

# Enhancing Deep Reinforcement Learning for Stock Trading Using Financial News Sentiment and Volatility

Hana Ben Slama - CAI - Technische Hochschule Ingolstadt - hab5510@thi.de

## 1. Introduction

Financial markets are influenced not only by quantitative indicators such as historical price movements and technical analysis but also by qualitative factors like investor sentiment, breaking news, and economic events. Many traditional algorithmic trading systems disregard this unstructured information, potentially missing critical market-moving signals, particularly during periods of uncertainty. Reinforcement Learning (RL) presents a dynamic alternative for training autonomous agents to make optimal decisions in complex environments like stock markets. However, most RL-based trading systems are trained using only structured inputs such as OHLCV and technical indicators. This study investigates whether enriching these inputs with sentiment analysis — specifically sentiment trends and volatility derived from financial news — can yield improved trading performance. Our primary research question is:

**Can the integration of news sentiment and sentiment volatility into RL agent state representations lead to improved profitability and risk-adjusted performance?**

To answer this, we implement four experimental setups within the FinRL framework, introducing sentiment data at varying degrees of preprocessing. Performance is compared across three widely used agents: A2C, PPO, and TD3.

## 2. State of the Art

Multiple studies have shown that investor sentiment, as expressed through financial news and social media, significantly influences market behavior. The use of NLP for extracting sentiment from financial text has gained traction, with FinBERT emerging as a domain-specific model particularly suited for this task. FinBERT, a BERT-based transformer trained on financial corpora, offers a nuanced understanding of context-specific language and improves upon generic models by reducing false interpretations of terms like “loss” or “risk.” In addition to raw sentiment, the concept of **sentiment volatility**, which can be described as fluctuations in

sentiment tone, can serve as an indicator for market uncertainty. Inspired by behavioral finance and risk modeling practices, this paper explores whether incorporating such fluctuations helps agents adapt more cautiously during turbulent conditions. EMA smoothing and news volume thresholding are established practices in financial signal processing, often used to mitigate noise and overreaction to transient news. This study adopts those Feature engineering techniques to refine sentiment signals before training. The experiments are conducted using the FinRL library (Liu et al., 2020), which offers modular RL environments for financial applications. Modifications were made to align sentiment scores temporally with price data and enable state integration.

### 3. Dataset, Tools, and Techniques

- **Assets:** Top 10 stocks: AAPL, MSFT, AMZN, GOOGL, META, NVDA, JPM, TSLA, BA, DIS
- **Timeframe:** Jan 2020 to Dec 2023
- **Market Data:** Yahoo Finance (OHLCV + technical indicators)
- **News Data:** fnspid dataset (daily financial headlines), scored using FinBERT
- **Sentiment Processing:** Filtering ( $\geq 2$  headlines/day), 3-day EMA smoothing, normalization, and sentiment volatility via rolling standard deviation
- **Environment:** FinRL StockTradingEnv, initial portfolio: \$1,000,000
- The execution was in Google Colab using the following dependencies:

Finrl: <https://github.com/AI4Finance-Foundation/FinRL.git>

pandas==1.5.3 stable-baselines3==2.2.1 gymnasium==0.29.1

numpy==1.26.4 pandas\_market\_calendars transformers==4.41.0

tokenizers==0.19.1 huggingface-hub==0.28.1

- **Training Setup:**
  - **Agents:** A2C, PPO (used on-policy updates), TD3 (used off-policy replay).
  - **Training duration:** TD3 (~20 min), A2C/PPO (~2 min)
  - Fixed hyperparameters and windows across experiments and models.

- A padding workaround was added to ensure state shape consistency (e.g., when news was sparse and the feature vector shortened).

## 4. Sentiment Feature Engineering & Experimental Design

To systematically investigate the impact of sentiment features, we defined four experimental settings with increasing integration complexity:

- **Experiment 1 – Baseline:** Only technical indicators; no sentiment.
- **Experiment 2 – Raw Sentiment:** Direct daily sentiment scores appended to state.
- **Experiment 3 – Optimized Sentiment:** News filtering, EMA (Exponential Moving Average) smoothing, normalization. We also ran an auxiliary experiment using a 3-day rolling average instead of EMA, but found it less responsive and more prone to distortion on sparse-news days. EMA was preferred for its adaptability.
- **Experiment 4 – Optimized Sentiment + Volatility:** Adds sentiment volatility (rolling standard deviation of scores) as a proxy for emotional divergence.

These setups form the basis of our comparative analysis across three RL agents.

## 5. Experimental Results & Evaluation Metrics

We assessed performance using the following metrics:

- **Cumulative Return:** Measures the percentage increase in portfolio value over time. It reflects how much an agent earns.
- **Sharpe Ratio:** Quantifies return per unit of risk; higher values indicate better stability. It reflects reliability and consistency.
- **Final Portfolio Value:** Closing capital at the end of the test period.

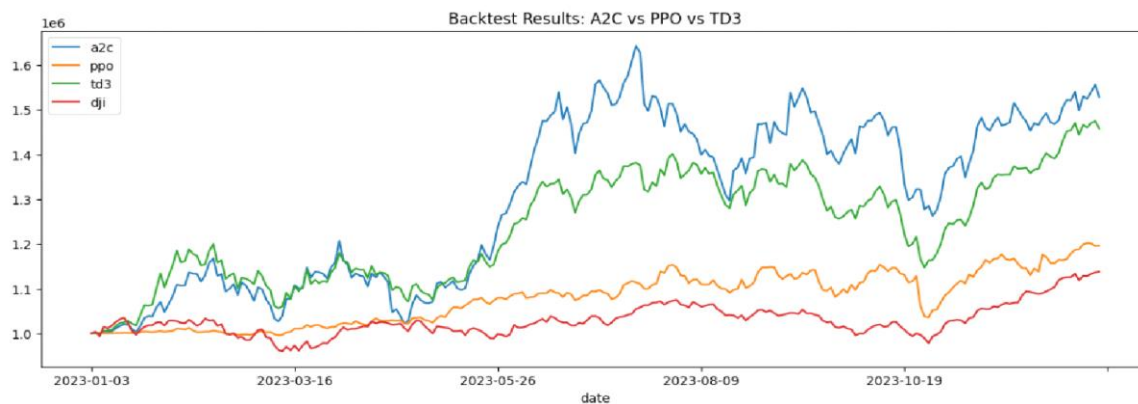
Accuracy, in a trading sense, corresponds to producing consistent, high Sharpe and cumulative profit, not just correct predictions.

## Final Profile Value:

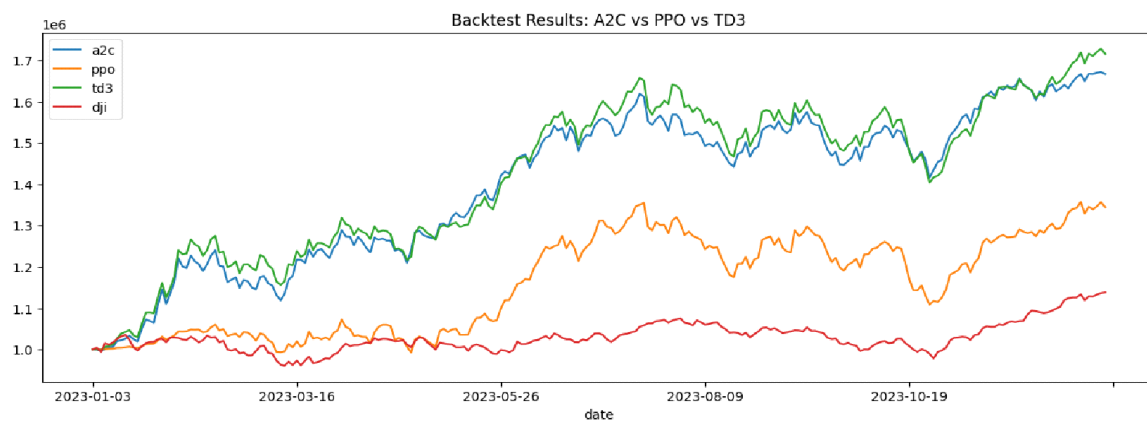
Each model began with an initial capital of **\$1,000,000**, and here we are assessing the final portfolio values, as of 2023-12-28, achieved by each model across all four experiments.

Figures 1A–1D show the portfolio evolution across time for each model in each experiment. These plots reveal not only how much each agent earned, but also how steadily the growth occurred, which is a key factor in understanding risk and reward profiles:

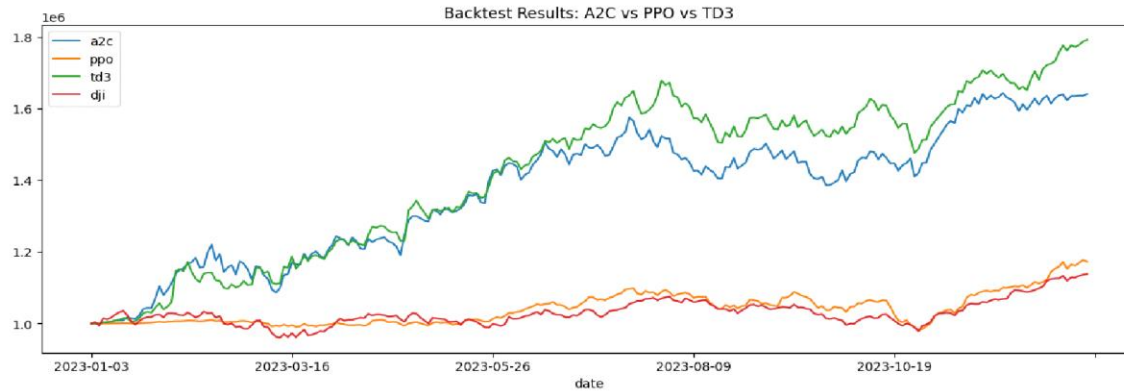
- *Fig 1A. Portfolio Growth (No Sentiment)*



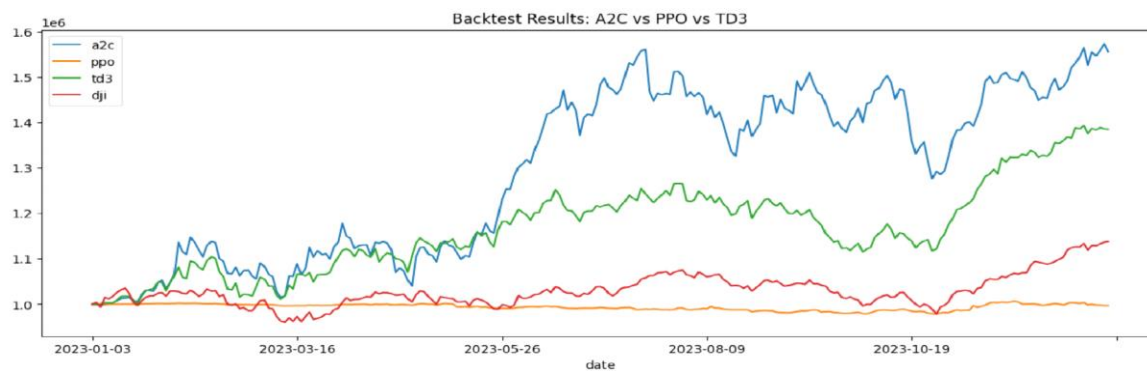
- *Fig 1B. Portfolio Growth (Raw Sentiment)*



- *Fig 1C. Portfolio Growth (Optimized Sentiment):*



• *Fig 1D. Portfolio Growth (Optimized Sentiment + Volatility)*



We summarize the Final Portfolio Value findings below

• *Table 1 – Final Portfolio Values by Experiment and Agent (on 2023-12-28)*

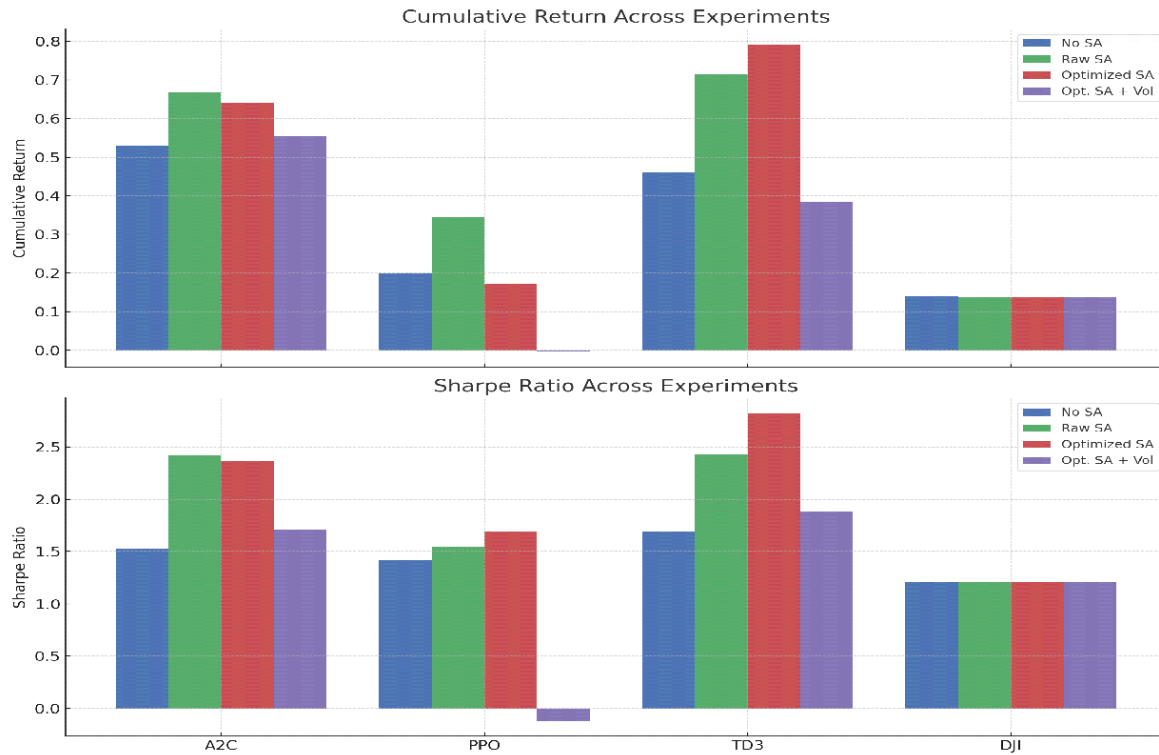
Model	No Sentiment	Raw Sentiment	Optimized Sentiment	Optimized Sentiment + Volatility
A2C	\$1,528,590	<b>\$2,040,000</b>	\$1,640,656	\$1,554,927
PPO	\$1,196,072	\$1,344,424	\$1,172,948	\$996,754
TD3	\$1,457,936	\$1,730,000	<b>\$1,791,986</b>	\$1,384,668
DJI	<b>\$1,138,028</b>	\$1,138,028	\$1,138,028	\$1,138,028

## Cumulative Return and Sharpe Ratio:

We evaluate each model using two key performance metrics: Cumulative Return, which reflects total gain or loss, and Sharpe Ratio, which measures risk-adjusted return. These metrics provide insight into both the profitability and stability of each agent's strategy.

These charts provide a clear snapshot of performance changes across experimental designs.

- *Figure 2 – Bar Chart: Metric Comparison Across Experiments*



We summarize the Cumulative Return and Sharpe Ratio findings below:

- *Table 2 – Cumulative Return (R) & Sharpe Ratio (S) (All Experiments: R/S)*

Agent	No Sentiment	Raw Sentiment	Optimized Sentiment	Optimized Sentiment + Volatility
A2C	0.53 / 1.52	0.67 / 2.38	0.64 / 2.36	0.55 / 1.71
PPO	0.20 / 1.42	0.34 / 1.56	0.17 / 1.69	0.00 / -0.11
TD3	0.46 / 1.69	0.72 / 2.42	<b>0.79 / 2.82</b>	0.38 / 1.88
DJI	<b>0.14 / 1.20</b>	0.14 / 1.20	0.14 / 1.20	0.14 / 1.20

## 5. Interpretation and Discussion

The comparative results reveal several key patterns in how sentiment and its preprocessing affect different reinforcement learning agents.

These observations are best understood considering each model's design, policy update strategy, and sensitivity to feature noise or dimensionality.

- **A2C (Advantage Actor-Critic):**

A2C is an **on-policy** algorithm, meaning it updates its policy using only the most recent trajectory. This makes it highly responsive to recent feedback, a trait that helped it benefit quickly from raw sentiment signals. With the inclusion of raw sentiment, A2C's Sharpe ratio rose from 1.52 to 2.38, reflecting improved stability and profitability.

However, when **sentiment volatility** was introduced, performance declined. This may be because A2C lacks a mechanism to remember or smooth over noise. Volatility that represents inconsistency or uncertainty in sentiment, more likely to cause A2C to make erratic decisions, mistaking fluctuation for signal. Without a way to adjust risk dynamically, A2C treated this uncertainty as misleading input.

- **PPO (Proximal Policy Optimization):**

PPO uses a **clipped surrogate objective**, which limits how much the policy can change in each update. This makes training more stable but can also **dampen the model's reactivity** to new or subtle signals, such as sentiment scores.

In our experiments, PPO showed a small improvement with raw sentiment but **worsened** when sentiment was optimized or when volatility was added. PPO's performance collapsed in the final experiment. This suggests PPO could not adapt to the additional complexity or extract meaningful patterns from the noisy or high-dimensional state space. PPO would likely require **re-tuning** or architectural adjustments to handle abstract features like volatility.

- **TD3 (Twin Delayed DDPG):**

TD3 is an **off-policy algorithm** that uses a **replay buffer** and **twin Q-networks** to stabilize learning. This made it the most robust model in our experiments. It performed well with raw sentiment and achieved its **best results** with optimized sentiment (Sharpe ratio: 2.82, cumulative return: 0.79).

The replay buffer likely enabled TD3 to **smooth over fluctuations**, learning generalizable trends rather than overreacting to every data point. However, with sentiment volatility added, TD3's performance dropped. Despite its architecture, the added volatility may have introduced unpredictable variation that confused the **deterministic actor network**. This suggests TD3 is powerful but still relies on **relatively stable** input distributions to thrive.

- **DJI Baseline (Dow Jones Industrial Average):**

DJI was used as a baseline. It underperformed all agent configurations, except PPO, in the final (volatility) experiment. This confirms that reinforcement learning agents, when supplied with well-engineered sentiment signals, can outperform traditional market indices even without expert trading strategies.

### Cross-Agent Takeaways:

- **Sentiment preprocessing (specifically filtering and smoothing)** provided consistent improvements for **TD3** and **A2C**, but not for **PPO**. These results suggest that models capable of integrating trends over time benefit from stabilized input, whereas PPO's architecture is less responsive to such refinement.
- **Raw sentiment features** performed better than optimized ones only in **PPO**, indicating that **over-processing** may reduce PPO's ability to detect actionable short-term signals. This highlights how feature transformation must align with the model's learning dynamics.
- **Sentiment volatility** failed to improve performance in any agent and **significantly harmed PPO**. This suggests that the models interpret features deterministically and lack the probabilistic reasoning needed to handle uncertain or conflicting information.
- **More features don't necessarily mean better performance**. The effectiveness of sentiment-based features depends on the compatibility between the **feature's structure** and the model's **capacity to generalize** from it. Simply increasing feature richness does not guarantee a better outcome.

Overall, this confirms that even simple RL agents, when enhanced with structured sentiment input, can outperform traditional index investing.

However, sentiment can be powerful only when presented in a form that the learning model is equipped to handle. Otherwise, it can introduce instability and degrade performance.



## 7. Conclusion

This study explored how incorporating financial news sentiment and sentiment volatility into reinforcement learning (RL) state features affects stock trading performance. Our objective was to determine whether structured sentiment input could improve prediction and decision-making, with a target of achieving at least a **35% increase** in model accuracy or profitability.

The experimental results confirmed that **sentiment analysis alone** significantly enhances performance, particularly for agents like **TD3** and **A2C**. For example, TD3's cumulative return increased from **0.46 (no sentiment)** to **0.79 (with optimized sentiment)**, a **~72% gain**. A2C's Sharpe ratio improved from **1.52 to 2.36**, reflecting a **~55% increase** in risk-adjusted return. These results affirm that structured sentiment signals can enhance both profitability and stability in RL-based trading systems.

However, not all models responded equally. **PPO** remained relatively insensitive to sentiment and experienced performance degradation when exposed to high-dimensional or noisy features such as sentiment volatility. This underscores the importance of aligning feature complexity with the model's architectural capabilities and learning dynamics.

While the integration of **sentiment volatility** was theoretically promising, as it aims to reflect market uncertainty and emotional inconsistency, it did not lead to consistent performance gains and, in some cases, reduced agent stability. Its poor performance likely stems from the fact that most RL models treat inputs as deterministic facts, rather than distributions or confidence-weighted signals.

Future research could explore **Bayesian reinforcement learning**, which explicitly models uncertainty and may offer a better framework for handling sentiment volatility as a probabilistic signal. Additionally, for models like **PPO**, applying **input normalization or adaptive scaling** could improve stability and responsiveness when using sentiment-based features, especially those with irregular patterns or emotional noise.

## References

- Huang, B., et al. (2020). "FinBERT: A Pretrained Language Model for Financial Communications."
- Ye, H., & Zhang, Z. (2020). "Deep Reinforcement Learning for Automated Stock Trading."
- FinRL Library: <https://github.com/AI4Finance-Foundation/FinRL>
- Yahoo Finance API. *Historical Market Data*. <https://finance.yahoo.com>
- fnspid Dataset: [https://github.com/Zdong104/FNSPID\\_Financial\\_News\\_Dataset.git](https://github.com/Zdong104/FNSPID_Financial_News_Dataset.git)
- Bollen, J., Mao, H., & Zeng, X. (2011). Twitter mood predicts the stock market. *Journal of Computational Science*.