# Week 2 Attempt Summary

## Introduction

In this week we were advised to create a load and process data function to a dataset with multiple features, with the following functional requirements

- Should be able to specify and start and end date for the whole dataset
- Deal with Nan values in the dataset
- Split train/test data according to a specified ratio (Date)
- Should be able to download the data to the local machine and store them to future use
- Scaling in the data structure

## Load_and_process_data Function

```python
import pandas as pd
import yfinance as yf
import os
from sklearn.preprocessing import MinMaxScaler

def load_and_process_data(ticker, start_date, end_date,
nan_method='drop', split_method='ratio',
                          train_size=0.8, test_size=0.2,
split_date=None, random_state=None,
                          save_data=False, local_path=None,
scale_features=False, save_scaler=False):

    # Load data locally if available
    if local_path and os.path.exists(local_path):
        data = pd.read_csv(local_path, index_col='Date',
parse_dates=True)
    else:
        # Download the data
        data = yf.download(ticker, start=start_date, end=end_date)
```

```python
        # Save the data locally
        if save_data and local_path:
            data.to_csv(local_path)

    # Handle NaN values
    if nan_method == 'drop':
        data = data.dropna()
    elif nan_method == 'fill':
        data = data.ffill().bfill()

    # Ensure data is not empty
    if data.empty:
        raise ValueError(f"No data found for {ticker} between
{start_date} and {end_date}. Please check the date range and ticker
symbol.")

    # Split the data
    if split_method == 'ratio':
        train_data, test_data = train_test_split(data,
train_size=train_size, test_size=test_size, random_state=random_state,
shuffle=True)
    elif split_method == 'date' and split_date:
        train_data = data[:split_date]
        test_data = data[split_date:]

    # Scale features if required
    scalers = {}
    if scale_features:
        for column in train_data.columns:
            scaler = MinMaxScaler(feature_range=(0, 1))
            train_data[column] =
scaler.fit_transform(train_data[[column]])
            if not test_data.empty:
                test_data[column] =
scaler.transform(test_data[[column]])
            if save_scaler:
                scalers[column] = scaler

    return train_data, test_data, scalers
```

This function is in a different python file (load_process_data.py). I am importing the function with a import statement at the top

```python
import numpy as np
import matplotlib.pyplot as plt
import pandas as pd
import yfinance as yf
import os
import tensorflow as tf
import load_process_data as data_function
from sklearn.preprocessing import MinMaxScaler
from tensorflow.keras.models import Sequential
from tensorflow.keras.layers import Dense, Dropout, LSTM
```

And this to instantiate it in the main python script

```python
    # Load and process data
    train_data, test_data, scalers =
data_function.load_and_process_data(
        ticker=company,
        start_date=start_date,
        end_date=end_date,
        nan_method='fill',
        split_method='date',
        split_date='2020-01-01',
        save_data=True,
        local_path=local_path,
        scale_features=scale_features,
        save_scaler=save_scaler
    )
```

This week 2 code base allows the user to enter the start date and the end date manually through the terminal before training the model

```python
start_date = input("Enter Start date (YYYY-MM-DD): ")
    end_date = input("Enter End date (YYYY-MM-DD): ")
```

And It handles the Nan values with predefined Nan value handlers that comes with the pandas library

Dropna :- Removes data Frames that contains Nan values

Fillna :- Fills missing data Frames with a specified value

And with with some exception handling we have this code

```python
# Handle NaN values
    if nan_method == 'drop':
        data = data.dropna()
    elif nan_method == 'fill':
        data = data.ffill().bfill()

    # Ensure data is not empty
    if data.empty:
        raise ValueError(f"No data found for {ticker} between
{start_date} and {end_date}. Please check the date range and ticker
symbol.")
```

it raises a Error massage when ever the data array might be empty (for debugging pruposes)

And this code facilitates the functionality to split the data into 2, Training and testing data
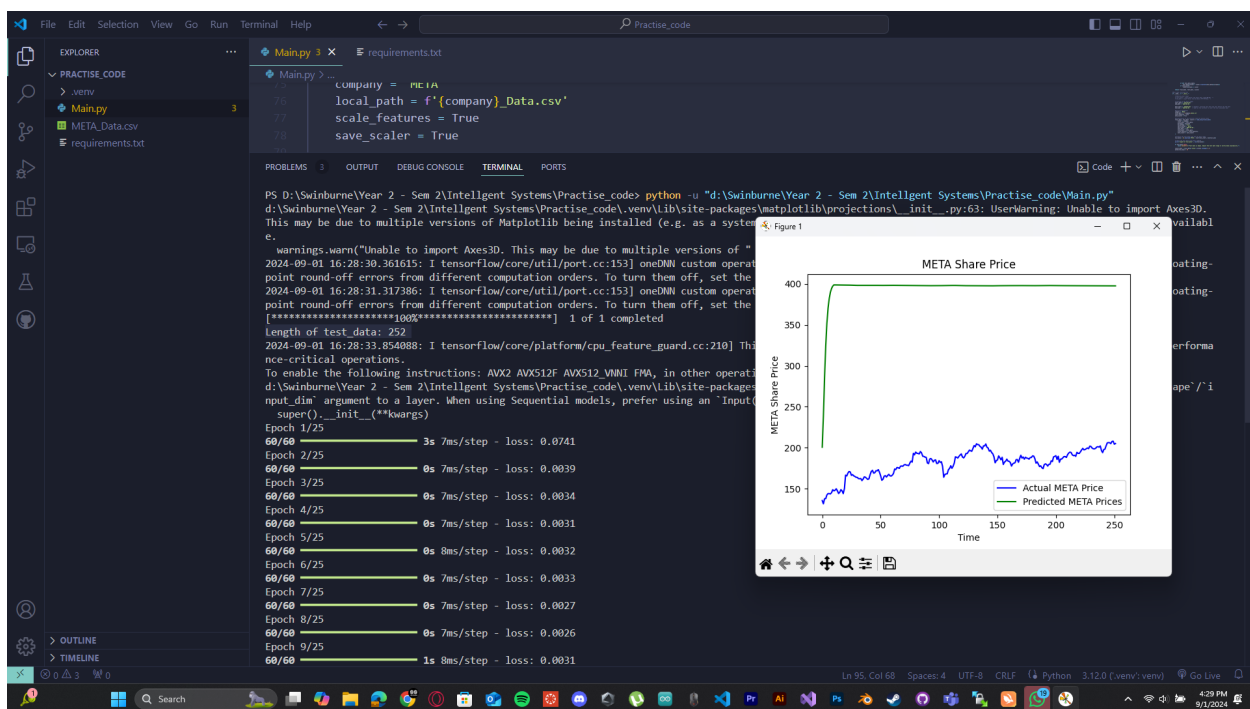
```python
# Split the data
    if split_method == 'ratio':
        train_data, test_data = train_test_split(data,
train_size=train_size, test_size=test_size, random_state=random_state,
shuffle=True)
    elif split_method == 'date' and split_date:
        train_data = data[:split_date]
        test_data = data[split_date:]
```
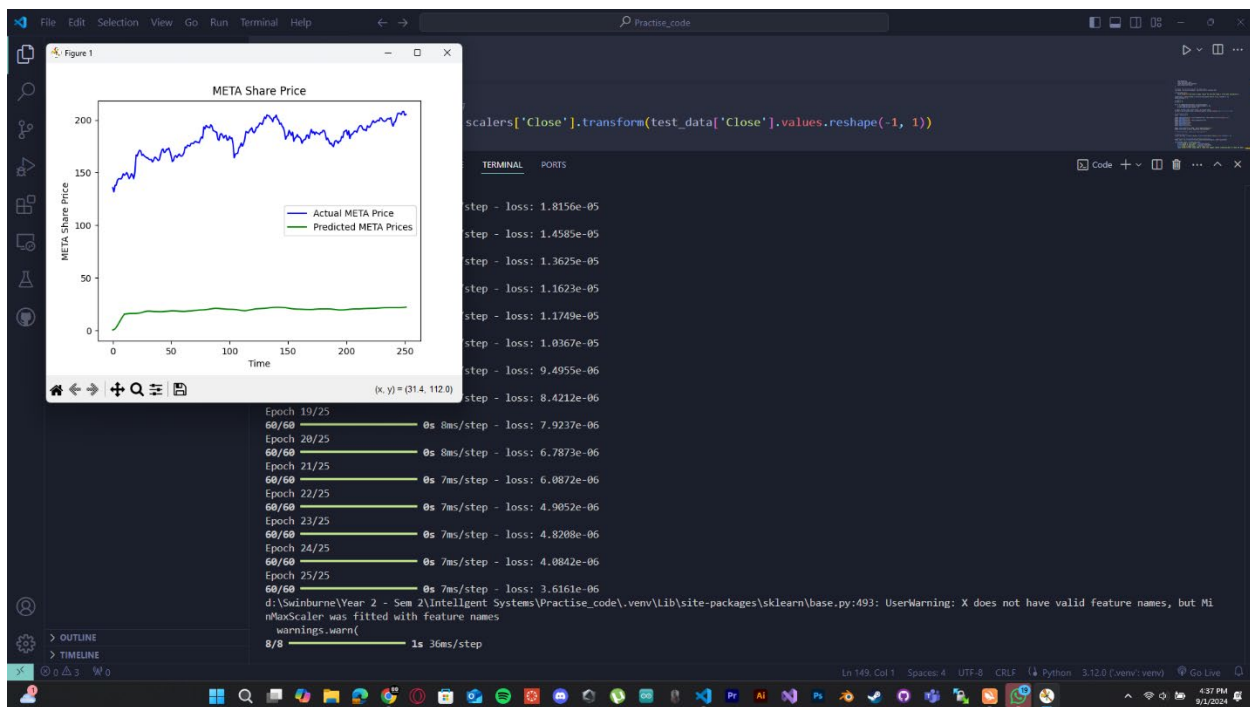
Yet the train_test_split function is still under development, for this week the date method will be used.

And then the last requirement was for adding scalers to the code base

```python
# Scale features if required
    scalers = {}
    if scale_features:
        for column in train_data.columns:
            scaler = MinMaxScaler(feature_range=(0, 1))
            train_data[column] =
scaler.fit_transform(train_data[[column]])
            if not test_data.empty:
                test_data[column] =
scaler.transform(test_data[[column]])
            if save_scaler:
                scalers[column] = scaler
```

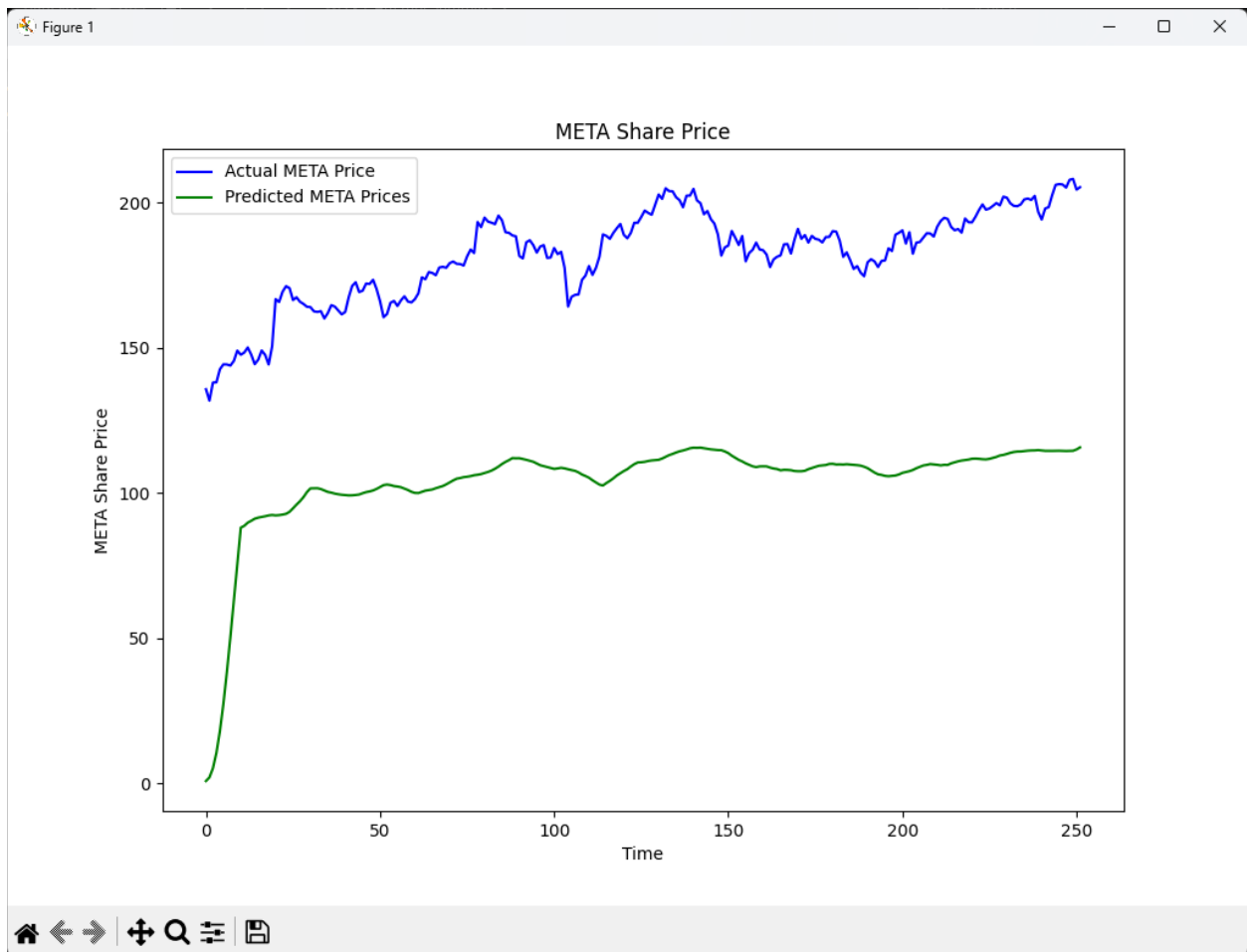I had problems with using the "Adam" optimizer when compiling the model

Then I changed the 'Adam' optimizer and used tensor flow keras adam optimizer instead, with 100 epochs

```python
    model.compile(optimizer =
tf.keras.optimizers.Adam(learning_rate=0.001),
loss='mean_squared_error')
    model.fit(x_train, y_train, epochs=50, batch_size=32)
```

By using that I was able to take a reasonable output



But this is very far from being perfect but this is still a work in progress