

Stock Trading using Reinforcement Learning

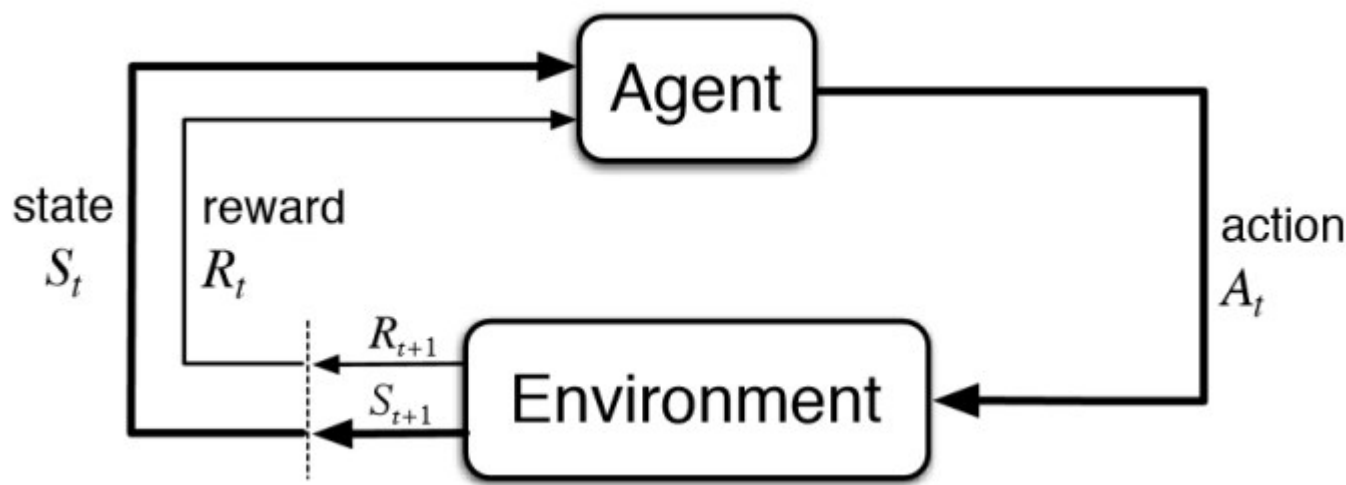
Reinforcement Learning Project

Naman Tejaswi and Arshabh Semwal



Introduction

Reinforcement Learning is one of the machine learning paradigm, supervised, unsupervised and reinforcement learning. A typical RL environment follows the Markov property and is modeled using a Markov Decision Process (MDP). The Markov property states that the future is independent of the past given present.



Objective

We aim to maximize the long-term cumulative trading returns of our Agent by trading a particular stock. This seems a daunting task given all the hedge funds in the world (with much more sophistication) are competing for the same thing.

Traditionally deep neural networks with LSTM which capture long term relationship have been used for stock trading we have tried to apply reinforcement learning and have used **DDPG** (Deep Deterministic Policy Gradient) and **TD3** (Twin Delayed Deep Deterministic Policy Gradient)

Environment Setting

We have selected a random time period of 100 trading days to evaluate our model and we have used the publicly available yahoo finance data. The model is given the input as the open close low and high prices along with the daily volume for 3000 trading days ~12 years.

Environment Setting

- The observation space is the daily price and volume of the stock, we have taken it as a continuous action space by considering each trading day as a timestep.
- The reward function is proportional to the absolute return and takes into account the drawdown periods.
- We have assumed the efficient market hypothesis meaning that we are able to buy at the quoted price without any slippages. We also haven't taken transaction cost into account for now.

DDPG Deep Deterministic Policy Gradient

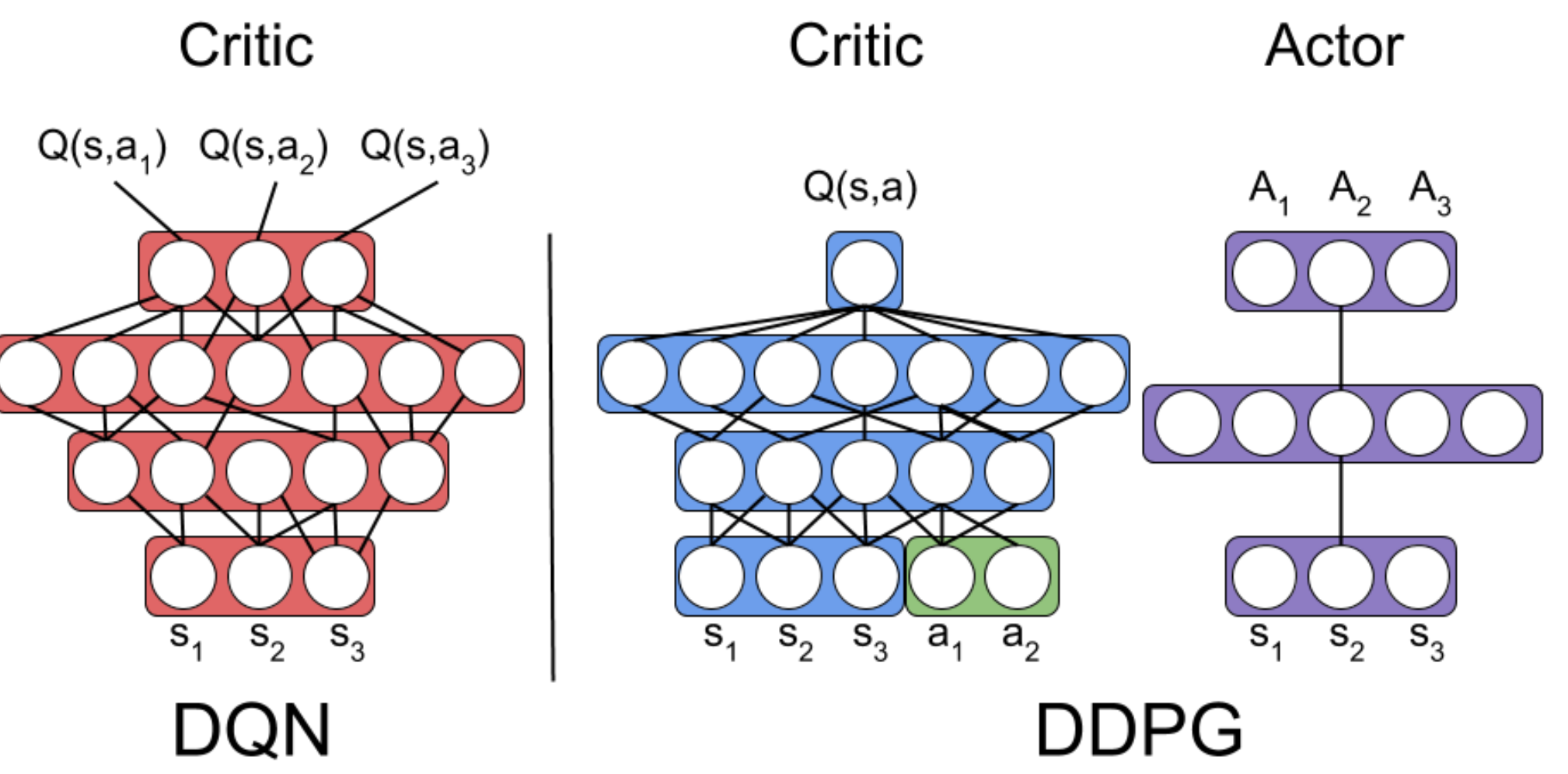
DDPG is similar to DQN but for continuous spaces instead of discrete spaces. DDPG combines the policy gradient method with DQN.

DDPG at a time learns both the Q function using bellman equation and a policy using policy gradient. DDPG uses 2 target networks like Actor Critic.

DDPG is an off-policy algorithm which can be used only in a continuous action space in the vanilla version.

DDPG does soft updates (“conservative policy iteration”) on the parameters of both actor and critic $\theta' \leftarrow \tau\theta + (1 - \tau)\theta'$

In this way, the target network values are constrained to change slowly, different from the design in DQN that the target network stays frozen for some period of time.



```
Algorithm 1 DDPG algorithm
Randomly initialize critic network Q(s, a|θ^Q) and actor μ(s|θ^μ) with weights θ^Q and θ^μ.
Initialize target network Q' and μ' with weights θ'^Q, θ'^μ ← θ^Q, θ^μ
Initialize replay buffer R
for episode = 1, M do
  Initialize a random process N' for action exploration
  Receive initial observation state s1
  for t = 1, T do
    Select action at = μ(st|θ^μ) + Nt according to the current policy and exploration noise
    Execute action at and observe reward rt and observe new state st+1
    Store transition (st, at, rt, st+1) in R
    Sample a random minibatch of N transitions (st, at, rt, st+1) from R
    Set yt = rt + γQ'(st+1, μ'(st+1|θ'^μ))|θ'^Q
    Update critic by minimizing the loss: L = 1/N ∑ (yt - Q(st, at|θ^Q))^2
    Update the actor policy using the sampled policy gradient:
      ∇_θ^μ J ≈ 1/N ∑ ∇_a Q(st, a|θ^Q)|_{a=μ(st|θ^μ)} ∇_θ^μ μ(st|θ^μ)|st
    Update the target networks:
      θ'^Q ← τθ^Q + (1 - τ)θ'^Q
      θ'^μ ← τθ^μ + (1 - τ)θ'^μ
  end for
end for
```

TD3 Twin Delayed DDPG

Twin Delayed DDPG is an extension to DDPG and has 3 additional techniques to DDPG which are:

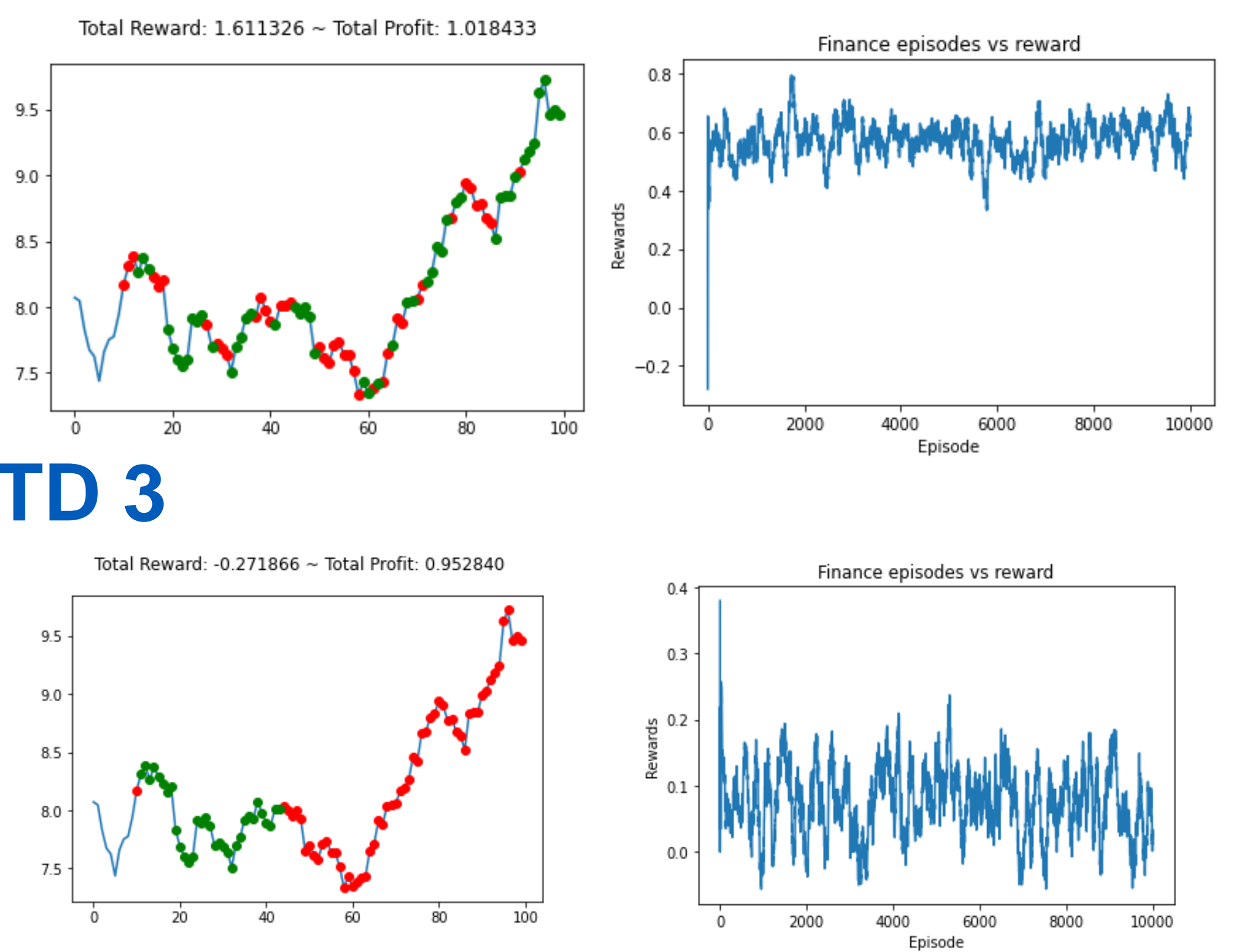
- Clipped Double-Q Learning. TD3 learns two Q-functions instead of one (hence “twin”) and uses the smaller of the two Q-values to form the targets in the Bellman error loss functions.
- “Delayed” Policy Updates. TD3 updates the policy (and target networks) less frequently than the Q-function
- Target Policy Smoothing. TD3 adds noise to the target action, to make it harder for the policy to exploit Q-function errors by smoothing out Q along changes in action.

Since TD3 uses 2 different critic networks it selects the smaller value of the two networks. TD3 reduces the overestimating bias and is more stable than DDPG.

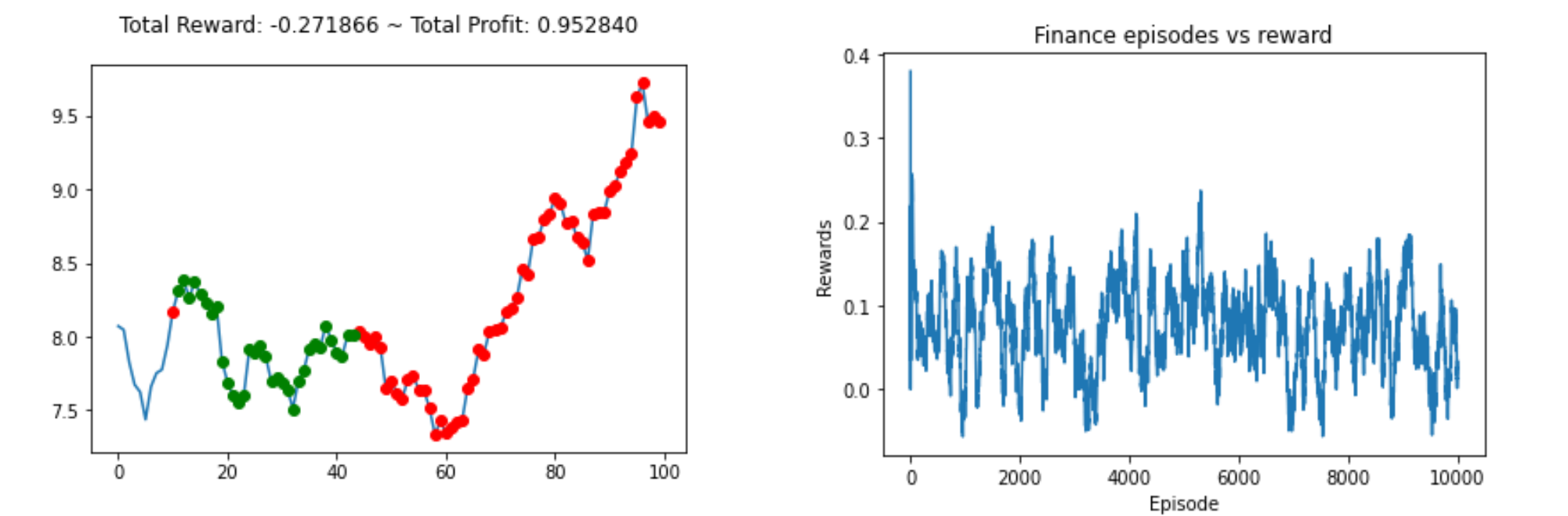
```
Algorithm 1 Twin Delayed DDPG
1: Input: initial policy parameters θ, Q-function parameters φ1, φ2, empty replay buffer D
2: Set target parameters equal to main parameters θ_tgt ← θ, φ_tgt1 ← φ1, φ_tgt2 ← φ2
3: repeat
4:   Observe state s and select action a = clip(μ(s) + ε, a_low, a_high), where ε ~ N
5:   Execute a in the environment
6:   Observe next state s', reward r, and done signal d to indicate whether s' is terminal
7:   Store (s, a, r, s', d) in replay buffer D
8:   If s' is terminal, reset environment state.
9:   if it's time to update then
10:    for j in range(however many updates) do
11:      Randomly sample a batch of transitions, B = {(s, a, r, s', d)} from D
12:      Compute target actions
        a'(s') = clip(μ_tgt(s') + clip(ε, -c, c), a_low, a_high), ε ~ N(0, σ)
13:      Compute targets
        y(r, s', d) = r + γ(1 - d) min_i Q_tgt(s', a'(s'))
14:      Update Q-functions by one step of gradient descent using
        ∇_φi 1/|B| ∑_{(s,a,r,s',d)∈B} (Q_φi(s, a) - y(r, s', d))^2 for i = 1, 2
15:    if j mod policy_delay = 0 then
16:      Update policy by one step of gradient ascent using
        ∇_θ 1/|B| ∑_{s∈B} Q_φ1(s, μ_θ(s))
17:      Update target networks with
        φ_tgt1 ← ρφ_tgt1 + (1 - ρ)φ1
        θ_tgt ← ρθ_tgt + (1 - ρ)θ for i = 1, 2
18:    end if
19:  end for
20: end if
21: until convergence
```

Results

The Red scatter indicates a sell action by our DDPG model, and the Green action indicates a Buy action by our agent on the AAPL and AMZN Stock over 100 trading days..



TD 3



Total Cumulative Returns on 100 trading Days
DDPG 1.84% TD3 -0.5%

Future Works and Improvements

We were expecting better results with TD3 but DDPG in this case DDPG performs better as the delayed update to the model policy is slowing down the policy learning. In addition, we need to further refine our models

We also tried working with options with black Scholes model, but we only had free data for 5-minute interval.

Our long-term goal is to use RL for portfolio optimization by hedging our positions using options and we hope to take this up as a project next semester.

References

- <https://www.coursera.org/learn/advanced-methods-reinforcement-learning-finance>
- <https://spinningup.openai.com/en/latest/>
- <https://github.com/AI4Finance-Foundation>
- <https://web.stanford.edu/class/psych209/Readings/SuttonBartoIPRLBook2ndEd.pdf>
- https://www.cs.princeton.edu/courses/archive/fall09/cos323/papers/b_lack_scholes73.pdf