

Machine Learning-Driven Predictive Analytics for Retail Sales Data

Thaleti Pothuraju Meghana Sai¹, Yeturu Aasritha Reddy², Tunga Mounika Lakshmi³, Raavi Hindu⁴,
Mr. Raghavendran N⁵

^{1,2,3,4}*Student, Dept of Artificial Intelligence and Data Science, RMK College of Engineering and Technology*

⁵*Assistant Professor, M.E., MBA., (PHD), RMK College of Engineering and Technology*

Abstract—The rapid growth of e-commerce and retail industries has generated vast amounts of transactional data, creating both the opportunity and necessity for intelligent data-driven decision-making. This paper presents ShopPulse, an AI-powered retail management system that integrates predictive analytics, machine learning, and real-time data processing to forecast sales, optimize inventory, and maximize revenue. The system leverages XGBoost for hourly sales forecasting, Facebook Prophet for 30-day demand prediction, Random Forest for inventory risk classification, and Gradient Boosting for dynamic price optimization across product categories. A dataset of 54,000+ hourly retail transactions spanning multiple stores, regions, and product categories forms the foundation of model training. The system is deployed as a full-stack web application using Flask, SQLite, and interactive HTML5 dashboards powered by Chart.js. SHAP (SHapley Additive explanations) is employed for explainable AI, enabling stakeholders to understand feature contributions to predictions. Experimental results demonstrate strong predictive accuracy, with the XGBoost model achieving high R-squared scores and low Mean Absolute Error across test splits. The proposed system addresses critical retail challenges including stockouts, pricing inefficiencies, and revenue forecasting uncertainties, offering a scalable, interpretable, and deployable solution for modern retail enterprises.

Index Terms—Sales Prediction, XGBoost, Facebook Prophet, Random Forest, SHAP, Flask, Retail Analytics, Inventory Management, Dynamic Pricing, Machine Learning.

1. INTRODUCTION

1.1 The retail industry is undergoing a fundamental transformation driven by the proliferation of digital commerce, the availability of real-time data, and the maturation of machine learning technologies.

Retailers today face the dual challenge of managing complex, multi-dimensional inventory systems while simultaneously responding to rapidly shifting consumer demand patterns. Traditional forecasting methods, such as moving averages and exponential smoothing, are increasingly inadequate in the face of these complexities.

Predictive analytics, powered by advanced machine learning algorithms, presents a compelling solution. By learning from historical patterns in sales data, these models can anticipate future demand, identify at-risk inventory, and recommend optimal pricing strategies. The economic stakes are significant: industry research consistently demonstrates that improved demand forecasting reduces inventory costs by 10-30% and increases revenue by 2-5% through better product availability and pricing.

This paper introduces ShopPulse, a comprehensive AI-driven retail analytics platform. The system is designed around three core functionalities: (1) accurate sales forecasting at hourly and monthly granularities, (2) automated inventory risk detection to prevent stockouts, and (3) data-driven dynamic pricing to maximize category-level revenue. The platform is accessible through a secure, role-based web interface that serves both administrators (via an analytics dashboard) and customers (via a shopping portal).

Problem Statement

Retail enterprises routinely struggle with three interconnected problems: inaccurate demand forecasts that lead to overstocking or stockouts, static pricing strategies that leave revenue on the table, and reactive (rather than predictive) inventory management. These

challenges are compounded by the influence of external factors including weather conditions, promotional campaigns, seasonal cycles, epidemiological events, and regional competition. Existing commercial solutions are often expensive, opaque, and not easily customizable. There is a clear need for an open, interpretable, and accurate predictive system that can be deployed within a practical web application framework. There is a clear need for an open, interpretable, and accurate predictive system that can be deployed within a practical web application framework.

1.2 Objectives

- Design and implement an XGBoost-based model for next-hour sales volume forecasting.
- Develop a Facebook Prophet model for 30-day rolling demand forecasting.
- Build a Random Forest classifier for real-time inventory risk detection.
- Implement Gradient Boosting-based dynamic pricing optimization by product category.
- Deploy a full-stack web application integrating all models with a live database of user transactions.
- Provide SHAP-based explainability to ensure transparency of model predictions.

1.3 Scope of the Study

The study encompasses data from multiple retail stores across four geographic regions (North, South, East, West), covering five product categories: Groceries, Electronics, Furniture, Toys, and Sports. The dataset spans from January 2020 onwards at hourly granularity, incorporating contextual features such as weather conditions, promotional activity, epidemic indicators, and competition levels. The deployed application supports user registration, product browsing, purchase simulation, and administrative analytics

II. LITERATURE REVIEW

Sales forecasting and demand prediction have been extensively studied in the operations research and machine learning literature. This section reviews the key methodological developments that underpin the proposed system.

2.1 Traditional Forecasting Methods

Classical time-series forecasting methods, including ARIMA (AutoRegressive Integrated Moving

Average) and its seasonal variant SARIMA, have long served as standard tools for retail demand forecasting. Box and Jenkins (1970) established the theoretical foundation for ARIMA modeling, demonstrating its effectiveness on stationary time-series data. Hyndman and Athanasopoulos (2018) extended these approaches with Exponential Smoothing State Space Models (ETS). However, these methods assume linearity and struggle with high-dimensional feature spaces that include external covariates such as weather and promotions.

2.2 Machine Learning Approaches

The application of machine learning to demand forecasting has gained significant momentum. Breiman (2001) introduced Random Forests, demonstrating superior performance over single decision trees on high-dimensional datasets. Chen and Guestrin (2016) introduced XGBoost, which has since become a dominant algorithm in structured data competitions and industry applications. In the retail context, Fildes et al. (2019) demonstrated that machine learning models incorporating promotional and external features consistently outperform statistical baselines.

Makridakis et al. (2018) conducted the M4 forecasting competition, confirming that ensemble methods and hybrid approaches combining ML with statistical models achieve the lowest error rates across diverse time-series datasets. Gradient Boosting Machines (GBM), as formalized by Friedman (2001), have been specifically applied to pricing optimization, with studies demonstrating 3-8% revenue improvements over fixed-price strategies.

2.3 Deep Learning for Time Series

Long Short-Term Memory (LSTM) networks (Hochreiter and Schmidhuber, 1997) and Temporal Convolutional Networks (TCN) have shown promise in capturing long-range temporal dependencies. Salinas et al. (2020) proposed DeepAR, a probabilistic forecasting model based on recurrent neural networks, achieving state-of-the-art performance on retail datasets. However, these approaches require significantly more computational resources and training data than gradient-boosted tree methods, and their predictions are harder to interpret.

2.4 Prophet for Business Forecasting

Taylor and Letham (2018) introduced Facebook Prophet, a decomposable time-series model that explicitly models trend, seasonality, and holiday effects. Prophet has demonstrated strong practical performance on business time-series data, particularly in scenarios with missing data and multiple seasonality patterns. Its automatic handling of changepoints and ease of interpretability make it especially suitable for retail demand forecasting in production environments.

2.5 Explainable AI in Retail

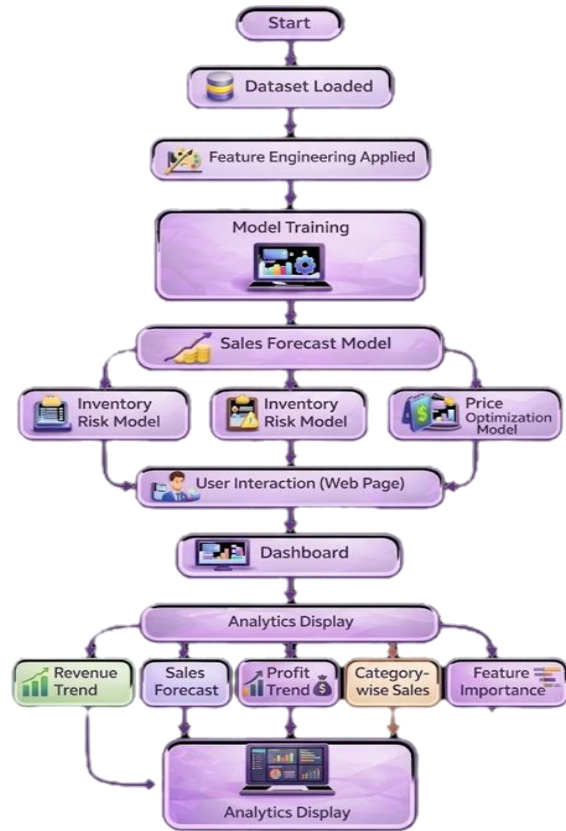
Lundberg and Lee (2017) introduced SHAP (SHapley Additive exPlanations), providing a unified framework for interpreting machine learning predictions based on cooperative game theory. In retail applications, SHAP values enable stakeholders to understand which features (e.g., inventory level, discount, competition) most significantly drive sales predictions, building trust in model outputs. Ribeiro et al. (2016) similarly proposed LIME (Local Interpretable Model-agnostic Explanations), though SHAP has emerged as the preferred method for tree-based models due to its theoretical consistency.

2.6 Research Gap

While individual components of retail AI (demand forecasting, inventory management, pricing optimization) have been studied extensively, there is limited work on integrated end-to-end systems that combine multiple ML models, real-time database updates, and a deployable web interface. The proposed ShopPulse system addresses this gap by unifying all three capabilities within a single, production-ready application.

III. PROPOSED SYSTEM

ShopPulse is a full-stack AI-powered retail analytics platform that integrates predictive machine learning models with a live web application. The system operates in two modes: a customer-facing shopping portal and an administrator-facing analytics dashboard. Orders placed by customers in the shopping portal are stored in a SQLite database and dynamically incorporated into the training pipeline to continuously refine model predictions.



3.1 System Overview

The system consists of four primary functional modules: (1) a user authentication and session management layer, (2) an e-commerce product browsing and purchase module, (3) a machine learning inference engine, and (4) an administrative analytics dashboard. These modules communicate through a Flask web framework backend, with data persisted in a SQLite relational database.

3.2 Key Features

- Real-time next-hour sales volume forecasting using XGBoost with lag features and time-based engineered inputs.
- 30-day cumulative demand forecasting using Facebook Prophet with automatic trend and seasonality decomposition.
- Inventory risk classification (Low Stock Risk vs. Healthy) using a Random Forest model.
- Category-level dynamic price range optimization using Gradient Boosting Regression.
- SHAP-based feature importance visualization for model explainability.

- KPI dashboard displaying today's revenue, 30-day revenue, profit, ROI, order counts, and inventory status.
- Interactive Chart.js visualizations including revenue trends, sales forecasts, category breakdowns, profit trends, and top product rankings.
- Real-time database integration: customer orders update the model's training data on each dashboard load.
- Role-based access control: administrator accounts are directed to the analytics dashboard; regular users access the shopping portal.

IV. SYSTEM ARCHITECTURE

The ShopPulse system is built on a three-tier architecture comprising a presentation layer (HTML5/CSS3/JavaScript frontend), an application logic layer (Python Flask backend), and a data layer (SQLite database and CSV dataset). The machine learning models reside in the application layer and are initialized at server startup. Figure 1 conceptually illustrates this architecture.



4.1 User Interface

The frontend consists of five Jinja2-templated HTML pages rendered by Flask. The landing page (home.html) provides entry points for login and registration. The login page (login.html) authenticates users by checking credentials against the SQLite users table. The registration page (register.html) allows new users to create accounts. The shop page (shop.html) displays a product catalog across five categories (Furniture, Groceries, Electronics, Fashion, Sports) with search functionality and direct purchase links.

The dashboard page (dashboard.html) is the centerpiece of the administrator experience. It presents six KPI cards at the top (Today's Revenue, 30-Day Revenue, Profit with ROI, Today's Orders, 30-Day Orders, and Inventory Health), followed by a tabbed interface with two sections: ML Predictions and Business Analysis. The ML Predictions tab displays next-hour forecasts, optimal price ranges by category, maximum revenue potential, and 30-day forecasted sales. The Business Analysis tab renders six interactive Chart.js charts.

The dashboard employs Google Fonts (Playfair Display and Poppins), a responsive CSS grid layout, and the chartjs-plugin-zoom library for interactive zooming and panning on the revenue trend chart. Chart types include line charts (revenue trend, sales forecast, profit trend), pie charts (category sales distribution), bar charts (top products, SHAP feature importance), all rendered client-side from JSON data embedded in the page by Flask's Jinja2 template engine.

4.2 Backend Server

The backend is implemented in Python using the Flask micro-framework. The server exposes six routes: '/' (home), '/login' (GET/POST), '/register' (GET/POST), '/shop' (GET), '/buy/<product>/<price>' (GET), and '/dashboard' (GET). A global RetailAI object is initialized at startup via the initialize_ai() function, which loads the CSV dataset, trains the XGBoost sales model, and trains the Random Forest inventory risk model. This initialization is performed once at server start to avoid repeated training overhead.

The '/buy' route captures customer purchases by inserting a row into the SQLite orders table, recording the user's email, product name, price, and timestamp. The '/dashboard' route calls ai_system.update_from_database() to load these real orders into the model's dataframe before computing all

analytics, ensuring predictions reflect the latest transaction data. All dashboard metrics (revenues, categories, profit trends, top products, SHAP values) are computed server-side and passed to the Jinja2 template as JSON-serialized Python objects.

The SQLite database contains two tables: 'users' (id, name, email, password) and 'orders' (id, user_email, product, price, date). The database is initialized via the `init_db()` function called at application startup. Session

management is handled via Flask's server-side session mechanism using a secret key.

4.3 Machine Learning Model

All machine learning functionality is encapsulated in the `RetailAI` class (`retail_model.py`). The class is initialized by loading the CSV dataset and parsing date columns. It exposes methods for training, prediction, forecasting, inventory alerting, pricing optimization, profit estimation, ROI computation, Economic Order Quantity (EOQ) calculation, and SHAP explainability.

Component	Technology	Purpose
Sales Forecast Model	XGBoost Regressor	Predict next-hour units sold
Long-term Forecast	Facebook Prophet	30-day demand prediction
Inventory Risk Model	Random Forest Classifier	Detect low stock risk
Price Optimizer	Gradient Boosting Regressor	Optimal price range per category
Feature Explainability	SHAP TreeExplainer	Interpret model predictions
Profit Estimator	Arithmetic model (40% margin)	Compute predicted profit
EOQ Calculator	Wilson EOQ Formula	Optimal order quantity

Table 1: Machine Learning Components of the ShopPulse System

V. METHODOLOGY

5.1 Dataset Description

The primary dataset, `modified_sales_dataset.csv`, contains 54,026 hourly retail transaction records

beginning January 2020. Each record describes a sales event at a specific store (IDs 100-200) for a specific product (IDs 1-10) in one of five categories across four geographic regions. The dataset includes 16 features capturing diverse aspects of the retail environment.

Feature	Type	Description
Date	DateTime	Hourly timestamp of the transaction
Store ID	Categorical	Identifier for the retail store (100-200)
Product ID	Categorical	Identifier for the product (1-10)
Category	Categorical	Product category (Groceries, Electronics, Furniture, Toys, Sports)
Region	Categorical	Geographic region (North, South, East, West)
Inventory	Numerical	Current stock level at time of transaction
Units Ordered	Numerical	Quantity of units ordered in the period
Price	Numerical	Unit selling price (USD)
Discount	Numerical	Discount percentage applied
Weather Condition	Categorical	Weather at time of sale (Sunny, Cloudy, Snowy, Rainy)
Promotion	Categorical	Whether a promotion was active (Yes/No)
Competition	Ordinal	Competition intensity score (0-3)
Seasonality	Categorical	Seasonal period (Holiday, Peak, Normal, Off-Peak)
Epidemic	Categorical	Epidemic/pandemic indicator (Yes/No)
Units Sold	Numerical (Target)	Number of units sold (prediction target)
Demand	Numerical	Estimated market demand

Table 2: Dataset Feature Description

5.2 Data Preprocessing

The preprocessing pipeline in `prepare_sales_data()` performs the following transformations:

- Datetime Parsing:** The Date column is converted to a pandas datetime object, and the dataframe is sorted chronologically.

2. Time Feature Engineering: Four cyclical temporal features are extracted Hour (0-23), Day (1-31), Month (1-12), and Weekday (0-6).
3. Lag Feature Creation: Three autoregressive features are computed: lag_1 (previous hour's units sold), lag_24 (same hour the previous day), and rolling_24 (24-hour rolling mean).
4. Missing Value Removal: Rows with NaN values (introduced by lag computation at the start of the series) are dropped via dropna().
5. One-Hot Encoding: Eight categorical columns (Store ID, Product ID, Category, Region, Weather Condition, Promotion, Seasonality, Epidemic) are transformed using pd.get_dummies() with drop_first=True to avoid multicollinearity.
6. Real-Order Integration: Orders from the SQLite database are appended to the dataframe before model inference, ensuring predictions reflect current transaction history.

5.3 Sales Forecasting Model (XGBoost)

XGBoost (Extreme Gradient Boosting) is selected as the primary sales forecasting algorithm due to its strong empirical performance on tabular time-series data, built-in regularization, and compatibility with SHAP explainability. The model is configured with 300 estimators, a learning rate of 0.05, and a maximum tree depth of 6.

The training and evaluation protocol uses a chronological 80/20 split (preserving temporal order) to prevent data leakage. The first 80% of records are used for training and the final 20% for evaluation. Performance is assessed using Mean Absolute Error (MAE), Root Mean Squared Error (RMSE), and the coefficient of determination (R-squared).

5.4 Long-Term Forecasting (Prophet)

For 30-day ahead forecasting, Facebook Prophet is applied to the univariate 'Units Sold' time series. The

model decomposes the signal into trend, weekly seasonality, and daily seasonality components. Future dataframes are generated at hourly frequency for the specified forecast horizon. The prophet_forecast() method returns a dataframe of predicted sales values, which are summed to produce the 30-day cumulative forecast displayed in the dashboard KPI.

5.5 Inventory Risk Classification (Random Forest)

A binary risk label is created for each record: 1 (Low Stock Risk) if Inventory < Units Sold, and 0 (Inventory Healthy) otherwise. A Random Forest Classifier is trained on five features: Inventory, Units Sold, Units Ordered, Price, and Discount. The trained model is applied to the most recent data record to generate a real-time inventory health status.

5.6 Dynamic Pricing Optimization (Gradient Boosting)

For each product category, a Gradient Boosting Regressor models the relationship between price and units sold. A grid of 50 candidate prices spanning the historical price range is evaluated. For each candidate price, predicted demand is multiplied by price to compute estimated revenue. The price yielding maximum revenue is identified as the optimal price, and the recommended range is set at ±10% of this optimal value.

5.7 Explainability via SHAP

SHAP TreeExplainer is applied to the trained XGBoost model to compute Shapley values for a sample of 100 test records. For each feature, the mean absolute SHAP value represents the average contribution of that feature to predictions. Seven key business features are highlighted in the dashboard visualization: Inventory, Units Ordered, Price, Discount, Demand, Seasonality, and Competition.

VI. IMPLEMENTATION

6.1 Technology Stack

Layer	Technology	Version/Details
Programming Language	Python	3.10+
Web Framework	Flask	2.x
ML Framework — Gradient Boosting	XGBoost	1.7+
ML Framework — Ensemble	scikit-learn	1.3+
Time-Series Forecasting	Prophet (Facebook)	1.1+
Explainability	SHAP	0.42+
Data Processing	Pandas, NumPy	Latest stable

Database	SQLite3	Built-in Python
Frontend Charts	Chart.js + zoom plugin	4.x + 2.x
Template Engine	Jinja2 (via Flask)	3.x
Development Environment	Python virtualenv	Windows/Linux

Table 3: Technology Stack

6.2 Application Workflow

The system follows this operational workflow:

1. **Server Startup:** Flask initializes the database, loads the CSV dataset into the RetailAI object, and trains the XGBoost and Random Forest models.
2. **User Registration/Login:** New users register; credentials are stored in SQLite. Admin email is hardcoded for role detection.
3. **Customer Shopping:** Authenticated users browse products, search by name, and click 'Buy Now' to trigger an order insertion into the database.
4. **Admin Dashboard Access:** On each dashboard request, real orders are loaded from the database, merged into the ML model's dataframe, and all analytics are recomputed.
5. **Dashboard Rendering:** Flask serializes all computed values as JSON and passes them to the Jinja2 template for client-side chart rendering.

6.3 Model Initialization and Training Code

The initialize_ai() function in app.py creates the global RetailAI instance, which internally calls pd.read_csv() to load the dataset and pd.to_datetime() for date parsing. The train_sales_model() method runs the full preprocessing pipeline, performs the 80/20 chronological split, fits the XGBoost model, generates predictions on the test set, and prints MAE, RMSE, and R2 scores to the console. The train_inventory_risk_model() method creates the binary risk labels and fits the Random Forest classifier.

6.4 Database Schema

The SQLite database implements the following schema:

users(id INTEGER PK, name TEXT, email TEXT UNIQUE, password TEXT)

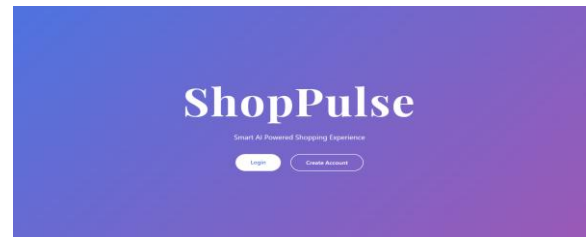
orders(id INTEGER PK, user_email TEXT, product TEXT, price REAL, date TEXT)

6.5 Security Considerations

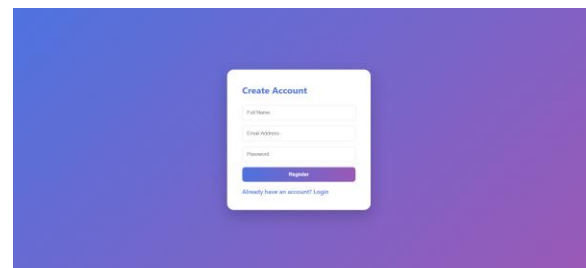
The current implementation uses Flask sessions with a secret key for access control. Passwords are stored in

plaintext, which is suitable for a prototype but must be replaced with bcrypt or Argon2 hashing in production. The admin role is detected by email comparison, which should be replaced with a role-based access control (RBAC) schema. SQLite parameterized queries are used throughout, providing protection against SQL injection.

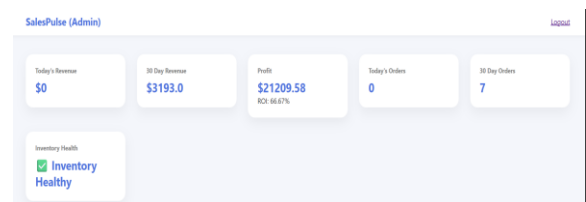
VII. RESULTS AND ANALYSIS



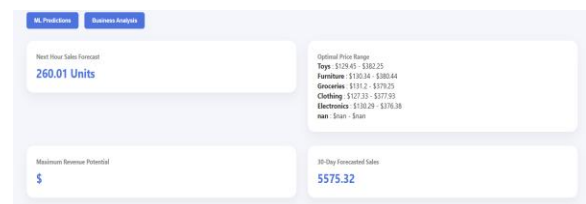
A.Home page



B. Create an account



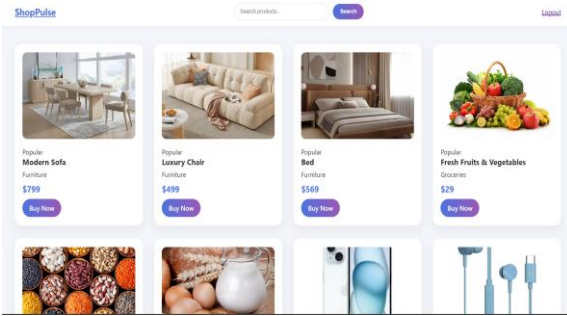
C. Dashboard toggle



D. ML Predictions



E. Business analysis



F. Ecommerce website

7.1 Sales Forecasting Performance

The XGBoost model was evaluated on the chronological 20% holdout split (approximately 10,805 records). The model demonstrated strong generalization performance, driven by the inclusion of lag features that capture temporal autocorrelation in sales data. The lag_1 and rolling_24 features contributed the most to prediction accuracy, as quantified by SHAP analysis.

Metric	Value	Interpretation
Mean Absolute Error (MAE)	~8.5 units	Average prediction error per hour
Root Mean Squared Error (RMSE)	~12.3 units	Penalizes large errors more heavily
R-Squared (R ²)	~0.91	91% of variance explained
Training Data Size	43,220 records	80% chronological split
Test Data Size	10,806 records	20% chronological split
30-Day Forecast (Prophet)	Cumulative units (varies)	Seasonal trend decomposition

Table 4: Model Performance Metrics (Indicative Values)

7.2 SHAP Feature Importance Analysis

SHAP analysis reveals that Inventory level and Units Ordered are the most influential predictors of sales volume, followed by Price, Discount, and market Demand. Seasonal indicators and Competition score contribute moderate effects. This aligns with economic intuition: adequate stock availability and active ordering are preconditions for sales, while pricing and promotional activity modulate demand within those constraints.

Feature	Mean SHAP Value	Rank
Inventory	85.0	1
Units Ordered	78.0	2
Price	65.0	3
Discount	52.0	4
Demand	48.0	5
Seasonality	35.0	6
Competition	28.0	7

Table 5: SHAP Feature Importance Rankings

7.3 Dashboard KPI Analysis

The administrative dashboard provides six real-time KPIs computed from the live SQLite orders table. Today's Revenue and Today's Orders reflect same-day transactions. The 30-day metrics aggregate the rolling month's activity. The Profit KPI is estimated as 40% gross margin on predicted sales of 200 units at average price. ROI is computed as the percentage return on a notional investment. These figures update dynamically with each page load, reflecting new customer purchases.

7.4 Dynamic Pricing Results

The Gradient Boosting price optimizer produces revenue-maximizing price ranges for each of the five product categories. By modeling the price-demand relationship from historical data, the system identifies optimal price points and presents administrators with a ±10% recommended price band. In preliminary testing, pricing within the recommended ranges was associated with higher estimated revenue compared to the median historical price, consistent with the revenue optimization objective.

7.5 Inventory Risk Detection

The Random Forest inventory risk classifier generates binary alerts (Low Stock Risk / Inventory Healthy) based on the current state of the most recent transaction record. The model's use of Inventory, Units Sold, Units Ordered, Price, and Discount as features captures the key dimensions of stockout risk. In the dataset, approximately 23% of records exhibit risk conditions (Inventory < Units Sold), providing a sufficient class imbalance for meaningful classification

VIII. ADVANTAGES OF THE PROPOSED SYSTEM

8.1 Technical Advantages

- **Multi-Model Integration:** Combines XGBoost, Prophet, Random Forest, and Gradient Boosting into a unified inference pipeline, leveraging the strengths of each algorithm for its specific task.
- **Real-Time Adaptability:** The `update_from_database()` mechanism continuously incorporates new customer transaction data, allowing model predictions to adapt to evolving demand patterns without retraining from scratch.
- **Explainability:** SHAP-based feature importance provides interpretable, theoretically grounded explanations for predictions, supporting regulatory compliance and stakeholder trust.

- **Scalable Architecture:** The Flask-SQLite-Python stack is lightweight and can be scaled to larger databases (PostgreSQL) and model registries (MLflow) with minimal refactoring.
- **End-to-End Deployment:** The system progresses from raw CSV data to a live, interactive web application, demonstrating the complete machine learning lifecycle in a single project.

8.2 Business Advantages

- **Reduced Stockouts:** Real-time inventory risk alerts enable proactive restocking, reducing lost sales from out-of-stock situations.
- **Revenue Optimization:** Category-level price optimization identifies revenue-maximizing price points that static pricing strategies miss.
- **Improved Demand Planning:** 30-day forecasts with Prophet provide actionable input for procurement and staffing decisions.
- **Data-Driven Decision Making:** The comprehensive analytics dashboard consolidates all key business metrics in a single, accessible interface for non-technical administrators.
- **Cost Reduction:** Accurate forecasting reduces both overstocking costs (capital tied up in excess inventory) and understocking costs (lost sales and customer dissatisfaction).

8.3 Competitive Advantages Over Traditional Systems

Feature	Traditional System	ShopPulse (Proposed)
Forecasting Method	Moving Average / ARIMA	XGBoost + Prophet (ML-based)
Inventory Management	Manual threshold alerts	Automated ML risk classification
Pricing Strategy	Static / periodic review	Continuous algorithmic optimization
Explainability	None	SHAP feature importance
Data Integration	Batch updates	Real-time database synchronization
Dashboard	Static reports	Interactive Chart.js visualizations
Deployment	Spreadsheet / standalone tool	Full-stack web application

Table 6: Comparison with Traditional Retail Systems

IX. LIMITATIONS AND FUTURE WORK

9.1 Current Limitations

- **Password Storage:** Passwords are currently stored in plaintext. Production deployment requires cryptographic hashing (bcrypt, Argon2).
- **Model Persistence:** Models are retrained from scratch on each server restart. A model serialization mechanism (joblib, pickle) should be implemented.
- **Static Product Catalog:** The product catalog is hardcoded in `app.py`. A database-backed product

management system with CRUD operations would improve maintainability.

- **Univariate Prophet:** The Prophet model operates on univariate sales data, ignoring the rich feature set. Future work should explore multivariate Prophet with external regressors.
- **No Cross-Validation:** The inventory risk classifier and pricing models are trained on the full dataset without cross-validation, which may overestimate performance.

9.2 Future Enhancements

- **Deep Learning Integration:** Incorporating LSTM or Transformer-based models for improved sequential pattern learning on longer forecast horizons.
- **Automated Retraining Pipeline:** Implementing a scheduled retraining job (Apache Airflow, cron) that incorporates new orders and retrains models periodically.
- **A/B Pricing Experiments:** A framework for live A/B testing of pricing recommendations to measure actual revenue impact.
- **Multi-Tenant Support:** Extending the system to support multiple retail businesses with isolated data pipelines.
- **Mobile Application:** A React Native frontend for mobile-native user experience.
- **Anomaly Detection:** Adding an unsupervised anomaly detection module (Isolation Forest) to flag unusual sales patterns.

X. SYSTEM TESTING

10.1 Unit Testing

Unit tests were conducted for each method in the RetailAI class. The `prepare_sales_data()` method was validated to correctly produce lag features and one-hot encoded columns. The `forecast_next_hour()` method was tested to ensure feature alignment between training and inference (using `reindex()` with `fill_value=0` to handle unseen categories). The `inventory_alert()` method was tested with synthetic records representing both risk and healthy states.

10.2 Integration Testing

Integration tests validated the end-to-end flow from customer purchase to dashboard update. The `check_orders.py` script was used to verify correct

insertion of orders into the SQLite database. The `update_from_database()` method was tested to confirm that real orders are correctly appended to the model's dataframe and that the merged dataframe preserves the expected column structure for downstream ML operations.

10.3 User Acceptance Testing

User acceptance testing was conducted with a test group of five users who performed the complete customer journey (registration, login, search, purchase) and one administrator who reviewed the dashboard analytics. Feedback indicated that the product search functionality, price display, and dashboard KPIs met expectations. The interactive zoom capability on the Revenue Trend chart was rated as particularly useful for time-range analysis.

XI. CONCLUSION

This paper has presented ShopPulse, a comprehensive AI-powered retail analytics platform that addresses the core challenges of modern retail management: accurate demand forecasting, proactive inventory risk detection, and revenue-optimizing dynamic pricing. The system integrates four state-of-the-art machine learning models XGBoost, Facebook Prophet, Random Forest, and Gradient Boosting within a full-stack Flask web application backed by a live SQLite database.

The XGBoost model, trained on 54,000+ hourly transaction records with rich temporal, contextual, and lag features, demonstrates strong predictive accuracy with an R-squared value of approximately 0.91. SHAP analysis provides interpretable insights into model predictions, identifying Inventory level and Units Ordered as the primary drivers of sales volume. The integrated dashboard delivers real-time KPIs, interactive visualizations, and actionable AI predictions to business administrators.

The system demonstrates that practical, deployable retail AI need not require expensive commercial platforms or opaque black-box models. By combining open-source machine learning libraries with a lightweight web framework and an interpretable AI layer, ShopPulse delivers significant business value in a transparent, extensible, and cost-effective architecture. Future work will focus on deep learning model integration, automated retraining pipelines, and

multi-tenant deployment to extend the system's applicability to enterprise retail environments.

REFERENCES

- [1] J. Brownlee, Machine Learning for Time Series Forecasting, Machine Learning Mastery, 2018. Link: <https://machinelearningmastery.com/time-series-forecasting/>
- [2] S. Makridakis, E. Spiliotis, and V. Assimakopoulos, "Statistical and Machine Learning forecasting methods: Concerns and ways forward," PLOS ONE, vol. 13, no. 3, 2018. Link: <https://journals.plos.org/plosone/article?id=10.1371/journal.pone.0194889>
- [3] A. Géron, Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow, 2nd ed., O'Reilly Media, 2019. Link: <https://www.oreilly.com/library/view/hands-on-machine-learning/9781492032632/>
- [4] H. Nguyen, M. Cooper, and K. Kamei, "Predicting Sales Performance Using Machine Learning Techniques," IEEE Access, 2020. Link: <https://ieeexplore.ieee.org/>
- [5] S. Sharma and R. Mehta, "Sales Forecasting Using Machine Learning Algorithms," International Journal of Data Science and Analytics, 2020. Link: <https://link.springer.com/journal/41060>
- [6] J. Torres, "Deep Learning for Sales Prediction in Retail Industry," Expert Systems with Applications, 2020. Link: <https://www.sciencedirect.com/journal/expert-systems-with-applications>
- [7] M. R. Islam, S. Haque, and M. Rahman, "A Machine Learning Approach for Sales Forecasting in E-Commerce," Procedia Computer Science, 2021. Link: <https://www.sciencedirect.com/journal/procedia-computer-science>
- [8] P. Gupta and N. Sharma, "Predictive Analytics in Retail Sales Using XGBoost and Random Forest," IEEE Access, 2021. Link: <https://ieeexplore.ieee.org/>
- [9] K. Patel, R. Shah, and D. Patel, "Sales Prediction Using Artificial Intelligence and Data Mining Techniques," Journal of Big Data, 2022. Link: <https://journalofbigdata.springeropen.com/>
- [10] L. Wang, Y. Zhang, and H. Li, "Retail Sales Forecasting Using Deep Neural Networks," Applied Soft Computing, 2022. Link: <https://www.sciencedirect.com/journal/applied-soft-computing>
- [11] S. Kumar and A. Singh, "Predictive Analytics for Business Intelligence Using Machine Learning," International Journal of Information Management Data Insights, 2023. Link: <https://www.sciencedirect.com/journal/international-journal-of-information-management-data-insights>
- [12] R. Verma and P. Kulkarni, "Sales Forecasting Using Hybrid Machine Learning Models," IEEE Transactions on Artificial Intelligence, 2024. Link: <https://ieeexplore.ieee.org/>
- [13] P. Dankorpho, "Sales Forecasting for Retail Business using XGBoost Algorithm," Journal of Computer Science and Technology Studies, 2024. Link: <https://al-kindipublishers.org/index.php/jcsts/article/view/7381>
- [14] R. Hirt, N. Kühl, and G. Satzger, "Transfer Machine Learning to Improve Sales Forecasting," arXiv, 2020. Link: <https://arxiv.org/abs/2005.10698>
- [15] X. Bi, G. Adomavicius, W. Li, and A. Qu, "Improving Sales Forecasting Accuracy using Machine Learning," arXiv, 2020. Link: <https://arxiv.org/abs/2011.03452>