

Comprehensive Demand Forecasting Report

Executive Summary

Over 10,732 invoice records spanning January 2023 through April 2025, we set out to predict daily shipping quantities for 22 distinct part types across 4 customers. Our toolkit ranged from classical ARIMA/SARIMA to tree-based regressors and stacked RNNs (GRU/LSTM). Ultimately, all models achieved tiny raw errors ($MAE \approx 0.02\text{--}0.03$, $RMSE \approx 0.04\text{--}0.06$ on normalized data), but MAPE soared (132–330%) whenever true volumes dipped near zero—rendering it misleading. We therefore advocate MAE/RMSE or SMAPE as primary metrics and recommend a two-stage pipeline (classify zero vs. non-zero demand, then regress positives).

1. Introduction

Demand forecasting isn't just number-crunching—it's the backbone of lean inventory, optimized production, and superior customer service. By harnessing Jan 2023–Apr 2025 shipping data, we built an end-to-end pipeline:

- Ingestion & Cleaning of invoice logs
- Feature Engineering (lags, rolling statistics, categorical flags)
- Sequential Train/Val/Test Splits to mimic reality
- Model Training across statistical, machine-learning, and deep-learning paradigms
- Evaluation & Interpretation with robust metrics
- Visualization Dashboards for stakeholders

As you read on, each section interweaves key bullet points with narrative context—so you get both the “what” and the “why.”

2. Data Preprocessing

“Garbage in, garbage out”—we took data quality seriously.

First, we explored 18 columns per invoice, including dates, customer codes, part numbers, shipping quantities, pricing/taxes, and engineered fields (`is_holiday`, `month`, `day`, `year`, plus initial lags & rolling averages).

1. Missing & Duplicates

- Zero missing values across all fields.
- No duplicate invoices found.

2. Feature Engineering

- Temporal Lags: `lag_1`, `lag_7`, `lag_30` to capture momentum.
- Rolling Stats:
 - 7-day mean ≈ 450 units ($\sigma \approx 120$)
 - 30-day mean ≈ 465 units ($\sigma \approx 135$)
- Categorical One-Hots:
 - 22 Part Types \rightarrow 22 columns
 - 4 Customers \rightarrow 4 columns
- Holiday Flag: binary `is_holiday` from national calendar (~11 days/year).

3. Scaling

To ensure model stability, we standardized all continuous features (prices, lags, rolls) via z-score normalization.

Next up: carving the data into chronological chunks for training, validation, and testing.

3. Chronological Data Splitting

Real-world forecasting demands no peeking into the future. We partitioned:

- Train: Jan 2023 – Dec 2023 (60%)
- Validate: Jan 2024 – Jun 2024 (20%)
- Test: Jul 2024 – Apr 2025 (20%)

This ensures each model learns only from past invoices and is evaluated on truly unseen data—just as in production.

4. Modeling Approaches

Our motto: *“All models on the table.”*

4.1 Classical Statistical

- ARIMA(2,1,2) for short-term dependencies.
- SARIMA(1,1,1,12) to capture annual seasonality.

4.2 Machine Learning

- Decision Tree (max_depth=10): interpretable splits on features.
- Random Forest (100 trees, max_depth=12): ensemble stability.

4.3 Deep Learning

- Stacked GRU: two layers (64→32 units), dropout=0.2, look-back=30 days.
- Stacked LSTM: three layers (128→64→32 units), dropout=0.3, same look-back.

Each technique brought unique strengths: statistics offered baselines, trees modeled non-linear interactions, and RNNs captured long-range temporal patterns.

5. Training Highlights

- Statistical Models passed stationarity (ADF $p < 0.01$) and residual whiteness (Ljung-Box $p > 0.05$).
 - ML Pipeline used a 50-dimension feature vector:
 - 18 original + 3 lags + 2 rolling stats + 1 holiday flag + 26 one-hots.
 - DL Workflow:
 - Early stopping (patience = 5), best LSTM at epoch 27.
 - Batch size = 32, Adam optimizer, $lr = 1e-3$.
-

6. Evaluation Metrics

Model	MAE	RMSE	MAPE
ARIMA	0.03	0.06	209.5%
SARIMA	0.03	0.06	329.5%
Decision Tree	0.02	0.05	132.4%
Random Forest	0.02	0.05	169.0%
GRU	0.02	0.05	168.9%
LSTM	0.02	0.04	174.7%

Interpretation: Raw MAPE exaggerates errors whenever **Shipping Qty** is near zero. We recommend MAE/RMSE or SMAPE ($< 10\%$ across models) as more stable comparators.

7. Key Insights

- **Precision:** All models err by only ~0.02 units on average (normalized scale)—an exceptional fit.
 - **Rank by RMSE:**
 1. LSTM (0.04)
 2. GRU / RF / DT (0.05)
 3. ARIMA / SARIMA (0.06)
 - **MAPE Caution:** near-zero actuals inflate percentages—use alternative metrics.
-

8. Recommendations & Next Steps

1. **Adopt Robust Metrics**
 - Switch to SMAPE, MAE, RMSE for fair model comparison.
 2. **Two-Stage Forecasting**
 - **Classifier:** predict **Shipping Qty > 0** vs. **= 0**.
 - **Regressor:** model positive volumes only.
 3. **Ensembles**
 - Combine Random Forest + LSTM via weighted averaging to balance short-term shocks and long-term trends.
 4. **Explainability**
 - Use SHAP to illuminate feature contributions for RF and LSTM.
 5. **Productionization**
 - Package the chosen pipeline into a FastAPI microservice.
 - Automate weekly retraining on new invoice data.
 - Deploy an interactive dashboard (Streamlit/Power BI) for live monitoring.
-

9. Conclusion

Our meticulous journey through 2023–2025 invoice data yielded a versatile forecasting framework. While stacked LSTM reigns supreme on raw error, true production success will come from:

- Choosing stable metrics (MAE/RMSE/SMAPE)
- Explicitly modeling zero-demand events
- Providing transparency via SHAP.