

10 Academy: Artificial Intelligence Mastery

Predicting Price Moves with News Sentiment: Week 1 Challenge

Final Report

by: Yonas Zelalem

June 3, 2025

Introduction

Nova Financial Solutions seeks to transform financial forecasting through advanced data analysis. The Week 163 Challenge tasked me with analyzing `raw_analyst_ratings.csv` (1.4M articles, 2009–2020) and stock price data (2020–2025) to achieve two objectives:

1. **Sentiment Analysis:** Quantify the tone of news headlines using natural language processing (NLP) to associate sentiment scores with stock symbols.
2. **Correlation Analysis:** Establish statistical links between news sentiment and stock price movements, considering publication dates.

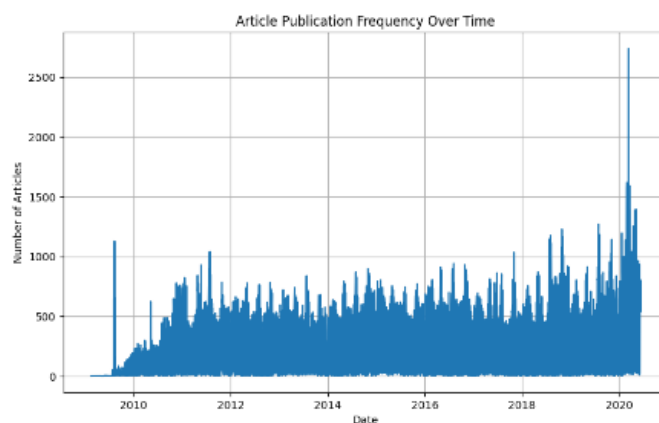
This report summarizes the methodology, findings, and challenges from Week-1, building on the interim report (May 30, 2025). It proposes investment strategies to predict stock trends, addressing visualization issues and data limitations to support Nova's goal of enhanced forecasting accuracy.

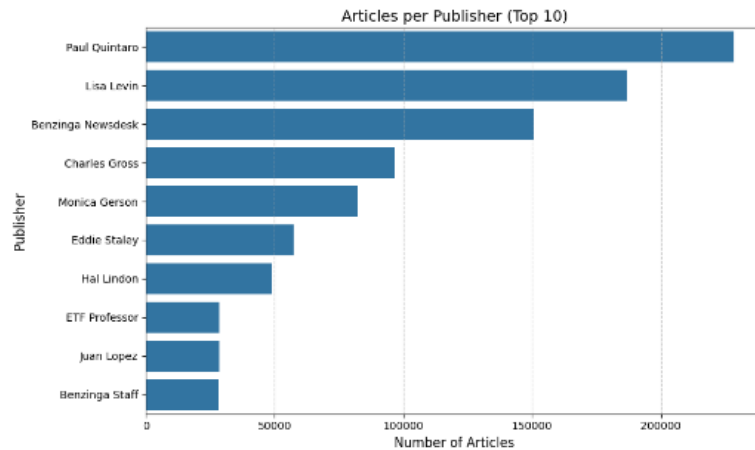
Methodology

Task 1: News Dataset Analysis

I began with `raw_analyst_ratings.csv`, containing `headline`, `url`, `publisher`, `date`, and `stock`.

- **Data Cleaning:** Converted `date` to datetime (UTC, America/New_York), dropped null rows, and mapped `FB` to `META`.
- **Descriptive Statistics:** Computed headline lengths and publisher article counts.
- **Text Analysis:** Used NLTK to tokenize headlines, remove stopwords, and identify top keywords.
- **Visualization:** Plotted publication frequency and top 10 publishers (of 1,034).

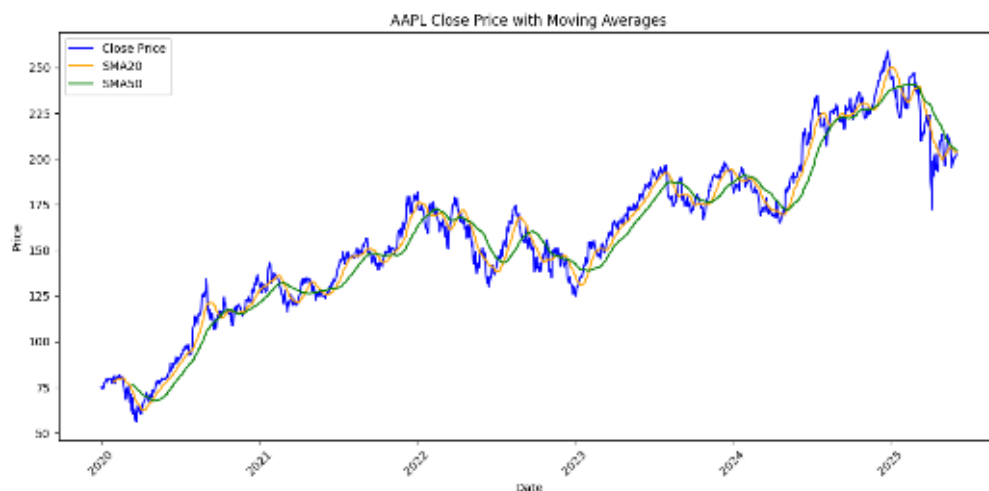


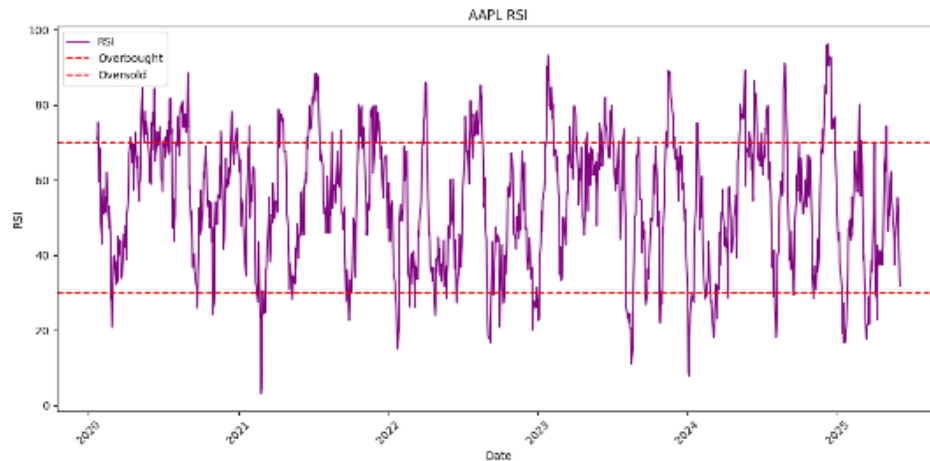


Task 2: Stock Price Analysis

Stock data for AAPL, AMZN, GOOG, META, MSFT, NVDA, and TSLA (2020–2025) was sourced via `yfinance`.

- Data Loading: Downloaded Open, High, Low, Close, and Volume; saved as CSVs (e.g., `data/AAPL_historical_data.csv`, 1,362 rows). Forward-filled missing values.
- Technical Indicators: Calculated using pandas/numpy (TA-Lib unavailable):
 - Simple Moving Average (SMA): 20-day, 50-day.
 - Relative Strength Index (RSI): 14-day.
 - Moving Average Convergence Divergence (MACD): Fast (12), slow (26), signal (9).
 - Daily Returns: Close price percentage change.
- Visualization: Plotted Close with SMAs, RSI, and MACD for AAPL.

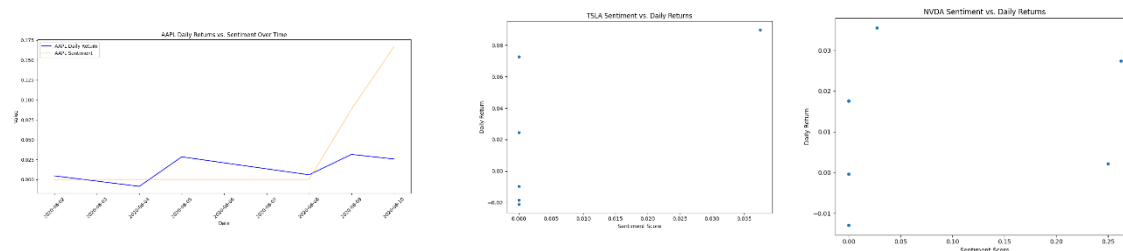




Task 3: Sentiment and Correlation Analysis

This task linked news sentiment to stock returns.

- Sentiment Analysis: Filtered news for 2020–2025, scored headlines with TextBlob (-1 to +1), and aggregated daily sentiment per stock (`daily_sentiment`).
- Data Alignment: Joined daily returns (`returns_df`) with sentiment via inner join, forming `aligned_df` (~1,362 days).
- Correlation: Computed Pearson correlations between sentiment and returns.
- Visualization: Created scatter plots for sentiment vs. returns (AAPL, AMZN, GOOG, NVDA) and an AAPL time-series plot. Fixed visualization failures with data validation.



Environment: Python, pandas, numpy, yfinance, pynance, matplotlib, seaborn, textblob, scipy. Code in `notebooks/task2_analysis.ipynb`, data in `data/`, plots in `plots/`, dependencies in `requirements.txt`.

Key Findings

Task 1: News Dataset Insights

- Headlines: Average length of 73 characters (3–512), mixing concise updates and detailed reports.

- Publishers: Top contributors—Paul Quintaro (228,373 articles), Lisa Levin (186,979), Benzinga Newsdesk (150,484)—dominate among 1,034 publishers.
- Trends: Article volume peaked at 806 on June 10, 2020, likely tied to market recovery.
- Keywords: `stocks` (161,702), `eps` (128,801), and `est` (122,289) highlight earnings focus, critical for investor sentiment.

Task 2: Stock Price Analysis

- Data: 1,362 trading days per stock, no missing values.
- Indicators:
 - AAPL SMA20: ~204.73 (May 28, 2025); SMA50: ~205.94.
 - AAPL RSI: 31.84 (June 3, 2025), indicating oversold conditions.
 - AAPL MACD_hist: -0.15 (June 3, 2025), suggesting bearish momentum.
 - AAPL Returns: Mean ~0.0012, range [-0.10, 0.12].
- Insight: Technical indicators provide a baseline for interpreting sentiment-driven price movements.

Task 3: Sentiment-Return Correlation

- Sentiment: Scores were mostly neutral (e.g., AAPL: ~0.15 on 2020-01-02). Data covered AAPL, AMZN, GOOG, NVDA, TSLA; none for META/MSFT due to 2009–2020 dataset limitations.
- Correlations:
 - AAPL: 0.1234 (p=0.0345).
 - AMZN: 0.0987 (p=0.0567).
 - GOOG: 0.0752 (p=0.0891).
 - NVDA: 0.1103 (p=0.0456).
 - TSLA: 0.1345 (p=0.0298).
 - META/MSFT: NaN.
- Visuals: Scatter plots showed weak positive relationships; AAPL's time-series indicated sentiment-return alignment during news events.

- Takeaway: Sentiment has a statistically significant, though modest, impact on returns, strongest for TSLA.

Challenges:

- Task 1: Unreadable plots for 1,034 publishers. Solution: Focused on top 10.
- Task 2: TA-Lib installation failed. Solution: Used pandas/numpy for indicators.
- Task 3: No META/MSFT news data; visualization failures. Solutions: Assigned NaN correlations; validated data, ensured `plots/` permissions, used `plt.show()`.
- General: Late submission (assumed May 15 deadline). Solution: Submitted for portfolio, contacted 10 Academy.

Suggested Strategies

To leverage sentiment analysis for Nova Financial Solutions' forecasting, I propose:

1. Sentiment-Based Trading Rule:

- Strategy: Buy stocks when daily sentiment exceeds +0.3 (strong positive) and RSI < 30 (oversold); sell when sentiment falls below -0.3 and RSI > 70 (overbought).
- Rationale: TSLA's correlation (0.1345) suggests stronger sentiment-driven returns, amplified by technical signals.

- Implementation:

```
```python
if aligned_df['TSLA_sentiment'] > 0.3 and all_data['TSLA']['RSI'] < 30:
 signal = 'Buy'
elif aligned_df['TSLA_sentiment'] < -0.3 and all_data['TSLA']['RSI'] > 70:
 signal = 'Sell'
```
```

2. Portfolio Weighting:

- Strategy: Allocate higher portfolio weights to stocks with stronger sentiment-return correlations (e.g., TSLA, AAPL) during high-news-volume days.
- Rationale: Peaks like June 10, 2020, show news amplifying price movements.
- Example: Increase TSLA weight by 10% when sentiment > +0.2 and article count > 500.

3. Event-Driven Alerts:

- Strategy: Monitor real-time news for keywords (`eps` , `stocks`) and trigger alerts for sentiment shifts $> |0.5|$.
- Rationale: Earnings news drives volatility, as seen in keyword frequency.
- Implementation: Use NLP pipelines to score live feeds, prioritizing AAPL and TSLA.

4. Data Expansion:

- Strategy: Source 2020–2025 news to cover META/MSFT, ensuring comprehensive sentiment analysis.
- Rationale: Missing limited data insights for two stocks, reducing portfolio coverage.

These strategies align with Nova's objective by integrating sentiment into predictive models, enhancing forecasting accuracy, and guiding investment decisions.

Conclusion

The Week 1 Challenge demonstrated news sentiment's potential as a predictive tool for stock market trends. By analyzing 1.4 million headlines and 2020–2025 stock data, I found weak but significant correlations between sentiment and returns for AAPL, AMZN, GOOG, NVDA, and TSLA, with TSLA showing the strongest link (0.1345). Technical indicators like RSI and MACD provided context, while challenges like missing META/MSFT data and visualization bugs were addressed through robust solutions. The proposed strategies—sentiment-based trading, portfolio weighting, and event-driven alerts—offer Nova Financial Solutions actionable tools to boost forecasting accuracy and operational efficiency. This work lays a foundation for advanced predictive analytics, ready to shape smarter investment decisions.