# Stock Data Streaming and Analysis

## Electric Vehicle Stock

**Emmitt Hasty**
Data Science
University of Missouri-Kansas City
Missouri
emhrv3@umsystem.edu

**Amulya Sara**
Computer Science
University of Missouri-Kansas City
Missouri
asc79@umsystem.edu

**Gautham Sidhardh Dommeti**
Computer Science
University of Missouri-Kansas City
Missouri
gsddbc@umsystem.edu

**Snehalatha Mallavarapu**
Computer Science
University of Missouri-Kansas City
Missouri
smfyp@umsystem.edu

**Naladala Sai Karthik**
Computer Science
University of Missouri-Kansas City
Missouri
sndzm@umsystem.edu

## ABSTRACT

Data analytics is the process of examining enormous amounts of data to find patterns, insights, and trends. Although data analytics technologies can be applied to a wide range of fields, including healthcare, politics, retail, banking, and government agencies, they are essential for survival in the competitive financial markets of the modern era. Data analytics are essential in finance to understand the ups and downs of capital markets. Great financial data science gives traders and investment advisors the assurance they need to decide for themselves whether to purchase, sell, or hold a certain security. They can use it to manage a portfolio according to short-, mid-, or long-term goals. This live streaming of data from a stock analysis is taken as source data. the data from Forex and Electric Vehicles trading API is streamed through Kafka. After streaming, the data stored in topics is analyzed. The analysis of data is visualized in Jupyter notebook for various use cases. Stocks for Electric vehicle brands are complicated, and insight can be gained from many things, as electric vehicles need special batteries, and valuable insight may be gained from looking into the stock of these companies that make the special batteries. These are companies like Amaraj, or HBL Power. As these stocks become more valuable it can become an indicator that the market for EV's is increasing and will attribute to an increase in stock value for EV's across the board in India. For this paper, we studied three company types, Auto Manufacturing, Battery, and Charging Point. In our study, we predict that increases in the volume of shares per day correlates to an increase in stock prices related to EV's and can be useful for stock trading and investing purposes. Using Kafka, daily streaming data was ingested from AlphaVantage, showing stock data relevant to the daily prices. Then using Spark, the data is loaded into a Jupyter Notebook, with which we can do data analysis on. Plotting the volume of stocks moving day to day can then be used for insight into long term or short-term goals.

## INTRODUCTION
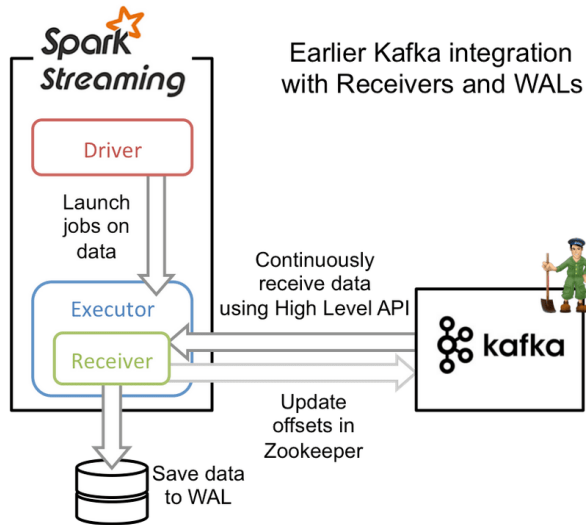
There's a Silent EV Revolution Coming in the world.

Over the world, electric vehicles are becoming more and more popular. The mass adoption of electric vehicles is being fueled by alarming emission rates and surging gasoline costs, making EV stocks the hottest trend.

It might be a business that manufactures lithium batteries, charging stations, associated components, and more. Due to the industry's youth and quick expansion, investing correctly in such businesses could result in a substantial return.

## TOOLS & TECHNOLOGIES

- Requests library,
- Pandas datareader library,
- API : Tiingo API, Alpha Vantage API.
- data visualization tool: matplotlib, Seaborn
- Spark
- Kafka
- Pyspark

## ARCHITECTURE

Earlier Kafka integration with Receivers and WALs

## KEYWORDS

EV: Electric Vehicle

## 1.1 KAFKA PRODUCER AND CONSUMER

Using Kafka, we have acquired daily stock reports detailing the price at market open, the high for the day, the low for the day, and the price at market close, as well as the time stamp of the day of the data. The data we gathered is the most recent 100 days, but the data can be altered to consider all data since 1999. The data is preprocessed to declare an index, as well as convert the date, which was the original index, to become a value of the data instead. Some minor changes are then made to the data for aesthetic purposes, and then the original columns are then removed from the data. This data is then loaded into the target cluster, to a topic declared 'mm.' Another topic, 'mm1' is also declared within the Kafka cluster, which is then used to load the data from the cluster into the target source,theJupyterNotebook.

```python
2 from kafka import KafkaProducer
3 import requests
4 from json import dumps
5 import time
6
7 def on_message1(message):
8     producer1.send('mm', message)
9     producer1.flush()
10
11 producer1 = KafkaProducer(value_serializer=lambda m: dumps(m).encode("utf-8"), bootstrap_servers=['localhost:9092'])
12
13
14 # url to collect data of Mahindra and Mahindra shares
15 url = 'Shttps://www.alphavantage.co/query?function=TIME_SERIES_DAILY&symbol=M&M.BSE&apikey=3650YNI4HJXT8L12&outputsize=full'
16 r = requests.get(url)
17 data=r.json()
18 #print(data)
19 #preprocessing
20 del(data['Meta Data'])
21 df=data['Time Series (Daily)']
22 #print(df)
23 jsonFile={}
24 j=0
25 for i in df.keys():
26     jsonFile[j]=df[i]
27     jsonFile[j]['0. date']=i
28     #print(i)
29     j+=1
30 print(jsonFile)
31 for i in jsonFile.keys():
32     jsonFile[i]['open']=jsonFile[i]['1. open']
33     jsonFile[i]['high']=jsonFile[i]['2. high']
34     jsonFile[i]['low']=jsonFile[i]['3. low']
35     jsonFile[i]['close']=jsonFile[i]['4. close']
36     jsonFile[i]['volume']=jsonFile[i]['5. volume']
37     jsonFile[i]['date']=jsonFile[i]['0. date']
38     del jsonFile[i]['1. open']
39     del jsonFile[i]['2. high']
40     del jsonFile[i]['3. low']
41     del jsonFile[i]['4. close']
42     del jsonFile[i]['5. volume']
```

**Figure 1.1: Source code from the Kafka Producer, to ingest the data and add it to the Kafka Cluster.**

## 1.2 SPARK PROCESSING

Spark is used to prepare the data for analyzing. Using readstream, Spark subscribes to the mm topic, then converts the entire data field to a string, and allowing the data to be structured. This is done by creating a data model type, and then casting the type onto the data. Once the data is prepared, it can then be sent back to the Kafka broker using writestream, and prepared for the Notebook.

```python
14 from pyspark.sql import SparkSession
15 from pyspark.sql.functions import *
16 from pyspark.sql.types import *
17
18 spark = SparkSession \
19     .builder \
20     .appName("Manufacturing_companies") \
21     .getOrCreate()
22
23 raw_df = spark \
24     .readStream \
25     .format("kafka") \
26     .option("kafka.bootstrap.servers", "localhost:9092") \
27     .option("subscribe", "mm") \
28     .option("startingOffsets", "latest") \
29     .load() \
30     .selectExpr("CAST(value AS STRING)")
31
32 # Created Schema for Structured Streaming
33
34 #timestamp,open,high,low,close,volume
35
36 schema = StructType(
37     [
38         StructField('open', StringType()),
39         StructField('high', StringType()),
40         StructField('low', StringType()),
41         StructField('close', StringType()),
42         StructField('volume', StringType()),
43         StructField('date', StringType())
44     ])
45 # Applied schema on data
46 data1= raw_df.select(from_json(raw_df.value, schema).alias('data'))
47
48 #Fetching required data for processing queries
49 output_df1 = data1.select(to_json(struct(col('data.volume'),col('data.date'))).alias('value'))
50
51 # Sending the data to kafka brocker
52 query = output_df1.writeStream.format("kafka").option("kafka.bootstrap.servers", "localhost:9092").option("checkpointLocation", "/tmp/checkpoint2").option("topic", "mm1").start()
```

**Figure 1.2: Code of the Spark Processing**

## 1.3   JUPYTER NOTEBOOK

Once the data is sent to the Notebook, it can be loaded into a dataframe, and then analysis can begin. The data is loaded from the Kafka cluster, and then cleaned to prepare a plot of the volume vs data. For convenience's sake, only the latest 30 days are plotted to have a legible graph. This allows the data to be analyzed for trends in volume change, which is beneficial knowledge due to its effect on stock price. Using matplotlib, and kafka-python allows for the gathering of data from the topic and moving it to the Jupyter Notebook.
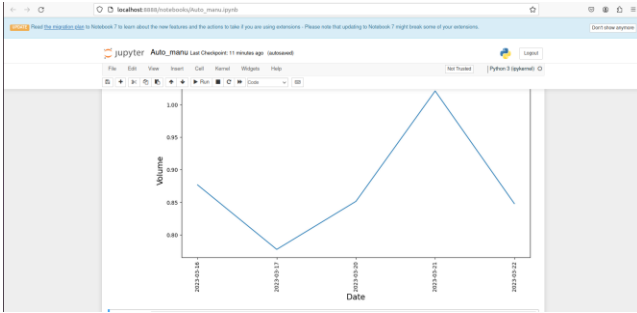
**Figure 1.3: The plot that comes from the notebook. This is a small selection of the data, only spanning six days but showing good trend perspective.**

## 2 INITIATION

Zookeeper is used as a cluster manager and must be initiated before the system starts. Starting it will begin the Kafka session, allowing for clusters to be maintained. After that, Kafka itself must then be initiated. Once Kafka is opened, the topics can then be created, mm and mm1, respectively. These topics are within the brokers, which are within the clusters. Running the mm file will then complete the process to ingest the data to the cluster, and then running mm1 will preprocess and send the data to the target source. The data is now ready for the notebook and can begin data analysis.
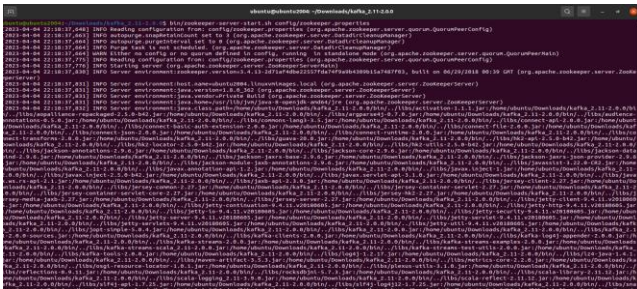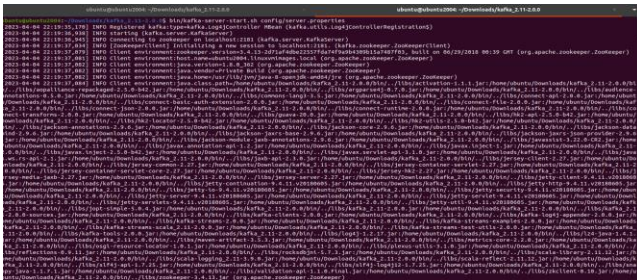


**Figure 2.1**



**Figure 2.2**



**Figure 2.3**

**Figure 2.1,2.2,2.3: These images show the code and processes of the Zookeeper initiation, the Kafka initiation, and then the topic declaration.**
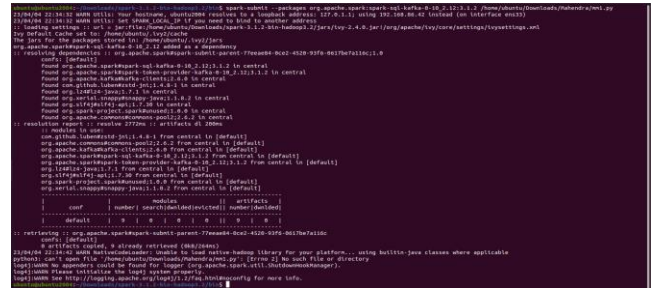


**Figure 2.4: Shows the code and processes of the Spark submittal process.**

## 3. PRELIMINARY RESULTS

Initial results have shown promising increases in the volume of shares being traded, which coinciding with the increase in stock prices, indicates the correlation between the two. Continued preparation of the data is expected to follow the trends of the volume plotting. This can become a valuable asset for stockbrokers to know when to buy and sell stocks, and for the average citizen to know when to invest in order to make capital gains.

## 4.CONTRIBUTIONS

### Emmitt Hasty
- Streaming Data Pipeline
- Install and configured Zookeeper, Kafka

### Snehalatha Mallavarapu
- Install and configured Jupyter Notebook
- Spark Streaming Jobs

### Amulya Sara
- Streaming Data Pipeline
- Install and Configured Zookeeper, Kafka

### Sai Karthik Naladala
- Spark Streaming Jobs
- Data Pre-processing
- Kafka Producer

### Gautham Sidhardh Dommeti
- Data Pre-processing
- Kafka Producer

## 5.REFERENCES

[1] https://www.digitalocean.com/community/tutorials/how-to-set-up-jupyter-notebook-with-python-3-on-ubuntu-20-04-and-connect-via-ssh-tunneling
[2] https://sparkbyexamples.com/kafka/apache-kafka-cluster-setup/
[3] https://www.alphavantage.co/documentation/