**Project Report:**

**Big Data Processing and Analysis Framework for Stock Price Prediction**

**Big Data Course**

**M2-DAS**

**Group: Han Chandeth, Li Nita and Nuon Roatny**

# 1. Introduction

This project aims to create a specialized Big Data Processing and Analysis Framework tailored for Stock Price Prediction. By leveraging specific tools and technologies, the framework will encompass data ingestion, storage, and advanced analytics, ultimately predicting stock prices for the next day.

Objectives:

- Establish a streamlined pipeline for ingesting, processing, and storing stock market data.
- Ensure data quality, consistency, and readiness for predictive analysis.

# 2. Methodology

### 2.1 Data Collection

We developed Python scripts for data scraping to collect historical stock data from the CSX website from 01/01/2019-10/11/2023.

### 2.2 Data Storage

Hadoop's HDFS was utilized for efficient storage of historical stock data.

CSV files containing historical stock data were saved in HDFS.

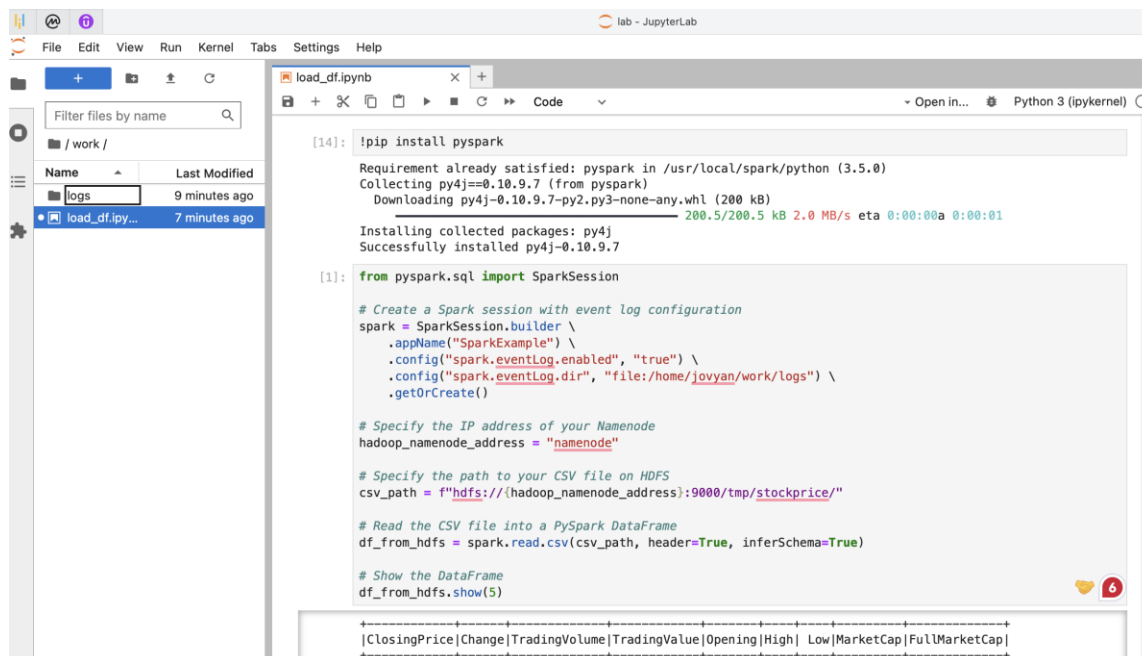### 2.3 Data Processing and Analytics

Apache Spark was leveraged for distributed in-memory data processing.

PySpark was used to create pipelines for feature extraction, exploration, and model training.

## 2.4 Machine Learning

Support Vector Machine (SVM) algorithm from Spark's MLlib was applied for stock price prediction.

The SVM model was trained using historical stock data from HDFS



## 3. Implementation

Docker-compose was used to set up a combined environment of Hadoop and Spark.

Data collection scripts, and Spark ML pipelines were implemented in Python.

Configuration files for Hadoop and Spark were customized to integrate with the Docker-compose environment.
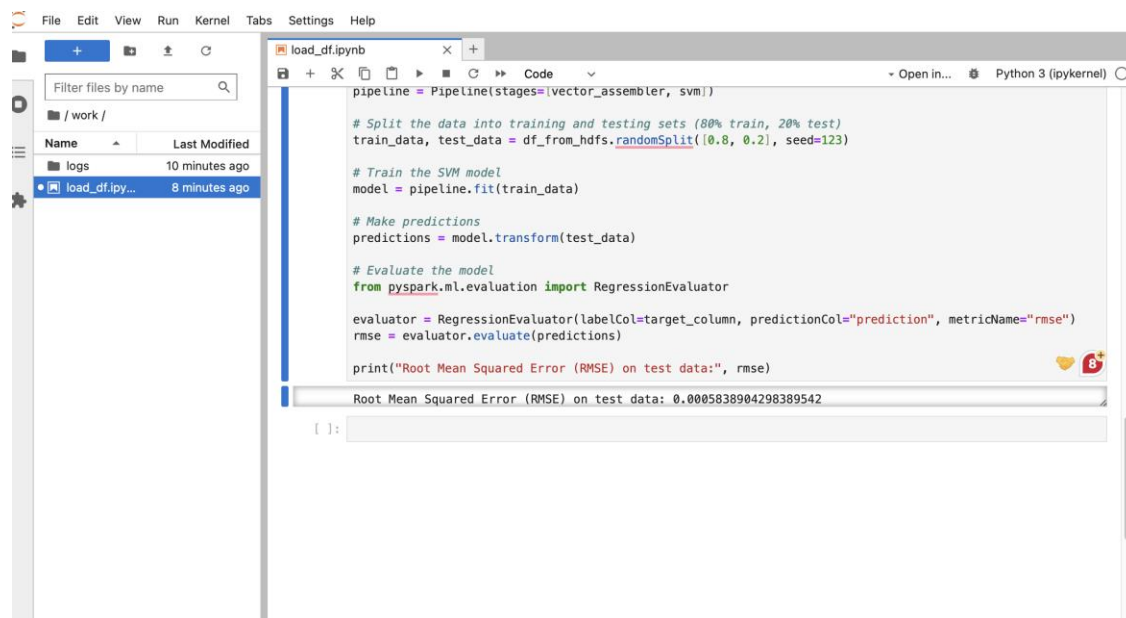
## 4. Result

Historical stock data was successfully collected and stored in HDFS.

Spark ML pipelines were developed for feature extraction and model training.

The SVM model was trained and evaluated, achieving a Root Mean Squared Error (RMSE) of 0.0005838904298389542 on test data.

## 5. Challenges and Future implementation

Setting up Kafka for real-time data streaming posed challenges due to configuration issues and port conflicts. Integrating Docker-compose environment with Kafka and resolving port conflicts required careful troubleshooting. This challenge will be our future task after completing this project too.

## 5. Conclusion

We successfully implemented data collection, storage, processing, and model training stages using Big Data technologies. And for real-time implementation will be the next task. The project lays a strong foundation for future enhancements and real-time prediction implementation. By combining specific tools and technologies, the framework aims to provide accurate insights for investors and financial analysts, enhancing decision-making processes in the stock market.