

HOMEWORK 2

Loïc Jannin

September 2023

1 Problem 1

Consider the Snakes and Ladders game (single-player version) that was covered in class. We conducted 5000 simulations to analyze the game and answer two key questions:

1.1 Probability Distribution Function (PDF) of Finishing the Game in X Dice Rolls

To understand the distribution of the number of dice rolls required to finish the game, we ran the simulation 5000 times. For each simulation, we recorded the number of dice rolls needed to reach the finish line at position 100. The histogram below illustrates the Probability Distribution Function (PDF) of finishing the game in X dice rolls:

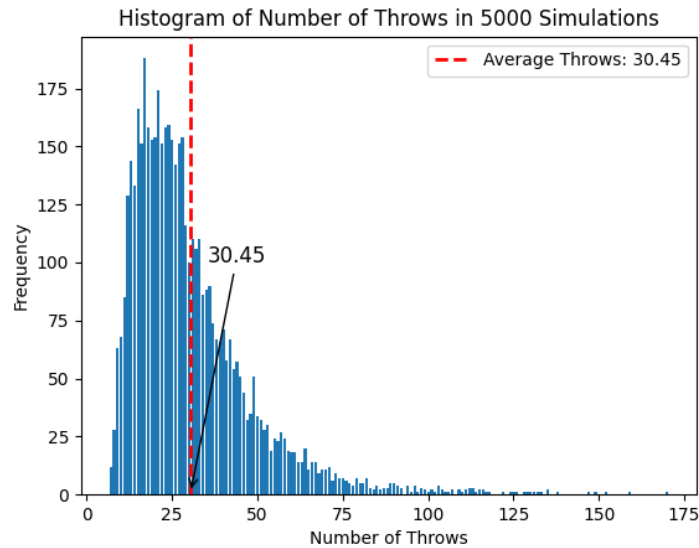


Figure 1: Histogram of Number of Throws in 5000 Simulations

In the histogram, the x-axis represents the number of throws (X), and the y-axis represents the frequency (i.e., how often that number of throws occurred) in the 5000 simulations. The red dashed line indicates the average number of throws required to finish the game, which is approximately 30.45

1.2 Expected Number of Dice Rolls Needed to Finish the Game

The average number of dice rolls needed to finish the game, also known as the expected value, was calculated from the simulation results. We found that the expected number of throws is approximately 30.45 based on the 5000 simulations. This means that, on average, it takes around 30.45 dice rolls to complete the game in our simulation.

2 Problem 2

Optimality of Policy in Markov Decision Processes

Consider two Markov Decision Processes, M1 and M2, with corresponding reward functions $R1$ and $R2$. We are given that M1 and M2 are identical except that the rewards for R2 are shifted by a constant from the rewards for R1, i.e., for all states s , $R2(s) = R1(s) + c$, where c does not depend on s .

We aim to prove that the optimal policy must be the same for both Markov Decision Processes.

Definition: Markov Decision Process (MDP).

A Markov decision process is a 4-tuple (S, A, P_a, R_a) , where:

- S is a set of states called the state space.
- A is a set of actions called the action space.
- $P_a(s, s')$ is the probability that action a in state s at time t will lead to state s' at time $t + 1$.
- $R_a(s, s')$ is the immediate reward (or expected immediate reward) received after transitioning from state s to state s' due to action a .

The state and action spaces may be finite or infinite, for example, the set of real numbers. Some processes with countably infinite state and action spaces can be reduced to ones with finite state and action spaces.

Definition: Policy Function.

A policy function π is a (potentially probabilistic) mapping from the state space S to the action space A .

Definition: Total cumulative rewards.

We define the return G_t from state S_t as :

$$G_t = \sum_{n=t+1}^{\infty} \gamma^{i-t-1} R_i \quad (1)$$

Definition: State-value function under a policy π .

$$V^{\pi}(s) = E_{\pi}(G_t | s_t = s) = E_{\pi} \left(\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \middle| s_t = s \right) \quad (2)$$

2.1 Proof

Let π_1^* be the optimal policy for M1 and π_2^* be the optimal policy for M2. We need to show that $\pi_1^* = \pi_2^*$. For a given state s , we have π_1^* and π_2^* respectively given by the following equations:

$$V^{\pi_1^*}(s) = \max_{\pi_1} V^{\pi_1}(s) \quad (3)$$

$$V^{\pi_2^*}(s) = \max_{\pi_2} V^{\pi_2}(s) \quad (4)$$

Since $R2(s) = R1(s) + c$ for all states s , we can substitute $R2$ into the equation for M2:

$$V^{\pi_2^*}(s) = \max_{\pi_2} \sum_{k=0}^{\infty} \gamma^k E_{\pi_2}[R1(s_{t+k+1}) + c | s_t = s]$$

We can split the sum into two parts:

$$V^{\pi_2^*}(s) = \max_{\pi_2} \sum_{t=0}^{\infty} \gamma^t E_{\pi_2}[R1(s_t) | s_0 = s] + \sum_{t=0}^{\infty} \gamma^t c$$

Notice that the second term, $\sum_{t=0}^{\infty} \gamma^t c$, is a constant that does not depend on the policy π_2 . Since the optimal policy is the one that maximizes the expected cumulative reward, this constant term does not affect the policy's optimality.

Hence, we have proven that the optimal policy must be the same for both Markov Decision Processes, M1 and M2.

3 Problem 3

Optimizing Frog's Escape in an MDP

3.1 Problem Statement:

Consider an array of $n + 1$ lilypads on a pond, numbered 0 to n . A frog sits on a lilypad other than the lilypad numbered 0 or n . When on lilypads i ($1 \leq i \leq n - 1$), the frog can croak one of two sounds, A or B.

- If it croaks A when on lilypad i ($1 \leq i \leq n - 1$), it is thrown to lilypad $i - 1$ with probability $\frac{i}{n}$ and is thrown to lilypad $i + 1$ with probability $\frac{n-i}{n}$.
- If it croaks B when on lilypad i ($1 \leq i \leq n - 1$), it is thrown to one of the lilypads $0, \dots, i - 1, i + 1, \dots, n$ with uniform probability $\frac{1}{n}$.

A snake, located on lilypad 0, will eat the frog if the frog lands on lilypad 0. The frog can escape the pond (and hence, escape the snake!) if it lands on lilypad n .

3.2 Modeling:

To model this problem as an MDP (Markov Decision Process) and derive the Optimal Action Value Function Q^* , we need to define the following components:

- **State Space (S):** The state space will represent the lilypad on which the frog is currently located. It ranges from 0 to n . States 0 and n are terminal states.
- **Action Space (A):** The action space corresponds to the frog's choice of croak (A or B) when on lilypads 1 to $n - 1$. $A = \{a, b\}$
- **Transition Probabilities (P):** We define the transition probabilities based on the frog's croak choice and the movement probabilities provided in the problem statement. $P_a(s, s')$ is the probability that action a in state s at time t will lead to state s' at time $t + 1$. For all $i \in \{1, 2, \dots, n - 1\}$ and for all $j \neq i \in \{0, 1, 2, \dots, n\}$, we have:

$$P_a(i, j) = \begin{cases} \frac{i}{n}, & \text{if } j = i - 1 \\ \frac{n-i}{n}, & \text{if } j = i + 1 \\ 0, & \text{else} \end{cases}$$

$$P_b(i, j) = \frac{1}{n}$$

- **Rewards (R):** The rewards assign values to different states and actions. In this case, we aim to maximize the probability of reaching lilypad n while avoiding lilypad 0 (being eaten by the snake). Different types of rewards can be considered. A natural approach would be to assign a reward of 0 to every lilypad from 1 to $n - 1$, a reward of 1 to lilypad n , and a reward of -1 to lilypad 0. But during the second part of the problem, we will experiment with different reward functions and analyze the results.

With these components, we can formulate the MDP and solve for the Optimal Action Value Function Q^* .

3.3 Results:

Below are the results of a Python code that models the MDP and calculates Q^* for $n = 3, 10$, and 25 , along with plots of $Q^*(s, a')$ as a function of the states of the MDP for each action a' . The code used is the code provided in class, but with some changes to play with the results. Additionally, this code plots the average cumulative reward of our trained agent. This code experiments with different hyperparameters to try and find the most efficient model for training our agent.

The parameters we're going to adjust are:

- The reward function.
- The value of alpha.
- The value of epsilon.
- The value of gamma.

As a reminder, the value of epsilon dictates the tendency of our agent to take a random action (here, to transition to a random state through action B) instead of following the best policy given the Q function. Alpha is the learning rate, generally set between 0 and 1. Setting the alpha value to 0 means that the Q-values are never updated, and nothing is learned. If we set alpha to a high value like 0.9, learning can occur quickly. Gamma is the discount factor set between 0 and 1. This models the fact that future rewards are worth less than immediate rewards.

3.3.1 Sensitivity to reward formulation.

First, we will examine the impact of the reward function. We propose four different reward functions:

1. **Reward 1:** State i gives $\frac{i}{n}$ as a reward.
2. **Reward 2:** Every state gives 0 except state n , which gives 1.
3. **Reward 3:** Every state gives 0 except state n , which gives 1, and state 0, which gives -1.
4. **Reward 4:** State i gives $1 - 1/i^2$ as a reward.

We will plot the optimal Q-values for each state and action, and then we will plot the average cumulative rewards over 10 simulations. We first use hyperparameter scheme 1, which means that the learning rate alpha linearly decreases from 1 to 0 during training with each episode. Then we use hyperparameter scheme 2 which means that the learning rate alpha linearly decreases from 1 to 0 during training with each episode.

Remarks:

1. The reward scheme 3 with hyperparameter scheme 2 seems to not have reached its optimal value in the previous simulations. We performed additional simulations with 2000 runs to investigate further; the results are far better (Figure 8).
2. We observe that the reward functions significantly impact the learning process of the agent. Reward scheme 2 led to the fastest convergence in both hyperparameter schemes (Figure 3 and 4).
3. Hyperparameter scheme 2, where the learning rate α linearly decreases from 1 to 0 during training with each episode, and epsilon keeps its value, shows promising results in terms of convergence (Figure 1 vs 2, Figure 7 vs. 8). The choice of epsilon, which affects the exploration vs. exploitation trade-off, plays a crucial role in the agent's learning. Further experimentation may help fine-tune this hyperparameter.
4. The results from the additional 2000 simulations for reward scheme 3 with hyperparameter scheme 2 show improved convergence, suggesting that more simulations can provide better insights into the agent's performance.
5. Future work may involve exploring different combinations of hyperparameters to optimize the training process further.
6. Understanding the impact of reward functions on the agent's behavior is essential for designing effective reinforcement learning systems.
7. It is evident that in reward function 3, where state 0 incurs a negative reward, the Q-value of action A at state 1 reaches its lowest point. This reflects the significant influence of the adverse consequence associated with state 0's negative reward and, consequently, the frog's cautious disposition under this reward scheme.
8. When examining scenarios where the "core states" receive positive rewards (as in reward functions 1 and 4), a noteworthy observation emerges. In such cases, the frog exhibits a greater propensity to employ croak A, as remaining in the pond yields favorable rewards, as illustrated in Figures 1, 2, 7, and 8.

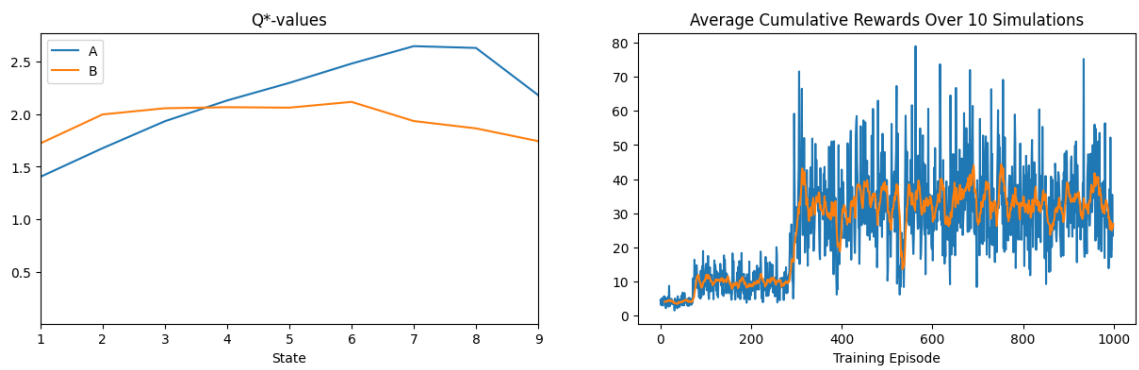


Figure 1: Results for reward 1, hyperparameters scheme 2.

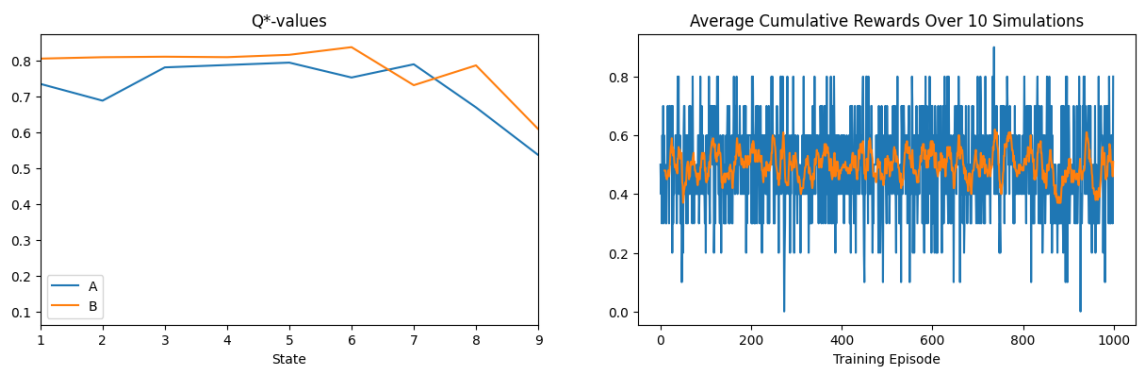


Figure 2: Results for reward 2, hyperparameters scheme 1.

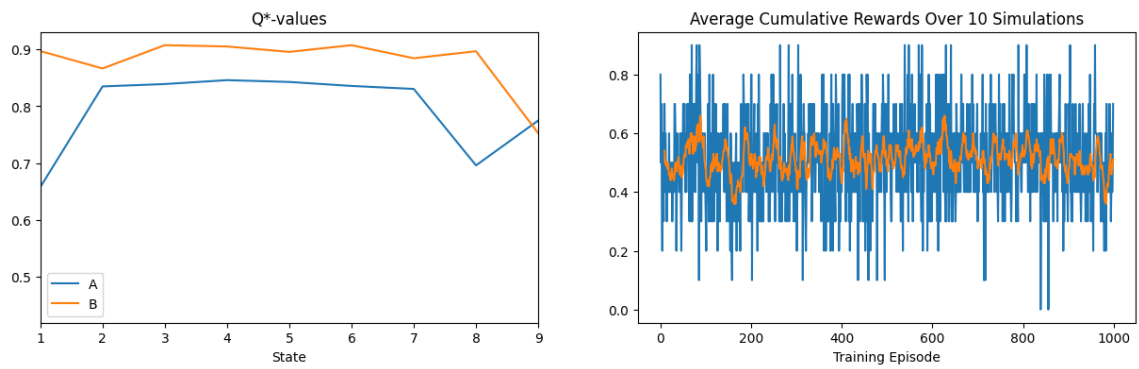


Figure 3: Results for reward 2, hyperparameters scheme 2.

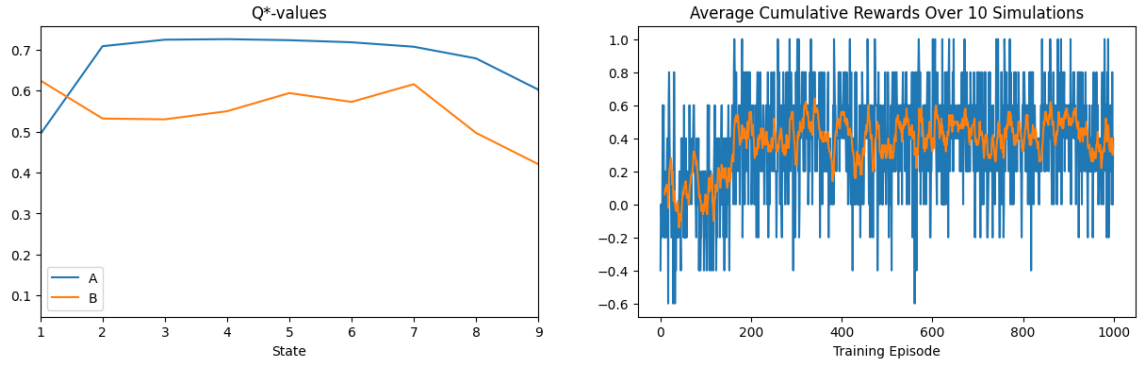


Figure 4: Results for reward 3, hyperparameters scheme 1.

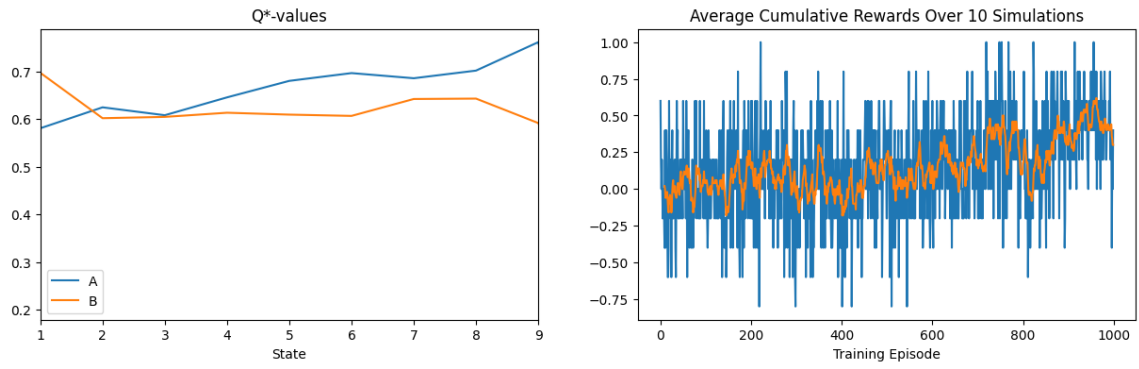


Figure 5: Results for reward 3, hyperparameters scheme 2.

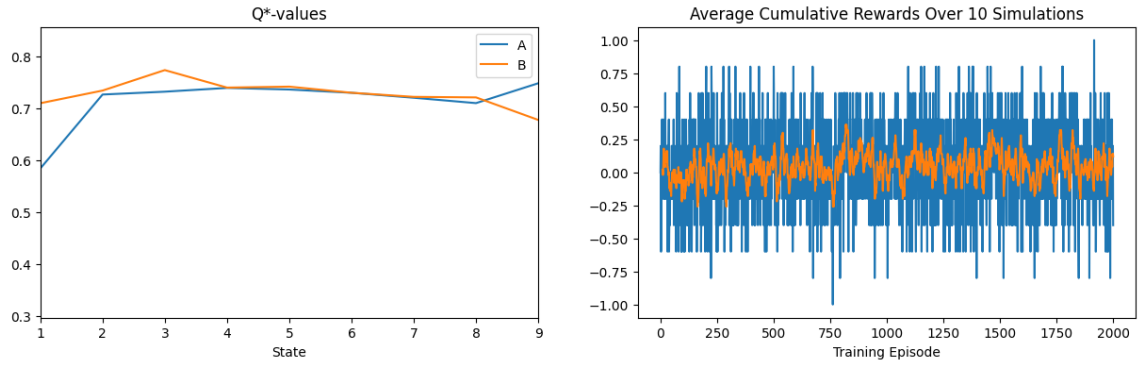


Figure 6: Results for reward scheme 3, hyperparameters scheme 2, and 2000 simulations.

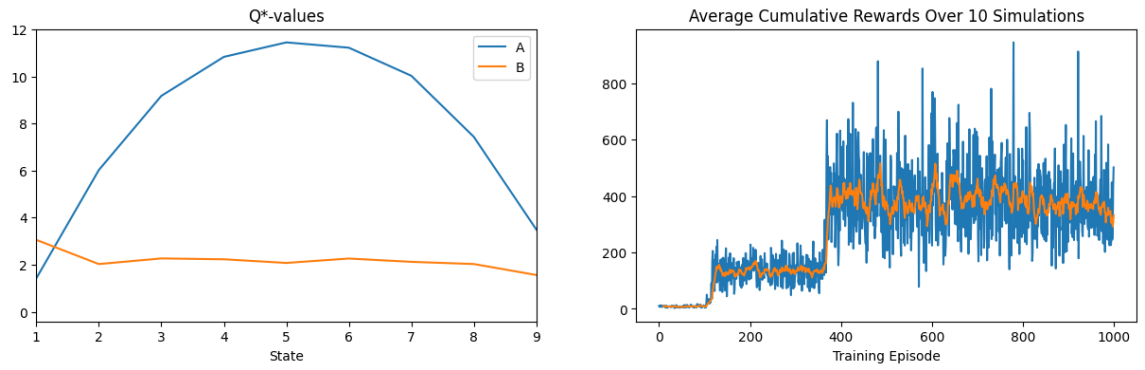


Figure 7: Results for reward scheme 4, hyperparameters scheme 1.

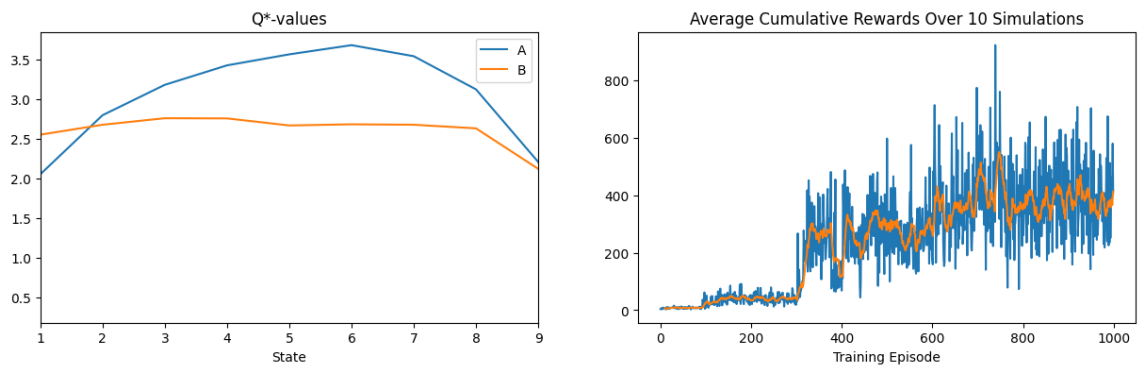


Figure 8: Results for reward scheme 4, hyperparameters scheme 2.

4 Problem 4

Using the depth 5 Amazon data from LOBSTER, we formulated an MDP (Markov Decision Process) to find an optimal execution strategy of a large client buy order of size X using a model-free RL (Reinforcement Learning) method. In this market replay regime, transition probabilities are implied by historical data with the assumption of no market impact. The execution agent makes decisions every 10 seconds. We provide clear mathematical notation for the proposed state space, action space, and rewards of the MDP.

4.1 State Space

Each state $s \in S$ is a vector of attributes (state variables) that describes the current configuration of our system. We use the term 'state' rather than 'observed state' due to the complexity of the environment, but we will assume that the Markov property is respected.

- State Variable 1: **Remain Inventory**. Noted "I". Gives the number of share left we have to buy.
- State Variable 2: **Imbalance**. Noted "Imb". Given by $Imb = \frac{BidVolume}{AskVolume}$
- State Variable 3: **Volatility**. Noted σ . Given by $\sigma = \sqrt{\text{Annualization Factor} \times \text{Var}(\text{Log Returns})}$ Computes the volatility for the last 100 ticks.
- State Variable 4: **Spread**. Noted Sp , given by the formula: $Sp = AskPrice - BidPrice$
- State Variable 5: **Traded Volumes**. Noted V , gives the sum of the traded volume during the last 10s.
- State Variable 6: **Best Bid Price**.
- State Variable 7: **Best Ask Price**.
- State Variable 8: **Last event** Dummy variables that can take values 1;2;3;4;5 given the last event, lobster notation.
- State Variable 9: **10s Returns**
- State Variable 10: **1s Returns**

4.2 Action Space

The action space defines the possible actions that the execution agent can take at each time step. Actions represent different execution strategies or trading decisions. The action space A includes:

- Action a: **Buy Limit order**. Action 'a' involves the placement of a limit order for any remaining unexecuted shares at a price equal to the current bid price plus 'a', with a fixed volume of 10 per cent the total share volume to buy. This action effectively cancels any existing outstanding limit orders we may have, aligning with the practices supported by actual exchanges. It's worth noting that 'a' can be either a positive or negative value, with 'a=0' indicating an entry at the current bid price. In reality, the action state is not limited to a single action 'a'; we can model a continuous action space with all possible values of 'a.' However, a more straightforward approach would be to discretize 'a', and allowing our agent to chose the different buy volumes (choice between 1 per cent, 10, 50) . Another important point is that there's no need for a separate 'market order' action, as a market order is equivalent to choosing a = spread.

4.3 Rewards

The rewards in the MDP play a critical role in providing immediate feedback to the agent based on its actions and the state of the environment. These rewards can be strategically designed to motivate the agent towards the ultimate goal of executing the large client buy order optimally. Since our primary objective is to successfully execute a large buy order, our immediate rewards essentially represent the cash outflows resulting from any (partial) execution of the limit order placed.

To facilitate meaningful comparisons across various training episodes with different buy volumes, it becomes imperative to 'normalize' the rewards for the purpose of comparison. Drawing inspiration from the work of Nevmyvaka, Feng, and Kearns, we will compare our total cash outflow to the ideal cash flow scenario, where we hypothetically purchase all the shares at the initial time midprice.

Let $CF(s, a, s')$ be the outflow realized by the transition between state s' and s given action a . Let $Id(V, t_0)$ be the ideal cash flow scenario. The reward function is defined by:

$$R(s, a, s') = \frac{CF(s, a, s')}{Id(V, t_0)}$$

4.4 Results

The selection of state variables is critical for training an effective RL agent. Among the state variables examined, imbalance, volatility, and spread are expected to be particularly useful. Imbalance provides insights into order book dynamics, while volatility and spread impact trading decisions and execution strategies. To address the question of the usefulness of each state variable, we could examine the correlation between each state variable and positive events for our agent, such as identifying favorable future stock prices for our agent to purchase shares. However, a more explicit method for assessing the utility of each variable would involve constructing a decision tree.

Further experimentation and fine-tuning of the MDP components will be necessary to achieve optimal execution results. Predicting which state variables will be the most useful for our MDP without running simulations can be quite challenging. Additionally, it's worth noting that, given the reward function, certain state variables may have more influence than others. (All codes are in annex)

4.5 Plots

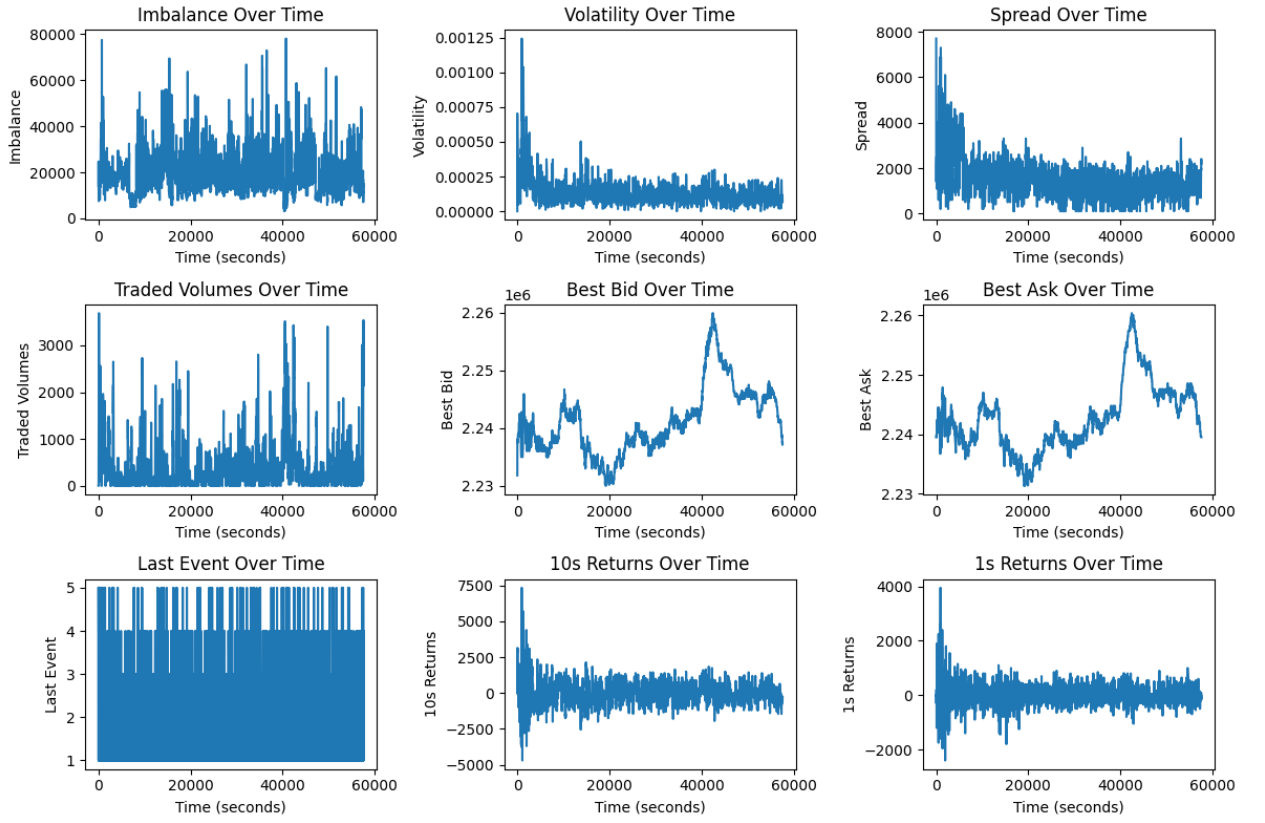


Figure 2: State Space variables as function of time

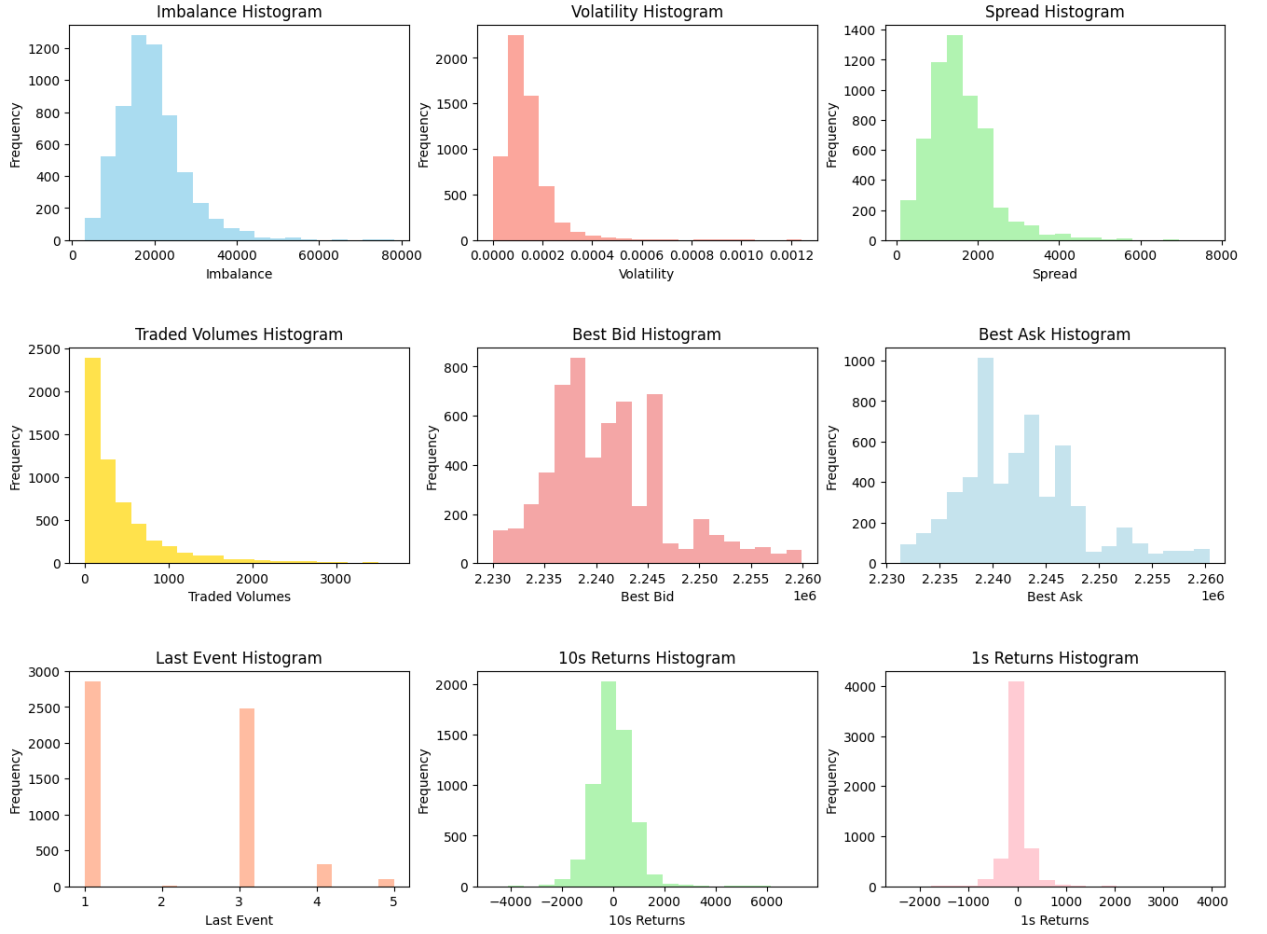


Figure 3: Histogram of States space variables

5 Problem 5

We propose a Markov Decision Process (MDP) formulation for designing an execution agent with deadlines. The goal is to find an optimal execution strategy for executing a large client buy order of size X by time T , ensuring that the entire order is executed by the specified deadline. In this problem, we are required to express the MDP with clear mathematical notation, including the state space, action space, rewards, and transition probabilities model. Once again, this work is inspired by the paper "Reinforcement Learning for Optimized Trade Execution", shared in class.

5.1 State Space

The state space for this MDP consists of 5 state variables that capture the relevant information about the market and the current state of the execution agent. We Limited the state space to 5 state variables for simplicity.

1. State Variable 1: **Remain Inventory**. Noted "I". Gives the number of share left we have to buy. Takes values between X and 0.
2. State Variable 2: **Imbalance**. Noted "Imb". Given by $Imb = \frac{BidVolume}{AskVolume}$
3. State Variable 3: **Volatility**. Noted σ . Given by $\sigma = \sqrt{\text{Annualization Factor} \times \text{Var}(\text{Log Returns})}$ Computes the volability for the last 10s.
4. State Variable 4: **Spread**. Noted Sp , given by the formula: $Sp = AskPrice - BidPrice$
5. State Variable 5: **Elapsed time**. Noted t. Takes values between 0 and T.

Let's mention that any state with a value of I=0 or t=T is a terminal state.

5.2 Action Space

The action space defines the set of possible actions that the execution agent can take at each time step. Actions represent different execution strategies or trading decisions.

1. Action 1: **Buy market order**
2. Action 2: **Buy limit order at price midprice(t) + a.**

For each of these actions, we attach a certain volume V. For complexity purposes, the volume V can take the values: $\frac{I}{2}, \frac{I}{5}, \frac{I}{10}, \frac{I}{100}$.

5.3 Rewards

The rewards in this MDP play a pivotal role in shaping the agent's behavior. When designing the reward function, we have several options to consider. One approach is to impose a significant penalty on any order executed after time T with a substantial negative reward. We can introduce a parameter β to model "how bad it is to exceed T". Re utilizing the former reward function notations (problem 4):

$$R(s, a, s') = \frac{CF(s, a, s')}{Id(V, t_0)} - \beta I_{t>T}$$

Where :

$$I_{t>T} = \begin{cases} 1, & \text{if } t > T \\ 0, & \text{if } t < T \end{cases}$$

Alternatively, we could implement a more nuanced strategy by gradually penalizing orders as the execution time progresses, encouraging the agent to execute a larger portion of the order early on. This reward function could introduce a risk aversion parameter, allowing us to control the rate at which penalization increases.

$$R(s, a, s') = \frac{CF(s, a, s')(1 - \alpha t)}{Id(V, t_0)}$$

To create a more accurate representation of our problem, we may also consider incorporating a constraint in the action space rather than the reward function. For instance, we can enforce that the agent must execute all remaining inventory at the market price when time reaches T . This constraint would ensure that the remaining shares are executed promptly without explicitly relying on the reward function.

5.4 Transition Probabilities Model

Given the complexity of the data we use, the transition probability model is not explicitly defined but rather inferred from the simulation results.

5.5 Market Impact

In addressing market impact within this problem, we face the choice of whether to explicitly model its effects on our execution strategy. If we opt not to model market impact, we should make assumptions, such as executing small orders in a highly liquid market, where our actions have minimal influence on prices. However, modeling market impact realistically can be intricate. It requires training the agent and other agents in the environment to respond to market dynamics and adapt as if trading in a real-world environment. Training the agent solely on the LOBSTER dataset, as seen in Problem 4, does not account for market impact accurately, as it lacks information on how the market reacts to our orders. Incorporating a temporary market impact model into the MDP is a more sophisticated approach but significantly increases complexity. Furthermore, defining the agent’s behavior in response to market impact becomes crucial—whether it should minimize impact by executing smaller orders over a more extended period or prioritize faster execution, even with higher price impact. These considerations highlight the complexities and choices involved in handling market impact effectively.

Listing 1: Python code for problem 1

```

1 import random as rd
2 import matplotlib.pyplot as plt
3
4 def next_position(position):
5     # Entry : Int position, position of the player in the game
6     # Output: Int position, position of the player after one play
7
8     # Create a dictionary to map positions after a dice roll
9     position_mapping = {
10         1: 38, 4: 14, 9: 31, 17: 6, 21: 42, 28: 84, 36: 44, 47: 26, 49: 11,
11         51: 67, 56: 53, 64: 60, 71: 91, 87: 24, 80: 100, 93: 73, 95: 75
12     }
13
14     # Compute the result of the dice
15     dice_result = rd.randint(1, 6)
16
17     # Calculate the new position
18     new_position = position + dice_result
19
20     # Check if the new position is in the position mapping
21     if new_position in position_mapping:
22         return position_mapping[new_position]
23
24     return new_position
25
26 def simulation():
27     # Entry - None
28     # Output - Int : The number of throws needed to arrive at 100
29
30
31     # Intiatlization
32     number_of_throws = 0
33     position = 0
34
35     # Play
36     while position < 100:
37         position = next_position(position)
38         number_of_throws += 1
39
40     # We arrived at 100
41     return number_of_throws
42
43 mean = 0
44 result_list = []
45 for i in range(5000):
46     result = simulation()
47     mean+=result
48     result_list.append(result)
49

```

```

50
51 # Plot the histogram of result_list
52 plt.hist(result_list, bins=range(min(result_list), max(result_list) + 1),
53          align='left', rwidth=0.8)
54
55 # Calculate the average number of throws
56 average_throws = mean / 5000
57
58 # Add a vertical line for the average number of throws
59 plt.axvline(x=average_throws, color='red', linestyle='dashed', linewidth=2,
60            label=f'Average_Throws: {average_throws:.2f}')
61
62 # Set labels and title
63 plt.xlabel('Number_of_Throws')
64 plt.ylabel('Frequency')
65 plt.title('Histogram_of_Number_of_Throws_in_5000_Simulations')
66
67 # Add a legend
68 plt.legend()
69
70 # Add the x value as text annotation
71 plt.annotate(f'{average_throws:.2f}', xy=(average_throws, 1), xytext=(
72     average_throws + 5, 100), fontsize=12,
73     arrowprops=dict(arrowstyle='->', color='black'))
74
75 # Show the plot
76 plt.show()

```

Listing 2: Python code for problem 3 (adapted from class)

```

1
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import pandas as pd
5
6 class Frog:
7     def __init__(self, n=10, gamma=0.99, initial_alpha=0.1, initial_eps
8         =0.9, num_simulations = 50,
9         hyperparameter_scheme=2, rewards_choice=2
10        ):
11         self.n = n
12         self.gamma = gamma
13         self.initial_alpha = initial_alpha
14         self.initial_eps = initial_eps
15         self.state_space = list(range(n+1))
16         self.terminal_space = [0, n]
17         self.action_space = ['A', 'B']
18
19         self.hyperparameter_scheme = hyperparameter_scheme
20         self.rewards_choice = rewards_choice

```

```

21     #initialization - might impact the speed of training (especially in
        early stages), but should not impact the outcome
22     #self.Q = np.zeros([len(self.state_space), len(self.action_space)])
23     self.Q = np.random.rand( len(self.state_space), len(self.
        action_space))

24
25     self.num_simulations = num_simulations
26     self.simulated_rewards = []
27
28
29     def eps_greedy_action(self, eps):
30         rv = np.random.uniform(0,1)
31         if rv < eps:
32             #random action
33             return np.random.choice(self.action_space)
34         else:
35             #follow the best action
36             return self.action_space[np.argmax(self.Q[self.state, :])]
37
38
39     def action_A(self, state):
40         rv = np.random.uniform(0,1)
41         if rv < float(state)/self.n:
42             return state - 1
43         else:
44             return state + 1
45
46     def action_B(self, state):
47         state_space = list(self.state_space)
48         state_space.pop(state)
49         return np.random.choice(state_space)
50
51     def choose_reward_function(self, state):
52         if self.rewards_choice == 1:
53             return self.get_reward(state)
54         elif self.rewards_choice == 2:
55             return self.get_reward2(state)
56         elif self.rewards_choice == 3:
57             return self.get_reward3(state)
58         elif self.rewards_choice == 4:
59             return self.get_reward4(state)
60         else:
61             print("ERROR!_Unknown_reward_function")
62
63     def get_reward(self, state):
64         #reward assignment - you can experiment with different rewards
65         return float(state)/self.n
66
67     def get_reward2(self, state):
68         #reward assignment - you can experiment with different rewards

```

```

69         if state == 0:
70             return 0
71         elif state == self.n:
72             return 1
73         else:
74             return 0
75
76     def get_reward3(self, state):
77         #reward assignment
78         if state == 0:
79             return -1
80         elif state == self.n:
81             return 1
82         else:
83             return 0
84
85     def get_reward4(self, state):
86         #reward assignment
87         if float(state) == 0:
88             return 0
89         else:
90             return 1-1/float(state)**2
91
92     def simulate(self, num_simulation):
93         simulated_rewards = []
94         for i in range(num_simulation):
95             state = np.random.randint(1, self.n - 1)
96             reward = 0
97             while state not in self.terminal_space:
98                 #follow the best policy
99                 action = self.action_space[np.argmax(self.Q[state, :])]
100                 if action == 'A':
101                     state_new = self.action_A(state)
102                 else:
103                     state_new = self.action_B(state)
104                 reward += self.choose_reward_function(state_new)
105                 state = state_new
106             simulated_rewards.append(reward)
107         #return cumlated rewards over num_episode simulations for a given
108         #policy
109         return np.mean(simulated_rewards)
110
111     def simulate_for_large_n(self, num_simulation):
112         #this function might be useful for testing/debugging your code for
113         #large n
114         simulated_rewards = []
115         for i in range(num_simulation):
116             state = np.random.randint(1, self.n - 1)
117             reward = 0
118             num_iter = 0

```

```

117         while state not in self.terminal_space and (num_iter<0.5e7):
118             #with a small probability pick action B not to be stuck in
            #the infinite loop traversing the lilypads,
119             #otherwise follow the best policy
120             rv = np.random.uniform(0,1)
121             if rv < 1e-2:
122                 action = 'B'
123             else:
124                 action = self.action_space[np.argmax(self.Q[state, :])]
125             if action == 'A':
126                 state_new = self.action_A(state)
127             else:
128                 state_new = self.action_B(state)
129             reward += self.choose_reward_function(state_new)
130             state = state_new
131             num_iter +=1
132         if (num_iter<0.5e7):
133             simulated_rewards.append(reward)
134         else:
135             print("Dropped_rewards_due_to_large_time_needed_to_simulate
                ")
136         #return cumulated rewards over num_episode simulations for a given
            #policy (this value is calibrated to n=25)
137         return np.mean(simulated_rewards)
138
139     def choose_hyperparameter_scheme(self, i, num_episode):
140         if self.hyperparameter_scheme == 1:
141             return self.my_hyperparameter_scheme_1(i, num_episode)
142         elif self.hyperparameter_scheme == 2:
143             return self.my_hyperparameter_scheme_2(i, num_episode)
144         else:
145             print("ERROR!_Unknown_hyperparameter_scheme")
146
147     def my_hyperparameter_scheme_1(self, i, num_episode):
148         if i<500:
149             eps = self.initial_eps
150             alpha = float(self.initial_alpha*(num_episode - i))/num_episode
151         else:
152             eps = 0
153             alpha = float(self.initial_alpha*(num_episode - i))/num_episode
154                     /10
155         return [eps, alpha]
156
157     def my_hyperparameter_scheme_2(self, i, num_episode):
158         eps = self.initial_eps
159         alpha = float(self.initial_alpha*(num_episode - i))/num_episode
160         return [eps, alpha]
161
162     def q_learning(self, num_episode):
163         for i in range(num_episode):

```

```

163         print("training_episode", i)
164         self.state = np.random.randint(1, self.n - 1)
165
166         #my hyperparameter scheme - feel free to implement your own
167         [eps, alpha] = self.choose_hyperparameter_scheme(i, num_episode
168         )
169
170         while self.state not in self.terminal_space:
171             #epsilon-greedy action selection
172             action = self.eps_greedy_action(eps)
173             #follow action to a new state
174             if action == 'A':
175                 state_new = self.action_A(self.state)
176             else:
177                 state_new = self.action_B(self.state)
178             #get reward at a new state
179             reward = self.choose_reward_function(state_new)
180             #Q-update
181             self.Q[self.state, self.action_space.index(action)] +=
182                 alpha * (reward + self.gamma * np.max(self.Q[state_new,
183                 :]) - self.Q[self.state, self.action_space.index(action)
184                 ])
185             self.state = state_new
186
187         #now simulated rewards for the fixed Q table
188         self.simulated_rewards.append(self.simulate(self.
189         num_simulations))
190
191     def all_plots(self):
192         plt.figure(figsize=(15, 4))
193         plt.subplot(121)
194         plt.title("Q*-values")
195         plt.plot(self.Q[:,0], label='A')
196         plt.plot(self.Q[:,1], label='B')
197         plt.xlabel('State')
198         plt.xlim((1, self.n-1))
199         plt.xticks(range(1, self.n))
200         plt.legend()
201         plt.subplot(122)
202         plt.title("Average_Cumulative_Rewards_Over_10_Simulations")
203         plt.plot(pd.Series(self.simulated_rewards))
204         plt.plot(pd.Series(self.simulated_rewards).rolling(10).mean())
205         plt.xlabel('Training_Episode')
206         plt.show()

```

Listing 3: Python code for problem 4

```

1 import numpy as np
2 import matplotlib.pyplot as plt
3 import pandas as pd

```

```

4
5
6 # Initialization of the orderbook and message file:
7 message_file = 'AMZN_2012-06-21_34200000_57600000_message_5.csv'
8 orderbook_file = 'AMZN_2012-06-21_34200000_57600000_orderbook_5.csv'
9 message_data = pd.read_csv(message_file,header=None)
10 orderbook_data = pd.read_csv(orderbook_file,header=None)
11
12 # Initialization of inventory, lets say we want to buy 100 shares
13 I = 100
14
15 def get_mean_for_tens():
16     # This function returns the average number of tick between every 10s
17     # It avoid the 10s and 10s volatility function to have a tremendous
18     # complexity by looking for the last
19     # 10s tick time among all the data
20     n=0
21     tick_list = []
22     previous_row = None
23
24     for index, row in message_data.iterrows():
25         if previous_row is not None:
26             if row[0]>34200+n*10 and previous_row[0]<34200+n*10:
27                 tick_list.append(index)
28                 n+=1
29             else :
30                 tick_list.append(index)
31                 n+=1
32             previous_row=row
33
34     for i in range(len(tick_list)-1):
35         tick_list[-i]=tick_list[-i]-tick_list[-i-1]
36     tick_list.pop(0)
37     mean = np.mean(tick_list)
38     print(mean)
39
40 # Retreive all the states variables
41 def get_remain_inventory():
42     return I
43
44 def get_midprice(time, orderbook_data):
45     # Entry:
46     #     The orderbook lobster file from a given stock on df format:
47     #         orderbook_data
48     #     The ticker time at format integer X: t
49     # Output:
50     #     The midprice at ticker time t
51     midprice = (float(orderbook_data.iloc[time,0]) + float(orderbook_data.
52         iloc[time,2])) / 2
53     return midprice

```

```

51
52 def get_midprice_physical(t, orderbook_data, message_data):
53     # --- midprice_physical ---
54     # Entry:
55     #     The orderbook lobster file from a given stock on df format:
56     #         orderbookFile
57     #     The message lobster file from a given stock on df format:
58     #         messageFile
59     #     The physical time at format XXXX.XXXXXXXX: t
60     # Output:
61     #     The midprice at time t
62
63     # Finds the index of the last transaction made before time = t
64     filtered_data = message_data[message_data.iloc[:, 0] < t]
65     tickerT = len(filtered_data)
66
67     # To avoid an indice out of bounds error
68     if tickerT==len(message_data):
69         return (float(orderbook_data.iloc[tickerT-1,0]) + float(
70             orderbook_data.iloc[tickerT-1,2])) / 2
71
72     # This index gives us the best bid and best ask at the given time
73     midprice = (float(orderbook_data.iloc[tickerT,0]) + float(
74         orderbook_data.iloc[tickerT,2])) / 2
75     return midprice
76
77 def get_imbalance(time, orderbook_data):
78     # Entry:
79     #     The orderbook lobster file from a given stock on CVS format:
80     #         orderbookFile
81     #     The ticker time at format integer X: t
82     # Output:
83     #     The imbalance at ticker time t
84     askVol = 0
85     bidVol = 0
86     for i in range(10):
87         askVol += orderbook_data.iloc[time,2*i+1]
88         bidVol += orderbook_data.iloc[time,2*i]
89     return bidVol/askVol
90
91 def get_volatility(time, orderbook_data):
92     # Entry : tick Time
93     # This function computes the volability of 10 tick returns during the
94     # last 100 ticks
95     # During the first min, lets say at time = 10 the function
96     # returns the volability during the first 10s
97
98     # Init

```



```

95     start_time = max(0 ,time - 100)
96     t = start_time
97     midprice_list =[]
98
99     # Computes the 1s return
100    while t <= time:
101        row = orderbook_data.iloc[t]
102        midprice = (row[0]+row[2])/2
103        midprice_list.append(midprice)
104        t+=10
105
106    # Compute the returns
107    midprice_return_list = []
108    for i in range(1, len(midprice_list),10):
109        midprice_return_list.append(midprice_list[i]-midprice_list[i-1])
110
111    # Computes the var and volatility
112    sigma = np.std(np.log(midprice_list))
113    return sigma
114
115    def get_spread(time,orderbook_data):
116        # --- spread_physical ---
117        # Entry:
118        #     The orderbook lobster file from a given stock on df format:
119        #         orderbookFile
120        #     The ticker time at format n: t
121        # Output:
122        #     The spread at ticker time t
123
124        # This index gives us the best bid and best ask at the given time
125        spread = abs(float(orderbook_data.iloc[time,0]) - float(orderbook_data.
126            iloc[time,2]))
127        return spread
128
129    def get_traded_volumes(time,message_data):
130        # Entry : Tick time
131        # This function returns the cumulative traded volumes durnig the last
132        # 100 ticks
133
134
135        start_time = max(0 ,time - 79)
136        t = start_time
137        cumulated_volume = 0
138
139        while t <= time:
140            row = message_data.iloc[t]
141            if row[1] in [4, 5]:
142                cumulated_volume+=row[3]
143            t+=1
144
145        return cumulated_volume

```

```

142
143 def get_best_bid(time, orderbook_data):
144     return orderbook_data.iloc[time, 2]
145
146 def get_best_ask(time, orderbook_data):
147     return orderbook_data.iloc[time, 0]
148
149 def get_last_event(time, message_data):
150     return message_data.iloc[time, 1]
151
152 def get_10s_returns(time, orderbook_data):
153     start_time = max(0, time - 79)
154     return ( get_midprice(time, orderbook_data) - get_midprice(start_time,
155                                     orderbook_data) )
156
157 def get_1s_returns(time, orderbook_data):
158     start_time = max(0, time - 8)
159     return ( get_midprice(time, orderbook_data) - get_midprice(start_time,
160                                     orderbook_data) )
161
162 def get_stateSpace(time, orderbook_data, message_data):
163     # returning the value of the state space at time t
164     # The value of I is supposed to be dynamically changed so here we will
165     # just use the initial
166     # value of I as we dont trade.
167     return I, get_midprice(t, orderbook_data), get_imbalance(t, orderbook_data
168                                     ), get_volatility(t, orderbook_data), get_spread(t, orderbook_data),
169     get_traded_volumes(t, message_data), get_best_bid(t, orderbook_data),
170     get_best_ask(t, orderbook_data), get_last_event(t, message_data),
171     get_10s_returns(t, orderbook_data)
172
173 # Retrieve functions
174 def retrieve_values(time, orderbook_data, message_data):
175     # Entry : Same as usual
176     # Output: diferents list containing the values of each states variables
177     # every 10s for t < time
178
179     # Init :
180     midprice_list = []
181     imbalance_list = []
182     volatility_list = []
183     spread_list = []
184     traded_volumes_list = []
185     bid_list = []
186     ask_list = []
187     last_event_list = []
188     tenS_return_list = []
189     oneS_return_list = []
190
191     for t in range(0, time, 10):

```

```

184     print(t)
185     midprice_list.append(get_midprice(t, orderbook_data))
186
187     imbalance_list.append(get_imbalance(t, orderbook_data))
188
189     volatility_list.append(get_volatility(t, orderbook_data))
190
191     spread_list.append(get_spread(t, orderbook_data))
192
193     traded_volumes_list.append(get_traded_volumes(t, message_data))
194
195     bid_list.append(get_best_bid(t, orderbook_data))
196
197     ask_list.append(get_best_ask(t, orderbook_data))
198
199     last_event_list.append(get_last_event(t, message_data))
200
201     tenS_return_list.append(get_10s_returns(t, orderbook_data))
202
203     oneS_return_list.append(get_1s_returns(t, orderbook_data))
204     return imbalance_list, volatility_list, spread_list,
205         traded_volumes_list, bid_list, ask_list, last_event_list,
206         tenS_return_list, oneS_return_list
207
208 # Plot function
209 def plot_values_over_time(time, orderbook_data, message_data):
210     # Call the retrieve_values function to get the lists of state variables
211     imbalance_list, volatility_list, spread_list, traded_volumes_list,
212     bid_list, ask_list, last_event_list, tenS_return_list,
213     oneS_return_list = retrieve_values(time, orderbook_data,
214     message_data)
215
216     # Create a time vector for plotting (assuming time increments of 10
217     seconds)
218     time_vector = list(range(0, time, 10))
219
220     # Plot each state variable
221     plt.figure(figsize=(12, 8))
222
223     plt.subplot(3, 3, 1)
224     plt.plot(time_vector, imbalance_list, label='Imbalance')
225     plt.title('Imbalance Over Time')
226     plt.xlabel('Time (seconds)')
227     plt.ylabel('Imbalance')
228
229     plt.subplot(3, 3, 2)
230     plt.plot(time_vector, volatility_list, label='Volatility')
231     plt.title('Volatility Over Time')
232     plt.xlabel('Time (seconds)')
233     plt.ylabel('Volatility')

```

```

228
229 plt.subplot(3, 3, 3)
230 plt.plot(time_vector, spread_list, label='Spread')
231 plt.title('Spread Over Time')
232 plt.xlabel('Time (seconds)')
233 plt.ylabel('Spread')
234
235 plt.subplot(3, 3, 4)
236 plt.plot(time_vector, traded_volumes_list, label='Traded Volumes')
237 plt.title('Traded Volumes Over Time')
238 plt.xlabel('Time (seconds)')
239 plt.ylabel('Traded Volumes')
240
241 plt.subplot(3, 3, 5)
242 plt.plot(time_vector, bid_list, label='Best Bid')
243 plt.title('Best Bid Over Time')
244 plt.xlabel('Time (seconds)')
245 plt.ylabel('Best Bid')
246
247 plt.subplot(3, 3, 6)
248 plt.plot(time_vector, ask_list, label='Best Ask')
249 plt.title('Best Ask Over Time')
250 plt.xlabel('Time (seconds)')
251 plt.ylabel('Best Ask')
252
253 plt.subplot(3, 3, 7)
254 plt.plot(time_vector, last_event_list, label='Last Event')
255 plt.title('Last Event Over Time')
256 plt.xlabel('Time (seconds)')
257 plt.ylabel('Last Event')
258
259 plt.subplot(3, 3, 8)
260 plt.plot(time_vector, tenS_return_list, label='10s Returns')
261 plt.title('10s Returns Over Time')
262 plt.xlabel('Time (seconds)')
263 plt.ylabel('10s Returns')
264
265 plt.subplot(3, 3, 9)
266 plt.plot(time_vector, oneS_return_list, label='1s Returns')
267 plt.title('1s Returns Over Time')
268 plt.xlabel('Time (seconds)')
269 plt.ylabel('1s Returns')
270
271 plt.tight_layout()
272 plt.show()
273
274 def plot_histograms(time, orderbook_data, message_data):
275     # Call the retrieve_values function to get the lists of state variables
276     imbalance_list, volatility_list, spread_list, traded_volumes_list,
        bid_list, ask_list, last_event_list, tenS_return_list,

```

```

    oneS_return_list = retrieve_values(time, orderbook_data,
    message_data)
277
278 # Create a time vector for plotting (assuming time increments of 10
    seconds)
279 time_vector = list(range(0, time, 10))
280
281 # Create a larger figure for improved readability
282 plt.figure(figsize=(16, 12))
283
284 # Adjust the subplot layout for better formatting
285 plt.subplots_adjust(hspace=0.5)
286
287 plt.subplot(3, 3, 1)
288 plt.hist(imbalance_list, bins=20, color='skyblue', alpha=0.7)
289 plt.title('Imbalance_Histogram')
290 plt.xlabel('Imbalance')
291 plt.ylabel('Frequency')
292
293 plt.subplot(3, 3, 2)
294 plt.hist(volatility_list, bins=20, color='salmon', alpha=0.7)
295 plt.title('Volatility_Histogram')
296 plt.xlabel('Volatility')
297 plt.ylabel('Frequency')
298
299 plt.subplot(3, 3, 3)
300 plt.hist(spread_list, bins=20, color='lightgreen', alpha=0.7)
301 plt.title('Spread_Histogram')
302 plt.xlabel('Spread')
303 plt.ylabel('Frequency')
304
305 plt.subplot(3, 3, 4)
306 plt.hist(traded_volumes_list, bins=20, color='gold', alpha=0.7)
307 plt.title('Traded_Volumes_Histogram')
308 plt.xlabel('Traded_Volumes')
309 plt.ylabel('Frequency')
310
311 plt.subplot(3, 3, 5)
312 plt.hist(bid_list, bins=20, color='lightcoral', alpha=0.7)
313 plt.title('Best_Bid_Histogram')
314 plt.xlabel('Best_Bid')
315 plt.ylabel('Frequency')
316
317 plt.subplot(3, 3, 6)
318 plt.hist(ask_list, bins=20, color='lightblue', alpha=0.7)
319 plt.title('Best_Ask_Histogram')
320 plt.xlabel('Best_Ask')
321 plt.ylabel('Frequency')
322
323 plt.subplot(3, 3, 7)

```

```

324 plt.hist(last_event_list, bins=20, color='lightsalmon', alpha=0.7)
325 plt.title('Last_Event_Histogram')
326 plt.xlabel('Last_Event')
327 plt.ylabel('Frequency')
328
329 plt.subplot(3, 3, 8)
330 plt.hist(tenS_return_list, bins=20, color='lightgreen', alpha=0.7)
331 plt.title('10s_Returns_Histogram')
332 plt.xlabel('10s_Returns')
333 plt.ylabel('Frequency')
334
335 plt.subplot(3, 3, 9)
336 plt.hist(oneS_return_list, bins=20, color='lightpink', alpha=0.7)
337 plt.title('1s_Returns_Histogram')
338 plt.xlabel('1s_Returns')
339 plt.ylabel('Frequency')
340
341 # Show the plots
342 plt.show()

```