# From Checkout to Forecast:

## A Retail Data Deep Dive

Dido De Boodt
June 2025

# Project Overview

**02**

**01**

**03**

**Business Goal**

Improve stock and promo planning by predicting sales per store–item–date

**Dataset**

Ecuadorian retailer with over 120 million rows.

**Focus**

Guayas region and top 3 item families

# Project Steps

| Data Preparation | Features | EDA | Models | App |
|---|---|---|---|---|

**Data Preparation**

Missing values and dates

Outliers and anomalies

**Features**

Lags

Rolling Windows

Time based Features

**EDA**

Time Series Behavior Analysis

Seasonal, holiday, promotion, oil, store and location

**Models**

XGBoost

LSTM

ARIMA and SARIMA
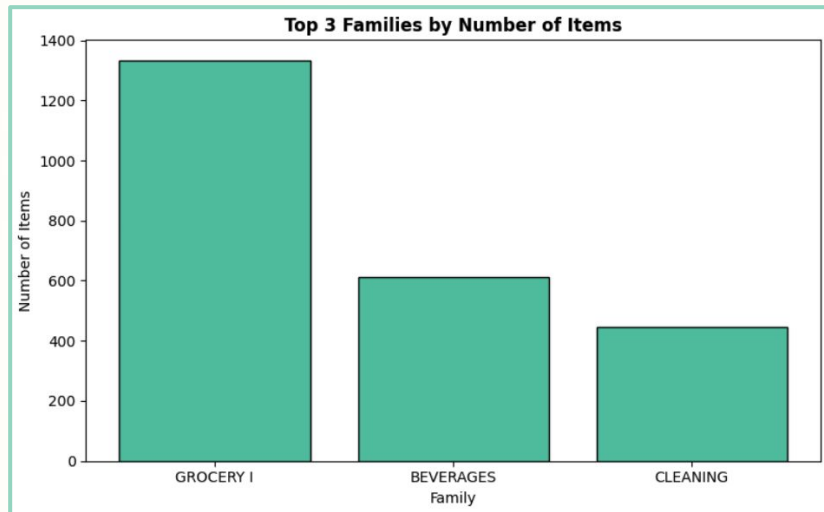
**App**

Streamlit App Deployment

Live Forecast Output

10–Day Forecast and download

# About the Data

## Key Steps

- *Merged: items, stores, and holidays*

- *Cleaned missing values & outliers*

- *Created lag, rolling avg, calendar features*



Top 3 Families by Number of Items

**Shape**: 3.39M rows × 41 cols

**Date range**: Jan 2013 − Mar 2014

**Stores**: 10

**Items**: 1,127

**Avg daily sales**: 3.97

**Zero sales**: 47.25%

# From Patterns to Predictions

## RMSE 7.77

### XGBoost

XGBoost performs well on large-scale tabular time series data, especially when engineered features include lag and rolling metrics.

- Delivered the best RMSE
- Performed well due to extensive feature engineering.
- Handles missing values and complex interactions effectively.

## RMSE 7.80

### LSTM

The LSTM model captures trend and seasonality reasonably well.

- LSTM did not outperforming XGBoost dramatically.
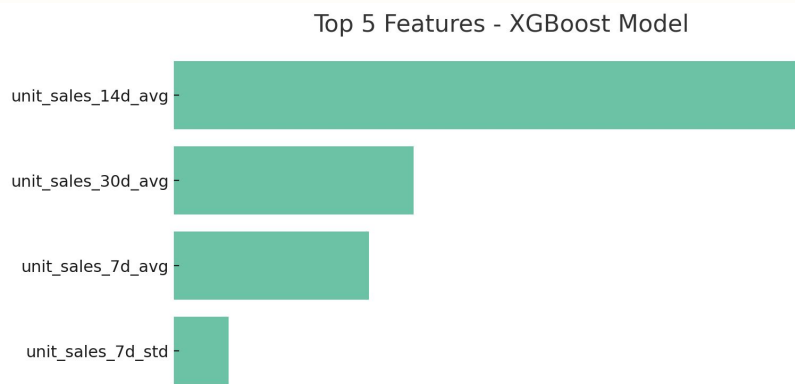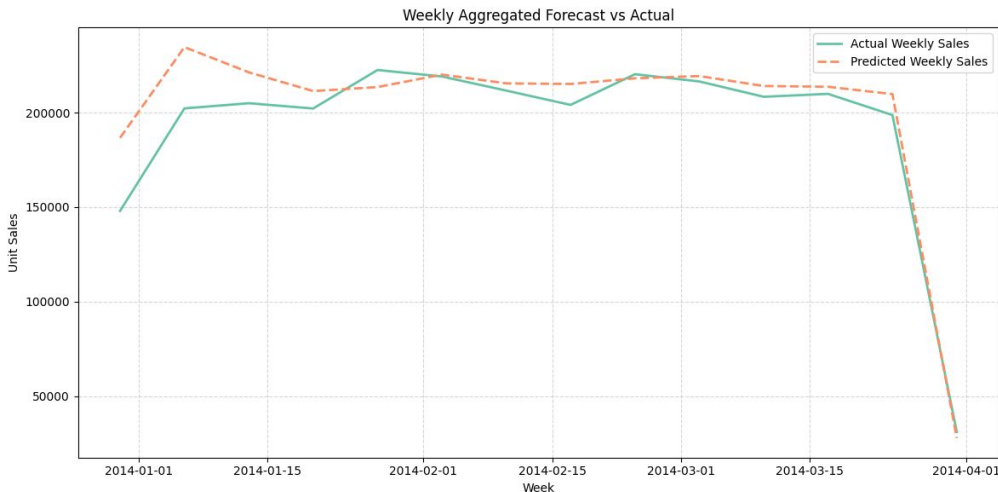- Improvements could include tuning sequence length and feature selection.

## RMSE 7.88

### SARIMA

The SARIMA model successfully captures both short-term and weekly seasonal dependencies in the sales data

- Captured weekly seasonality and trend components well.
- Higher RMSE due to limitations in handling features
- Strong interpretability, but less flexible for complex data.

# XGBoost Deep Dive



Weekly Aggregated Forecast vs Actual



Top 5 Features - XGBoost Model

## Why XGBoost Stood Out

- Lowest RMSE (7.77) and lowest bias among all models

- Effectively leveraged lag-based and rolling window features

- Adapted well to sparse & skewed data (47% zero sales)

- Weekly predictions aligned closely with true values

- Transparent model: top predictors are interpretable

# Evaluation Metrics

| Model | RMSE | MAD | SMAPE | Bias |
|---|---|---|---|---|
| XGBoost | 7.77 ✅ | 3.52 | 120.69% | 0.173 ✅ |
| LSTM | 7.80 | 3.49 ✅ | 121.55% ✅ | 0.195 |
| SARIMA | 7.89 | 3.53 | 122.40% | 0.189 |

Key Takeaways:

- **XGBoost** had the best RMSE and Bias → most accurate and least skewed.

- **LSTM** achieved lowest MAD and SMAPE → most consistent error.

- **SARIMA** performed well capturing seasonality, but had slightly higher errors overall.

# App Demo (Click here for the app)



**User Inputs**

**Introduction**

**Prediction**

## Retail Sales Forecasting App

### Predict daily unit sales for a selected store, item, and date.

Welcome to the Retail Sales Forecasting App!
This tool uses machine learning to predict daily unit sales for a specific store-item-date combination, based on historical patterns and calendar events.

Project by **Dido De Boodt**
Built using Python, XGBoost, and Streamlit.
Special thanks to **Kaggle** for the dataset! ❤️

ℹ️ About this model

**Forecasted Sales: 149 units**

### Historical Sales Trend

**Forecast Input**

Select Store
24

Select Item
257847

Select Forecast Date
2014/01/01

# Final Thought & Recommendations

## Key Takeaways

- Time-based features significantly improved results

- ML models outperform classical ones

- LSTM requires tuning & memory usage

- XGBoost has best trade-off accuracy vs explainability

## Recommendations

- Deploy app to support demand planning

- Automate updates and retraining

- Explore Prophet or attention-based models like Transformers

Thank you for your attention!

LinkedIn
GitHub Repo