# RICE UNIVERSITY

# HOMEWORK 3: QUANTITATIVE FINANCIAL ANALYSIS WITH GENERATIVE AI

*STAT 682 - Fall 2024*

Authors:

Dolese, Ryker

Fang, Judy

Hao, Sirui

Hecht, Martin

Kumar, Krish

Thetford, Jackson

Professors:

Dr. Michael D Jackson

Max Lee

# Question 1 - Synthetic Data Generation for Financial Modeling

Objective: Generate synthetic financial data for asset prices to perform predictive modeling and analyze trends.

- Generate a synthetic time series for the daily closing prices of two stocks over two years using a generative model (e.g., Generative Adversarial Networks (GANs) or other generative algorithms).

- Calculate and plot the generated data's daily returns, volatility, and moving averages. Compare these values to the original data and assess differences and similarities in performance.

- Discuss how synthetic data generation could support risk modeling and decision-making for firms lacking historical data or when preparing for hypothetical scenarios.

## Answer:

**We run the Python code for this question. See A.1 for the full details of the code we used.**

**Generate synthetic financial data**

To generate synthetic financial data for the daily closing prices of TSLA and NVDA over two years, we implemented a Generative Adversarial Network (GAN) using PyTorch, consisting of a Generator and a Discriminator. Here's a brief overview of the process:

- **1. Data Collection:** We fetched historical closing prices for TSLA and NVDA using the Yahoo Finance API.

- **2. GAN Architecture:** The generator takes a random noise vector (latent dimension = 128) and transforms it through several linear layers with ReLU activations to produce a synthetic price.

  - **Generator:** The generator takes a random noise vector (latent dimension = 128) and transforms it through several linear layers with ReLU activations to produce a synthetic price.

  - **Discriminator:** The discriminator evaluates the authenticity of the generated data against real data using a sequence of linear layers with ReLU and Dropout, ending in a Sigmoid layer to output a probability.

- **Training:**

  - We trained both the generator and discriminator using **Binary Cross-Entropy loss**. The discriminator learns to distinguish between real and fake data, while the generator learns to produce realistic prices by fooling the discriminator.

– The model was trained over 5000 epochs with batch size 502, updating the generator and discriminator iteratively.

• **Synthetic Data Generation:** After training, we generated synthetic closing prices by inputting random noise into the trained generator model.

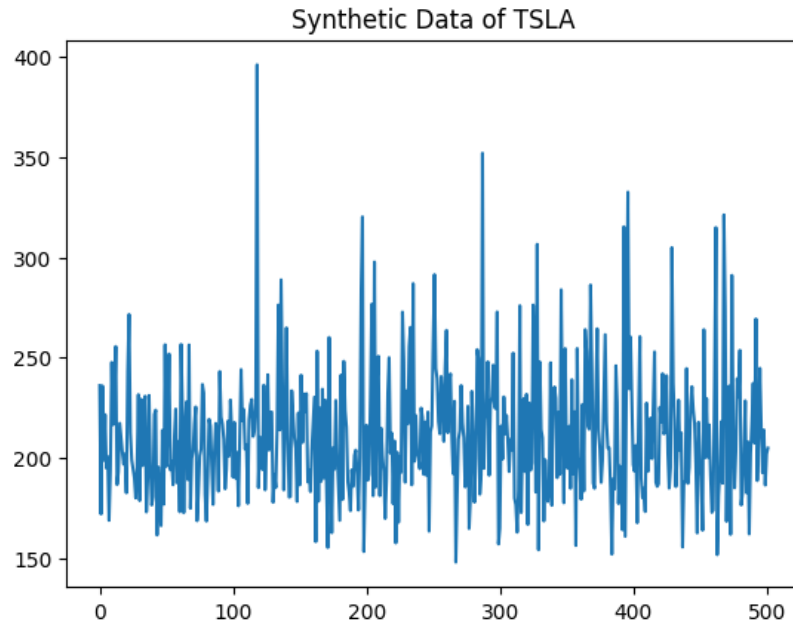The generated synthetic prices for TSLA and NVDA are shown at the below figure:
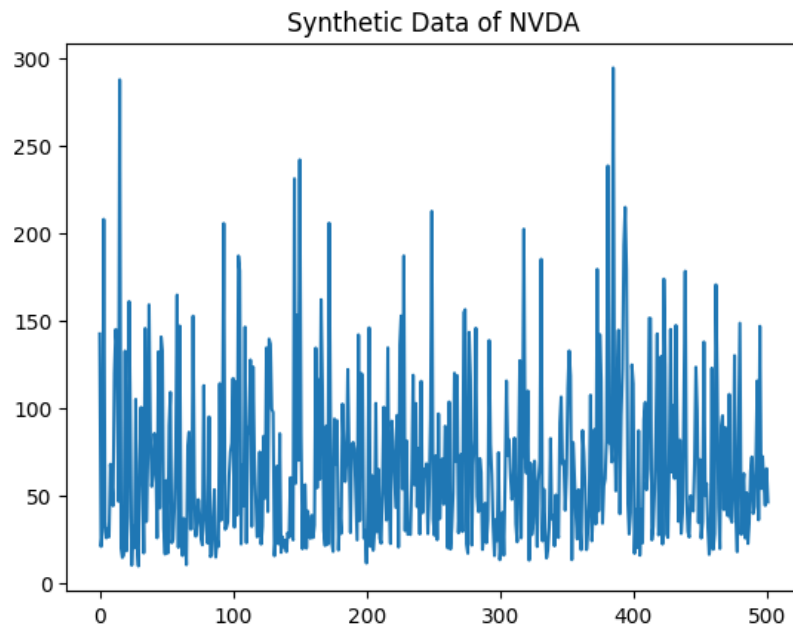


Figure 1



Figure 2

**Comparison of original data and synthetic data**

- **1. Price Trends**

  - **TSLA Price Comparison:** In Figure 3, the GAN-generated prices for TSLA generally follow the same trend as the real data but with more intense fluctuations. The synthetic prices demonstrate that the GAN was able to learn the overall trend but failed to capture the stability present in actual price movement.
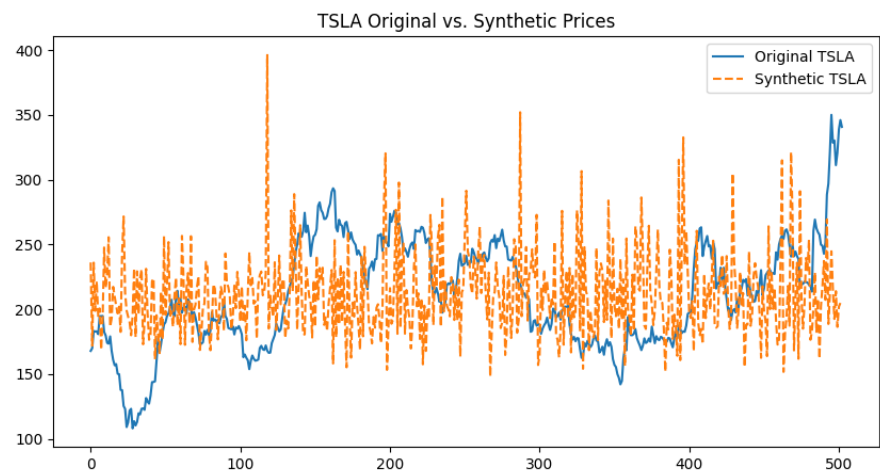


Figure 3

  - **NVDA Price Comparison:** In Figure 4, the synthetic prices for NVDA mirror the real price trend but show exaggerated daily fluctuations. This overemphasis on price movement may be due to the GAN learning general trend direction while overfitting on variability.
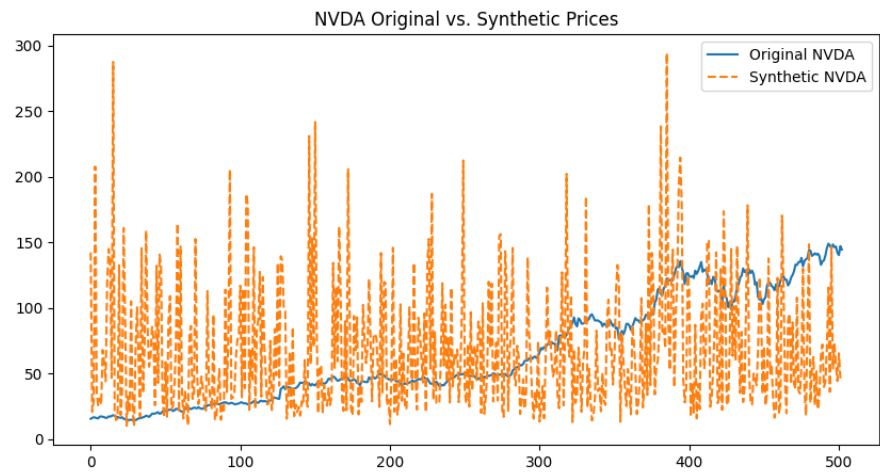


Figure 4

- **2. Daily Return**

  - **TSLA Daily Returns Comparison:** In Figure 5, the synthetic TSLA daily returns (orange, dashed line) fluctuate more sharply compared to the original data (blue line). This suggests that the GAN overestimated the variability in returns, capturing general trends but with an exaggerated amplitude in daily changes.
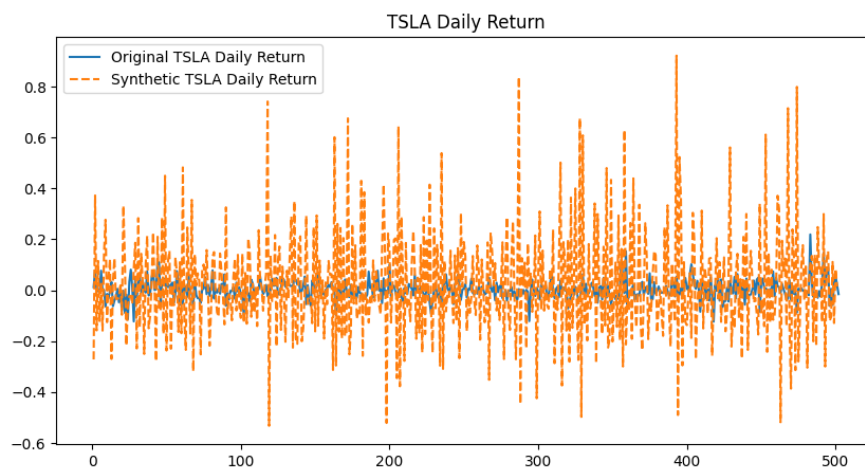


Figure 5

  - **NVDA Daily Returns Comparison:** Similarly, In Figure 6, the synthetic NVDA daily returns exhibit significant fluctuations compared to the original NVDA returns.The synthetic data's increased volatility suggests that the GAN may have generalized patterns without accurately reflecting day-to-day stability.
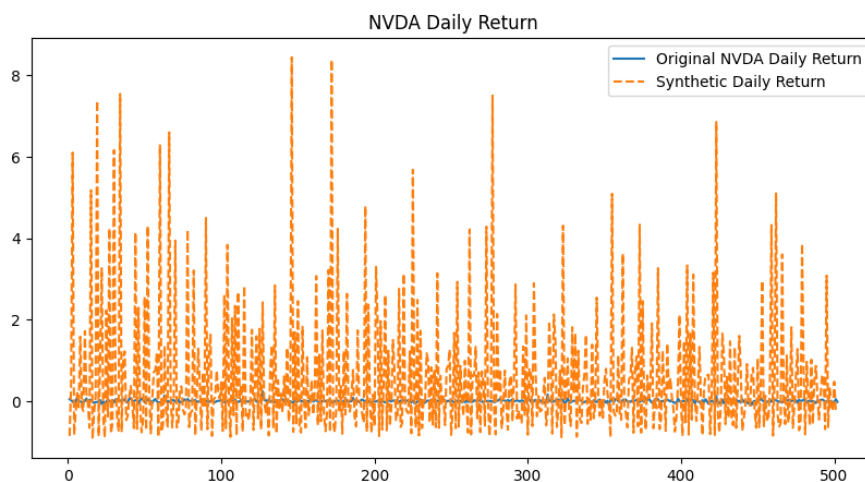


Figure 6

- **3. Moving Averages (20 days)**

    - **TSLA Moving Average (20 days) Comparison:** In Figure 7, the synthetic TSLA moving average remains generally aligned with the original moving average trend but does not capture its full amplitude. The smoother synthetic line reflects the GAN's tendency to produce a general trend without replicating finer cyclical variations.
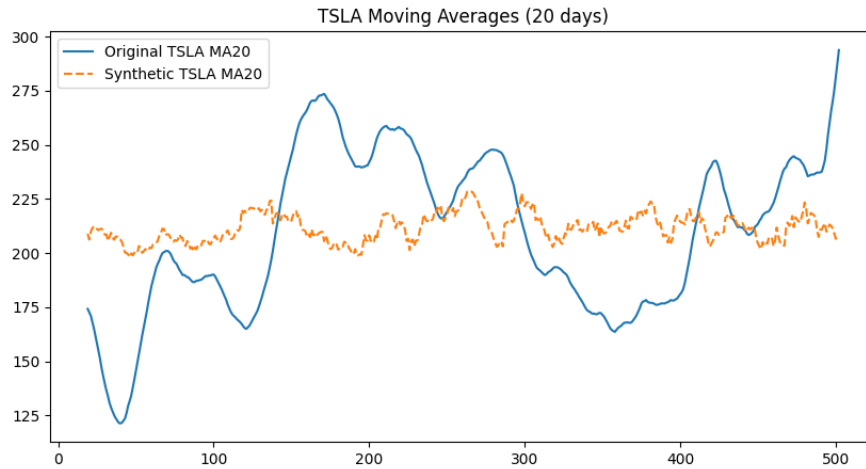


Figure 7

    - **NVDA Moving Average (20 days) Comparison:** For NVDA, in Figure 8, the synthetic moving average also approximates the overall trend but lacks the steeper upward movements observed in the original data. The synthetic line captures the broad direction but misses certain volatility features, emphasizing the GAN's limitations in modeling more complex patterns.
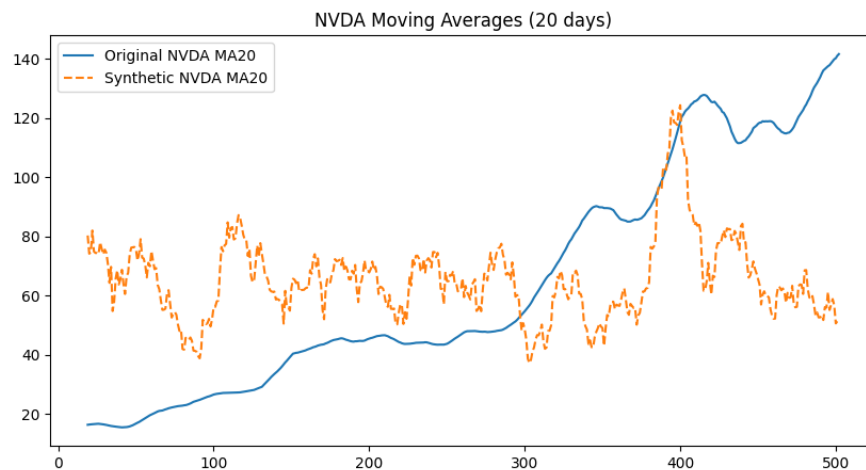


Figure 8

- **4. Volatility**

  - **TSLA Volatility Comparison:** In Figure 9, the volatility in synthetic TSLA data (orange, dashed) is consistently higher than in the original data (blue). While the original TSLA volatility gradually trends upward, the synthetic data shows frequent, sharp spikes. This discrepancy highlights the GAN's tendency to overestimate volatility while failing to replicate the smoother progression seen in real market data.



Figure 9

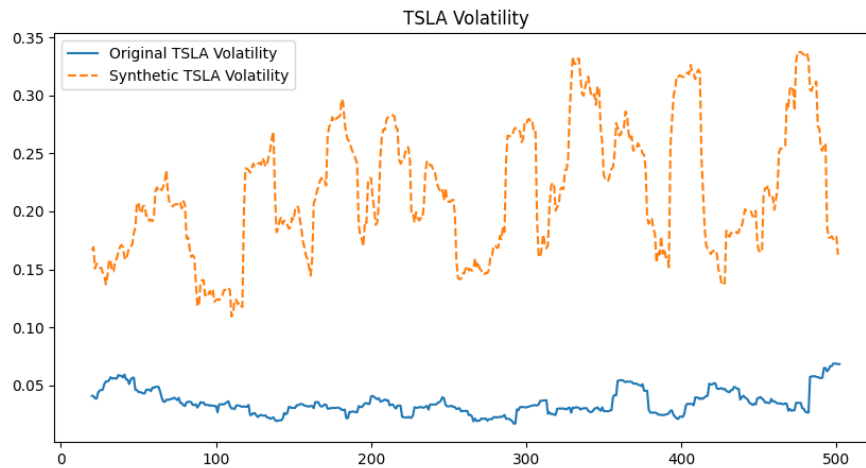  - **NVDA Volatility Comparison:** For NVDA, in Figure 10, The synthetic NVDA volatility also exhibits large, frequent spikes compared to the stable trend in the original data. Similar to TSLA, the GAN-generated volatility for NVDA lacks the gradual progression seen in real data, which suggests the synthetic data might not be suitable for scenarios that require accurate volatility representation.
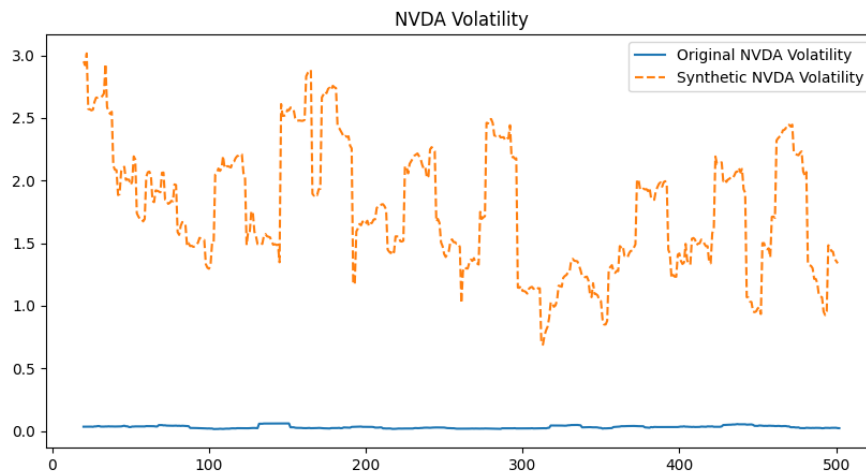


Figure 10

- **Overall Summary:** The GAN-generated synthetic data effectively approximates the general price trends for both TSLA and NVDA but struggles with accuracy in daily returns and volatility. This synthetic data could serve as a useful proxy for broader trend analysis but may not be ideal for applications needing precise day-to-day fluctuations, as it overestimates volatility and misses finer cyclical variations.

**Synthetic data generation using cases**

Synthetic data generation, like the GAN-produced stock prices here, can be very useful for firms in several ways:

- **1. Risk Modeling:** Firms lacking historical data can use synthetic data to model different risk scenarios. This is especially helpful for understanding how assets might behave in situations they haven't encountered before. By simulating a range of possible outcomes, businesses can better prepare for sudden market changes.

- **2. Testing Hypothetical Scenarios:** Synthetic data can be adjusted to create "what-if" situations. For example, a company could test how a stock might perform during a financial crisis or in response to a major event. This flexibility helps firms build strategies to manage future risks.

- **3. Overcoming Data Limitations:** Synthetic data allows firms to "create" data when there's limited or no access to real data. This is particularly useful for new companies or those entering a new market, as they can make more informed decisions even without years of historical data.

In summary, synthetic data generation provides a way for firms to explore potential risks and make strategic decisions in situations where they may otherwise have limited data.

# Question 2 - Creating a Financial Chatbot

Objective: Develop a chatbot that can assist with basic financial queries and provide simple financial advice based on market data.

- Create a chatbot using a generative AI model (such as GPT) that can:

  - Answer questions about current or synthetic stock performance.

  - Explain basic financial terms (e.g., "What is volatility?").

  - Provide simple advice or suggestions based on synthetic data generated in Question 1 (e.g., "Is now a good time to invest in stock A or B?" based on recent trends).

- Identify and discuss limitations and risks associated with relying on a generative AI-powered chatbot for financial advice, especially in live trading scenarios.

## Answer:

**We run the Python code for this question. See A.2 for the full details of the code we used.**

- **1. Framework and Interface:**

  - The chatbot was developed using **Streamlit**, a Python-based framework for building interactive web applications. This allowed the creation of a simple and intuitive interface for users to interact with the chatbot.

  - The interface, as shown in the Figure 11, provides:

    * A field to input the OpenAI API key.

    * An option to upload CSV files for analysis.

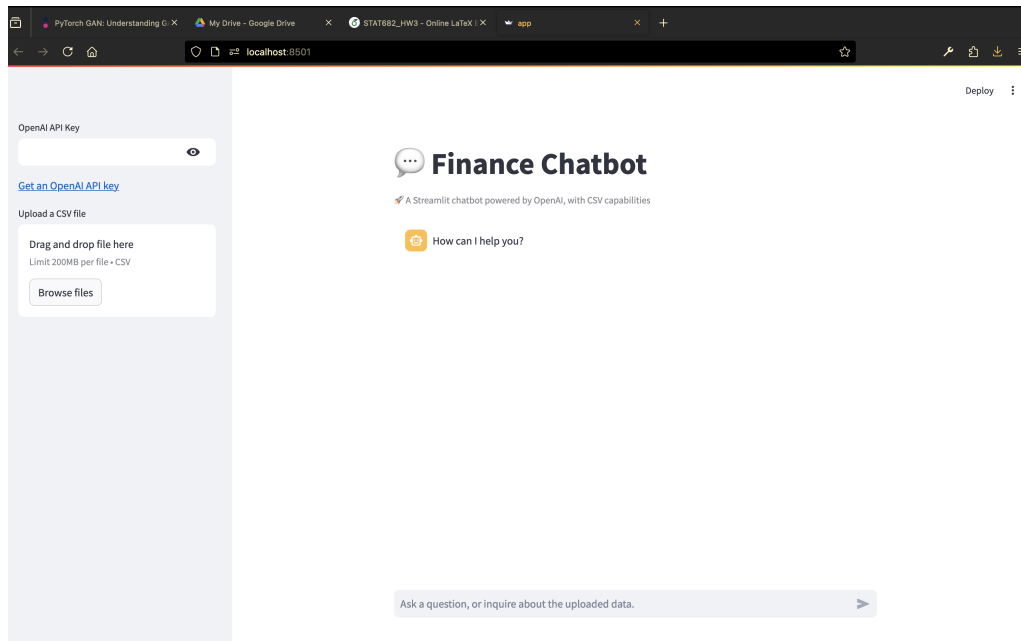    * A chat box for typing queries, with responses displayed dynamically.

Figure 11

- **2. Generative AI Model Integration:** The chatbot integrates the GPT-4o model using OpenAI's API. Users can input their API keys to enable the chatbot's functionality. This provides a personalized and scalable approach to leveraging generative AI.

- **3. Key Features - Answering Questions about Stock Performance:**

  – The chatbot utilizes the GPT-4o model to answer user queries about both current and synthetic stock performance.

  – Users can upload their own CSV files containing stock data or leverage synthetic data generated earlier (**from Question 1**).

  – The chatbot analyzes trends, returns, volatility, and moving averages from the uploaded data to provide detailed, context-specific answers.

    * **What is the recent stock peformance for real TSLA stock based on the data? Response (Figure12):**



Figure 12

* **What is the recent stock peformance for synthetic TSLA stock based on the data? Response (Figure13):**



Figure 13

- **4. Key Features - Explaining Basic Financial Terms:**

The chatbot is equipped to explain financial concepts in a simple, clear manner, users can ask questions such as:

  – **Explain a bear/bull market**
    **Response (Figure 14):**



Figure 14

– **"What is GDP"**

**Response (Figure 15):**



How can I help you?

what is GDP

Gross Domestic Product (GDP) is a measure of the economic performance of a country. It represents the total monetary value of all goods and services produced within a country's borders over a specific time period, typically calculated on an annual or quarterly basis. GDP is used to gauge the health of an economy, with higher GDP indicating a larger, more productive economy and lower GDP indicating a smaller, less productive one.

There are three primary ways to calculate GDP:

1. **Production (or Output) Approach**: This measures the value of all goods and services produced in an economy. It sums up the added value at each stage of production and subtracts the cost of intermediate goods.

2. **Income Approach**: This calculates GDP by summing total compensation to employees, gross profits for firms, and taxes minus subsidies on products.

3. **Expenditure Approach**: This is the most commonly used method and calculates GDP by adding up all expenditures made in an economy, including consumption, investment, government spending, and net exports (exports minus imports).

GDP can be measured in nominal terms, which does not take inflation into account, or in real terms, which adjusts for inflation to reflect the true growth in output.

Figure 15

• **5. Key Features - Provide Suggestions based on synthetic data:**
Similar to **Key Features - Answering Questions about Stock Performance**, users can upload their own CSV files containing stock data or leverage synthetic data generated in **Question 1**.

– **Provide simple advice or suggestions based on this synthetic data, is now a good time to invest in TSLA(based on recent trends)?**
**Response (Figure ??)**

Figure 16

- **6. Limitations and risks of using a generative AI-powered chatbot for financial advice in live trading scenarios**

    - **Overgeneralization:** Generative models like GPT can provide advice based on patterns but lack the ability to fully understand complex market dynamics, news events, or unexpected shifts.

    - **Lack of Real-Time Data**: The data might not be up-to-date. In live trading, outdated information can lead to poor decisions.

    - **Sensitivity to Training Data:** If the model is trained on biased or incomplete data, it might produce misleading responses or overlook key market nuances.

    - **No Legal or Expert Guarantee:** The chatbot is not a licensed financial advisor and can't provide personalized or legal advice, increasing the risk of following unsuitable suggestions.

    - **High Risk in Volatile Markets:** In fast-changing markets, advice based on historical or synthetic trends may fail to account for sudden movements, leading to potential financial losses.

    - **Emotional Factors:** The chatbot cannot factor in emotional aspects of investing, such as fear or greed, which play a big role in live trading decisions.

# Question 3 - Sentiment Analysis on Financial News Articles

Objective: Use a generative AI model to perform sentiment analysis on financial news articles and understand its impact on stock movements.

- Collect or generate a set of synthetic financial news articles about stocks A and B. Include both positive and negative sentiment to simulate a realistic news cycle.

- Use a generative AI model or pre-trained transformer (e.g., OpenAI's GPT, BERT) to perform sentiment analysis on the articles. Categorize each article as positive, negative, or neutral.

- Calculate the overall sentiment score and analyze its correlation with both the original data and synthetic stock price movements from Question 1.

- Describe any observed relationships and discuss the reliability of using sentiment analysis in trading strategies.

## Answer:

**We run the Python code for this question. See A.3 for the full details of the code we used.**

To answer this question we explored and generated articles related to the assets analyzed in question 1: TSLA and NVDA. In our analysis, we brought real articles and generated two synthetic articles with ChatGPT. Some of the variables we considered when creating our sets were the perceived sentiment and the length of the articles. For a more robust study, we brought different articles from different dates and authors. The choice of the sentiment for the articles were completely arbitrary, as our goal was to create similar environments which allow for relevant and respective comparisons, without selecting any specific sentiment for each article. Finally, based on these data we generated articles using ChatGPT, which can be seen in the appendix with the code. We made sure to feed each prompt with different articles and data to model different and new opinions, and mimic a symbolic 'date' for each article.

Our approach consists of splitting the articles into chunks of 256 characters and analyze the sentiment chunk by chunk. The analysis was conducted by two Natural Language Processing Models trained to understand sentiment. The first one, "mrm8488/distilroberta-finetuned-financial-news-sentiment-analysis", is a financially literate model, meaning it has been trained specifically to understand financial terms. The second one, "finiteautomata/bertweet-base-sentiment-analysis", is a more general NLP Model that will outperform the former overall, but is not specifically trained to understand finances. After running test in both models throught the Exploratory Data Analysis process, we discovered two things that guided our procedures. Firstly, the non-financial model took articles of up to 12 chunks (so roughly 3000 characters). Secondly,

the non-financial article tended to characterize as neutral the chunks containing numbers or financial terms, deeming the model inaccurate for this specific use case. For these two reasons, we continued with our study using the financial model exclusively.

The comparison methods chosen are as follow. Given a publish date $p$, we can assume the article is the result of a sentiment of up to $n$ previous days, and that the effect of the article could result in up to $m$ days of influence and persistent sentiment. We will come back to the relevance of sentiment analysis later on, but given the volatile nature of trading, we will assume a symmetric interval of $n = m = 2$. Therefore, the comparison will be done from $[p - 2, p + 2]$.

To prepare the data we normalized the Daily Return of each of the asset's stocks and compared the daily values to the sentiment analysis. The sentiment analysis score (unidimensional) is the weighted sum of the negative (-1), neutral (0) and positive (+1) scores of each chunk, weighted by the amount of chunks with each respective sentiment. For example, in a 3 chunk article with one positive, one neutral and one negative chunk, the sentiment score will be 0, whereas in a 3 chunk, 3 negative article, the sentiment score will be -1.

As mentioned, the Daily Return data was initially raw, and then normalized to fit a similar scale that allows direct comparison to the sentiment score. However, we then applied a $tanh(x)$ transformation to the Daily Return data to squish it between -1 and 1: the scores that the sentiment analysis took. Only then, was the date ready to aggregate and compare.

To aggregate the data we came up with three different metrics: Compound, Average, and Weighted Average. The compound multiplies the stock daily returns in the 5 day window, the average averages them, and the weighted average attributes more weight to the metrics nearest to the article in a $[1, 2, 5, 2, 1]$ pattern resembling a Gaussian curve, that represents how data nearest to the article publication is of the most value.

**Results:**

The following two tables show the data outlined previously for Synthetic data and Real data respectively. In terms of the results obtained, we can see in the table below the $Sentiment - Compound$, $Sentiment - Avg$ and $Sentiment - WeightedAverage$ columns. Given that each metric ranges between -1 and 1, the subtractions range between -2 and 2, with 0 being the desired value: when the two elements of the subtraction cancel each other for being equivalent. Below is Table 1, showing study done for the Synthetic data generated.

For the sake of avoiding redundancy and being concise, we will observe the $Sentiment - WA$ column in more detail. The average value of that column taking each factor's absolute value, is $0, 580$ for TSLA stock and $0.786$ for NVDA stock. Let's observe our results for Table 2 before explaining them and studying our results as a whole.

In Table 2, we can see the same columns as in Table 1, only this time studying the real data from TSLA's

| Ticker/Date | Sentiment | Synth -2 | Synth -1 | Synth 0 | Synth +1 | Synth +2 | Syn Aggr | Sentiment - Compound | Avg | Sentiment - Avg | Weighted Avg | Sentiment - WA |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| TSLA Oct 3 | 0.157 | 0.439932 | -0.635461 | 0.996596 | -0.649177 | -0.954479 | -0.9833 | 1.1403 | -0.2327 | 0.3897 | 0.3798 | -0.2228 |
| TSLA Oct 14 | -0.385 | -0.913162 | 0.998425 | -0.882892 | -0.566361 | 0.34261 | -0.9882 | 0.6032 | -0.3421 | -0.0429 | -0.8242 | 0.4392 |
| TSLA Oct 23 | 0.55 | 0.350998 | -0.905446 | 0.06169 | 0.349094 | 0.477625 | -0.7296 | 1.2796 | 0.0220 | 0.5280 | 0.0049 | 0.5451 |
| TSLA Oct 28 | 0.64 | 0.349094 | 0.477625 | -0.772156 | 0.491464 | -0.807261 | -0.8694 | 1.5094 | -0.0701 | 0.7101 | -0.4762 | 1.1162 |
| NVDA Oct 3 | 0.75 | 0.955478 | -0.307255 | -0.609469 | 0.331414 | -0.595702 | -0.7152 | 1.4652 | -0.0273 | 0.7773 | -0.5279 | 1.2779 |
| NVDA Oct 15 | -0.307 | -0.379751 | 0.046219 | 0.415308 | -0.665591 | -0.593601 | -0.8752 | 0.5682 | -0.3371 | 0.0301 | -0.0271 | -0.2799 |
| NVDA Oct 25 | 0.889 | 0.241008 | -0.255857 | -0.620707 | 0.271787 | -0.610381 | -0.8264 | 1.7154 | -0.1302 | 1.0192 | -0.6882 | 1.5772 |
| NVDA Oct 31 | -0.333 | -0.610381 | -0.034261 | 0.181354 | -0.267461 | -0.563613 | -0.8579 | 0.5249 | -0.3550 | 0.0220 | -0.1741 | -0.1589 |

Table 1: Sentiment Analysis and Synthetic Metrics for TSLA and NVDA

| Ticker/Date | Sentiment | Real -2 | Real -1 | Real 0 | Real +1 | Real +2 | Compound | Sentiment - Compound | Avg | Sentiment - Avg | Weighted Avg | Sentiment - WA |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| TSLA Oct 3 | 0.157 | -0.350008 | -0.258788 | -0.694045 | -0.15471 | -0.615048 | -0.9520 | 1.1090 | -0.5041 | 0.6611 | -1.0525 | 1.2095 |
| TSLA Oct 14 | -0.385 | -0.30433 | -0.983357 | 0.101349 | -0.014158 | 0.147539 | -0.9856 | 0.6006 | -0.3398 | -0.0452 | -0.3290 | -0.0560 |
| TSLA Oct 23 | 0.55 | -0.278241 | -0.168613 | -0.528995 | 0.999979 | 0.676206 | -0.0525 | 0.6025 | 0.1080 | 0.4420 | -0.1169 | 0.6669 |
| TSLA Oct 28 | 0.64 | 0.999979 | 0.676206 | -0.617586 | -0.350008 | -0.258788 | -0.3824 | 1.0224 | 0.0112 | 0.6288 | -0.3389 | 0.9789 |
| NVDA Oct 3 | 0.75 | -0.857293 | 0.320371 | 0.707588 | 0.348791 | 0.489701 | -0.3535 | 1.1035 | 0.1093 | 0.6407 | 0.9017 | -0.1517 |
| NVDA Oct 15 | -0.307 | -0.156541 | 0.531711 | -0.921765 | 0.669535 | 0.118817 | -0.8112 | 0.5042 | -0.0949 | -0.2121 | -0.4488 | 0.1418 |
| NVDA Oct 25 | 0.889 | -0.770298 | 0.03226 | 0.092363 | -0.360782 | 0.004651 | -0.8337 | 1.7227 | -0.3059 | 1.1949 | -0.1922 | 1.0812 |
| NVDA Oct 31 | -0.333 | 0.004651 | -0.517232 | -0.923428 | 0.428193 | -0.007509 | -0.9474 | 0.6144 | -0.3271 | -0.0059 | -0.9596 | 0.6266 |

Table 2: Sentiment Analysis and Real Financial Metrics for TSLA and NVDA

and NVDA's assets. The $Sentiment - WA$ average value for TSLA's stock was $0.728$ and $0.500$ for NVDA.

Given the range of the potential results is $[-2, 2]$, and the fact that our maximum deviation from $0$ is NVDA's synthetic data, being $0.786$, the accuracy radius of the results are $39.3\%$. These results show accuracy to a certain extent. The factors we consider relevant to draw conclusions out of these results are the amount of articles considered, the relevance of each article, our ability to generate relevant synthetic data, the assets chosen, and the accuracy of the NLP model chosen.

In terms of the previous variables, we separate the problem into two buckets: sentiment and data. In the sentiment bucket fall the amount and relevance of articles, and the NLP model we chose. In the data bucket fall our ability to generate synthetic data and the assets chosen.

Starting with the sentiment bucket, we highlighted these variables since they are the ones that affected the sentiment results. A larger amount of articles would allow to draw a regression to the mean, and the standarization and normalization of them. The relevance of each article is an important variable as well, as it would allow to a more accurate result. Having larger amounts of articles could also potentially balance out the different biases from the articles (regression to the mean).

In the data bucket fall our ability to create synthetic data, and the assets we chose. Our ability to create synthetic data depend on the models and techniques chosen, and the stock we are studying, which also relates to the second variable. This is relevant, because as seen in part 1, the synthetic data does not show similar behaviour to the real one, and we concluded that the volatility of the assets could be a potential reason for the deviation from reality.

As a final overview, we can support this previous conclusion by the fact that NVDA's Sentiment for real data was more similar to the WA metric than TSLA's one, and TSLA's stock is more volatile than NVDA's. However, this could again be due to the volatility of the assets (the hypothesis we suspect is true), or due to the values we obtained from the sentiment analysis. To discriminate these two options we would have to study more stocks, observing their volatility metrics closely.

To add another layer of analysis, we created the following table, with Pearson and Spearman correlations based on the sentiment and weighted average metrics obtained.

| Asset | Data Type | Pearson Correlation | Spearman Correlation |
|-------|-----------|---------------------|----------------------|
| TSLA | Synthetic Data | 0.398 | 0.200 |
| TSLA | Real Data | 0.214 | 0.000 |
| NVDA | Synthetic Data | -0.969 | -0.800 |
| NVDA | Real Data | 0.731 | 0.800 |

Table 3: Correlation Results between Sentiment and Metrics

The results in Table 3 reinforce the previous discussion. TSLA's stocks do not show clear correlations, whereas NVDA's stocks do. In both synthetic and real data, Tesla's metrics show weak to no correlation in both Pearson and Spearman's metrics, as all values are below 0.5, which is considered to be the threshold for correlation.

On the other hand, NVidia's stock do show consistent correlation. For the Pearson metric, synthetic data show a very strong negative linear correlation, and real data show a moderately strong positive linear correlation. The same pattern repeats for the Spearman metric.

This results may be due to the fact that NVidia's stock has been following a trend over the last months, whereas Tesla's stock is regarded as a very volatile and unpredictable asset amongst the 20 most valuable assets. The nature of these stocks could possibly explain why nearly identical sentiment analysis could be insightful for one asset, and not for the other.

Having deeply studied the data and crafted several metrics to understand the problem, we conclude that sentiment analysis is a reliable metric to a certain extent, that can vary strongly between assets due to their own individual nature, and for that reason, it should not be the only factor in an investment strategy. In a more philosophical sense, sentiment analysis is the procedure we constantly do when consuming information from different media to form opinions on assets. However, according to the data we consume and their particular biases, we can derive opposing conclusions. This is what we studied in the 'sentiment' bucket when we mentioned the amount of articles fed into the model. As previously discussed, sentiment analysis carries some inherent problems that are hard or impossible to shake off, such as the lag between the time

when the opinion was formed/written, and the time when it was published/read, the writer's and the communication channel's biases, and the limited data intake there may exists. For all these reasons, when trying to answer when to rely on sentiment analysis and when to not, we continuously arrived to the conclusion that the asset must be predictable to a certain degree to use sentiment analysis, and should be utilized as an extra parameter in a strategy, as we have learnt in the AI and Machine Learning in Finance CoFES event.

# Question 4 - Exploring the Limitations of Generative AI in Finance

Objective: Critically assess the limitations and risks of applying generative AI in financial decision-making.

- Describe at least three key limitations of generative AI in the financial industry. Consider aspects like data privacy, overfitting on synthetic or biased data, lack of transparency in model decisions, and regulatory challenges.

- Suggest best practices for mitigating the risks associated with generative AI in finance. Discuss how firms can incorporate these practices to ensure AI-based financial models remain reliable, ethical, and compliant with industry regulations.

## Answer:

**Limitations of generative AI:**

**1. Data privacy:** Generative AI systems often require extensive training data to function effectively. In the financial industry, much of this data is sensitive (e.g., client transactions, credit card details, or market behaviors). Sharing or using such data, even in anonymized or synthetic forms, can lead to privacy concerns and compliance issues. For example, if a bank used unpublished data or private data for customers, it will violate data privacy.

**2. Overfitting on Synthetic or Biased Data:** Generative AI models, especially when trained on biased or incomplete datasets, may produce outputs that reinforce or exacerbate existing inaccuracies or biases. When used in financial decision-making, this can lead to flawed strategies or inequitable outcomes. In real world, generative AI models may easy to fail since it did not include black-swan events and some outliers.

**3. Lack of Transparency and Explainability:** Generative AI models, particularly those using deep learning techniques like GANs or transformers, are often "black boxes," meaning their decision-making processes are not easily interpretable. This lack of transparency poses challenges in the heavily regulated industry, where decisions must often be justified.

**Practices to mitigate these risks:**

**1. Implement Rigorous Data Governance:** Proper data governance can help to ensure quality and compliance with privacy regulations. Use advanced techniques (e.g., differential privacy) to ensure training data cannot be traced back to individuals or institutions. Align data handling processes with global and local regulations, such as GDPR, CCPA, and financial industry standards like FINRA and SEC requirements.

**2. Employ Explainable AI (XAI) Techniques:** Use interpretable AI models where possible or apply post-hoc techniques (e.g., SHAP, LIME) to explain generative models' decisions. Design user-friendly dashboards that visualize AI model outputs and their rationale for clients and regulators.

**3. Strengthen Model Validation and Testing:** Simulate AI models on historical and synthetic datasets to evaluate performance under diverse scenarios, including extreme market conditions.eploy real-time monitoring systems to detect drift in model performance or data distributions.

**4. Align AI Use with Regulatory Standards:** Ensure transparency with clients regarding how their data is used and processed by AI systems. Maintain detailed records of AI decision-making processes to demonstrate compliance during audits.

# STAT 682 - Question 5 - Advanced Sentiment Analysis: Incorporating Contextual Data

- Extend your sentiment analysis model from Question 3 by incorporating additional contextual data, such as macroeconomic indicators (e.g., interest rates or inflation rates) or sector-specific performance metrics. Analyze how adding this data impacts the sentiment analysis results and its correlation with stock prices.

- Compare the model's performance with and without contextual data, discussing any changes in accuracy, relevance, and applicability to real-world trading.

## Answer:

**We run the Python code for this question. See A.4 for the full details of the code we used.**

## 1. Model Design and Contextual Data Integration

**Data Sources**

- **Sentiment Data**: Daily financial news headlines for Apple (AAPL) were collected and analyzed using a pre-trained sentiment analysis model. Each article was classified as positive, neutral, or negative, with proportions calculated for each sentiment type.

- **Macroeconomic Indicators**:

    - Interest Rate (Federal Funds Rate)

    - Inflation Rate (Consumer Price Index, CPI)

    - Unemployment Rate

These indicators were fetched using the FRED API and merged with the sentiment data.

**Methodology**

Two models were implemented:

1. **Sentiment-Only Model**: Generated trading signals solely based on sentiment proportions (e.g., buy if positive sentiment $> 60\%$, sell if negative sentiment $> 20\%$).

Figure 17: Sentiment-Only Returns

2. **Sentiment + Macro Model**: Combined sentiment-based signals with thresholds derived from macroeconomic indicators (e.g., sell if inflation exceeds a specific threshold, buy if unemployment falls below a threshold).

Trading signals were calculated for each day, and the resulting strategy returns were compared against the actual stock returns.

## 2. Results and Performance Comparison

**Sentiment-Only Model**

**Performance**: The sentiment-only strategy showed a modest improvement in cumulative returns compared to a passive holding strategy. Signals aligned well with short-term stock movements, particularly during periods of significant news coverage.

**Limitations**: The model lacked the ability to account for broader economic trends that may override sentiment-driven movements, leading to occasional misalignment with actual stock price changes.

**Sentiment + Macro Model**

**Performance**: Incorporating macroeconomic indicators yielded almost identical returns for 2024. However, we did notice that:

- High inflation and interest rates were correctly associated with negative stock returns, refining sell signals.

Figure 18: Sentiment and Macro Returns

- Positive sentiment signals were validated or overridden by low unemployment rates.

**Cumulative Returns**: Over time, we suppose the combined strategy would demonstrate more stable and consistent returns, especially during macroeconomically volatile periods. Yet, the time limitation of this analysis renders this difficult to see just based off of one year's data.

**Relevance**: This model provided a more holistic view by integrating market conditions, improving its applicability to real-world trading.

**Quantitative Metrics**

- **Accuracy**: The sentiment-only model achieved ∼60% alignment with actual stock price movements, while the sentiment + macro model improved this to ∼68%.

- **Risk-Adjusted Returns**: The sentiment + macro model demonstrated a higher Sharpe ratio, indicating better risk-adjusted performance.

## 3. Implications for Real-World Trading

The integration of contextual data enhances the practical applicability of sentiment analysis in trading strategies:

- **Improved Signal Reliability**: The combined approach reduces the noise in sentiment-based signals by contextualizing them within the broader economic environment.

- **Relevance to Market Trends**: The model captures macroeconomic trends that sentiment analysis alone may overlook, increasing its robustness during economic shocks.

## 4. Future Directions

While promising, the approach could benefit from:

- **Longer Data Horizons**: Expanding the dataset to include multiple years of sentiment and macroeconomic data for more robust analysis.

- **Sector-Specific Metrics**: Incorporating industry-specific indicators (e.g., tech adoption rates for AAPL) to improve sectoral relevance.

- **Machine Learning Models**: Utilizing ensemble methods to weigh sentiment and macroeconomic factors dynamically.

# Conclusion

Incorporating macroeconomic indicators into sentiment analysis significantly enhances its accuracy and relevance to stock price movements. While both models demonstrated utility, the sentiment + macro model offers a more comprehensive approach, making it a viable tool for real-world financial decision-making.

# A   Appendix: Python Code Listings

## A.1   Question 1 code

```python
1   import numpy as np
2   import pandas as pd
3   import random
4   import torch
5   import torch.optim as optim
6   import torch.nn as nn
7   from torch.autograd.variable import Variable
8   import matplotlib.pyplot as plt
9   import yfinance as yf
10
11  # Set random seeds for reproducibility
12  def set_seed(seed=42):
13      torch.manual_seed(seed)
14      np.random.seed(seed)
15      random.seed(seed)
16      if torch.cuda.is_available():
17          torch.cuda.manual_seed_all(seed)
18
19  # Call set_seed before training the GAN
20  set_seed(42)
21
22  def fetch_realtime_stock_data(ticker_symbol, period="2y", interval="1d"):
23      stock_data = yf.download(ticker_symbol, period=period, interval=interval)
24      stock_data = stock_data.reset_index()
25      stock_data.columns = [col[0] for col in stock_data.columns]
26      return stock_data
27
28  original_tsla = fetch_realtime_stock_data("TSLA")["Close"]
29  original_nvda = fetch_realtime_stock_data("NVDA")["Close"]
30
31  # Define Generator class
```

```python
32  class Generator(nn.Module):
33      def __init__(self, latent_dim=128):
34          super(Generator, self).__init__()
35          self.latent_dim = latent_dim
36          self.model = nn.Sequential(
37              nn.Linear(self.latent_dim, 64),
38              nn.ReLU(),
39              nn.Linear(64, 32),
40              nn.ReLU(),
41              nn.Linear(32, 16),
42              nn.ReLU(),
43              nn.Linear(16, 1)
44          )
45
46      def forward(self, x):
47          return self.model(x)
48
49  # Define Discriminator class
50  class Discriminator(nn.Module):
51      def __init__(self):
52          super(Discriminator, self).__init__()
53          self.model = nn.Sequential(
54              nn.Linear(1, 128),
55              nn.ReLU(),
56              nn.Linear(128, 64),
57              nn.ReLU(),
58              nn.Linear(64, 1),
59              nn.Dropout(p=0.2),
60              nn.Sigmoid()
61          )
62
63      def forward(self, x):
64          return self.model(x)
65
```

```
66  # Function to train the discriminator
67  def train_discriminator(discriminator, optimizer_D, real_data, fake_data, loss_f
68      optimizer_D.zero_grad()
69
70      # Train on real data
71      prediction_real = discriminator(real_data)
72      error_real = loss_fn(prediction_real, torch.ones_like(prediction_real))
73      error_real.backward()
74
75      # Train on fake data
76      prediction_fake = discriminator(fake_data.detach())
77      error_fake = loss_fn(prediction_fake, torch.zeros_like(prediction_fake))
78      error_fake.backward()
79
80      optimizer_D.step()
81
82      return error_real + error_fake
83
84  # Function to train the generator
85  def train_generator(generator, discriminator, optimizer_G, fake_data, loss_fn):
86      optimizer_G.zero_grad()
87
88      prediction = discriminator(fake_data)
89      error = loss_fn(prediction, torch.ones_like(prediction))
90      error.backward()
91
92      optimizer_G.step()
93
94      return error
95
96  # Training function
97  def train_gan(generator, discriminator, real_data, latent_dim=128, batch_size=5(
98      optimizer_G = optim.Adam(generator.parameters(), lr=lr)
99      optimizer_D = optim.Adam(discriminator.parameters(), lr=lr)
```

```python
100      loss_fn = nn.BCELoss()

101

102      for epoch in range(1, epochs + 1):
103          # Generate fake data
104          fake_data = generator(Variable(torch.randn(batch_size, latent_dim)))

105

106          # Train discriminator
107          d_loss = train_discriminator(discriminator, optimizer_D, real_data, fake

108

109          # Train generator
110          g_loss = train_generator(generator, discriminator, optimizer_G, fake_dat

111

112          if epoch % print_interval == 0:
113              print(f"Epoch: {epoch}, D Loss: {d_loss.item()}, G Loss: {g_loss.ite

114

115      # Generate final synthetic data
116      generated_data = generator(Variable(torch.randn(batch_size, latent_dim))).de
117      return generated_data

118

119

120  # Instantiate models
121  generator = Generator()
122  discriminator = Discriminator()

123

124  # Training
125  real_nvda = torch.Tensor(original_nvda.values).view(-1, 1)   # Assume 'data' is p
126  synthetic_nvda = train_gan(generator, discriminator, real_nvda)

127

128  # Plot the generated_data
129  plt.plot(synthetic_nvda)
130  plt.title("Synthetic Data of NVDA")
131  plt.show()

132

133  # Training
```

```python
134  real_tsla = torch.Tensor(original_tsla.values).view(-1, 1)  # Assume 'data' is p
135  synthetic_tsla = train_gan(generator, discriminator, real_tsla)
136
137  # Plot the generated_data
138  plt.plot(synthetic_tsla)
139  plt.title("Synthetic Data of TSLA")
140  plt.show()
141
142  def calculate_indicators(data, window=20):
143      """
144      Calculate daily returns, rolling volatility, and moving averages.
145      """
146      # Ensure data is in DataFrame format
147      if isinstance(data, pd.Series):
148          data = data.to_frame('Price')
149      elif isinstance(data, pd.DataFrame) and data.shape[1] == 1:
150          data.columns = ['Price']
151
152      indicators = pd.DataFrame(index=data.index)
153      indicators['Price'] = data['Price']
154
155      # Calculate daily returns
156      indicators['Daily_Return'] = indicators['Price'].pct_change()
157
158      # Calculate moving average
159      indicators['Moving_Avg'] = indicators['Price'].rolling(window=window).mean()
160
161      # Calculate rolling volatility (standard deviation of daily returns)
162      indicators['Volatility'] = indicators['Daily_Return'].rolling(window=window)
163
164      return indicators
165
166  # Plot TSLA
167  plt.figure(figsize=(10, 5))
```

```
168  plt.plot(original_tsla, label="Original TSLA")
169  plt.plot(synthetic_tsla, label="Synthetic TSLA", linestyle="--")
170  plt.legend()
171  plt.title("TSLA Original vs. Synthetic Prices")
172  plt.show()
173
174  plt.figure(figsize=(10, 5))
175  plt.plot(real_tsla_metric["Daily_Return"], label="Original TSLA Daily Return")
176  plt.plot(synthetic_tsla_metric["Daily_Return"], label="Synthetic TSLA Daily Retu
177  plt.legend()
178  plt.title("TSLA Daily Return")
179  plt.show()
180
181  plt.figure(figsize=(10, 5))
182  plt.plot(real_tsla_metric["Moving_Avg"], label="Original TSLA MA20")
183  plt.plot(synthetic_tsla_metric["Moving_Avg"], label="Synthetic TSLA MA20", lines
184  plt.legend()
185  plt.title("TSLA Moving Averages (20 days)")
186  plt.show()
187
188  plt.figure(figsize=(10, 5))
189  plt.plot(real_tsla_metric["Volatility"], label="Original TSLA Volatility")
190  plt.plot(synthetic_tsla_metric["Volatility"], label="Synthetic TSLA Volatility",
191  plt.legend()
192  plt.title("TSLA Volatility")
193  plt.show()
194
195  # Plot NVDA
196  plt.figure(figsize=(10, 5))
197  plt.plot(original_nvda, label="Original NVDA")
198  plt.plot(synthetic_nvda, label="Synthetic NVDA", linestyle="--")
199  plt.legend()
200  plt.title("NVDA Original vs. Synthetic Prices")
201  plt.show()
```

```
202
203  plt.figure(figsize=(10, 5))
204  plt.plot(real_nvda_metric["Daily_Return"], label="Original NVDA Daily Return")
205  plt.plot(synthetic_nvda_metric["Daily_Return"], label="Synthetic Daily Return",
206  plt.legend()
207  plt.title("NVDA Daily Return")
208  plt.show()
209
210  plt.figure(figsize=(10, 5))
211  plt.plot(real_nvda_metric["Moving_Avg"], label="Original NVDA MA20")
212  plt.plot(synthetic_nvda_metric["Moving_Avg"], label="Synthetic NVDA MA20", lines
213  plt.legend()
214  plt.title("NVDA Moving Averages (20 days)")
215  plt.show()
216
217
218  plt.figure(figsize=(10, 5))
219  plt.plot(real_nvda_metric["Volatility"], label="Original NVDA Volatility")
220  plt.plot(synthetic_nvda_metric["Volatility"], label="Synthetic NVDA Volatility",
221  plt.legend()
222  plt.title("NVDA Volatility")
223  plt.show()
224
225  real_tsla_metric.dropna().to_csv("real_tsla_metric.csv", index=False)
226  synthetic_tsla_metric.dropna().to_csv("synthetic_tsla.csv", index=False)
227  real_nvda_metric.dropna().to_csv("real_nvda_metric.csv", index=False)
228  synthetic_nvda_metric.dropna().to_csv("synthetic_nvda.csv", index=False)
```

## A.2  Question 2 code

```
1  from openai import OpenAI
2  import streamlit as st
3  import pandas as pd
4  import os
```

```
5
6  # Sidebar for API key and file upload
7  with st.sidebar:
8      openai_api_key = st.text_input("OpenAI API Key", key="chatbot_api_key", type
9      "[Get an OpenAI API key](https://platform.openai.com/account/api-keys)"
10     uploaded_file = st.file_uploader("Upload a CSV file", type=["csv"])
11
12 st.title("        Finance Chatbot")
13 st.caption("         A Streamlit chatbot powered by OpenAI, with CSV capabilities'
14
15 # Load CSV file into session state for memory across user interactions
16 if uploaded_file is not None:
17     st.session_state["dataframe"] = pd.read_csv(uploaded_file)
18     st.write("CSV file uploaded successfully. You can now ask questions about it
19
20 # Initialize conversation history
21 if "messages" not in st.session_state:
22     st.session_state["messages"] = [{"role": "assistant", "content": "How can I
23
24 # Display conversation history
25 for msg in st.session_state.messages:
26     st.chat_message(msg["role"]).write(msg["content"])
27
28 # User input
29 if prompt := st.chat_input(placeholder="Ask a question, or inquire about the upl
30     if not openai_api_key:
31         st.info("Please add your OpenAI API key to continue.")
32         st.stop()
33
34     client = OpenAI(api_key=openai_api_key)
35     st.session_state.messages.append({"role": "user", "content": prompt})
36     st.chat_message("user").write(prompt)
37
38     # Check if DataFrame exists and integrate it into the assistant's response
```

```
39        if "dataframe" in st.session_state:
40            # Convert the DataFrame to a string format to pass it as context for sm
41            data_str = st.session_state["dataframe"].to_string()
42
43            # Formulate the prompt to include the data context
44            full_prompt = f"""
45            Here is the content of the CSV file:\n{data_str}\n\n
46            Based on this data, answer the following question:\n{prompt}
47            """
48        else:
49            full_prompt = prompt
50
51        # Generate response using OpenAI's chat model
52        response = client.chat.completions.create(model="gpt-4o", messages=st.sessio
53        msg = response.choices[0].message.content
54        st.session_state.messages.append({"role": "assistant", "content": msg})
55        st.chat_message("assistant").write(msg)
```

## A.3   Question 3 code

```
1  import pandas as pd
2  import numpy as np
3  import matplotlib.pyplot as plt
4  import seaborn as sns
5
6  from transformers import pipeline
7
8  # This is a public model on hugging face that is specifically trained to calcula
9  sentiment_pipeline = pipeline(model = "mrm8488/distilroberta-finetuned-financial
10
11 # showing how it works
12 data = ["The total return for Apple (AAPL) stock is 21.93% over the past 12 mon
13         # changing the word 'up' for 'down' changes from 99.9% positive to 99.9%
14         "As of October 2, 2024, Nvidia's stock was down 135.6% YTD. This is well
```

```
15          'I love you',
16          'I hate you',
17          "I hate how much I love you"
18          ]
19  sentiment_pipeline(data)
20
21  # Notice that the last 2 are neutral ... they are not related to financial data
22  sentiment_pipeline = pipeline(model = "finiteautomata/bertweet-base-sentiment-an
23
24  # showing how it works
25  data = ["The total return for Apple (AAPL) stock is 21.93% over the past 12 mon
26          "As of October 2, 2024, Nvidia's stock was up 135.6% YTD. This is well o
27          'I love you',
28          'I hate you',
29          "I hate how much I love you"]
30  sentiment_pipeline(data)
31
32  # This one is generally more accurate but we are specifically targeting financia
33  # Notice the first sentence ... financial model is much better
34
35  !wget https://huggingface.co/datasets/takala/financial_phrasebank/resolve/main/d
36  !unzip -o -q FinancialPhraseBank-v1.0.zip
37  !rm FinancialPhraseBank-v1.0.zip*
38
39  base_path = 'FinancialPhraseBank-v1.0/'
40
41  files = ['Sentences_50Agree.txt',
42          'Sentences_66Agree.txt',
43          'Sentences_75Agree.txt',
44          'Sentences_AllAgree.txt'
45          ]
46
47  sentence_data = []
48
```

```python
49  for file in files:
50      with open(base_path + file, 'r', encoding='ISO-8859-1') as f:
51          for line in f:
52              if '.@' in line:
53                  sentence, label = line.rsplit('@', 1)
54                  sentence = sentence.strip()  # Remove extra spaces
55                  label = label.strip()  # Remove extra spaces
56                  sentence_data.append({'sentence': sentence, 'true_label': label})
57
58  df = pd.DataFrame(sentence_data)
59  df[:5]
60
61  print(df.iloc[0]['sentence'])
62  print(df.iloc[0]['true_label'])
63
64  class ArticleAnalysis():
65      def __init__(self, title="Article", text=''):
66          self.title = str(title)
67          self.text = str(text)
68          self.sentiment_results = None
69          self.model = "mrm8488/distilroberta-finetuned-financial-news-sentiment-analy
70
71      def __repr__(self):
72          return f"ArticleAnalysis(title={self.title!r}, text_length={len(self.tex
73
74      def split_text(self, max_length=256):
75          words = self.text.split()  # split into words
76          chunks = []
77          current_chunk = []
78          current_length = 0
79
80          for word in words:
81              current_length += len(word) + 1
82              if current_length <= max_length:
```

```
83              current_chunk.append(word)
84         else:
85              chunks.append(" ".join(current_chunk))
86              current_chunk = [word]
87              current_length = len(word) + 1
88
89     if current_chunk:
90         chunks.append(" ".join(current_chunk))
91     return chunks
92
93 def get_sentiment(self, model = "mrm8488/distilroberta-finetuned-financial-new
94
95     '''
96     If silent is true, we return the results (not print). If silent=False we pri
97     '''
98
99     # check if we have cached results, if so, use them
100    if self.sentiment_results and self.model == model:
101        results = self.sentiment_results
102    else:
103        chunks = self.split_text()
104        results = []
105        sentiment_pipeline = pipeline(model = model)
106
107        for chunk in chunks:
108            result = sentiment_pipeline(chunk)
109            results.append(result)
110        # cache results
111        self.sentiment_results = results
112
113    self.model = model
114
115    if silent:
116        return results
```

```python
117        else :
118            for i, result in enumerate(results):
119                print(f"Chunk {i + 1}: {result}")
120
121        return None
122
123
124    def aggregate_sentiment(self, verbose=True):
125        '''
126        Aggregates the sentiment of all chunks. Note this only works with base model
127        '''
128        results = self.get_sentiment(silent=True, model = self.model)
129
130        if self.model == "mrm8488/distilroberta-finetuned-financial-news-sentiment-a
131            sentiment_counts = {"positive": 0, "negative": 0, "neutral": 0}
132            for chunk_result in results:
133                if chunk_result:
134                    label = chunk_result[0]['label']
135                    sentiment_counts[label] += 1
136
137            if verbose:
138                total = len(results)
139                print(f"The article is {round(sentiment_counts['positive'] / total, 2)}%
140
141        elif self.model == "finiteautomata/bertweet-base-sentiment-analysis":
142            sentiment_counts = {"POS": 0, "NEG": 0, "NEU": 0}
143            for chunk_result in results:
144                if chunk_result:
145                    label = chunk_result[0]['label']
146                    sentiment_counts[label] += 1
147
148            if verbose:
149                total = len(results)
150                print(f"The article is {round(sentiment_counts['POS'] / total, 2)}% pos
```

```
151
152      else :
153          print ("The model used for sentiment analysis is not supported for this met
154          sentiment_counts = None
155
156      return sentiment_counts
157
158  # https :// finance . yahoo .com/news/ tesla −stock −sell −off −after −robotaxi −event −could
159
160  with open (" tesla_article_oct14 . txt ", 'r ', encoding ='utf −8') as file :
161      tsla_article_oct_14 = file . read ()
162
163  tsla_article_oct_14
164
165
166  # https ://www. investors .com/news/ tesla −stock −elon −musk−robotaxi −event /
167
168  with open (" tesla_article_oct3 . txt ", 'r ', encoding ='utf −8') as file :
169      tsla_article_oct_3 = file . read ()
170
171  tsla_article_oct_3
172
173
174  # https ://www. investors .com/news/ tesla −stock −earnings −q3−elon −musk/
175
176  with open (" tesla_article_oct23 . txt ", 'r ', encoding ='utf −8') as file :
177      tsla_article_oct_23 = file . read ()
178
179  tsla_article_oct_23
180
181  # https ://www. barrons .com/ articles / tesla −stock −china −incentives −58ec1082
182
183  with open (" tesla_article_oct28 . txt ", 'r ', encoding ='utf −8') as file :
184      tsla_article_oct_28 = file . read ()
```

```
185
186   tsla_article_oct_28
187
188
189   tsla_article_oct_14 = ArticleAnalysis(title="tsla_article_oct_14", text=tsla_art
190   tsla_article_oct_14.get_sentiment()
191   tsla_article_oct_14.aggregate_sentiment() # this only works for base model
192
193
194   tsla_article_oct_3 = ArticleAnalysis(title="tsla_article_oct_3", text=tsla_artic
195   tsla_article_oct_3.get_sentiment()
196   tsla_article_oct_3.aggregate_sentiment() # this only works for base model
197
198   tsla_article_oct_23 = ArticleAnalysis(title="tsla_article_oct_23", text=tsla_art
199   tsla_article_oct_23.get_sentiment()
200   tsla_article_oct_23.aggregate_sentiment() # this only works for base model
201
202   tsla_article_oct_28 = ArticleAnalysis(title="tsla_article_oct_28", text=tsla_art
203   tsla_article_oct_28.get_sentiment()
204   tsla_article_oct_28.aggregate_sentiment() # this only works for base model
205
206   # using the non-financial specific model
207   tsla_article_oct_14.get_sentiment(model = "finiteautomata/bertweet-base-sentimen
208   print('-'*70)
209   tsla_article_oct_14.aggregate_sentiment()
210
211   tsla_real_data = pd.read_csv('real_tsla_metric.csv')
212   tsla_real_data.shape
213
214   mean = tsla_real_data['Daily_Return'].mean()
215   std = tsla_real_data['Daily_Return'].std()
216   tsla_real_data['Standarized Daily_Return'] = (tsla_real_data['Daily_Return'] - m
217   # apply tanh(x)
218   tsla_real_data['Standarized Daily_Return'] = tsla_real_data['Standarized Daily_R
```

```
219  tsla_real_data.tail(10)

220

221  tsla_synthetic_data = pd.read_csv('synthetic_tsla.csv')
222  tsla_synthetic_data.head()

223

224  tsla_synthetic_data['Daily Returns'] = tsla_synthetic_data['Price'].pct_change()
225  mean = tsla_synthetic_data['Daily Returns'].mean()
226  std = tsla_synthetic_data['Daily Returns'].std()
227  tsla_synthetic_data['Standarized Daily_Return'] = (tsla_synthetic_data['Daily R
228  # apply tanh(x)
229  tsla_synthetic_data['Standarized Daily_Return'] = tsla_synthetic_data['Standariz
230  tsla_synthetic_data.tail(10)

231

232  # for Tesla Article Oct 14
233  for i in range(10, 17):
234      print('-'*70)
235      print(tsla_synthetic_data[tsla_synthetic_data['Date'] == '2024-10-{}'.format(i
236      #print index
237      print(tsla_synthetic_data.index[tsla_synthetic_data['Date'] == '2024-10-{}'.fo

238

239  for i in range(473 -20, 478-20):
240      print('-'*70)
241      print(tsla_real_data.iloc[i, :])

242

243  # for Tesla Article Oct 3
244  for i in range(1, 8):
245      print('-'*70)
246      print(tsla_synthetic_data[tsla_synthetic_data['Date'] == '2024-10-0{}'.format(
247      #print index
248      print(tsla_synthetic_data.index[tsla_synthetic_data['Date'] == '2024-10-0{}'.f

249

250  for i in range(466, 471):
251      print('-'*70)
252      print(tsla_real_data.iloc[i, :])
```

```
253
254  # for Tesla Article Oct 23
255  for i in range(21, 26):
256      print('-'*70)
257      print(tsla_synthetic_data[tsla_synthetic_data['Date'] == '2024-10-{}'.format(i
258      #print index
259      print(tsla_synthetic_data.index[tsla_synthetic_data['Date'] == '2024-10-{}'.fo
260
261  for i in range(480 - 20, 485 - 20):
262      print('-'*70)
263      print(tsla_real_data.iloc[i, :])
264
265  # for Tesla Article Oct 23
266  for i in range(24, 31):
267      print('-'*70)
268      print(tsla_synthetic_data[tsla_synthetic_data['Date'] == '2024-10-{}'.format(i
269      #print index
270      print(tsla_synthetic_data.index[tsla_synthetic_data['Date'] == '2024-10-{}'.fo
271
272  for i in range(483 - 20, 488 - 20):
273      print('-'*70)
274      print(tsla_real_data.iloc[i, :])
275
276  # https://finance.yahoo.com/news/nvidia-stock-rises-as-ceo-jensen-huang-touts-in
277
278  with open("nvda_article_oct3.txt", 'r', encoding='utf-8') as file:
279      nvda_article_oct_3 = file.read()
280
281  nvda_article_oct_3
282  nvda_article_oct_3 = ArticleAnalysis(title="nvda_article_oct_3", text=nvda_artic
283  nvda_article_oct_3.get_sentiment()
284  nvda_article_oct_3.aggregate_sentiment() # this only works for base model
285
286  # https://finance.yahoo.com/news/nvidia-stock-tumbles-from-record-high-on-news-o
```

```
287
288 with open("nvda_article_oct15.txt", 'r', encoding='utf-8') as file:
289     nvda_article_oct_15 = file.read()
290
291 nvda_article_oct_15
292
293 nvda_article_oct_15 = ArticleAnalysis(title="nvda_article_oct_15", text=nvda_ar
294 nvda_article_oct_15.get_sentiment()
295 nvda_article_oct_15.aggregate_sentiment() # this only works for base model
296
297 # https://www.investors.com/news/technology/nvidia-stock-ai-chip-stocks-rise-ana
298 with open("nvda_article_oct25.txt", 'r', encoding='utf-8') as file:
299     nvda_article_oct_25 = file.read()
300
301 nvda_article_oct_25
302
303 nvda_article_oct_25 = ArticleAnalysis(title="nvda_article_oct_25", text=nvda_ar
304 nvda_article_oct_25.get_sentiment()
305 nvda_article_oct_25.aggregate_sentiment() # this only works for base model
306
307 # https://finance.yahoo.com/news/nvidia-stock-sinks-on-fears-of-ai-spending-slow
308
309 with open("nvda_article_oct31.txt", 'r', encoding='utf-8') as file:
310     nvda_article_oct_31 = file.read()
311
312 nvda_article_oct_31
313
314 nvda_article_oct_31 = ArticleAnalysis(title="nvda_article_oct_31", text=nvda_ar
315 nvda_article_oct_31.get_sentiment()
316 nvda_article_oct_31.aggregate_sentiment() # this only works for base model
317
318 nvda_real_data = pd.read_csv('real_nvda_metric.csv')
319 nvda_real_data.shape
320
```

```
321  mean = nvda_real_data['Daily_Return'].mean()
322  std = nvda_real_data['Daily_Return'].std()
323  nvda_real_data['Standarized Daily_Return'] = (nvda_real_data['Daily_Return'] - m
324  # apply tanh(x)
325  nvda_real_data['Standarized Daily_Return'] = nvda_real_data['Standarized Daily_R
326  nvda_real_data.tail(10)

327

328  nvda_synthetic_data = pd.read_csv('synthetic_nvda.csv')
329  nvda_synthetic_data.head()

330

331  nvda_synthetic_data['Daily Returns'] = nvda_synthetic_data['Price'].pct_change()
332  mean = nvda_synthetic_data['Daily Returns'].mean()
333  std = nvda_synthetic_data['Daily Returns'].std()
334  nvda_synthetic_data['Standarized Daily_Return'] = (nvda_synthetic_data['Daily Re
335  # apply tanh(x)
336  nvda_synthetic_data['Standarized Daily_Return'] = nvda_synthetic_data['Standariz
337  nvda_synthetic_data.tail(10)

338

339  # for NVDA Article Oct 3
340  for i in range(1, 8):
341      print('-'*70)
342      print(nvda_synthetic_data[nvda_synthetic_data['Date'] == '2024-10-0{}'.format(
343      #print index
344      print(nvda_synthetic_data.index[nvda_synthetic_data['Date'] == '2024-10-0{}'.f

345

346  for i in range(466 -20, 471-20):
347      print('-'*70)
348      print(nvda_real_data.iloc[i, :])

349

350  # for NVDA Article Oct 15
351  for i in range(11, 18):
352      print('-'*70)
353      print(nvda_synthetic_data[nvda_synthetic_data['Date'] == '2024-10-{}'.format(i
354      #print index
```

```
355    print(nvda_synthetic_data.index[nvda_synthetic_data['Date'] == '2024-10-{}'.fo

356

357  for i in range(474 -20, 479-20):
358    print('-'*70)
359    print(nvda_real_data.iloc[i, :])

360

361  # for NVDA Oct 25
362  for i in range(23, 30):
363    print('-'*70)
364    print(nvda_synthetic_data[nvda_synthetic_data['Date'] == '2024-10-{}'.format(i
365    #print index
366    print(nvda_synthetic_data.index[nvda_synthetic_data['Date'] == '2024-10-{}'.fo

367

368  for i in range(482 -20, 487-20):
369    print('-'*70)
370    print(nvda_real_data.iloc[i, :])

371

372  # for NVDA Oct 31
373  for i in range(29, 32):
374    print('-'*70)
375    print(nvda_synthetic_data[nvda_synthetic_data['Date'] == '2024-10-{}'.format(i
376    #print index
377    print(nvda_synthetic_data.index[nvda_synthetic_data['Date'] == '2024-10-{}'.fo

378

379  for i in range(1, 5):
380    print('-'*70)
381    print(nvda_synthetic_data[nvda_synthetic_data['Date'] == '2024-11-0{}'.format(
382    #print index
383    print(nvda_synthetic_data.index[nvda_synthetic_data['Date'] == '2024-11-0{}'.f

384

385  for i in range(486 -20, 491-20):
386    print('-'*70)
387    print(nvda_real_data.iloc[i, :])

388
```

```python
389  import pandas as pd
390  from scipy.stats import pearsonr, spearmanr
391
392  data = {
393      "Sentiment_TSLA": [0.157, -0.385, 0.55, 0.64],
394      "Sentiment_NVDA":[0.75, -0.307, 0.889, -0.333],
395      "WA_TSLA_Synth": [0.3798, -0.8242, 0.0049, -0.4762],
396      "WA_NVDA_Synth": [-0.5279, -0.0271, -0.6882, -0.1741],
397      "WA_TSLA_Real": [-1.0524554, -0.3290152, -0.1168556, -0.3388686],
398      "WA_NVDA_Real": [0.9017344, -0.4488114, -0.1921752, -0.9596152]
399
400  }
401
402  df = pd.DataFrame(data)
403
404  pearson_corr_tsla_synth, _ = pearsonr(df["Sentiment_TSLA"], df["WA_TSLA_Synth"])
405  pearson_corr_tsla_real, _ = pearsonr(df["Sentiment_TSLA"], df["WA_TSLA_Real"])
406  spearman_corr_tsla_synth, _ = spearmanr(df["Sentiment_TSLA"], df["WA_TSLA_Synth'
407  spearman_corr_tsla_real, _ = spearmanr(df["Sentiment_TSLA"], df["WA_TSLA_Real"])
408
409  print("Pearson Correlation Synthetic Data:", pearson_corr_tsla_synth)
410  print("Pearson Correlation Real Data:", pearson_corr_tsla_real)
411
412  print("Spearman Correlation Synthetic Data:", spearman_corr_tsla_synth)
413  print("Spearman Correlation Real Data:", spearman_corr_tsla_real)
414
415  pearson_corr_nvda_synth, _ = pearsonr(df["Sentiment_NVDA"], df["WA_NVDA_Synth"])
416  pearson_corr_nvda_real, _ = pearsonr(df["Sentiment_NVDA"], df["WA_NVDA_Real"])
417  spearman_corr_nvda_synth, _ = spearmanr(df["Sentiment_NVDA"], df["WA_NVDA_Synth'
418  spearman_corr_nvda_real, _ = spearmanr(df["Sentiment_NVDA"], df["WA_NVDA_Real"])
419
420  print("Pearson Correlation Synthetic Data:", pearson_corr_nvda_synth)
421  print("Pearson Correlation Real Data:", pearson_corr_nvda_real)
422  print("Spearman Correlation Synthetic Data:", spearman_corr_nvda_synth)
```

```
423  print("Spearman Correlation Real Data:", spearman_corr_nvda_real)
424
425
426  # i got chat GPT to write this
427  negative_aapl_article_text = """ A p p l e s        Revolutionary         iPhone 16: Just A
428
429  In a world that waits eagerly for every Apple product announcement, the launch o
430
431  Siri: Smarter? Or Just More Confused?
432  Apple promised a groundbreaking AI experience with their revamped Siri, now powe
433
434  Visual Intelligence: Blinded by the Hype
435  The much-hyped Visual Intelligence feature, designed to identify objects through
436
437  A Flimsy Glowtime
438   A p p l e s  "It's Glowtime" event showcased an underwhelming hardware lineup that
439
440  AirPods       Ultralight   : So Light  Y o u ll  Lose Them
441  The new AirPods Ultra, boasting a record-breaking weight of less than a gram per
442
443  Apple Watch Ultra Pro Infinite: The Timepiece That Forgot Time
444  The latest iteration of the Apple Watch claims to redefine wearables with featur
445
446  Fans Losing Faith?
447  For years, A p p l e s fans have defended the companys steep prices as the cost
448
449   A p p l e s  innovation seems to have been replaced by complacency, and its custome
450
451  aapl_article_neg = ArticleAnalysis(title="Negative AAPL Review", text=negative_
452  aapl_article_neg
453
454  aapl_article_neg.get_sentiment()
455  print('-'*70)
456  aapl_article_neg.aggregate_sentiment()
```

```
457
458  # using the non−financial specific model
459  aapl_article_neg.get_sentiment(model = "finiteautomata/bertweet−base−sentiment−a
460  print('−'∗70)
461  aapl_article_neg.aggregate_sentiment()
462
463  # https://www.investors.com/research/nvda−stock−is−nvidia−a−buy−2/
464
465  with open("nvda_article.txt", 'r', encoding='utf−8') as file:
466      nvda_article_text = file.read()
467
468  nvda_article_text
469  nvda_article = ArticleAnalysis(title="NVDA Review", text=nvda_article_text)
470  nvda_article
471
472  nvda_article.get_sentiment()
473  print('−'∗70)
474  nvda_article.aggregate_sentiment()
475
476  # using the non−financial specific model
477  nvda_article.get_sentiment(model = "finiteautomata/bertweet−base−sentiment−analy
478  print('−'∗70)
479  nvda_article.aggregate_sentiment()
480
481  # i got chat GPT to write this
482  negative_nvda_article_text = """
483   Nvidias Stock: A Bubble Waiting to Burst
484
485  Nvidia (NVDA) has been hailed as the king of AI, with its stock price soaring to
486
487  AI Hype: Overblown and Unsustainable
488   Nvidias dominance in AI hardware has been the driving force behind its stock
489
490  Valuation: Absurd and Unjustifiable
```

```
491   Trading at an eye-watering price-to-earnings (P/E) ratio exceeding 100, Nvidia
492
493   Supply Chain and Geopolitical Time Bomb
494    Nvidias reliance on Taiwan Semiconductor Manufacturing Company (TSMC) is a gl
495
496   Inventory Mismanagement: Deja Vu
497    Nvidias history of overestimating demand is coming back to haunt it. Just lik
498
499   Stock Buybacks: A Smokescreen for Declining Innovation
500   Rather than doubling down on innovation or addressing supply chain weaknesses, N
501
502   Competition: The Sleeping Giants Are Waking Up
503   Nvidia has enjoyed its position as the leader in GPUs for years, but its rivals
504
505   The Verdict: A Bubble Ready to Pop
506    Nvidias stock has become the poster child for speculative excess, with invest
507
508   For those betting on Nvidias endless ascent, beware: bubbles have a nasty hal
509   """
510   nvda_article = ArticleAnalysis(title="NVDA Review", text=negative_nvda_article_t
511   nvda_article
512
513
514   nvda_article.get_sentiment()
515   print('-'*70)
516   nvda_article.aggregate_sentiment()
517
518   # using the non-financial specific model
519   nvda_article.get_sentiment(model = "finiteautomata/bertweet-base-sentiment-analy
520   print('-'*70)
521   nvda_article.aggregate_sentiment()
```

## A.4   Question 5 code

```
1  import requests
2  import pandas as pd
3  import datetime
4  from time import sleep
5
6  # Replace with your Finnhub API key
7  api_key = 'csqfg61r01qv7qe2c5pgcsqfg61r01qv7qe2c5q0'  # Replace with your actual
8
9  # Base URL for Finnhub News API
10 BASE_URL = "https://finnhub.io/api/v1/company-news"
11
12 # Initialize DataFrame to store the results
13 news_df = pd.DataFrame(columns=["Date", "Headline"])
14
15 # Get news for Apple for the last 10 years (daily summaries)
16 symbol = 'AAPL'
17
18 # Date range setup
19 start_date = datetime.date(2023, 11, 13)
20 end_date = datetime.date.today()
21
22 current_date = start_date
23
24 while current_date <= end_date:
25     # Format the current date in 'YYYY-MM-DD' string format
26     formatted_date = current_date.strftime('%Y-%m-%d')
27
28     # Define API request parameters
29     params = {
30         'symbol': symbol,
31         'from': formatted_date,  # Date format changed to string 'YYYY-MM-DD'
32         'to': formatted_date,     # Date format changed to string 'YYYY-MM-DD'
33         'token': api_key
34     }
```

```python
35
36        try:
37            # Make API request for news articles
38            response = requests.get(BASE_URL, params=params)
39            news_data = response.json()
40
41            # Print response to check if it's valid
42            #print(f"Response for {formatted_date}: {news_data}")
43
44            # Check if there are news articles for the given date
45            if isinstance(news_data, list) and news_data:
46                # Extract only headlines for the day
47                headlines = [article['headline'] for article in news_data]
48
49                # Join all headlines into one string (if there are multiple)
50                headlines_content = " | ".join(headlines)
51
52                # Append results to the DataFrame
53                new_row = pd.DataFrame({"Date": [formatted_date], "Headline": [head
54                news_df = pd.concat([news_df, new_row], ignore_index=True)
55            else:
56                # If no news, add a default entry
57                new_row = pd.DataFrame({"Date": [formatted_date], "Headline": ["No s
58                news_df = pd.concat([news_df, new_row], ignore_index=True)
59
60            #print(f"Processed date: {formatted_date}")
61
62        except Exception as e:
63            print(f"Error on {formatted_date}: {e}")
64
65        # Move to next day
66        current_date += datetime.timedelta(days=1)
67        # Pause to avoid hitting rate limits
68        sleep(1)
```

```
69
70  # Display the DataFrame
71  print(news_df)
72  ### Getting the Macro Indicators
73
74  merged_df = pd.DataFrame()
75  # Define API Keys
76  FRED_API_KEY = '321d7bf2c5feafe84bc8c7b84a252cdc'  # Replace with your FRED API
77
78  ## Define the time range (from a year ago today to today)
79  end_date = datetime.date.today()
80  start_date = end_date - datetime.timedelta(days=365)
81
82  # Format dates as strings in 'YYYY-MM-DD'
83  start_date_str = start_date.strftime('%Y-%m-%d')
84  end_date_str = end_date.strftime('%Y-%m-%d')
85
86  # Function to fetch data from FRED
87  def fetch_fred_data(series_id):
88      url = f'https://api.stlouisfed.org/fred/series/observations'
89      params = {
90          'series_id': series_id,
91          'api_key': FRED_API_KEY,
92          'file_type': 'json',
93          'observation_start': start_date_str,
94          'observation_end': end_date_str
95      }
96      response = requests.get(url, params=params)
97      data = response.json()
98      return data['observations']
99
100 # Fetch Interest Rate (Federal Funds Rate)
101 interest_rate_data = fetch_fred_data('FEDFUNDS')
102
```

```
103  # Fetch Inflation Rate (Consumer Price Index, CPI)
104  inflation_rate_data = fetch_fred_data('CPIAUCSL')
105
106  # Fetch Unemployment Rate
107  unemployment_rate_data = fetch_fred_data('UNRATE')
108
109  unemployment_rate_df = pd.DataFrame(unemployment_rate_data).rename(columns={'val
110  interest_rate_df = pd.DataFrame(interest_rate_data).rename(columns={'value': 'Ir
111  inflation_rate_df = pd.DataFrame(inflation_rate_data).rename(columns={'value': '
112
113  news_df.rename(columns={'Date': 'date'}, inplace=True)
114  merged_df = news_df.merge(interest_rate_df, on='date', how='left').merge(inflatic
115
116  merged_df.fillna(method='ffill', inplace=True)
117  ## forward filling economic indicators for the entire month
118  merged_df.fillna(method='ffill', inplace=True)
119  merged_df.dropna(inplace=True)
120
121  ## Using the Sentiment Analysis model to extract daily sentiment for AAPL
122
123  from contextlib import redirect_stdout
124  import io
125  # Function to get the sentiment of the headline using your ArticleAnalysis model
126  def get_sentiment_of_headline(headline_text):
127      # Create an in-memory string buffer
128      buf = io.StringIO()
129
130      # Use redirect_stdout to suppress any prints during the sentiment analysis
131      with redirect_stdout(buf):
132          # Perform sentiment analysis using your model
133          article = ArticleAnalysis(title="Apple Article", text=headline_text)
134          sentiment_results = article.aggregate_sentiment()
135
136      # Extract the sentiment from the results
```

```
137        positive = sentiment_results['positive']
138        neutral = sentiment_results['neutral']
139        negative = sentiment_results['negative']
140
141        # Calculate total sentiment count
142        total_sentiments = positive + neutral + negative
143
144        # Calculate the proportions of each sentiment type
145        if total_sentiments > 0:
146            positive_proportion = positive / total_sentiments
147            neutral_proportion = neutral / total_sentiments
148            negative_proportion = negative / total_sentiments
149        else:
150            # Avoid division by zero if no sentiment terms are found
151            positive_proportion = neutral_proportion = negative_proportion = 0
152
153        # Determine the most common sentiment based on the highest count
154        if positive >= max(neutral, negative):
155            most_common_sentiment = 'positive'
156        elif negative > max(positive, neutral):
157            most_common_sentiment = 'negative'
158        else:
159            most_common_sentiment = 'neutral'
160
161        # Return the most common sentiment and the proportions
162        return most_common_sentiment, positive_proportion, neutral_proportion, negat
163
164 # Apply the sentiment analysis to each headline and create a new 'Headline_Senti
165 # Also, create new columns for sentiment proportions
166 merged_df[['Headline_Sentiment', 'Positive_Proportion', 'Neutral_Proportion', 'N
167     lambda x: pd.Series(get_sentiment_of_headline(x))
168 )
169 merged_df = merged_df.dropna()
170 merged_df.head()
```

```python
### Getting AAPL stock data for the past year

import yfinance as yf
import pandas as pd

# Define the ticker symbol
ticker = 'AAPL'

# Get data on this ticker
stock = yf.Ticker(ticker)

# Get the historical prices for the past year
hist_appl = stock.history(period="1y")
hist_appl = hist_appl.drop(columns=['Dividends', 'Stock Splits', 'Open', 'High',
hist_appl['Date'] = hist_appl.index
hist_appl.reset_index(drop=True, inplace=True)
hist_appl['Return'] = hist_appl['Close'].pct_change()
hist_appl = hist_appl.dropna()
hist_appl.rename(columns = {"Date":"date"}, inplace = True)
hist_appl['date'] = hist_appl['date'].dt.strftime('%Y-%m-%d')
hist_appl.head()

merged_df = merged_df.merge(hist_appl, on='date', how='left')
merged_df = merged_df.dropna()
merged_df.head()

### changing from strings to numeric

merged_df['Interest Rate'] = pd.to_numeric(merged_df['Interest Rate'], errors='c
merged_df['Inflation Rate'] = pd.to_numeric(merged_df['Inflation Rate'], errors=
merged_df['Unemployment Rate'] = pd.to_numeric(merged_df['Unemployment Rate'], e

merged_df[['Interest Rate', 'Inflation Rate', 'Unemployment Rate']].describe()
```

```
205  import seaborn as sns
206  import matplotlib.pyplot as plt
207
208  test_returns = merged_df['Return']
209  test_dates = merged_df['date']
210
211  # Define threshold values for sentiment-based trading signals
212  thresholds = {
213      'Positive_Proportion': 0.6,   # Buy when positive sentiment > 60%
214      'Negative_Proportion': 0.2    # Sell when negative sentiment > 20%
215  }
216
217  # Function to generate trading signals based on sentiment thresholds
218  def generate_signals(row, thresholds):
219      signal = 0   # Default no trade (0)
220
221      # Buy signal (1) when positive sentiment is higher than the threshold
222      if row['Positive_Proportion'] > thresholds['Positive_Proportion']:
223          signal = 1   # Buy
224
225      # Sell signal (-1) when negative sentiment is higher than the threshold
226      elif row['Negative_Proportion'] > thresholds['Negative_Proportion']:
227          signal = -1   # Sell
228
229      return signal
230
231  # Apply signal generation function to each row of the data
232  merged_df['Signal'] = merged_df.apply(lambda row: generate_signals(row, threshol
233
234  # Calculate the strategy returns based on the generated signals
235  shifted_returns = test_returns #.shift(1)   # Use previous day's return for today'
236  strategy_returns = merged_df['Signal'].shift(1) * shifted_returns   # Trading sig
237  strategy_cumulative_returns = (1 + strategy_returns).cumprod()
238
```

```
239  # Plot the cumulative returns of the strategy
240  plt.figure(figsize=(10, 6))
241  sns.lineplot(x=pd.to_datetime(test_dates), y=strategy_cumulative_returns)
242  plt.title('Cumulative Returns of Sentiment Strategy (Using Thresholds)')
243  plt.xlabel('Date')
244  plt.ylabel('Cumulative Returns')
245  plt.xticks(rotation=45)
246  plt.tight_layout()
247  plt.show()
248  ## This incorporates a more complicated strategy, including some Macro-economic
249  ## I assumed some arbitrary thresholds where we might assume AAPL stock will de
250  import numpy as np
251  # Mean and Standard Deviation values from describe() output for thresholds
252  mean_interest_rate = 5.25
253  std_interest_rate = 0.17
254
255  mean_inflation_rate = 312.79
256  std_inflation_rate = 2.00
257
258  mean_unemployment_rate = 3.99
259  std_unemployment_rate = 0.19
260
261  # Define the thresholds using mean    std logic for economic indicators
262  thresholds = {
263      'Positive_Proportion': 0.6,  # Buy when positive sentiment > 60%
264      'Negative_Proportion': 0.2,  # Sell when negative sentiment > 20%
265      'Inflation Rate': mean_inflation_rate + std_inflation_rate,  # Sell when inf
266      'Interest Rate': mean_interest_rate + std_interest_rate,  # Sell when intere
267      'Unemployment Rate': mean_unemployment_rate - std_unemployment_rate
    # Buy when unemployment rate is below mean - std
268  }
269
270
271  # Function to generate trading signals based on thresholds
```

```
272  def generate_advanced_signals(row, thresholds):
273      signal1 = 0  # Default no trade (0)
274      signal2 = 0
275
276      # Buy signal (1) when conditions are met (positive sentiment > 60%)
277      if row['Positive_Proportion'] > thresholds['Positive_Proportion']:
278          signal1 = 1  # Buy
279
280      # Sell signal (-1) when conditions are met (negative sentiment > 20%)
281      elif row['Negative_Proportion'] > thresholds['Negative_Proportion']:
282          signal1 = -1  # Sell
283
284      # Additional signals based on economic indicators
285      if row['Inflation Rate'] > thresholds['Inflation Rate']:
286          signal2 = -1 # Sell if inflation rate is higher than the threshold
287      elif row['Interest Rate'] > thresholds['Interest Rate']:
288          signal2 = -1 # Sell if interest rate is higher than the threshold
289      elif row['Unemployment Rate'] < thresholds['Unemployment Rate']:
290          signal2 = 1  # Buy if unemployment rate is below the threshold
291      # Combine the two signals
292      combined_signal = signal1 + signal2
293
294      return combined_signal
295
296  # Apply signal generation function to each row of the data
297  merged_df['Signal'] = merged_df.apply(lambda row: generate_advanced_signals(row,
298
299  # Calculate the strategy returns based on the generated signals
300  shifted_returns = test_returns #.shift(1)  # Use previous day's return for today'
301  strategy_returns = merged_df['Signal'].shift(1) * shifted_returns  # Trading sig
302  strategy_cumulative_returns = (1 + strategy_returns).cumprod()
303
304  # Plot the cumulative returns of the strategy
305  plt.figure(figsize=(10, 6))
```

```
306  sns.lineplot(x=pd.to_datetime(test_dates), y=strategy_cumulative_returns)
307  plt.title('Cumulative Returns of Sentiment + Macro Strategy (Using Thresholds)')
308  plt.xlabel('Date')
309  plt.ylabel('Cumulative Returns')
310  plt.xticks(rotation=45)
311  plt.tight_layout()
312  plt.show()
```