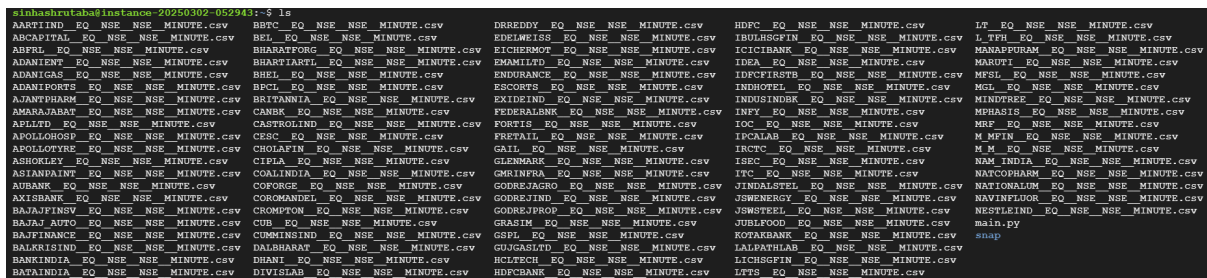Submitted by Shramana Sinha, 23f1002703

# Preparing the input files

## Downloading Data from GitHub

This step fetches the stock price data files from the specified GitHub repository into the vm machine.

```python
repo_url = "https://github.com/ShabbirHasan1/NSE-Data/tree/main/NSE%20Minute%20Data/NSE_Stocks_Data"
response = requests.get(repo_url)
soup = BeautifulSoup(response.content, "html.parser")
file_links = []
for link in soup.find_all("a"):
    href = link.get("href")
    if href and href.endswith(".csv"):
        file_links.append(
            f"https://raw.githubusercontent.com{href.replace('/blob', '')}"
        )
for file_link in file_links:
    try:
        filename = os.path.basename(file_link)
        response = requests.get(file_link)
        with open(filename, "wb") as f:
            f.write(response.content)
    except Exception as e:
        print(f"Error processing file {filename}: {e}")
```



*Screenshot: The files that were downloaded from github to the vm's local machine*

## Authenticating with Google Cloud

This step authenticates the vm machine to upload files in the gcs buckets.

```
Unset
gcloud auth login
```



Screenshot: The google cloud storage authentication

## Transferring Files to GCS Bucket

In this step, the files were transferred from the vm machine's local storage to the gcs bucket.

```
Unset
gcloud storage rsync . gs://oppe1_bucket
```



Screenshot: The files are uploaded from vm to gcs bucket

# Dataproc Cluster Setup

A Dataproc cluster on Compute Engine was created to run the PySpark job, with the following configuration:

**Manager Node:**

- Machine Series: E2
- Machine Type: e2-standard-2
- Primary Disk Size: 30 GB

**Worker Nodes:**
- Number of Nodes: 2
- Machine Series: E2
- Machine Type: e2-standard-2
- Primary Disk Size: 30 GB

| | |
|---|---|
| ← Cluster details   ➕ SUBMIT JOB   ↻ REFRESH   ▶ START   ■ STOP   🗑 DELETE   ☰ VIEW LOGS ↗ | |
| Advanced execution layer | Off |
| Google Cloud Storage caching | Off |
| Dataproc Metastore | None |
| Scheduled deletion | Off |
| Confidential computing enabled? | Disabled |
| Master node | Standard (1 master, N workers) |
| Machine type | e2-standard-2 |
| Number of GPUs | 0 |
| Primary disk type | pd-balanced |
| Primary disk size | 30GB |
| Local SSDs | 0 |
| Worker nodes | 2 |
| Machine type | e2-standard-2 |
| Number of GPUs | 0 |
| Primary disk type | pd-balanced |
| Primary disk size | 30GB |
| Local SSDs | 0 |
| Secondary worker nodes | 0 |
| Secure Boot | Disabled |
| VTPM | Disabled |
| Integrity Monitoring | Disabled |
| Cloud Storage staging bucket | dataproc-staging-us-central1-254833764668-fkb5pdnz |
| Network | default |

*Screenshot: Dataproc Cluster Configuration*

# File Upload

The required files were uploaded to the cluster's Master Node, using SSH on the cloud console.

- **Files:** `main.py` (PySpark script)

# PySpark Script Execution

The PySpark job was executed on the cluster's Master Node, using the following command:

```
Unset
spark-submit main.py
```

*Screenshot: The execution of the code*

# Code Explanation

## Import Required Libraries

```Python
from pyspark.sql import SparkSession
from pyspark.sql import Window
from pyspark.sql.functions import col, lag, abs, expr
import pyspark.sql.functions as F
```

## Create Spark Session

This creates a new Spark session or gets the existing one with the application name "Stock Price Analysis".

```Python
spark = SparkSession.builder.appName("Stock Price Analysis").getOrCreate()
```

## Load the data from GCS

The code reads all CSV files from the specified GCS bucket, automatically detecting headers and inferring the schema of the data.

```Python
gcs_path = "gs://oppe1_bucket/*.csv"
df = spark.read.option("header", "true").option("inferSchema",
"true").csv(gcs_path)
```

## Clean and Preprocess Data

This filters out rows with null values in any of the required columns to handle bad data. Data types are explicitly cast to ensure proper handling of calculations.

```Python
df_clean = df.filter(
    col("timestamp").isNotNull()
    & col("close").isNotNull()
    & col("open").isNotNull()
    & col("high").isNotNull()
    & col("low").isNotNull()
    & col("volume").isNotNull()
)
df_clean = (
    df_clean.withColumn("timestamp", col("timestamp").cast("timestamp"))
    .withColumn("close", col("close").cast("double"))
    .withColumn("open", col("open").cast("double"))
    .withColumn("high", col("high").cast("double"))
    .withColumn("low", col("low").cast("double"))
    .withColumn("volume", col("volume").cast("long"))
)
```

## Extract Stock Identifier

This extracts information about which file each row came from, serving as a stock identifier.

```Python
df_clean = df_clean.withColumn("stock_ticker", F.input_file_name())
```

## Sort Data for Time-Series Analysis

The data is sorted by stock ticker and timestamp to ensure correct calculation of price changes over time.

```Python
df_sorted = df_clean.orderBy("stock_ticker", "timestamp")
```

## Calculate Price Change Percentage

A window function creates partitions by stock ticker, then calculates the previous close price within each partition. The percentage change is calculated using the absolute value of the difference between current and previous close prices.

```Python
windowSpec = Window.partitionBy("stock_ticker").orderBy("timestamp")

df_with_change = (
    df_sorted.withColumn("prev_close", lag("close", 1).over(windowSpec))
    .filter(col("prev_close").isNotNull())
    .withColumn(
        "pct_change", abs((col("close") - col("prev_close")) /
col("prev_close") * 100)
    )
)
```

## Cache Data for Performance

The DataFrame is cached in memory to improve performance for subsequent operations.

```Python
df_with_change.cache()
```

## Calculate Percentile Values

This calculates multiple percentiles (95th, 99th, 99.5th, 99.95th, and 99.995th) in a single operation.

```Python
percentiles_expr = expr(
    "percentile(pct_change, array(0.95, 0.99, 0.995, 0.9995, 0.99995))"
)
percentile_values = df_with_change.select(
```

```python
    percentiles_expr.alias("percentiles")
).collect()[0][0]
```

## Extract Individual Percentile Values

The individual percentile values are extracted from the result array.

```python
p95_value = percentile_values[0]
p99_value = percentile_values[1]
p995_value = percentile_values[2]
p9995_value = percentile_values[3]
p99995_value = percentile_values[4]
```

## Count Trades Exceeding Thresholds

For each percentile threshold, the code counts how many trades exceed that threshold.

```python
count_exceeding_p95 = df_with_change.filter(col("pct_change") >
p95_value).count()
count_exceeding_p99 = df_with_change.filter(col("pct_change") >
p99_value).count()
count_exceeding_p995 = df_with_change.filter(col("pct_change") >
p995_value).count()
count_exceeding_p9995 = df_with_change.filter(col("pct_change") >
p9995_value).count()
count_exceeding_p99995 = df_with_change.filter(col("pct_change") >
p99995_value).count()
```

## Create Results DataFrame

The results are organized into a structured DataFrame with appropriate column names.

```python
result_data = [
    ("95th", p95_value, count_exceeding_p95),
    ("99th", p99_value, count_exceeding_p99),
```

```python
        ("99.5th", p995_value, count_exceeding_p995),
        ("99.95th", p9995_value, count_exceeding_p9995),
        ("99.995th", p99995_value, count_exceeding_p99995),
    ]

    result_df = spark.createDataFrame(
        result_data,
        [
            "Percentile of % change in stock price",
            "Value of % change in stock price",
            "Number of trades exceeding this value",
        ],)
```

## Display and Save Results

The final results are displayed in the console and saved as a CSV file with headers to an "output" location.

```python
result_df.show()
result_df.write.csv("output", header=True)
```

# Output

| Percentile of % change in stock price | Value of % change in stock price | Number of trades exceeding this value |
|---|---|---|
| 95th | 0.26856240126383024 | 1799468 |
| 99th | 0.5421184320266834 | 359897 |
| 99.5th | 0.7233747844977769 | 179949 |
| 99.95th | 1.674494638580514 | 17995 |
| 99.995th | 4.5015192422200085 | 1800 |



*Screenshot: The output of the execution*