

Project 2 - Reinforcement Learning

B(E)4M36SMU

Monday 24th March, 2025

In the second homework assignment, we will apply reinforcement learning techniques to Bitcoin trading. Your reinforcement learning agent will be given a set of training and testing data and will try to learn decisions so that it generates a profit. Again, we will use the *Gymnasium* library [1] to test our code, this time with the environment taken from the `gym-trading-env` library [2].

The data contain 1-hour windows with information about open, high, low, and close prices. Those numbers mean the price at the start of the window, the highest price over the window span, the lowest price, and the price at the end of the window. Your task will be then to calculate the position (i.e., the policy) that should be held in Bitcoin (or other cryptocurrency) to maximize your profit.

1 Warning — Comparison to Real-world

Important: Despite one can find many papers about applying reinforcement learning techniques to cryptocurrency trading, for example, in [3], it is not a good idea to apply such a simple algorithm we will develop in a real world situation. It might generate a profit for a while, but it is very likely to generate a loss as well. The real world does not live strictly in the MDP setting; it is slowly developing, and such a *concept drift* requires continuous adjustment of the algorithm decisions. At some points, events that radically change the behavior of the environment occur, and algorithms that worked before suddenly stop working - such events are called *black swans* [4] with examples as the Great Depression in 1930, or the *dotcom bubble crash* at the beginning of this century. One of the reasons for the change in the principles of trading can be seen in the fact that if one thing is profitable, many people try to copy it, thus making it less profitable for others, and at one point, the environment ends up in a new equilibrium [5]. To sum it up, in reinforcement learning language - the environment is not *stationary*.

Other reasons why the direct application of such a simple code is not a good idea can be seen in some simplifications we have made. For example, the official exchange rate being x does not mean you can trade the currency for x . If you are a large investor, there might simply be not enough people willing to buy or sell enough of the currency. We also expected that the environment is only a probability distribution, but in the real world, it is formed by many agents trying to maximize their gain. For all those reasons, please do not try to use this toy code in the real world. No warranties about its outcomes can be given.

2 Implementation

In the downloaded archive, you can find the `main.py` file that constructs the environment for you and allows you to train and test the provided data. Change this class to your own will, but keep in mind that all changes made to this class will not be reflected in BRUTE.

The other file, `tradingagent.py`, contains a draft of the class to implement. There, you might find several methods you will need to change. Change some or all of them; however,

keep the method signatures; otherwise, the automatic evaluation will not work. The contents of the methods are taken from the `gym-trading-env` documentation [2]. Changing the methods content is highly recommended.

- Method `reward_function` can be modified to define reward. Currently, it calculates the logarithm of the percentual gain of your method, but there are definitely better ones. See [6] for more details.
- Method `make_features` constructs features. We will use only static features. In this method, you will get a data frame with open, close, high, and low valuations, and your goal will be to create numeric features that will be good for learning. For example, Bollinger bands or moving averages are good examples of such features. Create new columns with names in the form of `feature_xy` to create a feature that will be returned by the `step` method. See [7] to find documentation.

Important: You are given the training and testing data in the data frame here. Please make sure that you do not create any look-ahead feature. To calculate the feature at time t , use only values from time $\leq t$. I will briefly look into your code to ensure that no such feature exists; creating a look-ahead feature might result in 0 points from the homework.

- The positions available to the agent can be customized in `get_position_list` method. In this method, you will specify the set of actions available to the agent; the values between 0 and 1 mean split of the portfolio into USD and Bitcoin, and values outside this interval represent borrowing either currency. There is a small cost for each exchange, as well as for a loan. More details can be found at [8].
- Method `get_test_position` will be your policy in test. Use trained model here to provide decisions.
- Method `train` will be used for training. You can run one or multiple episodes; there are 15 training episodes, each of 40,000 samples (5 currencies on 3 exchanges), the datasets are iterated in a random order. You can use any technique learned in this class. This time, I don't care whether you follow their assumptions as GLIE or Robbins-Munro theorem; you only need to achieve the required performance. Do not hardcode any strategy; use machine learning (only you can encode the inputs of the strategy into the features).

3 Problem Specification

Implement any reinforcement learning algorithm of your choice in `tradingagent.py` file. Do not hardcode any strategy, and while constructing the features, do not use any look-ahead features.

4 Submission and Evaluation

- All students must work individually.
- Upload the results to <https://cw.felk.cvut.cz/brute/>.
- Strict deadline is Wednesday 16th April, 2025 11:59 pm. I cannot guarantee answers to your queries placed three days before (or any time after) the deadline.
- A penalty for late submission is -1 points for each day of delay.
- Submit `tradingagent.py`.

- The python version and installed packages are listed at the end of the BRUTE submission report.
- The project is worth 10 points in total. On both test and validation datasets, your maximum score is 5 points. Both datasets contain 15 episodes, each with 10,000 samples. If you have positive percentual gain on 5 or less episodes, your score is 0, otherwise, each dataset is worth 0 to 0.5 points, scaled linearly between 0% and 10%. You cannot get more then 5 points on test (or validation) dataset.
- Do not use any look ahead features, do not hardcode any strategy. Use any reinforcement learning algorithm.
- Training time is set to 10 minutes; if you need a raise, email me. Similarly, if you miss a library that you think should be available to all submissions.
- Should you have any questions, or you find a bug in the code or project specification, feel free to email me and/or ask for a consultation.

5 Hints for a Successful Submission

- Start with simple features.
- Before designing the algorithm, think whether there is a connection between two consecutive time steps (the answer may depend on Your features). What are the implications? Should I change the discount factor accordingly? Should the value of a state be a sum of discounted rewards?
- Make sure not to overfit. In a hypothetical example, I might have decided to use tabular-based Q-learning. I have 400 states, half of them visited many times, half of them just once. Is this a good situation? Am I overfitting or not?
- Reward is also important. Keep in mind that the reward should be scale-invariant, an algorithm that has a gain of 5USD with 100USD investment is as good as the algorithm that has a gain of 50USD with 1000USD investment.
- Do not forget that in the `make_features` method, you can create features that contain moving averages, Bollinger bands, etc. It is always better to extrapolate from a window of, say, ten numbers than from only two data points.

References

- [1] <https://gymnasium.farama.org/index.html>
- [2] <https://gym-trading-env.readthedocs.io/en/latest/>
- [3] Otabek, S., Choi, J. Multi-level deep Q -networks for Bitcoin trading strategies. *Sci Rep* 14, 771 (2024). <https://doi.org/10.1038/s41598-024-51408-w>
- [4] https://en.wikipedia.org/wiki/Black_swan_theory
- [5] James Owen Weatherall. *The Physics of Wall Street: A Brief History of Predicting the Unpredictable*. (2013)
- [6] <https://gym-trading-env.readthedocs.io/en/latest/customization.html#custom-reward-function>

- [7] <https://gym-trading-env.readthedocs.io/en/latest/features.html>
- [8] https://gym-trading-env.readthedocs.io/en/latest/environment_desc.html#action-space