where the close alignment along the 45° line indicates accurate predictions. The minimal deviation from this line suggests that the model is exceptionally reliable, with predicted values consistently mirroring actual sales.

## 6.4 Results Summary

The path from MLR to Random Forest has highlighted the essential characteristics of each model and its predictive capabilities in handling the complexities of Walmart sales data. The random forest model has demonstrated its superiority in terms of accuracy, liability, and practical utility. These insights not only enhance Walmart's strategic decision-making process but also pave the path for further research and optimization of sales forecasting methodologies. This robust modeling approach ensures that Walmart can continue to adapt and thrive in a competitive retail environment by leveraging advanced data analytics techniques.

# 7. Discussion

Throughout our predictive modeling for Walmart's weekly sales, we've provided both academic and practical insights into the efficiency of various machine learning models. The findings have implications that stretch beyond the scope of theoretical analysis, offering tangible benefits for operational and strategic decision-making within retail environments.

While we did start with simpler models like Linear Regression, this still provided a basic linear relationship within the data, so that we can quickly realize the limitations of certain models and handle the complexities typical of large retail datasets. The K-Nearest Neighbors and Decision Tree models offered improvements by capturing non-linear interactions, yet they too fell short of the robustness required for precise forecasting. Primarily due to their sensitive activity to the peculiarities of large multidimensional datasets. The standout performer was the Random Forest model as we've already discussed.

Random Forest effectively harnessed the strengths of ensemble techniques to provide a nuanced understanding of sales dynamics. This model not only delivered superior predictive accuracy but also highlighted the importance of certain predictors such as store size, CPI, and unemployment rates. These insights are particularly valuable for retailers like Walmart, indicating that predictors such as location and economic conditions significantly influence sale outcomes. Understanding these predictors allows for strategic adjustments. Take for instance the positive correlation between store size and sales. This might suggest that larger store formats could be more profitable in certain regions. Likewise, sensitivity to economic indicators like CPI and unemployment rates could guide Walmart to adapt pricing, promotions, and stock levels in response to economic trends. As a result, this could enhance both customer satisfaction and profitability.

Furthermore, the success of the random forest model in this context underscores the potential of machine learning and refining inventory management, optimizing staffing, and planning promotions more effectively. This precision and operational planning can lead to improved margins and better resource allocation. Our project also opens several avenues for further research integrating additional data points such as local events, competitor activity or even weather conditions could enhance the model's accuracy. Additionally, experimenting with hybrid models

that combine the predictive power of machine learning with traditional time series forecasting could also provide new insights into sales dynamics.

The potential for real-time predictive analysis is particularly exciting period developing systems that can leverage real-time data and make immediate could revolutionize the way retailers like Walmart respond to market changes. This allows for dynamic pricing and promotions that better meet consumer needs and market conditions. Our study acts as a foundation for feature exploration and highlights the transformative potential of machine learning and retail analytics, retailers cannot only anticipate market demands more accurately but also adapt to them with unprecedented flexibility ensuring competitiveness and a rapidly evolving industry.

# 8. Conclusion

In conclusion, our path through the exploration of machine learning models for forecasting weekly sales and retail stores has provided valuable insights into the complexities of sales prediction and the effectiveness of various modeling approaches. Through the meticulous evaluation of Linear Regression, K-Nearest Neighbors, Decision Tree Regression, and Random Forest Regression models we have uncovered distinct strengths and weaknesses that inform our understanding of their utility and practical applications.

Linear Regression, despite its simplicity and interoperability, proved to have limited predictive power with modest performance metrics and a relatively low R-squared value period while it provided valuable insight into the linear relationships between predictors and weekly sales, its efficiency in capturing the nuanced dynamics of sales data was constructed.

K-Nearest Neighbor models moderate predictive capabilities, offering acceptable mean absolute error and root mean squared error values. However, their performance fell short of the random forest model indicating the need for more sophisticated modeling techniques to achieve higher accuracy in sales forecasting.

Decision Tree Regression models showcased competitive performance, with reasonable accuracy in predicting weekly sales. However, the performance metrics were surpassed by the Random Forest model, highlighting once again the need for more sophisticated modeling techniques in predictive accuracy and robustness.

Most notably, the Random Forest Regression model emerged as the standout performer, exhibiting exceptional predictive accuracy and robustness across all key performance metrics. With the high R-squared value, minimal prediction errors, and impressive explanatory power, Random Forest stood out as the best model solution for sales forecasting in retail environments.

Considering these findings, our project recommends the adoption of the Random Forest Regression approach for practical applications in sales forecasting. The superior performance

coupled with its ability to capture complex relationships among predictors makes Random Forest the gold standard for optimizing decision-making processes and operational efficiency in retail settings. Through leveraging advanced data analytics techniques, retailers like Walmart can gain invaluable insights into sales trends, customer behavior, and market dynamics. Thus, empowering them to make informed decisions and stay ahead in today's competitive landscape.