

19th AIAI 2023



A Sharpe Ratio Based Reward Scheme in Deep Reinforcement Learning for Financial Trading

Georgios Rodinos, Paraskevi Nousi, Nikolaos Passalis, and Anastasios Tefas

Presenter: G. Rodinos

E-mails: {grodinos, paranous, passalis, tefas}@csd.auth.gr

Computational Intelligence and Deep Learning Group (CIDL), AIIA Lab.

Department of Informatics, Aristotle University of Thessaloniki, Thessaloniki, Greece















Structure

- Introduction
- Proposed Method
- Experimental Evaluation
- Conclusion







Introduction

• Automated Financial Trading with Deep Learning

- The use of Deep Learning (DL) models aided the prediction of price movements
- Based on the forecast direction, a trader can make an informed decision about whether to take a long or short position

• Challenges

- Frequently, the generation of supervised labels is necessary
- This task might be a challenge due to the unpredictability of the financial markets







Introduction

• Deep Reinforcement Learning in Financial Trading

- Deep Reinforcement Learning (DRL) is an effective approach that addresses the challenges associated with supervised learning limitations
- The integration with DL has enabled the direct optimization of trading policies to maximize expected profits even in the presence of volatility and uncertainty







Introduction

Reward Approaches

- Profit and Loss (P&L): most DRL systems for financial trading use P&L as a reward function
- Sharpe ratio: captures the risk-related component of an agent's performance, and is used to evaluate a portfolio's risk-adjusted performance

• Drawbacks

- Profit and Loss (P&L): Doesn't take into account the risk associated with the returns
- Sharpe ratio: Limited samples are available during the training phase. To perform the calculation, it is necessary to consider the returns over a long period of time



ARTIFICIAL INTELLIGENCE APPLICATIONS & INNOVATIONS

Introduction

• **Goal:** Incorporating the Sharpe ratio into the reward function of a DRL agent aims to enhance the overall performance of the portfolio by mitigating the risk associated with the agent's decisions













- Our study outlines a method for integrating the Sharpe ratio into the training procedure of a DRL agent
- A dynamic window is proposed which adjusts its size based on the returns obtained within an RL episode. Consequently, an estimation of the Sharpe ratio can be incorporated into the reward function

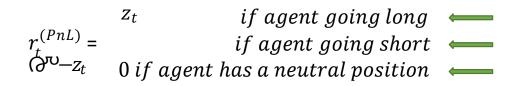






PnL reward

The profit-based reward is defined as:



where z_t is the <u>return change</u> and is defined as:

$$z_t = \frac{p_c(t) - p_c(t-1)}{p_c(t-1)}$$

which is also referred to as the change of the close price p_c .







With the return definition, the profit-based reward can be written as:

$$r_t^{(PnL)} = e_t \cdot z_t \quad \longleftarrow$$

Where e_t is the <u>current market position</u> and $e_t \in \{1, 0, -1\} = \{\text{long, neutral, short}\}$

When the agent changes position is obligated to pay an extra fee. That is called the <u>commission</u>, in which case an additional reward is formulated as:

$$r_t^{(fee)} = -c \cdot |e_t - e_{t-1}| \longleftarrow$$

where *c* denotes the <u>commission</u>. The total <u>PnL reward</u> can be defined as:

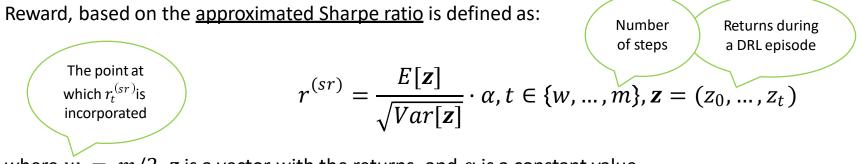
$$r_t^{(total)} = r_t^{(PnL)} + r_t^{(fee)} \quad \longleftarrow$$

Il DeepFinance





• PnL and Sharpe ratio reward



where w = m/2, z is a vector with the returns, and α is a constant value.

The total <u>PnL and Sharpe ratio reward</u> is defined as:

$$r_t^{(total)} = \begin{array}{l} r_t^{(PnL)} + r_t^{(fee)}, & t < w \\ r_t^{(PnL)} + r_t^{(fee)} + r_t^{(sr)}, & for t \ge w \end{array}$$

Il DeepFinance





• Proposed

The total reward of the proposed scheme is defined as:

$$r_{t}^{(total)} = \begin{cases} r_{t}^{(PnL)} + r_{t}^{(fee)}, & for t < w \\ r_{t}^{(PnL)} + r_{t}^{(fee)} + r_{t}^{(sr)}, & for t = w \\ r_{t}^{(PnL)} + r_{t}^{(fee)} + r_{t}^{(sr)}, & if r_{t}^{(sr)} > r_{t-1}^{(sr)}, & for t > w \\ r_{t}^{(PnL)} + r_{t}^{(fee)} - r_{t}^{(sr)}, & if r_{t}^{(sr)} < r_{t-1}^{(sr)}, & for t > w \end{cases}$$







Experimental Evaluation







Dataset

- Crypto trading data
- The dataset provided by SpeedLab AG from 2017-08-17 up to 2022-02-12.
- Training set from 2017-08-17 to 2021-03-15, Test set from 2021-03-15 to 2022-02-12
- In total: 439.737 candles, where the train/test candles are 327.596 and 112.141 candles, respectively.
- 14 currency pairs such as the BTC/BUSD, BTC/USDT, and ETH/USDT among others
- The Open-High-Low-Close (OHLC) price level technique was used to preprocess the data. More specifically, OHLC values are:
 - open price (the first traded price of the set interval)
 - highest and lowest traded prices within the interval
 - close price (the last price that a trade did occur during the interval)
- The minute-price candles are resampled to hour candles





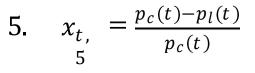
Input Features

1.
$$x_{t, j} = \frac{p_c(t) - p_c(t-1)}{p_c(t-1)}$$

2.
$$x_{t, j} = \frac{p_h(t) - p_h(t-1)}{p_h(t-1)}$$

3.
$$x_{t,} = \frac{p_l(t) - p_l(t-1)}{p_l(t-1)}$$

4.
$$x_{t,4} = \frac{p_h(t) - p_c(t)}{p(t)}$$



time-related features are created, including day, month, week, and year features.

The described features are concatenated into a feature vector $x_t \in \mathbb{R}^{13}$ for each time t.









DRL setup characteristics

- The DRL agent is trained using the Proximal Policy Optimization (PPO) approach
- Neural Network architecture is Long-Short Term Memory
- The number of steps that an DRL episode consists of is equal to 100
- Each experiment is executed 10 times, with each instance using a different random seed
- The PnLs presented, are averaged throughout the 10 experiments as well as the Annualized Sharpe ratios





Proposed Reward Evaluation

| Reward type | Annualized Sharpe ratio Monthly Returns | Annualized Sharpe ratio Hourly Returns |
|--------------------|--|---|
| PnL | 1.462 ± 0.055 | 2.374 ± 0.079 |
| PnL + Sharpe ratio | 1.499 ± 0.060 | 2.484 ± 0.090 |
| Proposed | 1.617 ± 0.056 | 2.641 ± 0.083 |





Proposed Reward Evaluation



Il DeepFinance





Conclusion





Conclusion



Developed and evaluated **a Sharpe ratio-based reward shaping scheme** for training DRL agents that are capable of decreasing the risk that often occurs in agents' trading decisions and improving the overall performance of a portfolio







Acknowledgements

This work has been co-financed by the European Union and Greek national funds through the Operational Program Competitiveness, Entrepreneurship and Innovation, under the call RESEARCH - CREATE - INNOVATE (project code: T2EDK-02094).





Co-funded by Greece and the European Union







Thank you!

Questions?

Il DeepFinance



Co-funded by Greece and the European Union