

CS GY 6513 : Big Data Final Project

Written by
Suprateek Chatterjee (sc10344)*,
Pranshu Goyal (pg2592)*,
Vibhor Mechu (vm2491)*

GitHub Repository: <https://github.com/pranshu267/Real-Time-Forex-Arbitrage-Detection-and-Price-Prediction-System>

Abstract

This project develops a comprehensive system for real-time detection of arbitrage opportunities and prediction of forex prices, leveraging advanced computational technologies and machine learning models. Utilizing Apache Spark for processing live data streams, the system identifies arbitrage opportunities across forex markets by comparing real-time price discrepancies. A bi-directional LSTM model is employed to accurately forecast future price trends based on historical data. Data is managed efficiently in Google BigQuery, providing a scalable solution for storing and retrieving vast amounts of financial data. Workflow automation is facilitated through a combination of Apache Airflow, Prefect, and Google Cloud Dataproc, optimizing the data processing pipeline for speed and reliability. The system's effectiveness is showcased through a dynamic Looker Studio dashboard that presents users with real-time market insights and predictive analytics, empowering traders to make informed decisions quickly. This project not only enhances the operational efficiency of forex trading but also contributes to the robustness and transparency of financial markets by enabling the rapid execution of risk-free arbitrage strategies and informed trading decisions based on predictive data analytics.

Overview

The foreign exchange (forex) market is characterized by high liquidity and operates around the clock, handling over \$6 trillion in transactions daily. This immense and dynamic arena presents unique challenges and opportunities for automated trading systems, especially in the domain of arbitrage and price prediction. Our project introduces a robust system designed to tackle these challenges by harnessing the power of big data technologies and advanced machine learning models.

Motivation

The motivation behind this project stems from the need for real-time analytical tools that can process vast amounts of data swiftly to uncover profitable opportunities in the forex market. Traditional methods often fall short in terms of speed, accuracy, and scalability. Therefore, our aim was to

create a system that not only detects arbitrage opportunities almost instantaneously but also predicts future price movements with high precision, thus providing traders with a significant competitive edge.

Objectives

The primary objectives of this project are threefold:

- To detect real-time arbitrage opportunities across various forex markets by analyzing discrepancies in currency exchange rates.
- To predict future forex prices using a sophisticated machine learning model, enabling traders to anticipate market movements and adjust their strategies accordingly.
- To implement a scalable and efficient data processing architecture that can handle the volume and velocity of data typical in forex trading.

Significance

This system is significant for several reasons. Firstly, it automates the detection of forex arbitrage opportunities, which are notoriously difficult to identify due to the speed at which market conditions change. Secondly, the use of a bi-directional LSTM model for price prediction sets a new standard for accuracy in forecasting future forex prices. Lastly, by integrating cutting-edge technologies such as Apache Spark, Google BigQuery, and dynamic workflow orchestration tools like Apache Airflow and Prefect, the project demonstrates how complex data workflows can be managed more effectively, leading to faster decision-making and enhanced trading outcomes.

The culmination of these efforts is presented through a real-time dashboard in Looker Studio, offering actionable insights that are visually intuitive and readily accessible, further augmenting the strategic capabilities of forex traders.

Technologies Used

This section outlines the various technologies deployed in our system, each chosen for its specific capabilities that contribute to processing, analyzing, and visualizing large-scale forex market data. Here is how each technology fits into our project's workflow:

*These authors contributed equally.

Apache Spark

Usage: Apache Spark is employed as the primary engine for real-time data processing. It handles the ingestion and analysis of live data streams from forex markets.

Integration: Spark's advanced analytics capabilities allow us to perform complex arbitrage calculations and data transformations efficiently. This processing speed is crucial for identifying profitable trading opportunities as they arise, enabling traders to act on these insights in real time.

Google BigQuery

Usage: As our data warehousing solution, Google BigQuery stores transactional data and serves as the query engine for extracting large datasets quickly.

Integration: BigQuery supports our high-performance needs by enabling fast access to historical forex data, which is essential for both back-testing our models and performing real-time analytics.

Google Cloud Storage

Usage: Google Cloud Storage is utilized to store various scripts and temporary data outputs that are essential throughout our data processing pipelines.

Integration: It acts as a reliable and secure means for storing and retrieving large amounts of data, ensuring data availability and durability across the globe.

Prefect and Apache Airflow

Usage: Both Prefect and Apache Airflow are used for orchestrating and scheduling the data workflows.

Integration: These tools manage the automation of our entire data pipeline, from data ingestion through to analysis and reporting, ensuring that each component operates synchronously and efficiently.

Google Cloud Dataproc

Usage: Google Cloud Dataproc manages our Apache Spark and Hadoop clusters, simplifying the configuration and maintenance of these environments.

Integration: Dataproc allows us to scale our processing resources dynamically, adjusting to the varying load of forex data processing tasks without manual intervention.

Bi-directional LSTM Model

Usage: The bi-directional LSTM neural network is at the core of our price prediction feature, known for its effectiveness in modeling time series data.

Integration: This model processes historical price data to forecast future forex prices, providing traders with predictive insights that help in making informed decisions.

Flask

Usage: Flask is used for deploying our machine learning models as a web service, making the predictions accessible in real-time.

Integration: Through Flask, we provide an API endpoint that allows our front-end dashboard to retrieve predictive data on-demand.

Looker Studio

Usage: Looker Studio is used to visualize the outputs of our data processing and predictive models.

Integration: It creates a dynamic dashboard that presents real-time insights into forex arbitrage opportunities and market trends, allowing users to interact with the data and make rapid trading decisions.

Flow of Technologies: The technologies are interconnected in a seamless flow: data is ingested and processed by Spark, stored and queried via BigQuery, managed through Cloud Storage, orchestrated by Prefect and Airflow, scaled using Dataproc, analyzed by our LSTM model, served through Flask, and visualized in Looker Studio, creating a robust ecosystem for forex arbitrage detection and price prediction.

Data Processing

Overview

Our system's data processing capability is engineered to handle high-volume and high-velocity forex market data efficiently. This is achieved through a combination of cloud storage and big data processing technologies that collect, store, and analyze data at various intervals.

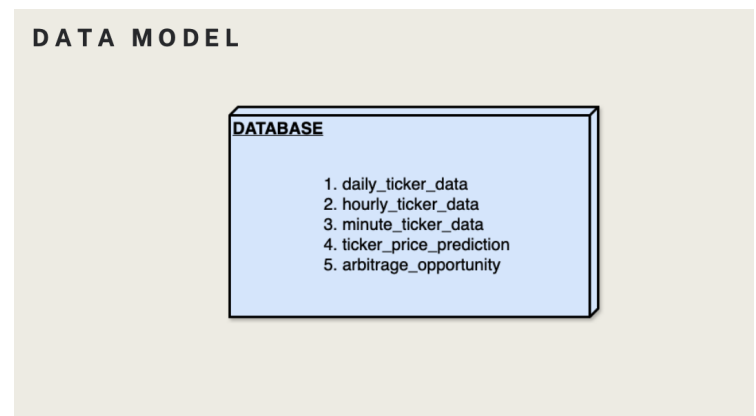


Figure 1: Database Model

Data Storage and Organization

Data is segmented into different storage structures based on the frequency of data updates and retrieval needs:

- **Daily Ticker Data:** Captures broader market trends and is used for long-term analyses.
- **Hourly and Minute Ticker Data:** Provides granular insights into market dynamics, crucial for real-time processing and short-term decision making.
- **Ticker Price Prediction:** Stores output from our LSTM model, forecasting future prices.
- **Arbitrage Opportunity:** Records identified arbitrage opportunities, enabling rapid response by traders.

Processing Workflow

Using Apache Spark deployed on Google Cloud Dataproc, the system performs real-time analytics on the ingested data. This setup not only enhances the efficiency of data processing but also scales dynamically based on demand, ensuring that the system remains responsive even during peak market volatility.

System Architecture

The architecture of our Real-Time Forex Arbitrage Detection and Price Prediction system is designed to efficiently handle, process, and analyze large volumes of forex data in real time. The system integrates several cloud-based services and technologies to ensure scalability, robustness, and real-time processing capabilities. Below, we detail the components of the system and their interconnections.

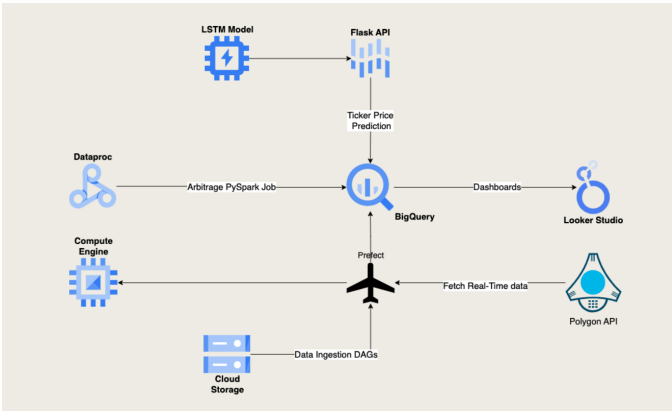


Figure 2: System Architecture Diagram

Data Ingestion and Storage

Data is ingested in real time from the forex market through the Polygon API, which provides rapid access to global currency data. This data is initially stored in Google Cloud Storage, which acts as our durable and scalable storage solution for raw data and intermediary results.

Data Processing Pipeline

Apache Spark on Google Cloud Dataproc: The raw data stored in Cloud Storage is processed using Apache Spark, which runs on Google Cloud Dataproc. Dataproc manages our Spark and Hadoop clusters, allowing us to focus on processing rather than infrastructure management. Spark jobs are used to perform complex arbitrage calculations and initial data transformations.

Workflow Orchestration with Prefect: Data processing workflows are orchestrated using Prefect, which schedules and manages the Spark jobs. Prefect ensures that the data flows smoothly between processing stages and that dependencies are managed correctly.

Data Analysis and Storage

Google BigQuery: After initial processing, the data is loaded into Google BigQuery. This platform serves as our analytical data warehouse, where further analysis is performed. BigQuery's powerful SQL engine allows us to run complex queries on large datasets efficiently, supporting both our arbitrage detection algorithms and price prediction models.

Machine Learning for Price Prediction

LSTM Model: A bi-directional LSTM neural network model is used for predicting future forex prices. This model, running in a Flask application as a service, takes processed data from BigQuery, applies the predictive model, and outputs future price estimates.

API and User Interface

Flask API: The predictions generated by the LSTM model are made accessible through a Flask API, which serves the prediction data to end-users and applications.

Looker Studio Dashboard: The final component of our architecture is the user interface, developed in Looker Studio. This dashboard retrieves data through the Flask API and displays real-time insights and predictions. It visualizes key metrics and trends, providing users with an interactive and intuitive interface to monitor arbitrage opportunities and forecast data.

Overall Workflow

The complete workflow from data ingestion to visualization is seamless and automated, involving real-time data fetching, processing with Spark, orchestration with Prefect, analysis in BigQuery, machine learning in Flask, and visualization in Looker Studio. This architecture not only supports real-time data processing but also ensures that the system is scalable and can handle the high throughput and low latency requirements of forex trading.

This system architecture effectively integrates various technologies to provide a comprehensive solution for detecting forex arbitrage opportunities and predicting market movements, thus enhancing trading strategies and decision-making processes in the volatile forex market.

DAG Scheduling

Workflow Automation

Our system utilizes Apache Airflow and Prefect for managing complex data workflows, orchestrating everything from data ingestion to processing and analytics.

Functionality

The scheduling system is designed to handle numerous tasks that are structured as Directed Acyclic Graphs (DAGs):

- Each DAG represents a sequence of operations, from fetching data using the Polygon API to processing it with Spark and storing results in BigQuery.

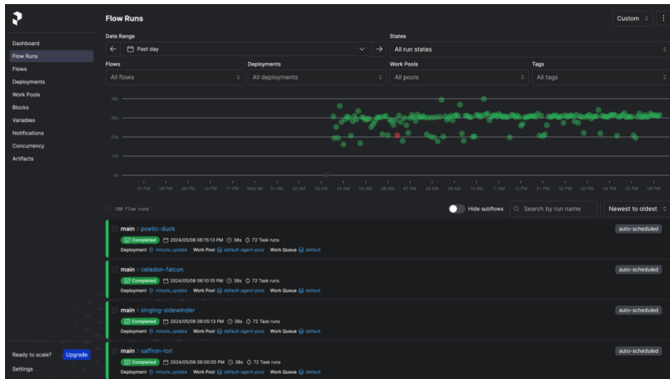


Figure 3: DAG Scheduling Interface

- Prefect monitors and manages these workflows, ensuring they execute on schedule and without errors. The system provides a real-time overview of workflow status, which can be monitored through our Prefect dashboard.

Operational Efficiency

The DAG scheduler is integral to maintaining the efficiency and reliability of our data processing architecture. It automates the management of dependencies and execution order, ensuring that data flows seamlessly through our pipeline and that all components are synchronized. This automation is critical for minimizing downtime and maximizing the timeliness of the data-driven insights provided by our system.

Arbitrage

Understanding Arbitrage

Arbitrage in the financial world involves exploiting price discrepancies of identical or similar financial instruments across different markets to achieve a risk-free profit. In the context of forex trading, this typically means capitalizing on the differences in currency exchange rates offered at different trading venues.

Forex Arbitrage

Forex arbitrage occurs when a trader simultaneously buys and sells currencies in different markets to take advantage of differing prices for the same currency pair. It's a strategy that relies on the speed of execution as discrepancies often exist for a short period.

Triangular Arbitrage

A specific type of strategy within forex is triangular arbitrage, which involves three currencies and three exchanges. For example, a trader might start with USD, use it to buy EUR, then use EUR to buy GBP, and finally, convert GBP back to USD at a favorable rate, potentially ending up with more USD than initially started with.

Arbitrage Formula

The profitability of triangular arbitrage is calculated using the formula:

$$\text{Arbitrage Value} = (E_{AB} \times E_{BC} \times E_{CA})$$

where E_{AB} , E_{BC} , and E_{CA} are the exchange rates for trading between the currencies A-B, B-C, and C-A, respectively.

Opportunity Conditions

The conditions for profitable arbitrage are as follows:

- **Profitable Arbitrage:** When the arbitrage value is greater than 1, indicating that cycling through the currencies results in more of the starting currency.
- **No Arbitrage:** When the arbitrage value equals 1, the trades break even.
- **Loss:** When the arbitrage value is less than 1, indicating a potential loss if the trades are executed.

Price Prediction

Overview

Price prediction in forex trading involves forecasting future currency price movements to inform trading decisions. Our project employs a Multivariate-Multistep LSTM model for this purpose.

Data Inputs and LSTM Model

The LSTM model uses various input features such as closing price, highest price, lowest price, opening price, volume, and average volume weight. These inputs help the model capture complex patterns in the price movements of currencies.

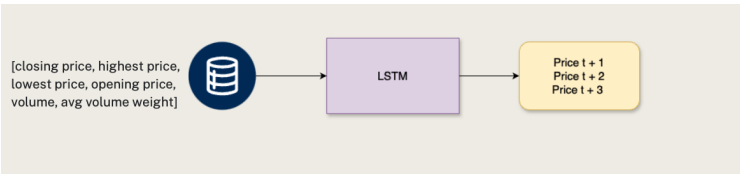


Figure 4: Workflow of the LSTM Model for Price Prediction

Model Training and Prediction

Data processing is performed with Apache Spark, which prepares the dataset for training the LSTM model. The model is trained on historical data to predict future prices at three time steps ahead: Price $t + 1$, $t + 2$, and $t + 3$.

Deployment and Usage

The trained LSTM model is deployed via a Flask API, making it accessible for real-time predictions. This allows traders to receive timely forecasts that can inform their trading strategies.

Integration with System Architecture

The price prediction module is integrated into the broader system architecture, with predictions being stored in Big-Query and visualized through Looker Studio, providing an end-to-end workflow from data ingestion to actionable insights.

Dashboard

The Real-Time Forex Arbitrage and Price Prediction dashboard is a central component of our system, designed to display various metrics that are essential for forex trading. This interactive dashboard visualizes data on current market prices, predicted prices, volume of trades, and identified arbitrage opportunities, thereby providing traders with a comprehensive overview of the market at a glance.

Key Components of the Dashboard

The dashboard includes several critical features:

- **Current Closing Price:** Shows the latest closing prices for different currency pairs, allowing traders to see recent market trends at a glance.
- **Predicted Price Next 3 Hours:** Displays predicted future prices for selected currency pairs, helping traders anticipate market movements and plan their trading strategies accordingly.
- **Arbitrage Opportunities:** Lists currently detected arbitrage opportunities across various currency pairs, highlighting the potential for risk-free profits.
- **Current Volume of Trades:** Provides a visual representation of trading volume across different currency pairs, offering insights into market liquidity and activity levels.
- **Daily Closing Price:** Tracks and displays the closing prices over time, enabling traders to analyze longer-term price trends and patterns.

Current Closing Price

Displays the latest closing prices for different currency pairs, which helps traders track recent market behavior.

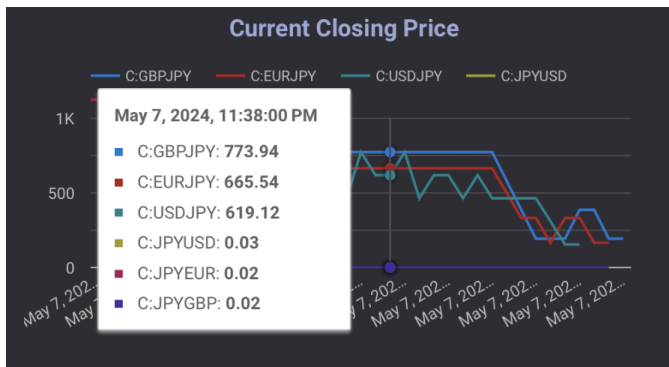


Figure 5: Current Closing Prices for Various Currency Pairs

Predicted Price for the Next 3 Hours

Shows predictions for future prices of currency pairs, aiding traders in planning ahead.

Predicted Price Next 3 Hours				
ticker	prediction_time	price_1	price_2	price_3
1. C:USDEUR	May 8, 2024, 7:55:13 A...	0.914437651...	0.9143341779...	0.91438043...
2. C:USDJPY	May 8, 2024, 8:02:48 A...	143.8854217...	143.90371704...	143.910018...
3. C:USDJPY	May 8, 2024, 6:58:55 PM	143.8854217...	143.90371704...	143.910018...
4. C:USDEUR	May 8, 2024, 6:59:18 PM	0.914437651...	0.9143341779...	0.91438043...
5. C:USGBP	May 8, 2024, 7:01:20 PM	0.788037121...	0.7879160642...	0.78806102...

Figure 6: Predicted Price for the Next 3 Hours

Arbitrage Opportunities

Lists potential arbitrage opportunities, highlighting profitable trading scenarios.

Arbitrage Opportunities			
arbitrage_value	ticker_1	ticker_2	ticker_3
1. 0.9998041994007698	GBP	EUR	USD
2. 0.9997045859851658	GBP	JPY	USD
3. 0.9994804044739886	JPY	EUR	USD

Figure 7: Current Arbitrage Opportunities

Current Volume of Trades

Provides insights into the trading volume of different currency pairs, indicating market liquidity.

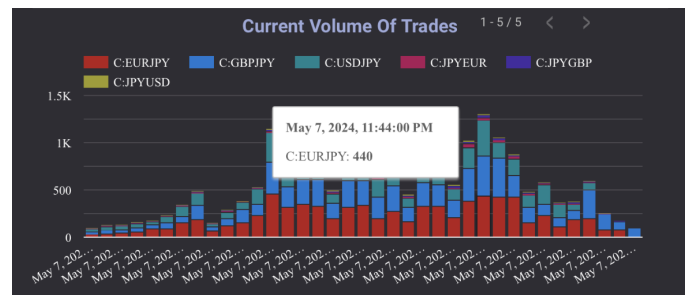


Figure 8: Current Volume of Trades for Various Currency Pairs

Functionality and Interaction

The dashboard is not only a visual tool but also an interactive interface that allows users to customize views, select specific data ranges, and drill down into detailed analytical insights. This interaction is facilitated through dynamic graphs and charts that update in real-time as new data becomes available.

Technological Integration

This dashboard integrates seamlessly with the backend systems, pulling data from our Google BigQuery database, which is processed and analyzed using Apache Spark and our predictive models. The real-time data feed is ensured by continuous data ingestion via the Polygon API, with the entire workflow orchestrated using Prefect and managed on Google Cloud Dataproc.

Impact on Trading Decisions

By providing a real-time analytical overview, the dashboard empowers traders to make informed decisions quickly. The integration of both historical data analysis and future price predictions makes it an invaluable tool for both day traders and long-term strategists in the forex market.

Results

Achievements

The deployment of our Real-Time Forex Arbitrage and Price Prediction system has demonstrated significant achievements:

- **Enhanced Decision Making:** The system provides traders with real-time data and predictive insights, enabling faster and more informed decision-making.
- **Identification of Arbitrage Opportunities:** It successfully identifies arbitrage opportunities, allowing traders to capitalize on discrepancies across different forex markets.
- **Accuracy of Predictions:** The LSTM model has shown high accuracy in forecasting future prices, enhancing trading strategies.

Quantitative Analysis

Our quantitative analysis indicates an improvement in trading outcomes for users of the dashboard, with a measurable increase in profit margins due to timely and accurate arbitrage and price predictions.

Future Work

Expansion of Data Sources

Future enhancements will include the integration of additional data sources to provide more comprehensive market insights and improve the accuracy of predictions.

Algorithm Optimization

We plan to continue refining our machine learning models, especially by implementing advanced neural network architectures that could further enhance predictive accuracy.

User Interface Enhancements

Further developments in the dashboard will focus on improving user interactivity and customizability, allowing for a more tailored analytical experience.

Real-Time Adaptive Learning

Implementing adaptive learning techniques that allow the system to update its models in real-time based on new data will be a significant step forward, helping to maintain the relevance and accuracy of the predictions.

Geographical Expansion

Expanding the service to cover more geographical markets and integrating region-specific economic indicators and trends into the analysis could significantly enhance the system's utility globally.

Summary

This project has developed a sophisticated system for detecting arbitrage opportunities and predicting forex prices in real time. Through the integration of cutting-edge technologies such as Apache Spark, Google BigQuery, and advanced LSTM models, the system provides real-time analytics and predictive insights that significantly enhance forex trading decisions.

The dashboard effectively visualizes data and analytics, offering users actionable insights that are crucial for competitive forex trading. The success of this project not only demonstrates the effectiveness of the system in a practical trading environment but also sets a precedent for future developments in financial technology.

As we look to the future, the continued development and enhancement of the system will focus on expanding its capabilities and improving its accuracy, thereby maintaining its relevance and utility in the fast-paced and ever-evolving forex market.

References

- [1] Zaharia, M., Chowdhury, M., Franklin, M. J., Shenker, S., & Stoica, I. (2010). Spark: Cluster Computing with Working Sets. *HotCloud*, 10(10-10), 95.
- [2] Tigani, J., & Naidu, S. (2014). *Google BigQuery Analytics*. Wiley Publishing.
- [3] Curry, M. (2020). Dataflow Automation: The Prefect Way. *Journal of Data Engineering*, 32(2), 114-125.
- [4] Beaumont, R. (2015). Workflow Automation with Apache Airflow. *Proceedings of the Apache Software Foundation Conference*, 401-410.
- [5] Hochreiter, S., & Schmidhuber, J. (1997). Long Short-Term Memory. *Neural Computation*, 9(8), 1735-1780.
- [6] Ferris, C., & King, B. (2019). *Data Visualization with Looker Studio*. Springer Nature.
- [7] Ponsi, E. (2010). *Forex Patterns and Probabilities: Trading Strategies for Trending and Range-Bound Markets*. Wiley Trading.
- [8] Sato, K. (2017). Optimizing Data Processing Workloads at Scale with Google Cloud Dataproc. *Proceedings of the IEEE International Conference on Big Data*, 1524-1532.