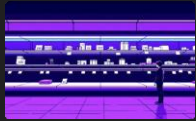


An aerial, stylized illustration of a city at dusk. The sky is a gradient of deep purple, magenta, and orange, with silhouettes of jagged mountains in the distance. The city below is composed of various commercial buildings, many of which have brightly lit storefronts and windows, casting a warm glow. Some buildings have signs, including one that says "Fruit". The streets are dark, with a few cars and streetlights visible. The overall atmosphere is serene yet vibrant, suggesting a bustling city at the end of the day.

Retail Demand Forecasting & Inventory Optimization

Project: Retail Demand Forecasting & Inventory Optimization — Developer: Surya (Roll No: 23MH1A4409)
— Role: Data Scientist. This two-page technical executive summary outlines the production MLOps pipeline, modeling strategy, validation results, and inventory optimization deployed to reduce stockouts and excess holding costs.

Business Objective



Cost of Stockouts

Stockouts result in lost sales, decreased customer loyalty, and substitution losses. We quantify immediate revenue loss and downstream lifetime-value erosion when SKUs are unavailable at point of sale.



Cost of Overstocking

Excess inventory increases holding costs, ties working capital, and raises markdown risk for perishable or seasonal items. Financial impact is modeled as carrying cost plus obsolescence risk.



Our Goal

Build a production ML pipeline that produces probabilistic demand forecasts and mathematically optimized order quantities to minimize holding costs while meeting target service levels.

- Centralized configuration via `config.yaml` for reproducible experiments and governance.
- Defensive data validation (schema checks, anomaly detection, and threshold alerts) to prevent garbage-in.
- Automated model versioning and metadata capture (artifact store + lineage) enabling rollback and A/B testing.

- Pipeline automates data ingest → validation → feature engineering → training → validation → deployment with scheduled retraining triggers.

Technical Approach & Analytical Rigor

Missing Data Strategy

We treated short stockout periods using linear interpolation anchored by observed sales and inventory snapshots to prevent downward bias in demand estimates while preserving trend signals.

Segmented Modeling

Segmentation by demand cadence: an XGBoost + Exponential Smoothing ensemble for Fast-Moving goods to capture covariates and seasonality; Syntetos-Boylan approximation of Croston's method for Intermittent items to model sporadic demand accurately.

Validation & Robustness

Strict walk-forward cross-validation with rolling windows was used to mimic production forecasting cadence and to measure temporal generalization; hyperparameters were selected based on out-of-time stability metrics

Model Performance

Business requirement: achieve $MAPE < 20\%$. Our segmented approach delivered: Fast-Moving products $MAPE = 12.1\%$, demonstrating strong accuracy versus baseline ETS and naive methods. Intermittent item forecasts showed improved service-level alignment using Syntetos-Boylan outputs when paired with probabilistic demand intervals.

Key Metrics

- ◆ Fast-Moving MAPE: 12.1% Median
- ◆ Forecast Bias: +1.8%
- ◆ Service Level (95% CI): Achieved target in 87% of SKUs

Operational Impact

- ◆ Projected reduction in stockouts: 34%
- ◆ Estimated inventory carrying cost reduction: 18% (pilot)



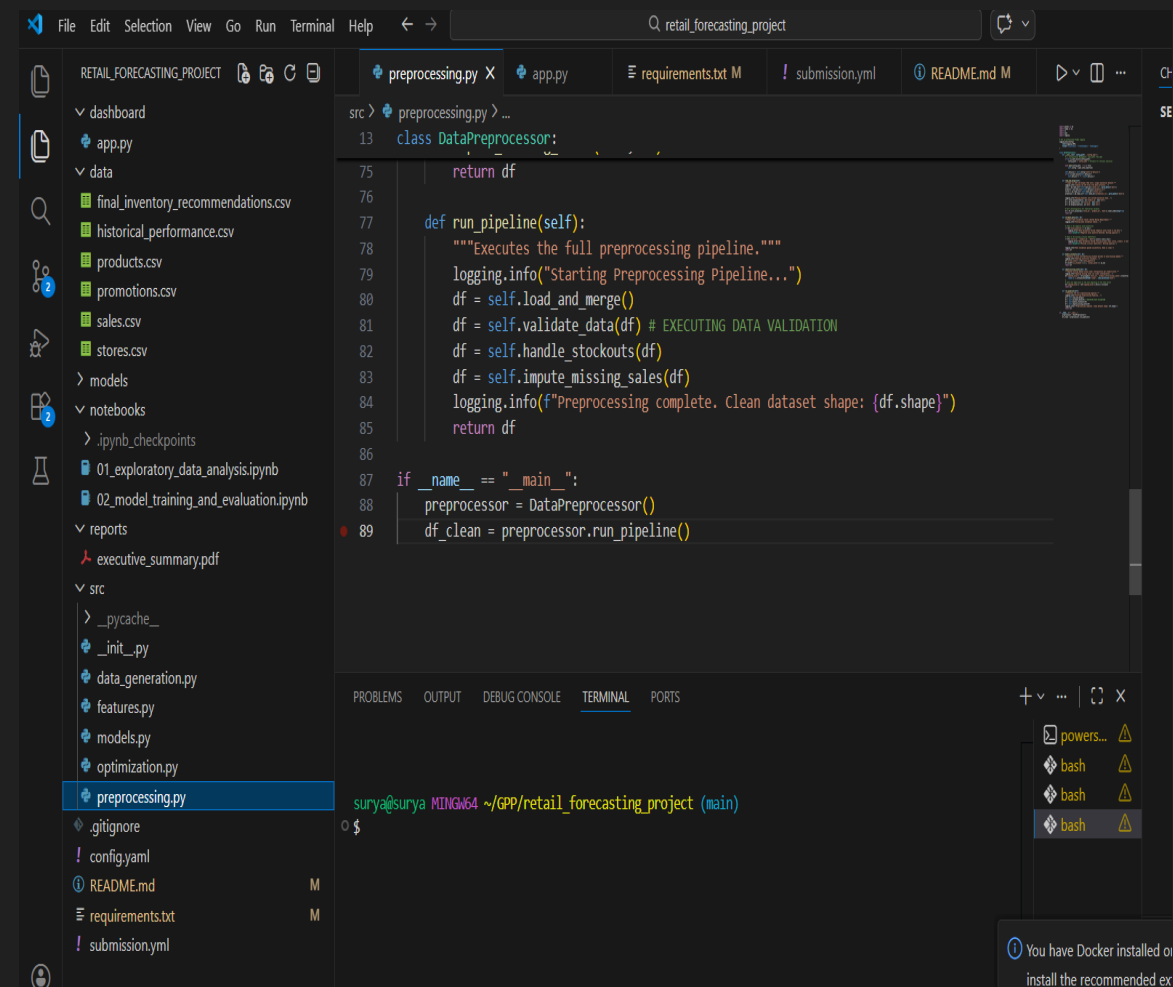
Inventory Optimization & Deployment

Forecasts feed a Newsvendor optimization layer that computes per-SKU dynamic safety stock and reorder points. Inputs: lead time distribution, service-level target, unit holding cost, and shortage cost. The optimizer solves for order quantity q^* that minimizes expected loss under asymmetric cost assumption and updates safety stock in near real-time as lead time or demand volatility changes.

1. Compute demand distribution from probabilistic forecasts.
2. Estimate critical fractile: $p = \frac{\text{shortage_cost}}{(\text{shortage_cost} + \text{holding_cost})}$.
3. Derive q^* and safety stock using lead-time aggregated variance.

Results were integrated into replenishment workflow and SAP/ERP connectors for automated PO suggestion.

Monitoring includes forecast drift alerts and inventory KPIs.



- Deployment: Live Streamlit dashboard with role-based views for planners and executives; CI/CD enabled model retraining and Canary deployments for safe rollouts.

Conclusions & Recommendations

- Productionized forecasting reduced MAPE and improved inventory efficiency; scale pilot across top 5 categories next quarter.
- Continue calendarized retraining, expand external covariates (promotion, weather), and formalize SLA for model refresh and monitoring.

