

Retail Sales Forecasting

*June 2025
Machine Learning Project*

Manas R. Das

Under the Supervision of
Mr. Sambit Subhasish Sahu

ABSTRACT

This document presents the **Advanced Retail Sales Forecasting Dashboard**, a comprehensive and interactive web application developed using **Python** and the **Streamlit** framework. Designed to empower retail businesses, this project provides a sophisticated platform for in-depth analysis of historical sales data, visualization of key performance indicators (KPIs) and intricate trends, and the generation of accurate future sales forecasts.

At its core, the dashboard leverages an **XGBoost Regressor** machine learning model, meticulously trained on diverse features including temporal attributes (date, month, year, season, holiday status) and store-specific data (store, department, type, size, temperature, fuel price, markdown events, CPI, unemployment). To ensure transparency and build user trust, the application incorporates **SHAP (SHapley Additive exPlanations)** for model interpretability, visually detailing how each feature contributes to the sales predictions.

Key functionalities include dynamic data filtering for granular analysis across various dimensions (store, department, type, season, month, holiday status), real-time KPI monitoring, and interactive visualizations (line charts, bar charts, pie charts, box plots, correlation heatmaps) powered by **Plotly**. Crucially, the dashboard features a robust mechanism for **uploading new sales data and automatically retraining the XGBoost model**, ensuring its continuous accuracy and adaptability to evolving market conditions. This holistic solution addresses real-world retail challenges such as sales volatility and inventory mismanagement, enabling data-driven decision-making, optimizing resource allocation, and ultimately driving operational efficiency and profitability.

Introduction

In the dynamic landscape of retail, understanding sales patterns and predicting future demand are critical for operational efficiency, inventory management, and strategic planning. Traditional methods often fall short in handling the complexity and volume of modern sales data. This project addresses these challenges by offering a sophisticated, interactive web-based dashboard. Leveraging the power of modern data science tools, it transforms raw sales data into actionable insights, providing a clear view of past performance and a reliable projection of future trends.

Objective

The primary objectives of the Advanced Retail Sales Forecasting Dashboard are:

1. **Comprehensive Data Visualization:** To provide an intuitive and interactive platform for visualizing historical retail sales data through various charts and graphs, enabling users to quickly identify trends, distributions, and outliers.
2. **Key Performance Indicator (KPI) Monitoring:** To present crucial sales metrics in an easily digestible format, allowing users to track overall sales performance.
3. **Advanced Sales Forecasting:** To implement a robust machine learning model (XGBoost) capable of predicting future weekly sales for specific stores and departments, aiding in proactive decision-making.
4. **Model Interpretability:** To offer insights into the model's predictions using SHAP (SHapley Additive exPlanations) values, helping users understand the contribution of each feature to the sales forecast.
5. **Dynamic Filtering and Customization:** To allow users to filter data based on various dimensions (Store, Department, Type, Season, Month, Holiday status) to focus on specific segments of interest.
6. **Model Retraining Capability:** To enable users to upload new sales data to continuously update and improve the accuracy of the forecasting model, ensuring its relevance over time.
7. **User-Friendly Interface:** To present a clean, responsive, and easy-to-navigate multi-page application suitable for business analysts, managers, and data scientists.

Technical Specification

Technologies & Libraries:

- **Python:** The core programming language for the entire application.
- **Streamlit:** The primary framework for building the interactive web dashboard. Its simplicity and speed of development were key factors in its selection.
- **Pandas:** Used extensively for data loading, manipulation, cleaning, and aggregation.
- **NumPy:** Essential for numerical operations and array handling, particularly within the machine learning pipeline.
- **Matplotlib & Seaborn:** Employed for generating static statistical plots, notably the correlation heatmap.
- **Plotly Express & Plotly Graph Objects:** Utilized for creating interactive and aesthetically pleasing visualizations (line charts, bar charts, pie charts, box plots) that enhance user engagement.
- **scikit-learn:** Provides tools for data splitting (`train_test_split`) and performance evaluation metrics (`r2_score`, `mean_absolute_error`, `mean_squared_error`).
- **XGBoost (eXtreme Gradient Boosting):** The chosen machine learning algorithm for sales forecasting. Its high performance, ability to handle various data types, and interpretability make it suitable for this task.
- **SHAP (SHapley Additive exPlanations):** A powerful library for explaining the output of machine learning models. It helps in understanding individual predictions and overall feature importance.
- **Joblib:** Used for serializing (saving) and deserializing (loading) the trained XGBoost model, allowing persistence across application sessions.
- **datetime module:** For handling date and time operations, crucial for extracting features like Month, Day, Year, and calculating future dates.

Architecture:

The application follows a multi-page Streamlit architecture, organized as follows:

- **app.py (Main Entry Point):**
 - Initializes global Streamlit page configuration (`st.set_page_config`).
 - Defines and stores global helper functions (e.g., `get_season`) in `st.session_state`.
 - Initializes and stores LabelEncoder instances in `st.session_state` to ensure consistent categorical encoding across the application.
 - Loads and preprocesses the primary dataset (`dataset/dataset.csv`) using `@st.cache_data` for performance optimization. The processed DataFrame is then stored in `st.session_state.df`.
 - Initializes a placeholder for the `xgboost_model` in `st.session_state` and attempts to load a pre-trained model from disk.
 - Provides a brief introduction and navigation guidance to the user.
- **pages:** Contains individual Python files, each representing a distinct page in the application. Streamlit automatically creates a sidebar navigation based on these files.
 - **1_Dashboard.py:**
 - Accesses the preprocessed `df` and encoders from `st.session_state`.
 - Implements data filtering logic based on user selections in the sidebar.
 - Displays Key Performance Indicators (KPIs) such as Total Sales, Avg. Weekly Sales, Unique Stores, and Departments.
 - Generates various interactive sales performance charts (e.g., Weekly Sales Over Time, Sales by Store, Sales by Store Type, Holiday vs. Non-Holiday Sales, Season-wise Sales, Monthly Sales Trends by Store Type).
 - Presents a correlation heatmap of numerical features.
 - Shows an aggregated data summary table. Allows downloading of filtered raw data.
 - **2_Forecasting.py:**
 - Accesses `df`, encoders, `model_path`, and `get_season` from `st.session_state`.
 - Loads the XGBoost model from `st.session_state` (or trains it if not available/outdated).
 - Evaluates the model's performance on historical data using R^2 , MAE, and RMSE.
 - Visualizes actual vs. predicted sales trends.
 - Displays feature importance and SHAP summary plots for model interpretability.
 - Provides interactive controls in the sidebar for selecting specific Store and Department for future sales forecasting.
 - Generates and displays the future sales forecast plot and a detailed table.
 - Enables downloading of the forecasted data.
 - **3_Retrain_Model.py:**
 - Accesses `df`, encoders, `model_path`, and `get_season` from `st.session_state`.
 - Provides a file uploader for users to upload new sales CSV data.
 - Performs validation and preprocessing on the uploaded data, ensuring compatibility.
 - Combines new data with existing historical data, handling duplicates.
 - Retrains the XGBoost model on the combined dataset.
 - Saves the retrained model to disk and updates `st.session_state` with the new data, encoders, and model, triggering a rerun of the application to reflect changes.

Data Flow:

1. **Initialization (app.py):** Raw dataset.csv is loaded, preprocessed (date parsing, season derivation, categorical encoding via LabelEncoder), and the resulting DataFrame (df) is stored in st.session_state. Encoders and the get_season helper function are also stored.
2. **Dashboard (1_Dashboard.py):** Filters are applied to st.session_state.df to create filtered_df. Visualizations and summary tables are generated from filtered_df.
3. **Forecasting (2_Forecasting.py):** The full st.session_state.df is used to prepare X and y for model training/prediction. The xgboost_model is loaded from st.session_state (or trained). Predictions are made, and visualizations are generated. Future data points are constructed using historical data patterns and st.session_state.get_season and st.session_state.season_encoder.
4. **Retraining (3_Retrain_Model.py):** New CSV data is uploaded, preprocessed (using updated st.session_state.type_encoder and st.session_state.season_encoder to handle new categories), and concatenated with st.session_state.df. The model is retrained on this combined dataset, saved to disk, and the updated df, encoders, and model are pushed back to st.session_state, forcing a full app rerun.

Output

The application provides the following key outputs:

- **Interactive Dashboard Pages:**
 - **"Sales Performance Dashboard":** Displays interactive line, bar, pie, and box plots illustrating weekly sales trends, sales by store/department/type/season, and holiday impact. Includes a correlation heatmap.
 - **"Advanced Sales Forecasting":** Presents the XGBoost model's performance on historical data (Actual vs. Predicted Sales plot, R-squared, MAE, RMSE metrics), feature importance bar chart, and detailed SHAP summary plots (bar and beeswarm).
- **Key Performance Indicators (KPIs):** Instant display of total sales, average weekly sales, unique stores, and unique departments based on filtered data.
- **Aggregated Data Summary Table:** A detailed table summarizing sales metrics (Total, Average, Max, Min) grouped by Store, Department, Type, and Season.
- **Future Sales Forecast Plot:** A line chart visualizing historical sales alongside projected future sales for a user-selected store and department.
- **Detailed Future Forecast Table:** A tabular breakdown of predicted weekly sales for each forecasted day.
- **Downloadable Data:** Users can download the filtered raw data and the generated future sales forecast as CSV files.
- **Model Retraining Confirmation:** Clear messages indicating the success or failure of new data processing and model retraining.

Outcome

This project empowers retail businesses to achieve several beneficial outcomes:

- **Enhanced Business Intelligence:** Provides a centralized and dynamic platform for exploring sales data, leading to a deeper understanding of sales drivers and patterns.
- **Improved Forecasting Accuracy:** The XGBoost model offers more reliable sales predictions compared to traditional methods, reducing guesswork in planning.
- **Optimized Inventory Management:** Accurate forecasts enable businesses to better manage inventory levels, reducing overstocking (and associated costs) and understocking (and lost sales).
- **Proactive Strategic Planning:** Insights from sales trends and forecasts allow for more informed decisions regarding promotions, staffing, marketing campaigns, and store operations.
- **Identification of Key Performance Drivers:** Feature importance and SHAP analysis help identify which factors (e.g., Markdown, Holiday, Season) have the most significant impact on sales, allowing for targeted interventions.
- **Continuous Improvement:** The model retraining feature ensures that the forecasting model remains accurate and adapts to evolving market conditions and new data.
- **Democratization of Data Insights:** The user-friendly Streamlit interface makes complex sales analytics and machine learning predictions accessible to a broader audience within the organization, not just data scientists.

Real-life Problem:

Retail businesses constantly grapple with several critical challenges:

1. **Sales Volatility:** Weekly sales can fluctuate dramatically due to seasonality, holidays, promotions, economic factors, and competitive pressures, making it difficult to predict future demand.
2. **Suboptimal Inventory:** Inaccurate sales forecasts lead to either excess inventory (tying up capital, increasing storage costs, risk of obsolescence) or insufficient stock (resulting in lost sales, customer dissatisfaction, and missed opportunities).
3. **Inefficient Resource Allocation:** Without clear sales projections, businesses struggle to allocate staff effectively, plan marketing budgets, or schedule supply chain logistics.
4. **Lack of Data-Driven Decision-Making:** Businesses often rely on intuition or simplistic historical averages, missing deeper insights hidden within their sales data.
5. **Model Obsolescence:** Predictive models built on old data can quickly become irrelevant as market dynamics change, requiring continuous updates.

How This Project Solves It:

This project directly addresses these problems:

- **Mitigating Sales Volatility:** The XGBoost model, trained on various features beyond just historical sales (like holidays, temperature, fuel price, markdowns, CPI, unemployment), learns complex non-linear relationships to provide more accurate predictions, accounting for diverse influencing factors.
- **Optimizing Inventory:** By providing specific future sales forecasts for individual stores and departments, the application enables precise inventory ordering, minimizing waste and maximizing product availability.
- **Enabling Efficient Resource Allocation:** Management can use forecasted sales to better plan staffing levels, optimize promotional calendars, and align supply chain operations with anticipated demand.
- **Facilitating Data-Driven Decisions:** The interactive dashboards, KPIs, and model interpretability (SHAP) tools empower users to explore data, understand "why" sales are happening, and base decisions on empirical evidence rather than guesswork.
- **Ensuring Model Relevance:** The "Retrain Model with New Data" feature ensures that the forecasting model can be regularly updated with the latest information, adapting to new trends and maintaining high predictive accuracy over time.

Future Scope

The project has significant potential for expansion and enhancement:

1. **More Advanced ML Models:** Explore other time series forecasting models (e.g., Prophet, ARIMA, LSTM neural networks) for comparison and potentially better performance. Implement ensemble methods.
2. **External API Integration:** Incorporate real-time external data sources like local weather forecasts, economic indicators, or social media sentiment analysis (where applicable) to enrich feature sets for prediction.
3. **User Authentication and Access Control:** Implement user login and role-based access to protect sensitive data and customize dashboards for different user types.
4. **Database Integration:** Replace CSV file loading with a connection to a robust database (e.g., SQL, NoSQL, or cloud data warehouses like BigQuery) for scalable data storage and faster retrieval, enabling truly real-time updates.

5. **What-If Analysis:** Allow users to simulate the impact of changing certain features (e.g., launching a new markdown campaign, predicting holiday sales for a specific date) on future sales.
6. **Uncertainty Quantification:** Provide confidence intervals or prediction ranges for forecasts, giving users a better sense of the prediction's reliability.
7. **Automated Retraining and Deployment:** Implement a backend system for automated model retraining on a schedule (e.g., weekly) and continuous deployment (CI/CD) of the updated model.
8. **Improved Holiday and Markdown Handling:** Develop more sophisticated features for future holidays (e.g., identifying upcoming holidays and their historical impact) and more granular markdown predictions instead of defaulting to zero.
9. **Anomaly Detection:** Add functionality to identify unusual sales spikes or dips in historical data for further investigation.
10. **Customizable Dashboards:** Allow users to save their preferred filter settings and chart layouts for quicker access.
11. **Scalability for Large Datasets:** Optimize data processing and model training for very large datasets, potentially using distributed computing frameworks.

Conclusion

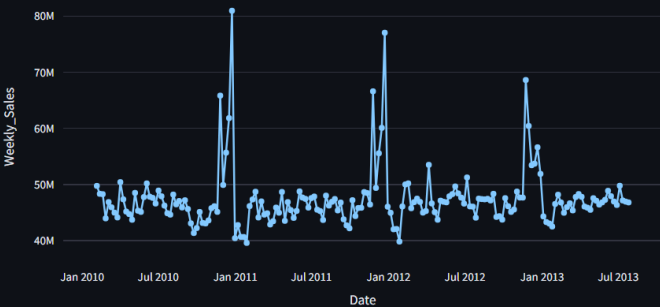
The Advanced Retail Sales Forecasting Dashboard is a powerful, user-centric application designed to revolutionize how retail businesses interact with their sales data. By integrating interactive data visualization, robust machine learning forecasting, and transparent model interpretation, it provides actionable insights that drive strategic decisions. This project not only addresses immediate operational challenges like inventory optimization but also lays the groundwork for continuous improvement through its adaptable and retraining-capable model, positioning businesses for sustained growth in a competitive market.



Sales Performance Analysis

Weekly Sales Over Time

Total Weekly Sales Trend

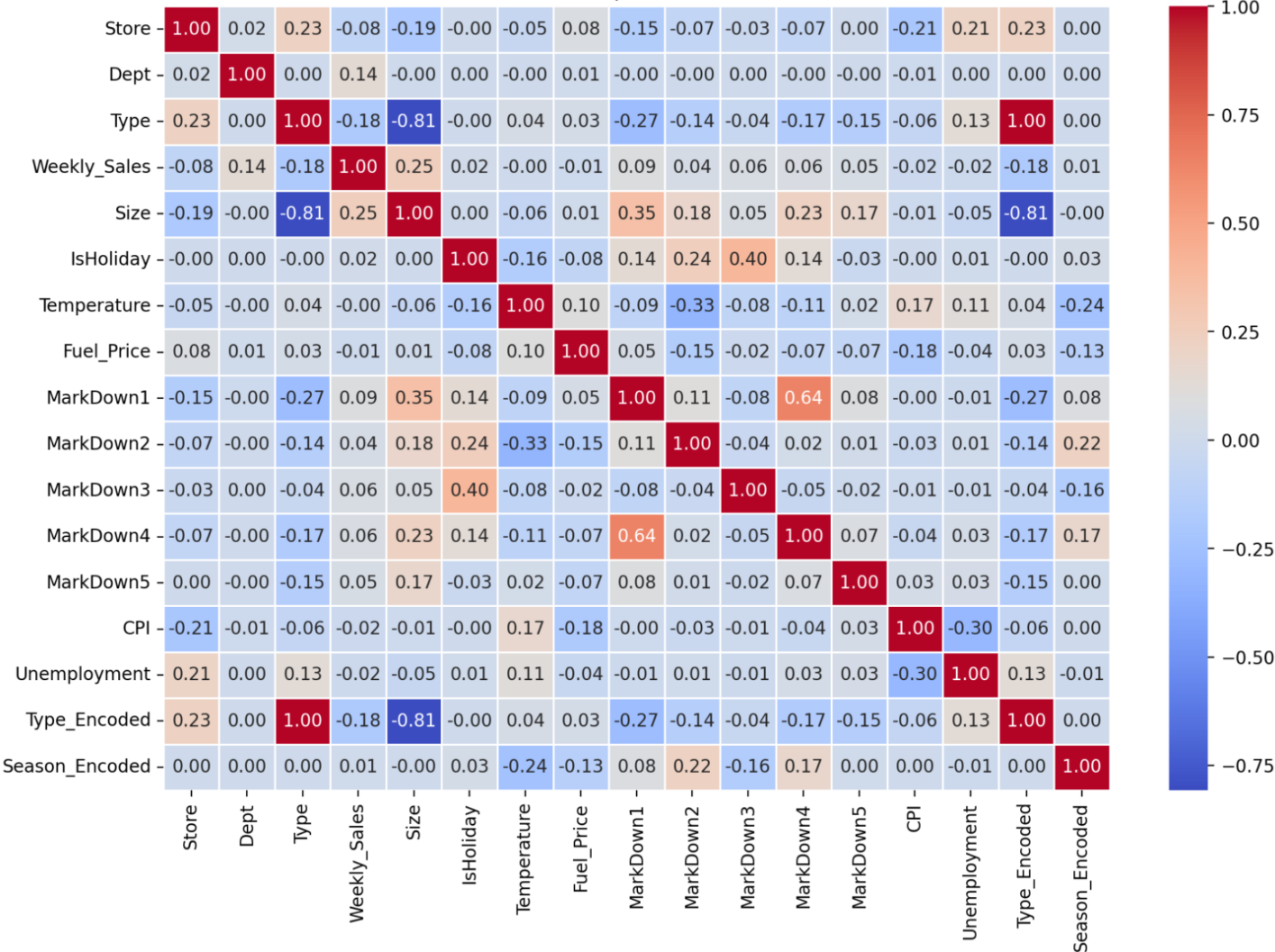


Sales by Store

Total Sales per Store

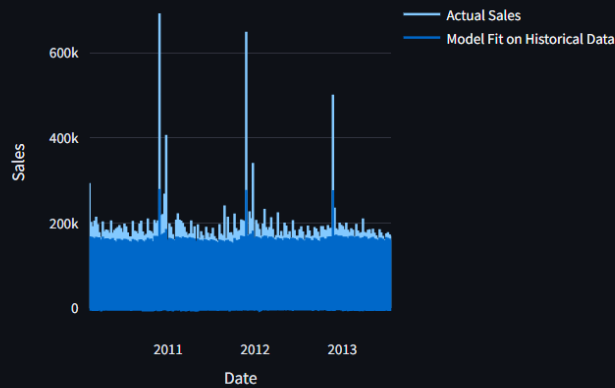


Correlation Heatmap of Numerical Features



Model Performance (XGBoost)

Weekly Sales: Actual vs. Model Fit



Evaluation Metrics (on Test Set)

	Metric	Value
0	R-squared (R^2)	0.8984
1	Mean Absolute Error (MAE)	4,252.81
2	Root Mean Squared Error (RMSE)	7,130.01

Feature Importance

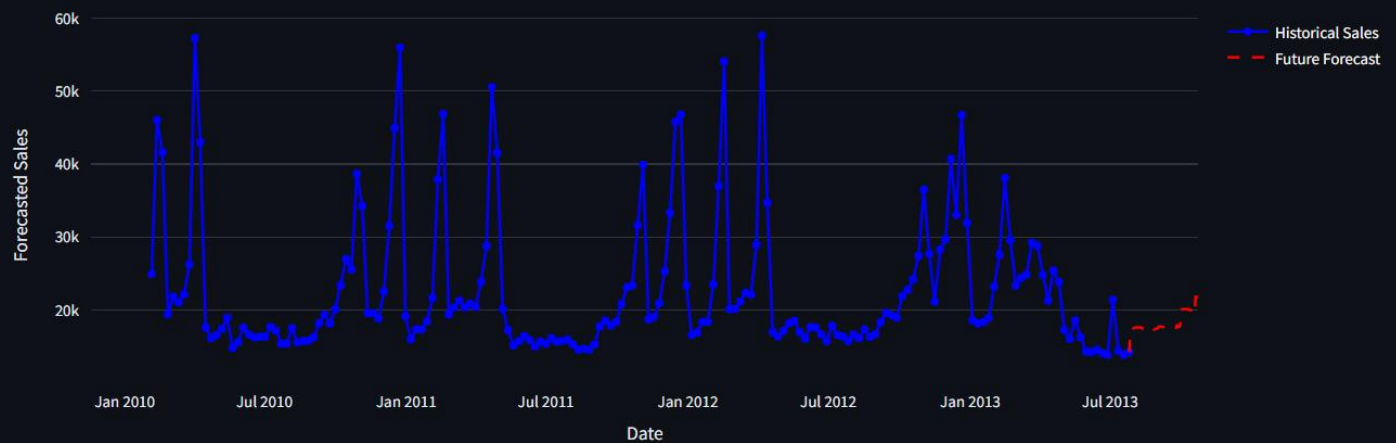
XGBoost Feature Importance



Future Sales Forecast for Specific Store/Department

Forecasting sales for Store: 1 and Department: 1 for the next 90 days.

Future Sales Forecast for Store 1, Dept 1



**Our greatest glory is
not in never failing,
but in rising every
time we fail.**

- Confucius