# Table of Contents

# 01

# Overview

# Overview of Project

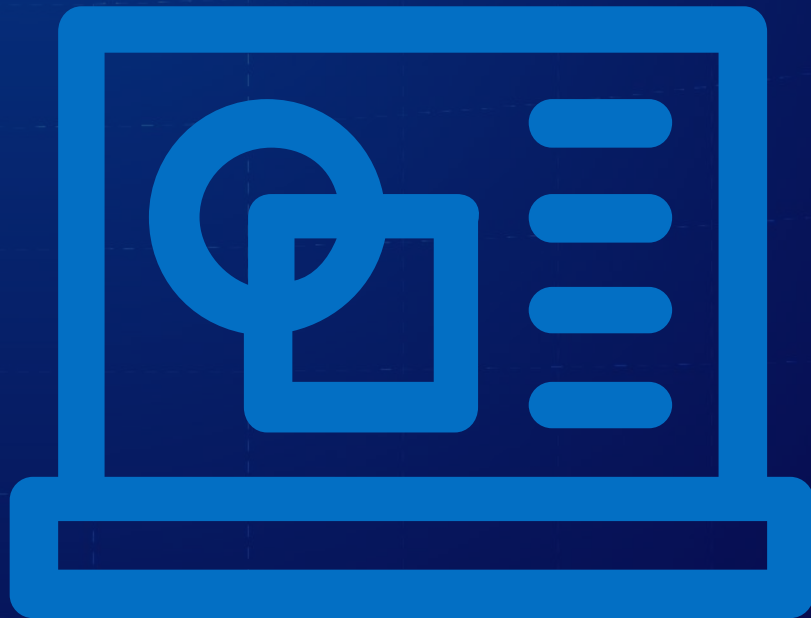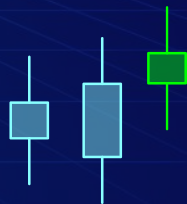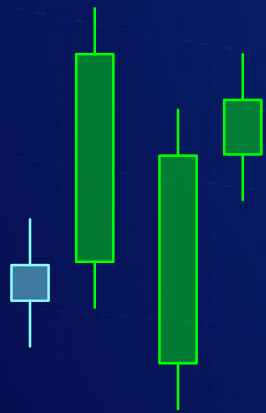By using semi-structured data, specifically about daily changes in Amazon stock prices, we will create a data pipeline locally, utilizing extraction, transforming, loading and analysis to achieve clean, structured data as well as insightful analysis.

The deliverable of this product will be a SQL database with transformed and cleaned data as well as statistical insights.

**02**

Checkpoint 1: Extraction

# Extraction

For extraction, we choose to go with a dataset from https://www.alphavantage.co/ using their stock market API.

Using the requests module in Python, we first retrieved the dataset from the URL, with some parameters.

```python
def api_call(output_size, api_key):
    return rq.get(f'https://www.alphavantage.co/query?function=TIME_SERIES_DAILY&symbol=AMZN&outputsize={output_size}&apikey={api_key}').json()
```
✓ 1.0s

This would efficiently return a JSON file of the dataset.
Later in the code, we can then make a call to api_call to retrieve the JSON and store it in a variable.

```python
data = api_call('full', 'KVUJT3KI90X21VGB')
```
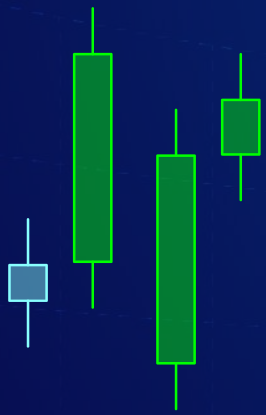
From this data, we can then further specify that we want the actual data from within the JSON.

```python
time_series_data = data["Time Series (Daily)"]
```

After this, we can then use this to finally convert into a dataframe.

```python
df = pd.DataFrame(time_series_data).T
```

**03**

Checkpoint 2: Transformation

Upon converting the data into json we converted all the values into floats.

This will detect and transform any values which are received as strings, ensuring minimal data loss while maintaining data quality and consistency.

```
stmnt1 = '''create table if not exists amzn (
    date date primary key,
    open float,
    high float,
    low float,
    close float,
    volume int
    )
'''
```

Although we found our data source to be very clean, we decided to quickly convert any leftover 'string' values in the dataset. This would not delete any irregular data but convert it and turn all of it into a more desired data type

```python
df = pd.DataFrame(time_series_data).T
df.rename(columns={'1. open': 'open', '2. high': 'high', '3. low': 'low', '4. close': 'close', '5. volume': 'volume'}, inplace=True)
df = df.reset_index().rename(columns={'index': 'date'})
df
```
✓  0.0s

Followed by the transformation code which set the new columns for our Dataframe

| | date | open | high | low | close | volume |
|---|---|---|---|---|---|---|
| 0 | 2024-02-16 | 168.7400 | 170.4200 | 167.1700 | 169.5100 | 48107744 |
| 1 | 2024-02-15 | 170.5800 | 171.1700 | 167.5900 | 169.8000 | 49855196 |
| 2 | 2024-02-14 | 169.2100 | 171.2100 | 168.2800 | 170.9800 | 42815544 |
| 3 | 2024-02-13 | 167.7300 | 170.9500 | 165.7500 | 168.6400 | 56345122 |
| 4 | 2024-02-12 | 174.8000 | 175.3900 | 171.5400 | 172.3400 | 51050440 |

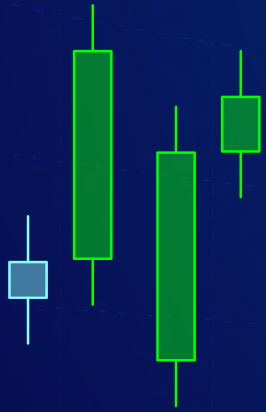| date | open | high | low | close | volume |
|---|---|---|---|---|---|
| 2024-02-16 | 168.7400 | 170.4200 | 167.1700 | 169.5100 | 48107744 |
| 2024-02-15 | 170.5800 | 171.1700 | 167.5900 | 169.8000 | 49855196 |
| 2024-02-14 | 169.2100 | 171.2100 | 168.2800 | 170.9800 | 42815544 |
| 2024-02-13 | 167.7300 | 170.9500 | 165.7500 | 168.6400 | 56345122 |
| 2024-02-12 | 174.8000 | 175.3900 | 171.5400 | 172.3400 | 51050440 |

Finally changing the dataframe from the index as primary key —> to the date, as each date is already unique.
Keeping our dataframe to essential information only.

# Adding a final column

Using the pandas .rolling() function we added a useful piece of data.

- Collects an average value of the closing average
- The rolling mean at each point is the average of the current value and the two preceding values within the window.
- Valuable information for traders

```python
# Calculate 30-day moving average
df2['30_MA'] = df2['close'].rolling(window=30).mean()
```

04

Checkpoint 3: Loading

# Data Loading

## Schema

Given the nature of the retrieved data, a simple schema was judged to be sufficient. Primary key would be the date, as that was guaranteed to be a unique value.

| AMZN_Stock | |
|---|---|
| PK | date_DATE_NOT NULL |
| | open FLOAT |
| | high FLOAT |
| | low FLOAT |
| | close FLOAT |
| | volume INT |

## Selective Insertion

Most of this checkpoint's complexity lay in the mechanics of only inserting non-duplicate data into the SQLite table every time the code was run to avoid the 'unique key' error.

Depending on the state of the table ie. empty/not empty, the parameters of the API call would change.

The output size is either 'full' (entire dataset) or 'compact' (last 100 days' worth):

```python
def api_call(output_size, api_key):
    return rq.get(f'https://www.alphavantage.co
```

If database table is empty, full dataset for AMZN stock is retrieved via api_call() function

If there is already data in the table, only last 100 rows of data will be retrieved

Either way, retrieved data is then converted to a dataframe.

```python
is_empty = True
if conn.execute(f"SELECT COUNT(*) FROM amzn;").fetchone()[0] == 0:
    data = api_call('full', 'KVUJT3KI90X21VGB')
else:
    data = api_call('compact', 'KVUJT3KI90X21VGB')
    is_empty = False


time_series_data = data["Time Series (Daily)"]
```

```python
if is_empty:
    df.to_sql('amzn', conn, index=False, if_exists='append')
    print("All data inserted.")
else:
    max_date = conn.execute("SELECT MAX(date) FROM amzn;").fetchone()[0]
    df_filtered = df[pd.to_datetime(df['date']) > pd.to_datetime(max_date)]
    if df_filtered.empty:
        print("No new data to insert.")
    else:
        df_filtered.to_sql('amzn', conn, index=False, if_exists='append')
        conn.commit()
        print("New data inserted.")
```

If table empty, entire dataframe will be inserted via SQLite

If there is already data in table, only the most recent, non-duplicate entries will be inserted

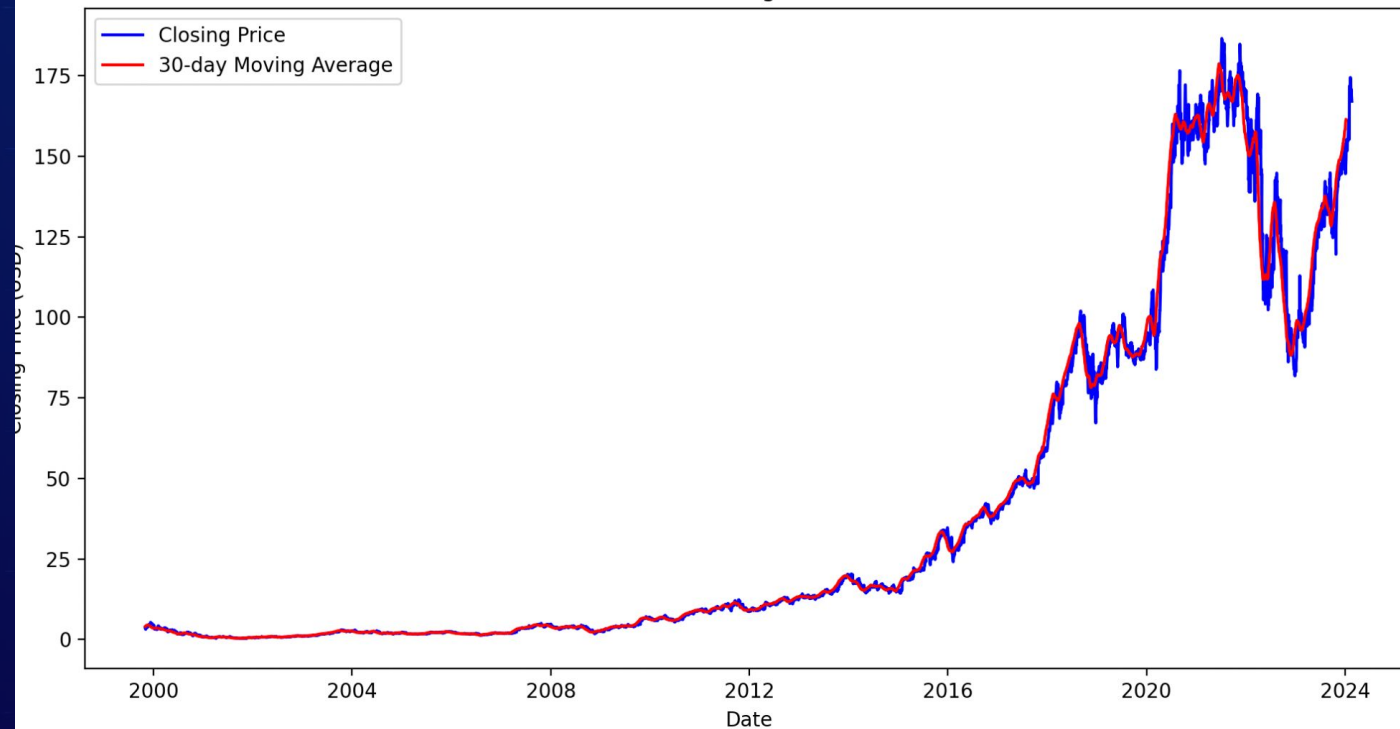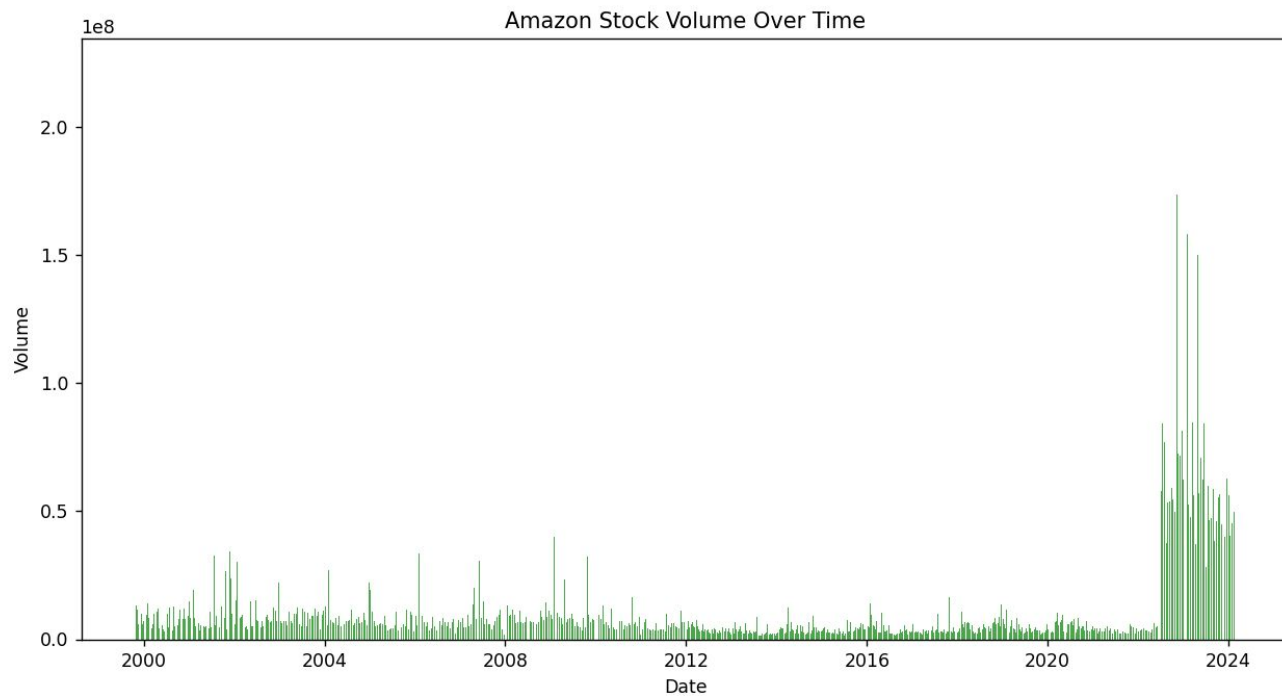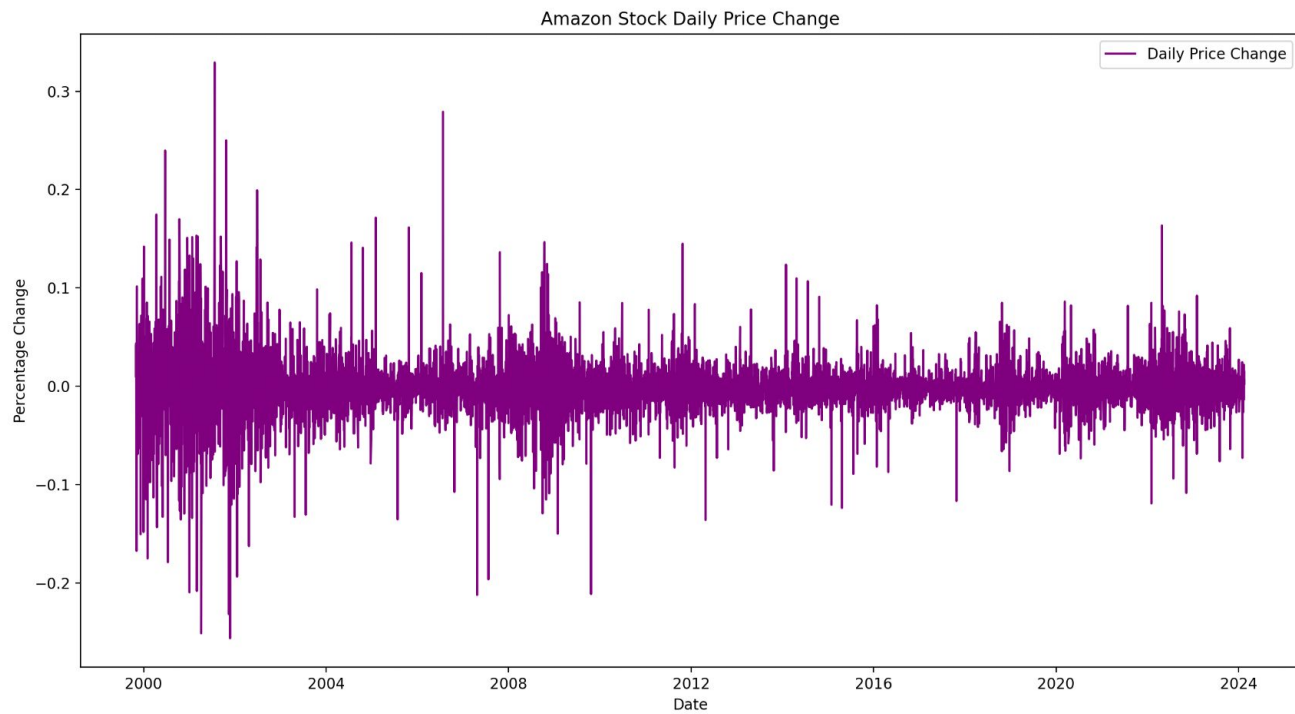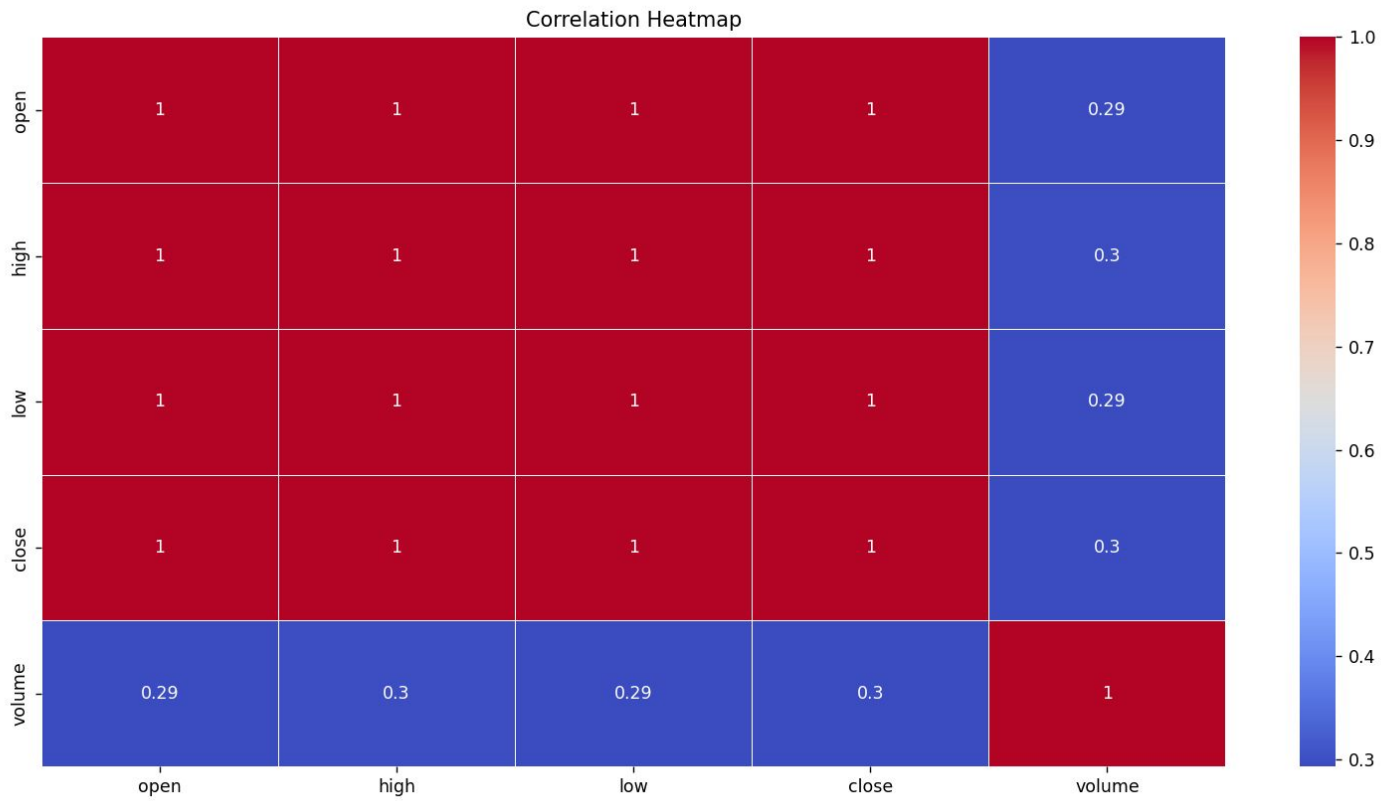By comparing every entry date to the latest date already present in the table, redundant data and unique key error are no longer issues
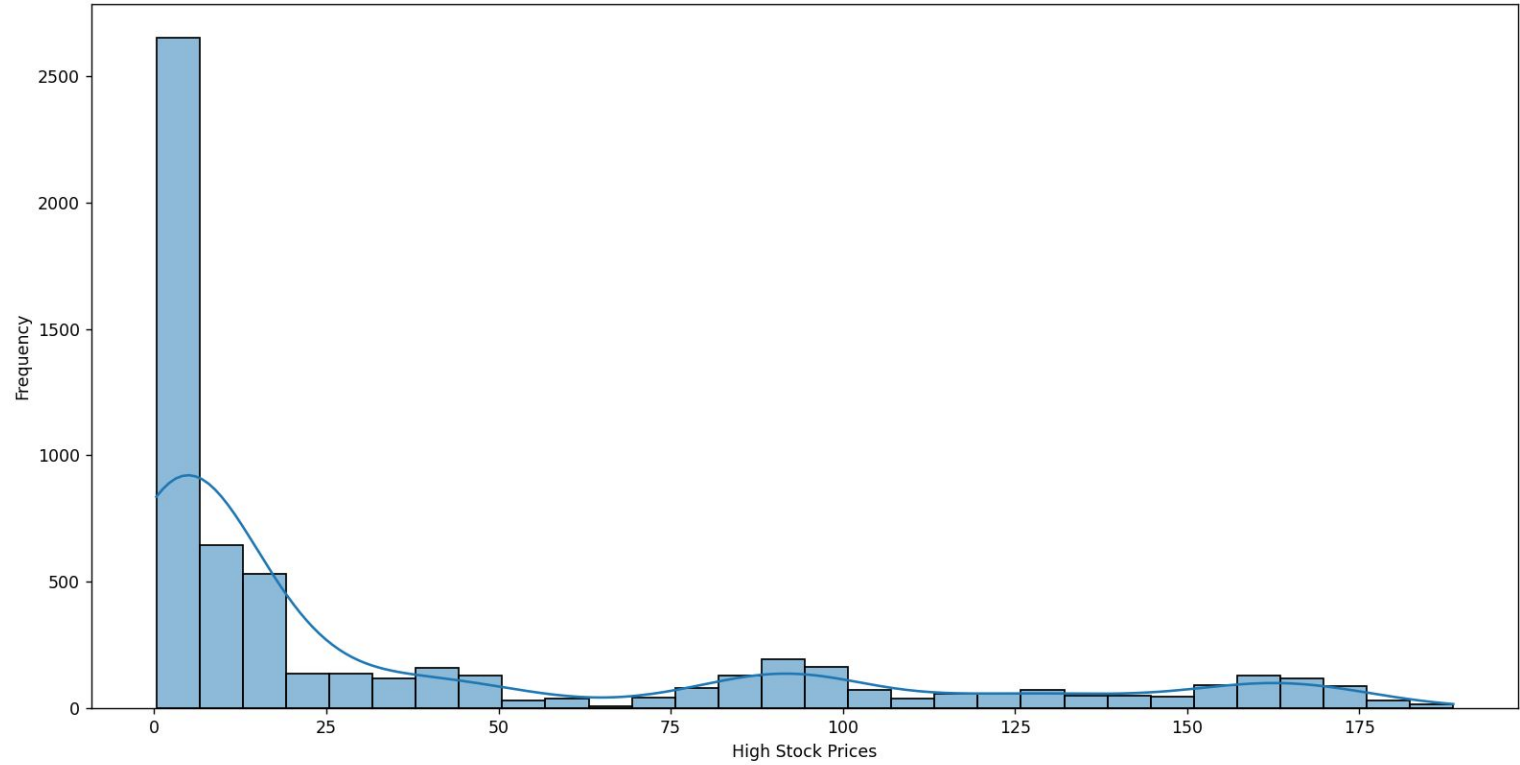
05

Checkpoint 4: Analysis

Amazon Stock Closing Prices Over Time
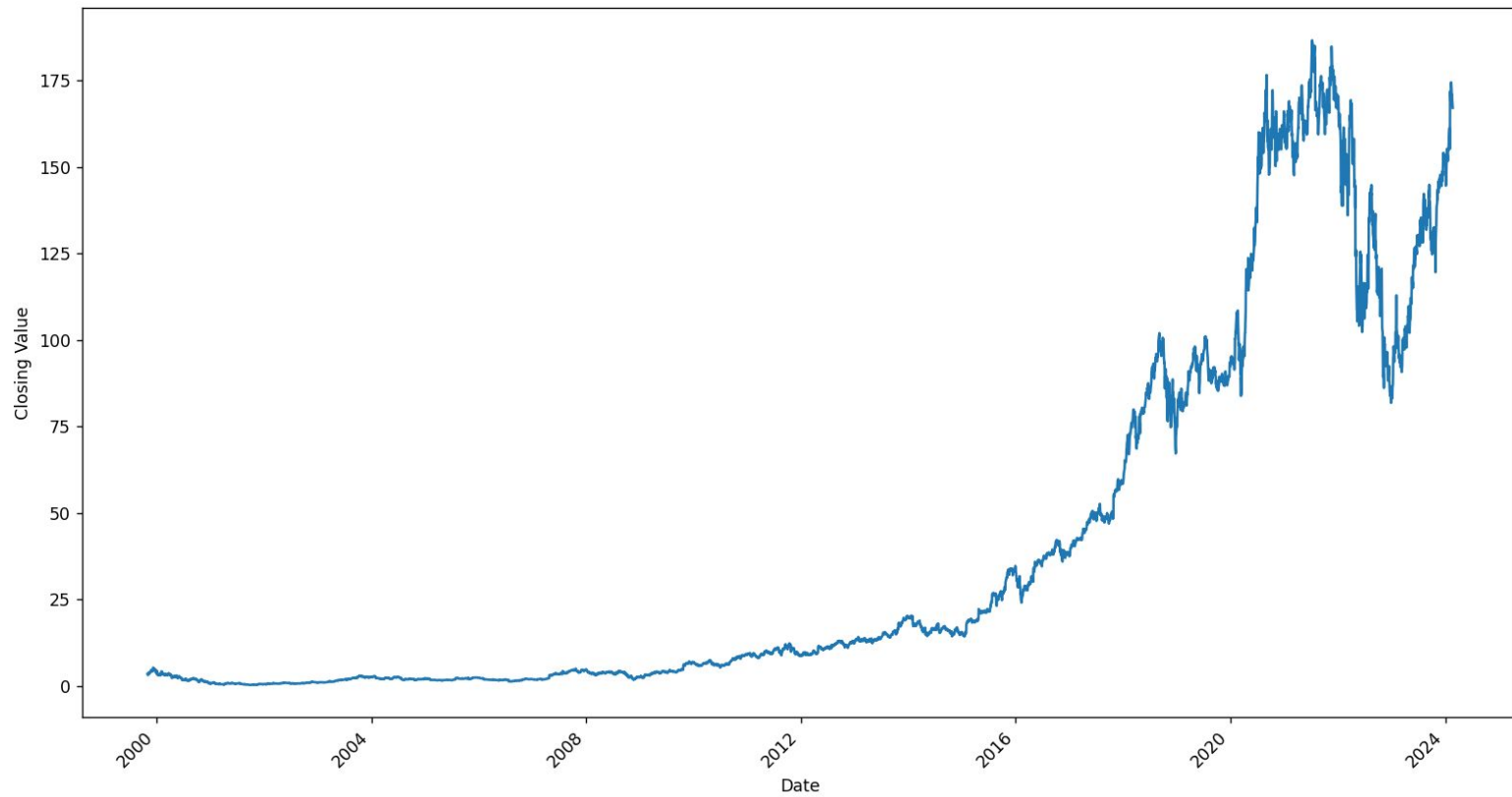
Amazon Stock Volume Over Time

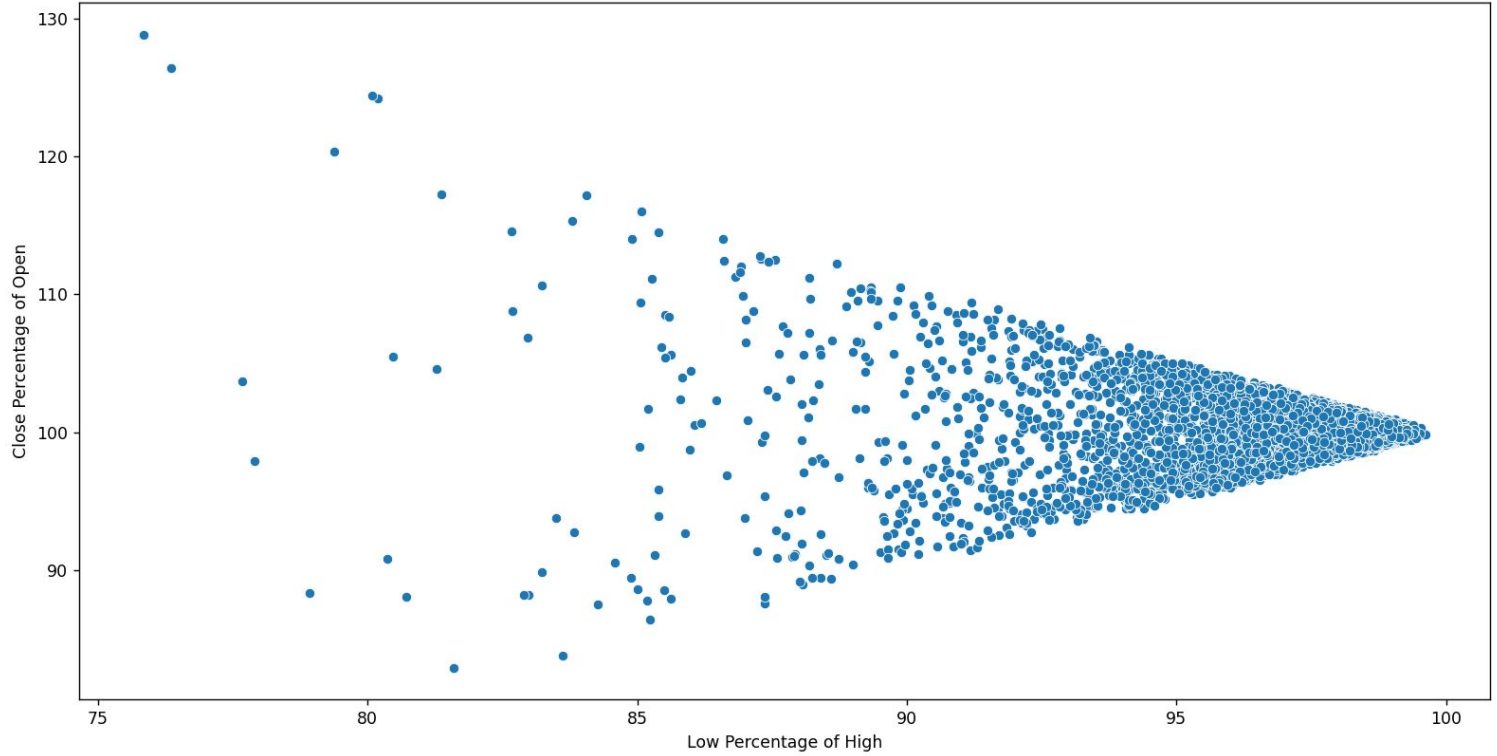Amazon Stock Daily Price Change
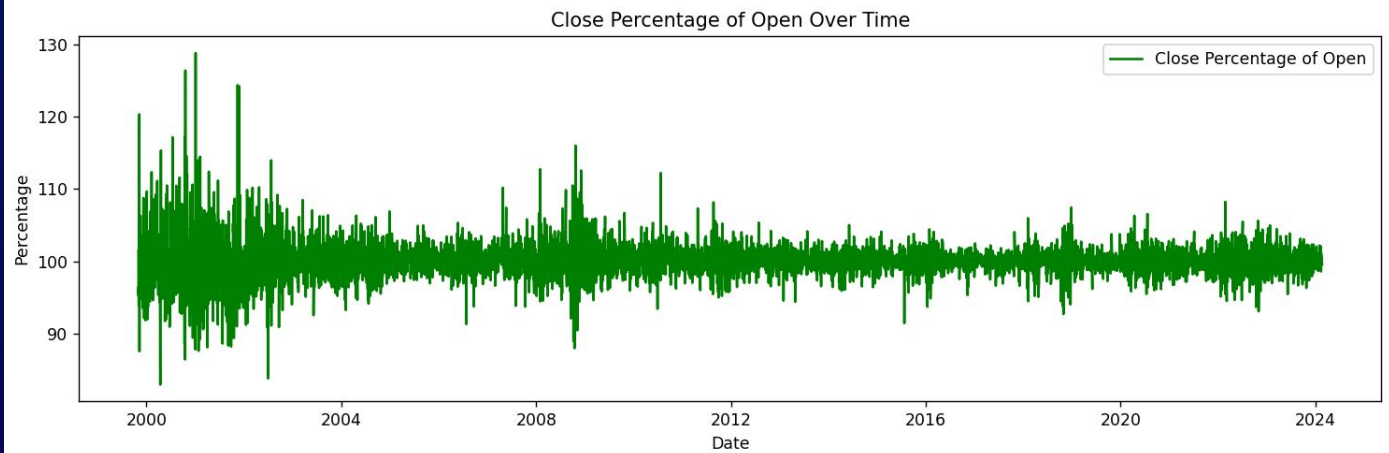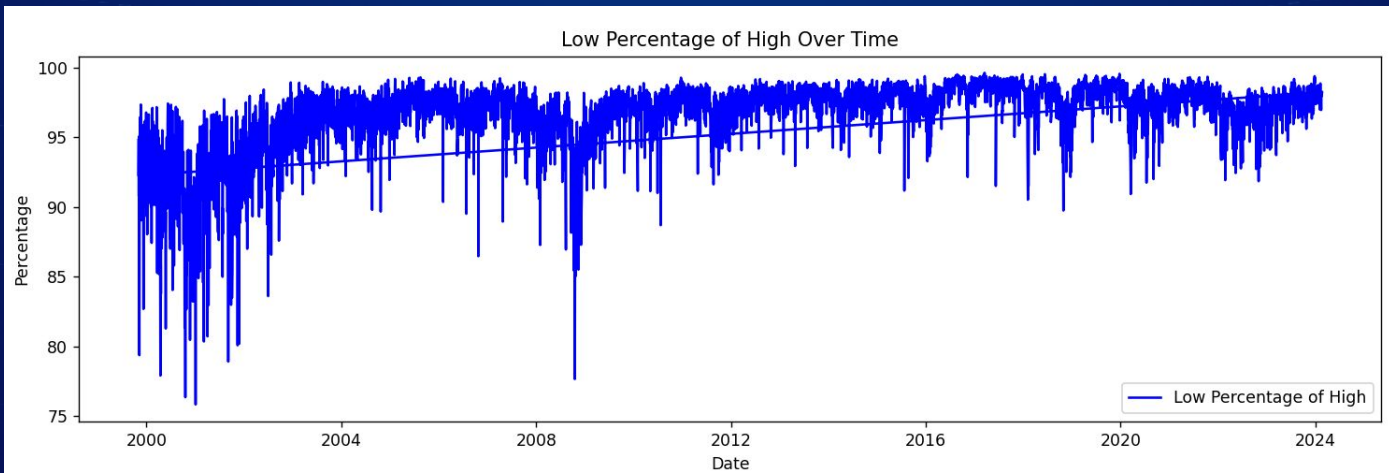
Correlation Heatmap

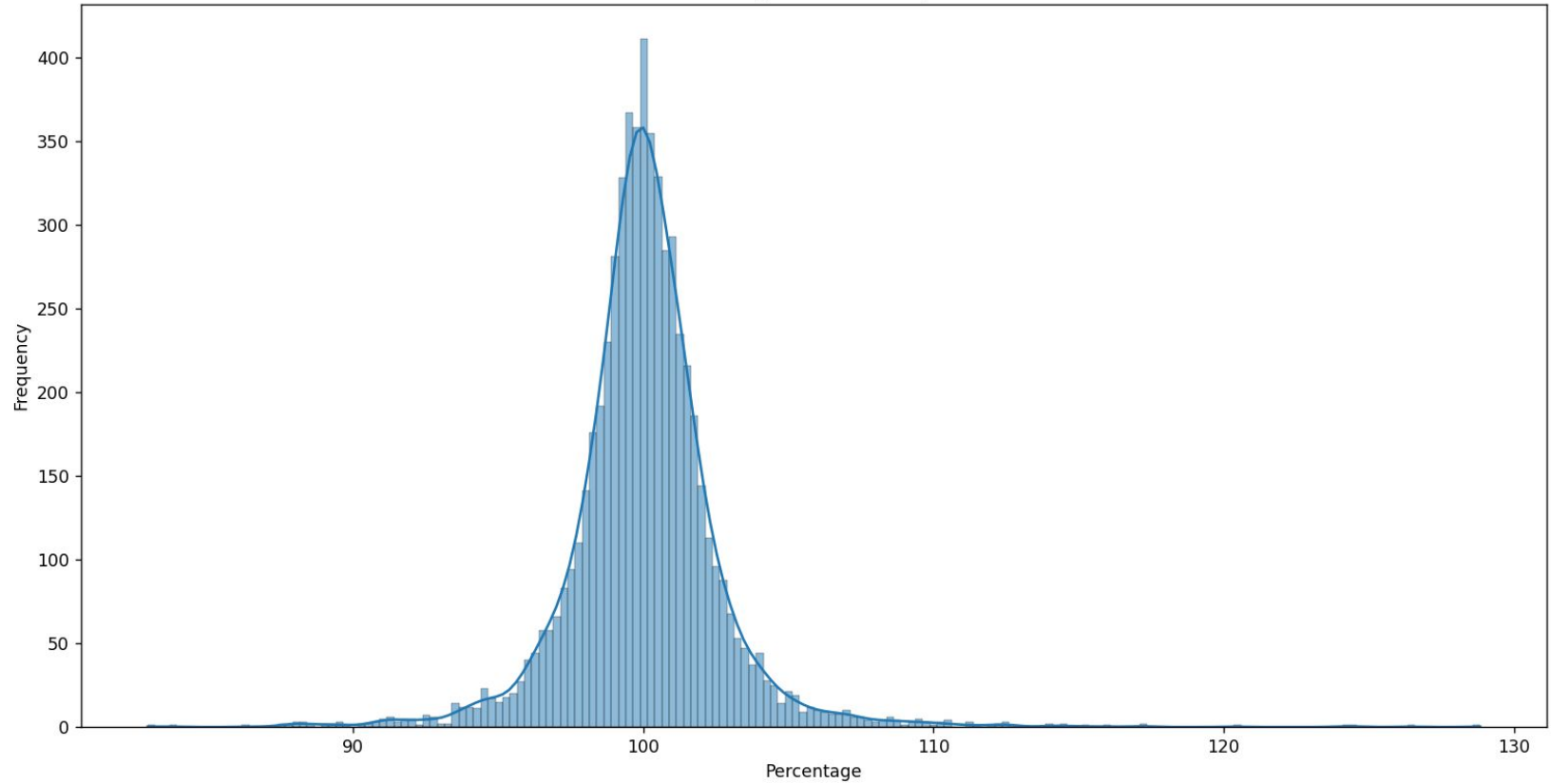Distribution of High Stock Prices

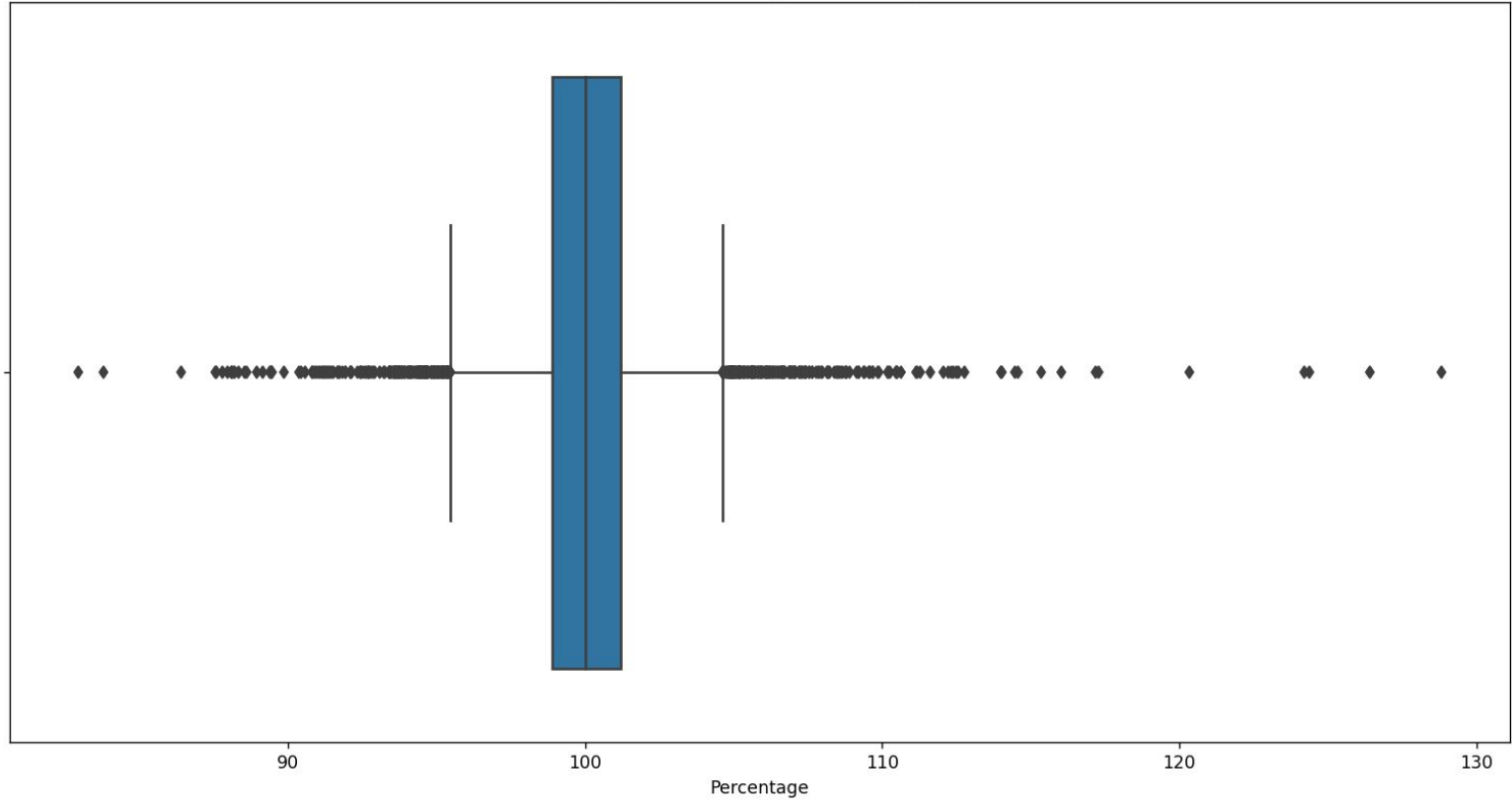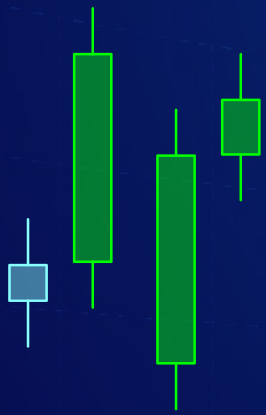Scatterplot: Correlation Between Price Percentage Relationships

Percentage of Close in Open

Boxplot of Percentage Close in Open

**06**

Lessons learned

# Lesson #1

## The Problem

Had to make the code infinitely reusable without flooding the database with duplicates or getting the 'unique key' error

Took current state of table into account, implemented checks to determine parameter of API call

Filtered dataframe to only latest date(s) not already present in table

## The Lesson

Take a goal-oriented approach to solve problems, don't be afraid to come at it from different angles

Thoroughly account for as many possible use-case scenarios as possible

# Lesson #2

## The Problem

Raw stock price data saw an enormous drop-off after 6/5/22

AMZN had underwent a stock split (by a factor of 20) at that time, dramatically skewing the data from that point onwards; full graphs from 1999-present showed a massive, misleading drop in price

We divided all data from before 6/5/22 by 20 prior to visualization and analysis

## The Lesson

Always look into significant discrepancies or outliers in raw data and account for them, otherwise the analysis could be rendered worthless

Many times, external research should be conducted

# Lesson #3

## The Problem

The nature of this project structure lent itself to team members finishing certain checkpoints very quickly. Leaving some team members done with their jobs in order for others to begin..

## The Lesson

Most people don't want to go and ask team members for help, and feel like their burdening others. Let them know you are bored! And want to help!

# Lesson #4

## The Problem

1. Initial challenge was deciding which aspects of the data to graph.
2. Once the visuals were created, extracting meaningful insights demanded a thorough understanding of the dataset
3. Selecting the right visuals was crucial.

## The Lesson

1. Deeper understanding of the dataset for better analysis
2. Experimenting with different graphs and refining them based on insights for overall quality of the visuals

# END