# Advanced Time Series Forecasting for Optimizing Retail Sales

# in Corporación Favorita Stores

Dare, Kunle

Data Science Capstone: Practicum (MSDS-699-A01) Dr. Eve Thullen October 13, 2024

#### **Advanced Time Series Forecasting for**

#### **Optimizing Retail Sales in Corporación Favorita Stores**

#### Abstract

Retail sales forecasting is a critical challenge in business management, especially for large organizations like Corporación Favorita, a leading grocery chain. Accurate demand predictions are vital for optimizing inventory management, reducing wastage, ensuring product availability, and enhancing customer satisfaction. This study addresses these challenges by developing a sophisticated, data-driven time series forecasting model tailored to predict storelevel sales across thousands of items in multiple locations.

This report presents a comprehensive analysis of advanced time series forecasting techniques aimed at improving sales forecasting accuracy for Corporación Favorita stores. The study leverages a rich dataset from the Kaggle "Store Sales - Time Series Forecasting" competition, containing daily sales, promotions, holidays, oil prices, and other external factors. The key objective of this project is to develop a robust forecasting solution that can handle the intricacies of retail sales, including seasonal patterns, economic factors, promotions, and store-level variations.

The methodology integrates an ensemble of four different models: ARIMA, Prophet, XGBoost, and Long Short-Term Memory (LSTM) neural networks. Each model brings distinct strengths—ARIMA captures linear trends and simple seasonality, Prophet handles multiple seasonalities and holidays, XGBoost captures complex interactions between features, and LSTM

models long-term dependencies. These models were meticulously tuned through hyperparameter optimization, and feature engineering was conducted to enhance their predictive performance. Time-based features such as Fourier terms, lag features, rolling statistics, and interaction terms were engineered to capture seasonality, promotional effects, and regional economic impacts. Oil prices, in particular, were used as a proxy for regional economic health.

A stacking ensemble method was employed to combine the predictions from these models, utilizing a Generalized Linear Model (GLM) as the meta-learner to produce a final, optimized forecast. The ensemble model outperformed the individual models in terms of accuracy and robustness, achieving a Root Mean Squared Logarithmic Error (RMSLE) of 0.512, a Mean Absolute Scaled Error (MASE) of 0.876, and a Mean Absolute Error (MAE) of 2.34.

The evaluation of the model was conducted through time-series cross-validation and holdout validation, ensuring the results were reliable and generalizable. This project highlights several key findings, including the importance of interaction features (such as promotions and holidays) and the role of external economic factors (like oil prices) in driving sales. The study also identified challenges, such as predicting sales for new products and handling unexpected events not captured in historical data.

In terms of practical implications, the project demonstrates the potential for improving inventory management, optimizing promotion planning, and enhancing store-level resource allocation through more accurate sales forecasts. The results provide actionable insights that could help Corporación Favorita reduce waste, improve customer satisfaction, and increase profitability. Future work could focus on incorporating more external data sources (e.g., weather, consumer sentiment), developing models for new products, and implementing real-time online learning for continuous model improvement. This study underscores the power of machine learning and data science in solving complex retail forecasting challenges and driving data-driven decision-making in the retail industry.

# **Table of Contents**

Advanced Time Series Forecasting for		
Optimizing Retail Sales in Corporación Favorita Stores		
Abstract	iii	
1. Introduction	1	
2. Background	2	
3. Problem Statement	2	
4. Project Objectives	3	
5. Significance of the Project	3	
6. Literature Review	4	
7. Methodology	7	
7.1 Data Collection and Preprocessing	7	
7.2 Feature Engineering	8	
7.3 Model Development	10	
7.4 Model Tuning	11	
7.5 Ensemble Method	11	
8. Data Analysis	12	
9. Model Development	14	
10.Evaluation	17	
10.1 Evaluation Metrics	17	
10.2 Validation Techniques	17	
11.Model Comparison	19	
12.Conclusion	20	
13.Future Work	21	
References	27	
Certification of Authorship	30	

# 1. Introduction

The retail industry operates in a highly dynamic and competitive environment, where predicting sales accurately can significantly impact a company's bottom line. For Corporación Favorita, a large grocery retailer operating across multiple stores, the ability to forecast sales at the store-item level is crucial for optimizing inventory management, enhancing promotion strategies, reducing wastage, and improving customer satisfaction. Sales forecasting also plays a central role in operational decision-making, influencing staffing, resource allocation, and supply chain management.

This project focuses on developing an advanced time series forecasting solution aimed at predicting daily sales for thousands of items across multiple stores. The primary objective is to create a robust and scalable model capable of handling the complexities of retail data, including seasonal variations, promotional effects, holidays, and external economic factors. The data used for this study is derived from the Kaggle competition "Store Sales - Time Series Forecasting," which provides a rich dataset of sales information, promotions, holidays, oil prices, and other factors that influence retail performance.

In this report, we explore a combination of traditional statistical models (ARIMA and Prophet), machine learning models (XGBoost), and deep learning models (LSTM) to create an ensemble forecasting approach. The study also delves into the importance of feature engineering, model tuning, and error analysis in improving forecasting accuracy.

#### 2. Background

In the dynamic and competitive retail industry, accurate sales forecasting is crucial for optimizing inventory management, preventing stockouts, reducing waste, and enhancing customer satisfaction and profitability (Fildes et al., 2019). Traditional forecasting methods often fall short in capturing complex patterns and adapting to rapidly changing market conditions. This project aims to address these challenges by developing an advanced machine learning model capable of accurately predicting unit sales for thousands of items across various Favorita stores.

Corporación Favorita, a major Ecuadorian-based grocery retailer, faces the ongoing challenge of balancing inventory levels with customer demand across its diverse product range and multiple store locations. By leveraging cutting-edge time series forecasting techniques and incorporating external factors, this project seeks to create a robust and adaptable forecasting system that can significantly improve Favorita's operational efficiency and strategic decision-making capabilities (Borovykh et al., 2017).

#### 3. Problem Statement

The primary challenge addressed in this project is the development of a highly accurate and scalable forecasting system for retail sales. Specific issues include:

- 1. Handling high-dimensional data with thousands of products across multiple stores.
- Incorporating external factors such as promotions, holidays, and economic indicators into the forecasting model.
- Capturing complex temporal patterns, including multiple seasonalities and long-term trends.

- 4. Developing a model that can adapt to changing market conditions and maintain accuracy over time.
- Balancing model complexity with interpretability to provide actionable insights for business decision-makers.

# 4. Project Objectives

The main objectives of this project are:

- Develop a state-of-the-art time series forecasting model that accurately predicts unit sales for Favorita's product range across all store locations.
- 2. Incorporate external factors such as promotions, holidays, and economic indicators to enhance prediction accuracy (Bandara et al., 2020).
- 3. Implement a hybrid approach that combines traditional time series models with advanced machine learning techniques, including deep learning (Makridakis et al., 2018).
- Achieve a lower Root Mean Squared Logarithmic Error (RMSLE) compared to Favorita's current forecasting methods.
- Provide actionable insights and recommendations for inventory management and promotional strategies based on the model's predictions and feature importance analysis (Lundberg & Lee, 2017).

# 5. Significance of the Project

This project has significant implications for Corporación Favorita and the broader retail industry:

- Improved Inventory Management: Accurate sales forecasts will enable Favorita to optimize stock levels, reducing both stockouts and excess inventory.
- Enhanced Customer Satisfaction: By ensuring product availability, the company can improve customer experience and loyalty.
- **Reduced Waste**: Particularly important for perishable goods, better forecasting can significantly reduce food waste.
- **Optimized Promotions**: Insights from the model can guide more effective promotional strategies.
- Data-Driven Decision Making: The project will provide a framework for incorporating advanced analytics into strategic business decisions.
- **Competitive Advantage**: Implementing cutting-edge forecasting techniques can give Favorita an edge in the competitive retail market.

# 6. Literature Review

Time series forecasting has been a crucial area of study in retail analytics, with significant advancements in recent years. Traditional methods such as ARIMA (AutoRegressive Integrated Moving Average) and exponential smoothing have been widely used in retail forecasting (Box et al., 2015). However, these methods often struggle with the complexity and scale of modern retail data.

Fildes et al. (2019) provided a comprehensive review of retail forecasting research and practice, highlighting the importance of incorporating domain knowledge and external factors

into forecasting models. They emphasized the need for methods that can handle the hierarchical nature of retail data and the ability to produce probabilistic forecasts.

The application of machine learning techniques to sales forecasting has gained significant traction in recent years. Breiman (2001) introduced Random Forests, which have shown strong performance in various forecasting tasks, including retail sales prediction. Gradient Boosting Machines, particularly XGBoost (Chen & Guestrin, 2016), have demonstrated superior performance in many forecasting competitions.

Bandara et al. (2020) proposed a clustering-based approach for forecasting across multiple related time series, which is particularly relevant for retail settings with numerous products and stores. Their method leverages similarities between different time series to improve overall forecast accuracy.

Deep learning models have shown promising results in capturing complex temporal patterns in sales data. Long Short-Term Memory (LSTM) networks, introduced by Hochreiter & Schmidhuber (1997), have been particularly effective in capturing long-term dependencies in time series data.

More recently, Salinas et al. (2020) introduced DeepAR, a probabilistic forecasting model with autoregressive recurrent networks. This model has shown strong performance in retail demand forecasting tasks, particularly when dealing with cold-start problems and intermittent demand patterns.

Temporal Convolutional Networks (TCN), as described by Bai et al. (2018), have emerged as a powerful alternative to recurrent models for sequence modeling tasks, including time series forecasting. TCNs offer advantages in terms of parallelization and the ability to capture long-range dependencies efficiently.

Ensemble methods, which combine predictions from multiple models, have consistently shown superior performance in forecasting tasks. Makridakis et al. (2018) conducted a comprehensive comparison of statistical and machine learning methods for time series forecasting, concluding that hybrid methods often outperform individual models.

Taylor & Letham (2018) introduced Prophet, a procedure for forecasting time series data based on an additive model where non-linear trends are fit with yearly, weekly, and daily seasonality, plus holiday effects. Prophet has shown promise in handling multiple seasonalities and incorporating domain knowledge into the forecasting process.

The importance of feature engineering and incorporating external factors in retail forecasting has been emphasized in numerous studies. Arunraj & Ahrens (2015) demonstrated the value of including weather data in food sales forecasting, while Gur Ali et al. (2009) showed how promotional information could significantly improve forecast accuracy.

Borovykh et al. (2017) proposed a method for conditional time series forecasting using convolutional neural networks, which can effectively incorporate external conditional information into the forecasting process.

The choice of appropriate evaluation metrics is crucial in assessing and comparing forecasting models. Hyndman & Koehler (2006) introduced the Mean Absolute Scaled Error (MASE) as a generally applicable measurement of forecast accuracy without the problems seen in other measures.

In terms of model interpretation, Lundberg & Lee (2017) introduced SHAP (SHapley Additive exPlanations) values, a unified approach to explaining the output of any machine learning model. This method has gained popularity in the forecasting community for its ability to provide insights into complex models.

The literature review reveals a trend towards hybrid approaches that combine the strengths of traditional statistical methods, machine learning techniques, and deep learning models. There is also a growing emphasis on incorporating external factors and domain knowledge into forecasting models. The challenge lies in developing scalable and interpretable models that can handle the complexity and high-dimensionality of retail sales data while providing actionable insights for business decision-makers.

# 7. Methodology

The methodology for this project follows a systematic approach to develop an advanced time series forecasting system for Corporación Favorita's retail sales. The process incorporates multiple stages, from data preprocessing to model evaluation, employing a hybrid approach that combines traditional time series models, machine learning techniques, and deep learning models.

## 7.1 Data Collection and Preprocessing

The dataset for this project was collected from the Kaggle competition "Store Sales -Time Series Forecasting." It contains daily sales data for multiple stores and items, along with additional information such as promotions, holidays, and oil prices. The dataset spans a period of several years, allowing for an in-depth analysis of sales patterns over time.

The preprocessing steps included:

- Handling Missing Values: Missing values in the dataset were addressed through imputation techniques. For numerical columns, missing values were replaced using the median or mean, while missing categorical data was filled using the most frequent value or mode.
- Categorical Encoding: Categorical variables, such as store type and product category, were encoded using label encoding and one-hot encoding. Label encoding was used for variables with ordinal relationships, while one-hot encoding was applied to nominal variables.
- Normalization of Numerical Features: To ensure that features were on a similar scale, numerical features were normalized using z-score normalization, ensuring that all features had a mean of zero and a standard deviation of one.
- Feature Generation: Time-based features were created to capture various seasonal and cyclical patterns. These included day of the week, day of the month, month of the year, and Fourier terms to capture cyclical effects. Lag features and rolling statistics (such as rolling means and standard deviations) were also introduced to capture temporal dependencies in the data.

# 7.2 Feature Engineering

Feature engineering played a pivotal role in improving the predictive accuracy of the models. The following feature sets were developed to capture different aspects of the data:

#### **Time-Based Features:**

**Fourier Terms:** Fourier terms with K=10K=10 were used to capture cyclical patterns in the sales data. These terms helped the model capture seasonality effects, which are critical in retail environments where sales are often driven by weekly and monthly patterns.

Lag Features: Lag features were created for 7, 14, and 28 days to capture the temporal dependencies in the sales data. These features allow the model to understand the impact of past sales on future sales.

**Rolling Means and Standard Deviations:** Rolling statistics were calculated over 7-day and 28-day windows, providing insights into short-term and long-term trends in sales.

#### **Interaction Features:**

**Promotion and Holiday Interactions:** Interaction terms between promotions and holidays were created to capture the compounded effects of these factors on sales. For instance, sales during a holiday period may be higher when a promotion is also active.

**Day of the Week and Month Interactions:** Sales patterns can vary significantly based on the day of the week and the time of the month. Interaction terms between these variables helped capture these variations more effectively.

#### **External Features:**

**Oil Prices:** Oil prices were included as a proxy for the overall economic conditions, particularly in regions where the economy is heavily influenced by

the oil industry. The hypothesis was that higher oil prices could lead to lower consumer spending, affecting retail sales.

**Local Events:** Information about local events, such as festivals or public gatherings, was incorporated to account for regional variations in sales patterns.

# 7.3 Model Development

The ensemble approach adopted for this project consists of four models, each of which was chosen for its ability to capture different aspects of the time series data:

- ARIMA (Autoregressive Integrated Moving Average): ARIMA models are wellsuited for capturing linear trends and seasonal components in time series data. The model was configured to handle seasonal sales patterns and linear dependencies between lagged sales values.
- **Prophet:** Facebook's Prophet model was employed for its ability to capture multiple seasonalities and holiday effects. Prophet is also well-suited for data with missing values and outliers, making it a valuable component of the ensemble.
- XGBoost (Extreme Gradient Boosting): XGBoost is a powerful machine learning model that captures complex interactions between features. It was particularly useful in understanding the importance of various time-based and external features in sales forecasting.
- LSTM (Long Short-Term Memory): LSTM networks are a type of recurrent neural network that excels at modeling long-term dependencies in sequential data. The

LSTM model was designed to capture the temporal dependencies in sales data that traditional models struggle with.

# 7.4 Model Tuning

Each model was meticulously tuned to optimize its performance:

- **ARIMA:** The ARIMA model's order parameters (p, d, q) were optimized using the Akaike Information Criterion (AIC), ensuring that the model captured the underlying seasonal patterns while avoiding overfitting.
- **Prophet:** The seasonality priors and changepoint sensitivity in the Prophet model were fine-tuned to ensure that the model adapted effectively to both long-term trends and short-term seasonal fluctuations.
- **XGBoost:** Hyperparameters for the XGBoost model were tuned using grid search. The final parameters included a learning rate of 0.01, a maximum tree depth of 7, a subsample ratio of 0.8, and a column sampling ratio of 0.8.
- **LSTM:** The LSTM model was tuned to have 128 units in the LSTM layer, with a dropout rate of 0.2 to prevent overfitting. Recurrent dropout was also set at 0.2, and the model was trained using the Adam optimizer

# 7.5 Ensemble Method

After developing and tuning the individual models, a stacking ensemble approach was implemented. This method combines the predictions of the base models (ARIMA, Prophet,

XGBoost, and LSTM) into a final prediction using a Generalized Linear Model (GLM) as the meta-learner. The purpose of the ensemble was to leverage the strengths of each base model and mitigate their individual weaknesses. The final ensemble was designed to capture both linear and non-linear patterns in the data, providing a more comprehensive forecast than any individual model.

The stacking ensemble integrates the strengths of:

- ARIMA, which effectively captures linear and seasonal trends.
- **Prophet**, which is adept at modeling holidays and multiple seasonalities.
- **XGBoost**, which excels at detecting complex feature interactions.
- LSTM, which models long-term dependencies in the data.

By combining these models, the ensemble delivers a more robust and accurate sales forecast, optimizing Corporación Favorita's ability to manage stock levels and respond to market conditions.

# 8. Data Analysis

#### **Exploratory Data Analysis**

Before building the models, an exploratory data analysis (EDA) was conducted to identify patterns, trends, and relationships in the sales data. Key findings from the EDA include:

• **Strong Seasonality:** Sales exhibited strong weekly and monthly seasonality, indicating that shopping behavior follows predictable patterns across different time periods.

- **Promotions:** Promotions had a significant impact on sales volumes, with sharp increases in sales during promotional periods. However, the promotional effect varied by product category and store location.
- Store-Specific Trends: Sales trends varied across different stores, reflecting regional differences in consumer behavior. Some stores showed strong periodicity, while others exhibited irregular sales patterns.
- External Factors: Oil prices, included as an external factor, showed a moderate correlation with sales, indicating that broader economic conditions could influence consumer spending.

## **Time Series Decomposition**

To gain a deeper understanding of the sales data, Seasonal-Trend decomposition using LOESS (STL) was performed. STL separates the time series into three components: trend, seasonal, and residual. This decomposition provided insights into:

- Trend: Long-term movements in sales, reflecting overall growth or decline in consumer demand.
- Seasonality: Regular, repeating patterns in sales driven by weekly, monthly, or annual cycles.
- **Residual:** Irregular fluctuations in sales that could not be explained by the trend or seasonality components.

The decomposition revealed clear seasonal patterns in the sales data, particularly around holidays and weekends. These insights were valuable for feature engineering and model development.

#### **Correlation Analysis**

A correlation analysis was conducted to examine relationships between sales and other variables. Correlation matrices and heatmaps were used to visualize these relationships. Key observations include:

**Positive Correlation between Promotions and Sales:** Promotions had a strong positive impact on sales, with promotional periods leading to higher sales volumes across most product categories.

**Moderate Correlation with Oil Prices:** Oil prices exhibited a moderate correlation with overall sales trends, particularly in regions where economic conditions are heavily influenced by the oil industry.

**Varying Correlations Across Product Categories:** The strength of correlations between sales and time-based features (e.g., day of the week, month) varied across different product categories, suggesting that some products are more sensitive to seasonal factors than others.

#### 9. Model Development

#### ARIMA Model

The ARIMA model was developed as a baseline model for forecasting sales. ARIMA is a traditional statistical model that captures linear relationships and seasonality in time series data. The model's parameters—pp (autoregressive), dd (differencing), and qq (moving average)—

were optimized using the Akaike Information Criterion (AIC), which balances model fit and complexity.

Although ARIMA provided reasonable forecasts for simpler sales patterns, it struggled to capture non-linear relationships and complex interactions between features, which led to the need for more advanced models.

#### **Prophet Model**

Facebook's Prophet model was selected for its ability to handle complex seasonal patterns and holidays. Prophet allows for multiple seasonalities and is robust to missing data, which made it suitable for the retail sales dataset. The model was fine-tuned by adjusting its seasonality priors and changepoint sensitivity, allowing it to adapt to long-term trends and short-term fluctuations.

Prophet performed well in capturing holiday effects and adapting to changing trends, particularly in stores where sales patterns followed distinct seasonal cycles.

#### **XGBoost Model**

XGBoost, a tree-based machine learning model, was chosen for its ability to handle nonlinear relationships and complex feature interactions. XGBoost excels at identifying the most important features that influence the target variable. In this case, XGBoost was able to capture the intricate relationships between promotions, holidays, oil prices, and sales.

The model was tuned using grid search to optimize its hyperparameters, including learning rate, maximum tree depth, subsample ratio, and column sampling ratio. XGBoost demonstrated strong performance in predicting sales for items with complex sales patterns and

provided insights into feature importance, with promotions, holidays, and store-specific factors emerging as key predictors.

#### LSTM Network

The LSTM neural network was designed to capture long-term dependencies in the sales data. LSTMs are a type of recurrent neural network (RNN) that excel at modeling sequential data. The LSTM model was particularly useful for identifying patterns that spanned longer time horizons, such as year-over-year trends and long-term promotional effects.

The LSTM architecture consisted of 128 units in the LSTM layer, with dropout and recurrent dropout rates of 0.2 to prevent overfitting. The model was trained using the Adam optimizer and performed well in capturing complex temporal patterns that were not easily modeled by traditional statistical methods.

#### **Ensemble Integration**

The final model combined the strengths of the individual models through a stacking ensemble approach. By using a Generalized Linear Model (GLM) as the meta-learner, the ensemble was able to produce more accurate and reliable forecasts than any single model. The GLM was trained on the predictions of the base models (ARIMA, Prophet, XGBoost, LSTM), effectively balancing their strengths and weaknesses.

The ensemble approach proved to be robust across different store locations and product categories, delivering superior performance in both accuracy and generalizability.

### **10.Evaluation**

#### **10.1 Evaluation Metrics**

To evaluate the performance of the models, several metrics were used:

- Root Mean Squared Logarithmic Error (RMSLE): RMSLE is a suitable metric for time series data with varying scales. It penalizes large errors more than smaller ones, making it effective for measuring sales forecasts.
- Mean Absolute Scaled Error (MASE): MASE provides a comparison of the model's performance relative to a naive forecast (such as predicting the previous day's sales). It is a scale-independent measure, allowing for easy comparison across models.
- Mean Absolute Error (MAE): MAE measures the average magnitude of errors between predicted and actual values.
- Root Mean Squared Error (RMSE): RMSE is sensitive to large errors and was used to measure the model's overall prediction accuracy.

# 10.2 Validation Techniques

Two primary validation techniques were employed to ensure the robustness and reliability of the model:

• **Time Series Cross-Validation:** A time-series-specific cross-validation approach was used, where the initial training window consisted of 365 days, and the model was tested on the following 30 days. The process was repeated for 12 folds, ensuring that the model was trained and tested on different time windows.

• Holdout Validation: The last three months of data were used as a holdout set, providing an additional measure of model performance on unseen data.

# **11. Model Comparison**

The table below summarizes the performance of the individual models and the final ensemble:

Model	RMSLE	MASE	Training Time
ARIMA	0.589	0.934	1.2 hours
Prophet	0.543	0.901	2.5 hours
XGBoost	0.528	0.892	3.7 hours
LSTM	0.535	0.897	5.2 hours
Ensemble	0.512	0.876	12.6 hours

The ensemble model achieved the lowest RMSLE and MASE, indicating its superiority over individual models in forecasting sales for Corporación Favorita.

#### **Error Analysis**

A detailed error analysis was conducted to identify areas where the model struggled:

- New Products: The model exhibited higher prediction errors for new products with limited historical data. This was a limitation of all models, as they relied heavily on past sales trends.
- **Special Events:** Sales during special events (e.g., one-time promotions or product launches) were more difficult to predict, as these events were not always captured in the training data.

• Irregular Sales Patterns: Some stores exhibited irregular sales patterns, leading to higher forecasting errors. This was particularly true for stores with significant seasonal fluctuations or those heavily influenced by regional events.

Residual analysis showed no significant autocorrelation in the residuals, indicating that the model effectively captured most of the underlying patterns in the data. However, slightly higher errors were observed during holiday periods, which could be attributed to the complex interplay between promotions, holidays, and consumer behavior.

# **12.**Conclusion

This project demonstrated the successful application of advanced time series forecasting techniques for optimizing retail sales predictions at Corporación Favorita stores. The combination of traditional models like ARIMA and Prophet, machine learning models like XGBoost, and deep learning architectures like LSTM yielded a robust and accurate forecasting system. By using a stacking ensemble approach, we leveraged the strengths of each model, resulting in improved accuracy and the ability to generalize well across different store locations and product categories.

Key findings from the project include:

• Seasonality and Promotions: Sales patterns exhibited strong seasonality, with weekly, monthly, and annual cycles. Promotions were a key driver of sales, and their effects varied across different products and locations.

- Feature Engineering: The inclusion of time-based features, external factors like oil prices, and interaction terms significantly improved the performance of the models, highlighting the importance of domain knowledge in retail forecasting.
- Ensemble Model Superiority: The ensemble model outperformed individual models in terms of accuracy and generalizability, demonstrating the value of combining diverse forecasting approaches.
- Challenges with New Products and Special Events: The models struggled with forecasting sales for new products and special events, underscoring the need for future enhancements to handle these scenarios.

Overall, the forecasting system developed in this project provides Corporación Favorita with valuable insights that can enhance inventory management, reduce stockouts and overstock, and improve overall business operations. The models can be integrated into existing business processes to support data-driven decision-making at the store level.

#### **13.Future Work**

While the ensemble model developed in this project demonstrated substantial improvements in sales forecasting accuracy for Corporación Favorita stores, several areas of potential enhancement could further optimize the results and expand the project's applicability in real-world retail operations.

1. **Incorporating Additional External Data:** Currently, the model incorporates economic indicators like oil prices to capture external market influences. However, retail sales are

also affected by various other external factors that could be incorporated into the forecasting system:

- Weather Data: Weather conditions, such as temperature, rainfall, and seasonal changes, significantly impact consumer behavior, especially in the grocery sector. For example, cold weather could drive sales of hot beverages, while heat waves may boost demand for cold drinks and ice cream. By integrating localized weather data into the model, the forecasts could become more responsive to short-term shifts in demand.
- **Competitor Data:** Competitor activity, such as changes in pricing, promotions, and store openings or closings, directly influences market share. Publicly available competitor pricing information or event data could be used as an external feature to account for these effects in the forecasting process.
- Social Media and Consumer Sentiment: Analyzing sentiment data from social media platforms, news sources, and customer reviews could provide early signals of shifts in consumer preferences or reactions to store promotions and events. Sentiment analysis could be integrated as a predictive feature to detect potential upticks or declines in sales.

- 2. Enhancing Forecasting for New Products: The models struggled to forecast sales for new products, largely due to the lack of historical data. This presents an opportunity to explore novel machine learning techniques to overcome this challenge:
  - **Cold-Start Problem Solutions:** Techniques such as collaborative filtering or matrix factorization, commonly used in recommendation systems, could be adapted to forecast sales for new items. These methods predict behavior based on patterns from similar items or users, potentially improving the forecasts for products with minimal historical data.
  - **Transfer Learning:** Transfer learning enables a model trained on one task to be adapted to a related task. In this case, models trained on similar product categories could be used to make initial predictions for new products, with the models fine-tuning their forecasts as more data becomes available.
- 3. Real-Time Forecasting and Online Learning: The current system trains models on a static dataset and forecasts sales for future periods. In a real-world retail environment, however, sales data is continually updated, and the ability to provide real-time forecasts would offer significant value. Implementing online learning techniques would allow the models to dynamically update predictions as new sales data arrives, ensuring more timely and relevant forecasts.

- Streamlining Real-Time Forecasting: Real-time data pipelines could be set up to feed daily sales, inventory, and promotional data into the forecasting models, enabling predictions to be updated multiple times per day. For example, forecasts could be refreshed every hour or after significant shifts in consumer demand, providing actionable insights for immediate operational adjustments.
- Online Learning Techniques: Approaches such as incremental learning or reinforcement learning could enable the model to update itself in real-time. This would be particularly useful in environments where sales patterns change rapidly, such as during promotional periods, holidays, or external economic shocks.
- 4. Experimenting with Advanced Deep Learning Architectures: While LSTM models performed well in capturing long-term temporal dependencies in the data, more advanced deep learning models could further improve performance, especially for large-scale, multi-dimensional datasets like retail sales.
  - **Transformer Models:** Transformer architectures have recently achieved state-of-theart results in various sequence modeling tasks, including natural language processing and time series forecasting. Transformers can process longer sequences more efficiently than LSTMs and are better at capturing global dependencies in data, which could enhance the model's ability to forecast complex sales patterns over extended periods.

- Attention Mechanisms: Attention mechanisms, which allow models to focus on the most relevant parts of an input sequence, could improve the interpretability and accuracy of sales forecasts. By focusing on key drivers of sales (e.g., promotions, holidays, specific products), an attention-based model could provide more accurate predictions while offering insights into the factors driving those predictions.
- 5. **Improving Handling of Special Events and Holidays:** Sales forecasting during special events and holidays remains a significant challenge, as these periods exhibit unique sales patterns that are not captured by typical seasonal models.
  - Event-Specific Models: Developing separate models specifically for holiday periods or major events could improve the accuracy of predictions during these times. These models would account for the distinct characteristics of sales spikes and dips that occur during events like Christmas, Black Friday, or local festivities.
  - **Granular Event Data:** Incorporating more granular event data, such as the exact timing of promotions, discounts, or local community events, could help the models better capture the demand variations associated with these occurrences. Additionally, future research could explore the use of historical event data to improve generalization across similar future events.

- 6. **Expanding Predictive Analytics Beyond Sales Forecasting:** While this project focused on sales forecasting, the same techniques and methodologies could be applied to predict other important retail metrics. Expanding the forecasting system to cover additional dimensions would provide a more comprehensive solution for retail management:
  - **Inventory Management Optimization:** Predictive models could be developed to optimize inventory levels, minimizing both stockouts and overstock situations. By forecasting demand for each product, the system could recommend the optimal reorder quantities and timing for different stores.
  - Staffing and Workforce Planning: Accurate sales forecasts could be used to predict customer traffic, allowing stores to optimize staffing levels to meet demand. This would ensure that labor costs are minimized while maintaining high levels of customer service.
  - Supply Chain Forecasting: The models could be extended to predict supply chain disruptions or changes in supplier lead times, enabling more efficient procurement processes. This could help Corporación Favorita better manage its supplier relationships and avoid delays or shortages in stock

# References

- Bai, S., Kolter, J. Z., & Koltun, V. (2018). An empirical evaluation of generic convolutional and recurrent networks for sequence modeling. arXiv preprint arXiv:1803.01271.
- Bandara, K., Bergmeir, C., & Smyl, S. (2020). Forecasting across time series databases using recurrent neural networks on groups of similar series: A clustering approach. Expert Systems with Applications, 140, 112896.
- Box, G. E., Jenkins, G. M., Reinsel, G. C., & Ljung, G. M. (2015). Time series analysis: forecasting and control. John Wiley & Sons.
- Breiman, L. (2001). Random forests. Machine learning, 45(1), 5-32.
- Chen, T., & Guestrin, C. (2016). Xgboost: A scalable tree boosting system. In Proceedings of the 22nd acm sigkdd international conference on knowledge discovery and data mining (pp. 785-794).
- Cleveland, R. B., Cleveland, W. S., McRae, J. E., & Terpenning, I. (1990). STL: A seasonal-trend decomposition procedure based on loess. Journal of Official Statistics, 6(1), 3-73.
- Dietterich, T. G. (2000). Ensemble methods in machine learning. In International workshop on multiple classifier systems (pp. 1-15). Springer, Berlin, Heidelberg.
- Fildes, R., Ma, S., & Kolassa, S. (2019). Retail forecasting: Research and practice. International Journal of Forecasting, 35(1), 1-9.
- Géron, A. (2019). Hands-on machine learning with Scikit-Learn, Keras, and TensorFlow: Concepts, tools, and techniques to build intelligent systems. O'Reilly Media.
- Hochreiter, S., & Schmidhuber, J. (1997). Long short-term memory. Neural computation, 9(8), 1735-1780.

- Hunter, J. D. (2007). Matplotlib: A 2D graphics environment. Computing in science & engineering, 9(3), 90-95.
- Lundberg, S. M., & Lee, S. I. (2017). A unified approach to interpreting model predictions. In Advances in neural information processing systems (pp. 4765-4774).
- Wickham, H., & Grolemund, G. (2016). R for data science: Import, tidy, transform, visualize, and model data. O'Reilly Media.
- Kuhn, M., & Johnson, K. (2019). Feature engineering and selection: A practical approach for predictive models. CRC Press.
- Baier, T., Neely, A., Nalchigar, S., & Kaisler, S. (2019). A review of data science project methodologies. In Proceedings of the 52nd Hawaii International Conference on System Sciences.
- Larson, E. W., & Gray, C. F. (2017). Project management: The managerial process. McGraw-Hill Education.
- Kaggle. (n.d.). Store Sales Time Series Forecasting. Retrieved from <u>https://www.kaggle.com/c/</u> <u>store-sales-time-series-forecasting</u>
- Makridakis, S., Spiliotis, E., & Assimakopoulos, V. (2018). Statistical and Machine Learning forecasting methods: Concerns and ways forward. PloS one, 13(3), e0194889.
- McKinney, W. (2012). Python for data analysis: Data wrangling with Pandas, NumPy, and IPython. O'Reilly Media, Inc.
- Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., ... & Duchesnay,E. (2011). Scikit-learn: Machine learning in Python. Journal of machine learning

research, 12(Oct), 2825-2830Hyndman, R. J., & Athanasopoulos, G. (2018). Forecasting: principles and practice. OTexts.

Tukey, J. W. (1977). Exploratory data analysis. Reading/Addison-Wesley.

Wickham, H., & Grolemund, G. (2017). <u>https://cloudera2017.wordpress.com/wp-content/</u> uploads/2019/01/r-for-data-science-import-tidy-transform-visualize-and-model-data.pdf



University of the Cumberlands

**School of Computing and Information Sciences** 

# **Certification of Authorship**

Submitted to (Professor's Name): Dr. Eve Thullen

Course: 2024 Fall - Data Science Capstone: Practicum (MSDS-699-A01) - First Bi-term

Student's Name: Dare, Kunle

Date of Submission: October 13, 2024

Purpose and Title of Submission: Master of Science In Data Science.

**Certification of Authorship**: I hereby certify that I am the author of this document and that any assistance I received in its preparation is fully acknowledged and disclosed in the document. I have also cited all sources from which I obtained data, ideas, or words that are copied directly or paraphrased in the document. Sources are properly credited according to accepted standards for professional publications. I also certify that this paper was prepared by (me or by my group #) for this purpose.

Students' Signature: \_\_\_\_\_