
EL2805 Reinforcement Learning

Computer Lab 1

Abdessamad Badaoui & Nasr Allah Aghelias
20011228-T118 & 20010616-T318

1 The Maze and the Random Minotaur

1.1 Basic Maze

1.1.1

We model the state space \mathcal{S} as a tuple of all possible positions of the player (x_p, y_p) (excluding walls), and all possible positions of the Minotaur (x_m, y_m) :

$$\mathcal{S} = \{((x_p, y_p), (x_m, y_m)) : \text{such that } (x_p, y_p) \text{ is not an obstacle and } 0 \leq x_p, x_m \leq 7 \text{ and } 0 \leq y_p, y_m \leq 6\} \quad (1)$$

We let the player choose any action from the action space \mathcal{A} defined by :

$$\mathcal{A} = \{\text{up, down, right, left, stay}\} \quad (2)$$

The objective of the player is to identify a strategy maximizing the probability of exiting the maze (reaching B) before time T :

$$r(s, a) = p_M * R_{mp} + (1 - p_M) * R_{step} \quad (3)$$

where p_M is the probability of the Minotaur being at the same position as the new deterministic state of the player after taking action a .

$R_{mp} = -\infty$: represents the reward of being eaten by the Minotaur.

R_{step} : equals to -1 when taking action a , leads the player to some position in the maze other than the exit, the walls or the Minotaur's position, 0 if it leads to the exit, and $-\infty$ if it leads to a wall.

14

The player transitions are deterministic, whereas the Minotaur's transitions are random, and we can express that formally as follows :

16

$$\mathbb{P}((p_t, m_t) | (p_{t-1}, m_{t-1}, a)) = \frac{1}{N} \quad (4)$$

where N equals to the number of Minotaur's possible actions ($2 \leq N \leq 4$) and p_t, m_t are the possible positions of the player and the Minotaur respectively after taking the action a from the state (p_{t-1}, m_{t-1}) , which means that the probability of the other states is 0.

1.1.2

To model this case where the player and the Minotaur do not move simultaneously, we need just to modify the reward function as follows:

22

- 23 • If at state s , taking action a , leads the player to a wall or to the Minotaur's position then
24 $r(s, a) = -\infty$
- 25 • If at state s , taking action a , leads the player to some position in the maze other than the
26 exit, the walls or the Minotaur's position then $r(s, a) = -1$
- 27 • if at state s , taking action a , leads the player to the exit then $r(s, a) = 0$

28 1.2 Dynamic Programming

29 1.2.1

30 For a time horizon $T = 20$, we compute the policy that maximizes the probability of leaving the
31 maze alive using dynamic programming. Then we simulate the Minotaur's behavior randomly and
32 we compute the player's next move according to the best policy. The following figures show the
33 player and Minotaur's movement in the maze.

34 1.2.2

35 We plotted the probability of the player exiting the maze alive for all time horizons $T \in \{1, \dots, 30\}$,
36 while taking into account that the Minotaur can or can not stand still during rounds. The result can be
37 observed in the following figure.

38 Allowing the Minotaur to remain stationary introduces an extra potential location for the Minotaur
39 at each time step. This additional option adds complexity to the player's decision-making process,
40 potentially reducing the likelihood of successfully navigating the maze and exiting alive.

41 1.3 Value Iteration

42 1.3.1

43 Now, in this scenario, we are dealing with a random time horizon following a geometric distribution.
44 Consequently, we discount the importance of future rewards using the parameter λ . This parameter
45 can be determined using the following relationship, which links the expectation of the time horizon
46 with λ :

$$\mathbb{E}(T) = \frac{1}{1-\lambda} = 30 \iff \lambda = \frac{1-30}{30} = \frac{29}{30} \quad (5)$$

Here we solve the discounted infinite-horizon MDP problem using value iteration, the objective here is to find a stationary policy π that minimizes the infinite horizon objective with a discount factor λ

$$\mathbb{E} \left[\sum_{t=0}^{\infty} \lambda^t r(s_t, \pi(s_t)) \right].$$

47 Bellman equations in the case of a stationary policy π :

$$\forall s \in \mathcal{S} \quad V^*(s) = \max_{\pi} \left\{ r(s, \pi(s)) + \lambda \sum_{s' \in \mathcal{S}} \mathbb{P}(s'|s, \pi(s)) V^*(s') \right\} \quad (6)$$

48 or equivalently in terms of the Bellman operator \mathcal{L}

$$V^* = \mathcal{L}(V^*) \quad (7)$$

49 where

$$\forall s \in \mathcal{S} \quad \mathcal{L}(V)(s) = \max_{\pi} \left\{ r(s, \pi(s)) + \lambda \sum_{s' \in \mathcal{S}} \mathbb{P}(s'|s, \pi(s)) V(s') \right\}. \quad (8)$$

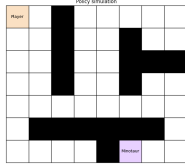


Figure 1: State at instant 0

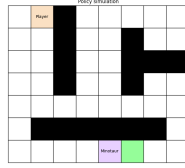


Figure 2: State at instant 1

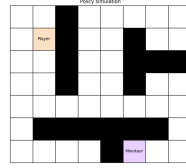


Figure 3: State at instant 2

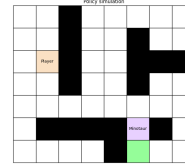


Figure 4: State at instant 3

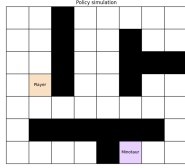


Figure 5: State at instant 4

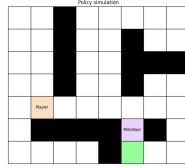


Figure 6: State at instant 5

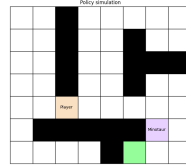


Figure 7: State at instant 6

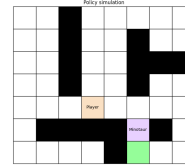


Figure 8: State at instant 7

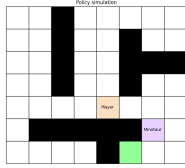


Figure 9: State at instant 8

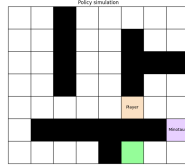


Figure 10: State at instant 9

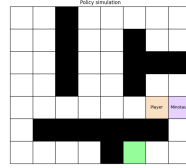


Figure 11: State at instant 10

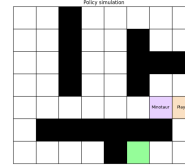


Figure 12: State at instant 11

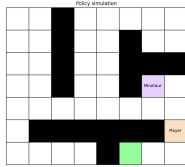


Figure 13: State at instant 12

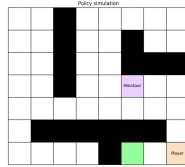


Figure 14: State at instant 13

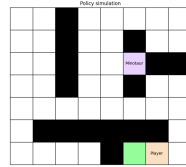


Figure 15: State at instant 14

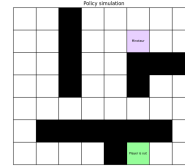


Figure 16: State at instant 15

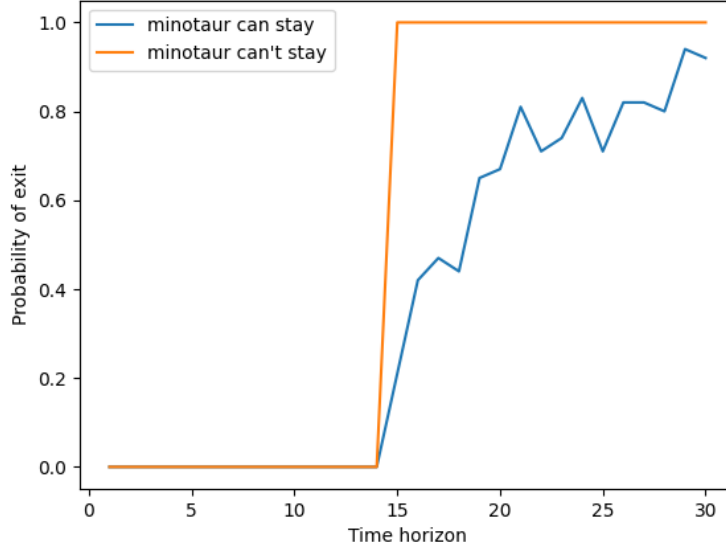


Figure 17: The probability of successful exit with respect to time horizon

1.3.2

Based on the graphics, we observe that the player needs at least $T=15$ to reach the exit. To calculate the probability of the player successfully reaching the exit alive, we just need to calculate the probability of T being greater than 15. This gives us:

$$\mathbb{P}(T \geq 15) = 1 - p(T \leq 14) = 1 - F_T(15) = 1 - 1 + \left(\frac{29}{30}\right)^{14} = 0.622 \quad (9)$$

By simulating 10,000 games, we obtained an estimate of the probability of successfully exiting using the optimal policy: 0.619. This value is close to the theoretical result.

1.4 Additional questions

1.4.1

In off-policy learning, the algorithm learns the value of the optimal policy independently of the actions taken by the agent. The policy employed by the agent is often referred to as the behavior policy, denoted as π_b . On the contrary, in on-policy learning, the algorithm learns the value of the policy being executed by the agent. This policy is determined from previously collected data, making it an active learning method where the gathered data is controlled and used for further learning.

1.4.2

For this new problem, the state space and action space remain the same as the previous MDP. We modify the reward function as follows :

$$r(s, a) = p_M * R_{mp} + (1 - p_M) * R'_{step} \quad (10)$$

where R'_{step} equals to -1 when taking action a , leads the player to some position in the maze other than the exit, the walls, the Minotaur's position or the first visit of C. 50 if it leads to the first visit of C, 0 if it leads to the exit and $-\infty$ if it leads to a wall.

69 $R_{mp} = -100$: represents the reward of being eaten by the Minotaur.
70 p_M is the probability of the Minotaur being at the same position as the new deterministic state of the
71 player after taking action a .
72

73 The Minotaur's transitions are no longer 100% random. There is 35% chance that the Minotaur
74 will move towards the player and 65% chance he will move randomly in all possible directions. Let
75 us consider d_{Man} the Manhattan distance between two positions. We can express the transition
76 probabilities as follows :

$$\mathbb{P}((p_t, m_t)|(p_{t-1}, m_{t-1}, a)) = \begin{cases} 0.35 \cdot \frac{1}{M} + 0.65 \cdot \frac{1}{N} & \text{if } d_{Man}(m_t, p_{t-1}) < d_{Man}(m_{t-1}, p_{t-1}) \\ 0.65 \cdot \frac{1}{N} & \text{if } d_{Man}(m_t, p_{t-1}) \geq d_{Man}(m_{t-1}, p_{t-1}) \end{cases} \quad (11)$$

77 where N equals to the number of Minotaur's possible actions ($2 \leq N \leq 4$) and p_t, m_t are the
78 possible positions of the player and the Minotaur respectively after taking the action a from the state
79 (p_{t-1}, m_{t-1}) .
80 M is equal to the number of the Minotaur's possible moves that minimize the Manhattan distance
81 between him and the player (i.e. the moves that take him to the player) ($1 \leq M \leq 2$).
82 M is equal to 1 when the player and the Minotaur are in the same column or row at instant $t - 1$.
83 M is equal to 2 when the player and the Minotaur are in a different column and a different row at
84 instant $t - 1$

85 1.5 Q-Learning and Sarsa [Bonus questions]

86 (i)

87 1.5.1

88 The pseudo-code is given in Algorithm 1 :

89 1.5.2

90 We plot figures (See Figures 20) of the value function for the initial states using dif-
91 ferent initial values for the Q-function. It is observed that our Q-function converges
92 to -50. Therefore, starting with an initial value of the Q-function that is close to -50
93 accelerates the convergence process.

94 1.5.3

95 When the exploration parameter is held constant, the algorithm converges more
96 rapidly as the value of alpha approaches 0.5, in contrast to a slower convergence
97 observed as alpha approaches 1 as shown in figure 21

Algorithm 1 Q-learning

```
1: procedure Q_LEARNING(env, number_episodes, initial_state, epsilon, gamma, alpha)
2:    $N \leftarrow \text{number\_episodes}$ 
3:    $n\_states \leftarrow \text{env.n\_states}$ 
4:    $n\_actions \leftarrow \text{env.n\_actions}$ 
5:    $Q \leftarrow \text{initialize\_Q}(n\_states, n\_actions)$  ▷ Initialize Q-values
6:    $\text{visits\_state\_action} \leftarrow \text{zeros}(n\_states, n\_actions)$  ▷ Initialize visit counts
7:    $V \leftarrow []$ 
8:   for  $k \leftarrow 1$  to  $N$  do
9:      $\text{state} \leftarrow \text{initial\_state}$ 
10:     $\text{cond} \leftarrow \text{False}$ 
11:    while not  $\text{cond}$  do
12:       $\text{action} \leftarrow \text{epsilon\_greedy}(Q, \epsilon, \text{state})$ 
13:       $\text{visits\_state\_action}[\text{state}, \text{action}] \leftarrow \text{visits\_state\_action}[\text{state}, \text{action}] + 1$ 
14:       $\text{next\_state}, \_ \leftarrow \text{env.move}(\text{state}, \text{action})$ 
15:       $\text{reward} \leftarrow \text{env.reward}(\text{state}, \text{action}, \text{next\_state})$ 
16:       $Q[\text{state}, \text{action}] \leftarrow Q[\text{state}, \text{action}] + \frac{1}{\text{visits\_state\_action}[\text{state}, \text{action}]^\alpha} \cdot (\text{reward} + \gamma \cdot \max(Q[\text{next\_state}]) - Q[\text{state}, \text{action}])$ 
17:       $\text{state} \leftarrow \text{next\_state}$ 
18:      if  $(\text{env.maze}[\text{env.states}[\text{next\_state}][0]] == 2 \text{ and } \text{env.states}[\text{next\_state}][2]) \text{ or } \text{env.states}[\text{next\_state}][1] == \text{env.states}[\text{next\_state}][0]$  then
19:         $\text{cond} \leftarrow \text{True}$ 
20:      end if
21:    end while
22:     $V.\text{append}(\max(Q[\text{initial\_state}]))$ 
23:  end for
24:   $\text{policy} \leftarrow \text{argmax}(Q, 1)$ 
25:  return  $V, \text{policy}$ 
26: end procedure
```

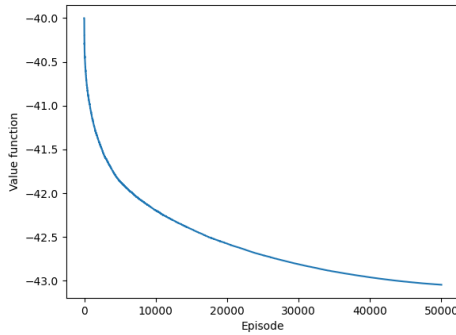


Figure 18: Q-values initialized to -40

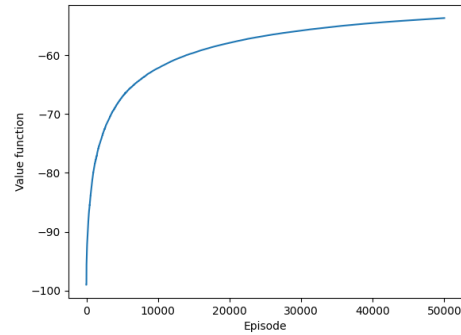


Figure 19: Q-values initialized to -100

Figure 20: Value Function Of The Initial State

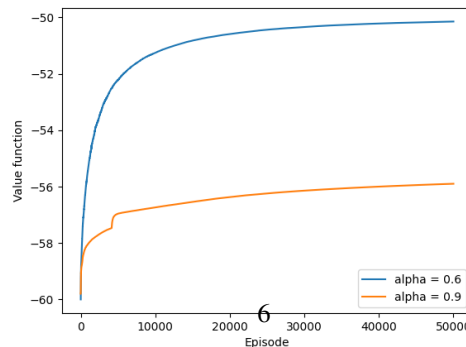


Figure 21: Value Function with Step Sizes $1/n(s, a)^\alpha$, $\alpha = 0.6$ and $\alpha = 0.9$

98 (j)

99 1.5.4

100 The pseudo-code of the SARSA algorithm is given in Algorithm 2. The main
 101 difference between Q-learning and SARSA is that SARSA uses an additional action.
 102 This is necessary to obtain an estimate of the state-action value function. Thanks to
 103 the policy improvement theorem, SARSA can approximate the optimal policy when
 104 epsilon is very small (epsilon decreasing over time towards 0). On the other hand,
 105 Q-learning attempts to learn the Q function directly, and ultimately converges to the
 106 true optimal Q function.

Algorithm 2 SARSA

```

1: function SARSA(env, number_episodes, initial_state,  $\epsilon$ ,  $\gamma$ ,  $\alpha$ ,  $\delta$ )
2:    $N \leftarrow \text{number\_episodes}$ 
3:    $n\_states \leftarrow \text{env.n\_states}$ 
4:    $n\_actions \leftarrow \text{env.n\_actions}$ 
5:    $Q \leftarrow Q_{init}$ 
6:    $visits\_state\_action \leftarrow \text{np.zeros}((n\_states, n\_actions))$ 
7:    $V \leftarrow []$ 
8:   for  $k \leftarrow 0$  to  $N$  do
9:      $state \leftarrow \text{initial\_state}$ 
10:     $cond \leftarrow \text{False}$ 
11:     $\epsilon_k \leftarrow \frac{1}{(k+1)^\delta}$ 
12:    while not  $cond$  do
13:       $action \leftarrow \text{epsilon\_greedy}(Q, \epsilon_k, state)$ 
14:       $visits\_state\_action[state, action] \leftarrow visits\_state\_action[state, action] + 1$ 
15:       $next\_state, \_ \leftarrow \text{env.move}(state, action)$ 
16:       $next\_action \leftarrow \text{epsilon\_greedy}(Q, \epsilon_k, next\_state)$ 
17:       $reward \leftarrow \text{env.reward}(state, action, next\_state)$ 
18:       $Q[state, action] \leftarrow Q[state, action] + \frac{1}{(visits\_state\_action[state, action])^\alpha} \times$ 
         $(reward + \gamma \times Q[next\_state, next\_action] - Q[state, action])$ 
19:       $state \leftarrow next\_state$ 
20:      if ( $\text{env.maze}[\text{env.states}[next\_state][0]] == 2$  and  $\text{env.states}[next\_state][2] == 0$ ) or
         $\text{env.states}[next\_state][1] == \text{env.states}[next\_state][0]$  then
21:         $cond \leftarrow \text{True}$ 
22:      end if
23:    end while
24:     $V.append(\text{np.max}(Q[\text{initial\_state}]))$ 
25:  end for
26:   $policy \leftarrow \text{np.argmax}(Q, 1)$ 
27:  return  $V, policy$ 
28: end function

```

107 1.5.5

108 For $\epsilon = 0.1$ and $\epsilon = 0.2$ and a step size of $1/n(s, a)$, with $\alpha = 2/3$. We obtain the
 109 figures (22) and (23)

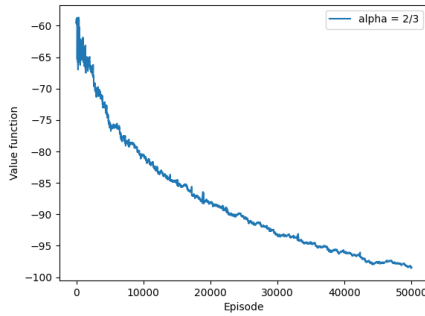


Figure 22: Value function of the SARSA algorithm for $\varepsilon = 0.1$ wto. episodes

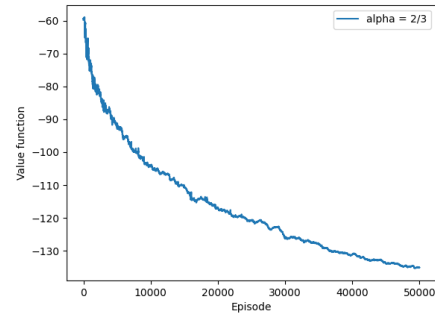


Figure 23: Value function of the SARSA algorithm for $\varepsilon = 0.2$ wto. episodes

SARSA will converge faster using an epsilon equal to 0.2 for exploration compared to epsilon equal to 0.1. This can be explained by the fact that setting ε to a relatively high value leads to discover the optimal policy faster in that case.

Proper initialization of Q-values in SARSA can contribute to faster convergence. Initializing Q-values appropriately sets the initial estimates for the expected cumulative rewards, influencing the exploration-exploitation trade-off. Careful initialization can help balance exploration and exploitation, potentially leading to quicker convergence to optimal policies. Initializing with small random values or prior knowledge can be effective, but it's crucial to consider the specific characteristics of the environment and the learning task.

1.6

For $\varepsilon = 0.1$ and $\varepsilon = 0.2$ and a step size of $1/n(s, a)$, with $\alpha = 2/3$. We obtain the figures (22) and (23)

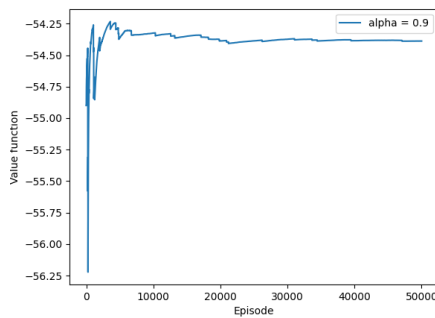


Figure 24: Value function of the SARSA algorithm for $\alpha = 0.9$ and $\delta = 0.6$ wto. episodes

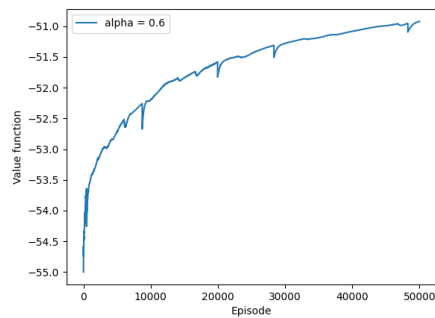


Figure 25: Value function of the SARSA algorithm for $\alpha = 0.6$ and $\delta = 0.8$ wto. episodes

A larger α implies a smaller step size and consequently a slower update.

A larger δ implies a smaller ε_k and consequently a smaller exploration.

$\alpha > \delta$ means we're going for more exploration and a slower update whereas $\alpha < \delta$

126 means we're going for less exploration and a faster update.
127 In our particular case, we witnessed a better performance when $\alpha < \delta$ i.e. when we
128 explored less and updated faster.

129 2 RL with linear approximators : MountainCar

130 2.1 Question c :

131 We have trained our model for 200 episodes using the following hyperparameter
132 values :

- 133 • $\lambda = 0.5$
- 134 • $\alpha = 0.01$
- 135 • $\epsilon = 0.1$
- 136 • $m = 0.01$

137 SARSA(λ) with Fourier basis and SGD with momentum equal to m for the Moun-
138 tainCar problem involves approximating the state-action value function using Fourier
139 basis functions (We have used $\eta_i = [k, l] \quad \forall k, l \in \{0, \dots, 2\}$, except $[0, 0]$). The algo-
140 rithm initializes weights of the linear function approximator and eligibility traces. It
141 iteratively selects actions (ϵ -greedy selection with respect to the current estimation of
142 Q_w), updates eligibility traces, and employs SGD with momentum to adjust weights
143 based on temporal difference errors. The exploration strategy guides action selection,
144 and the algorithm repeats this process for a number of episodes equal to 200. The
145 combination of function approximation, eligibility traces, and momentum-based
146 optimization allows SARSA() to efficiently learn and update the state-action value
147 function, balancing exploration and exploitation to solve the MountainCar problem.
148 Fine-tuning hyperparameters was crucial for achieving optimal performance.

149 2.2 Question d :

150 2.2.1

151 Here is the figure 26 showing the convergence of our algorithm. We observe an
152 initial phase during which the average total reward is -200, indicating scenarios
153 where the car fails to reach the flag. Once the car successfully reaches the flag (upon
154 exploration), the actual learning process begins. Subsequently, the learning algorithm
155 converges towards the optimal state-action value function, gradually improving the
156 average performance.

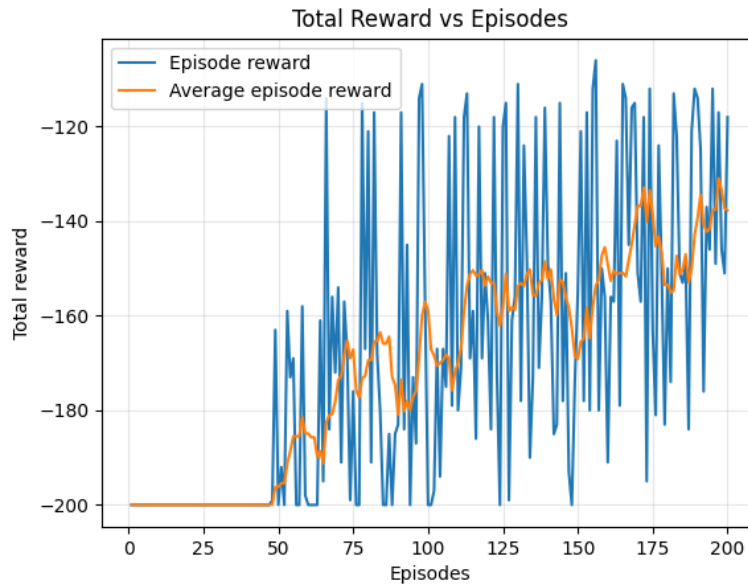


Figure 26: Average episode reward converging to ~ -135

157 2.2.2

158 Here is the 3D plot in figure 2.2.2 of the value function of the optimal policy over
 159 the state space domain :

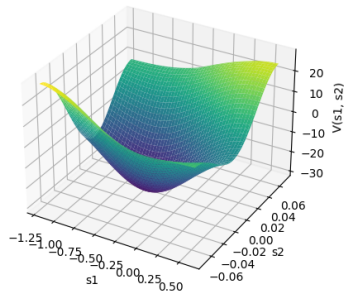


Figure 27: 3D plot of the optimal policy value function

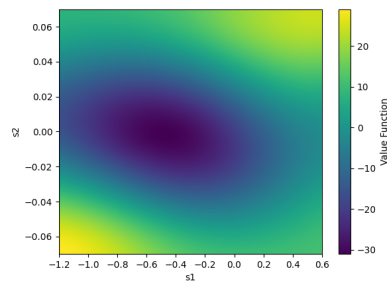


Figure 28: Temperature plot of the value function of the optimal policy

160 The plots 2.2.2 reveal that the value function, under the optimal policy that was
 161 identified, exhibits a global minimum around the state where the position is -0.4,
 162 and the velocity is equal to 0. This state corresponds to a location close to the initial
 163 state used during training and symbolizes the bottom of the hill. Consequently, it is
 164 evident that the value function for this state is very low, indicating that it requires
 165 numerous steps to reach the flag. As we increase both the position and the velocity,
 166 the value function starts to rise until we reach the position of the flag, where it attains
 167 its maximum value.

168 **2.3**

169 Here is the 3D plot in figure 29 and 30 of the optimal policy :

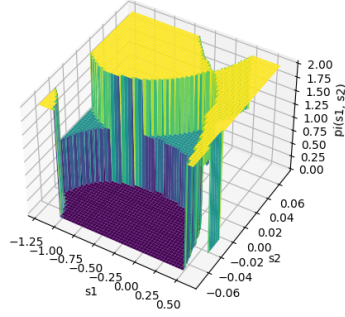


Figure 29: 3D plot of the optimal policy

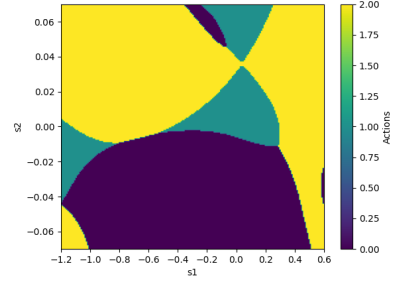


Figure 30: Temperature plot of the optimal policy

170 In examining these figures, a notable trend is apparent: generally, when the velocity
171 is negative, the predominant action is moving left. Conversely, the usual action is
172 to move right when the velocity is positive (though not always, as indicated in the
173 plots), with occasional instances of choosing not to move.

174 **2.4**

175 Figure 31 shows the evolution of both the episode reward and average episode reward
176 when adding $\eta = [0, 0]$ to our Fourier basis.

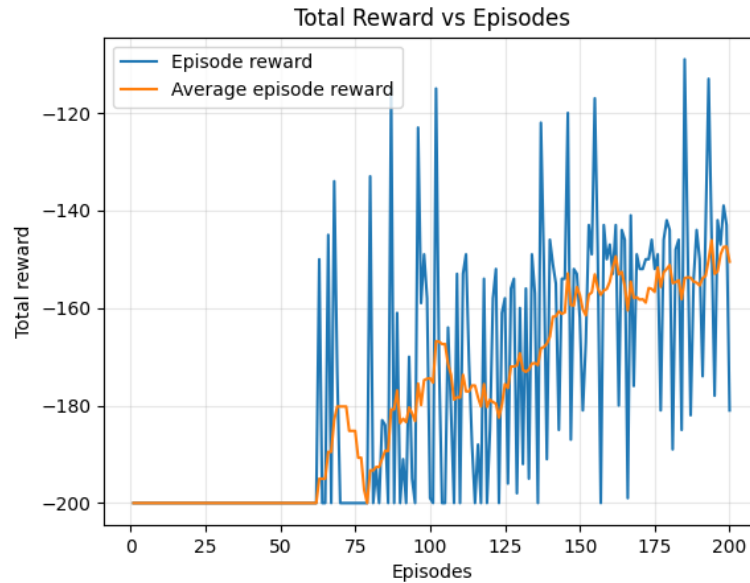


Figure 31: Average episode reward in the case of using $\eta = [0, 0]$

177 Compared to (26) where $[0, 0]$ is excluded from the basis, the average episodic
 178 reward plots are fairly similar.

179 The following plots 2.4 and ?? show the optimal policy and the value function in
 180 that case :

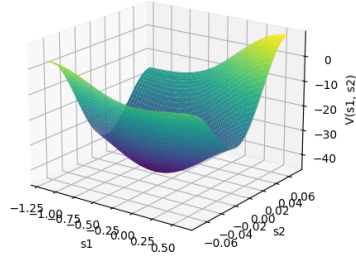


Figure 32: 3D plot of the optimal policy value function

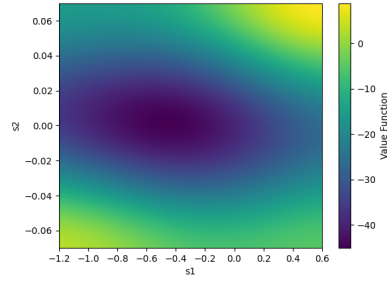


Figure 33: Temperature plot of the value function of the optimal policy

181 Same as before, the above figures are visually similar to the ones obtained where
 182 $[0, 0]$ (27 and 28). The quantitative difference is the order of magnitude of the value
 183 function that appears higher in the former case.

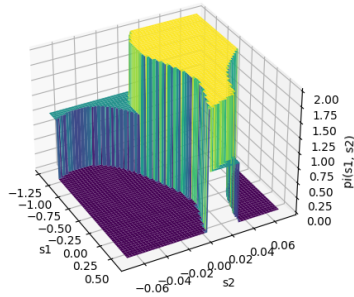


Figure 34: 3D plot of the optimal policy

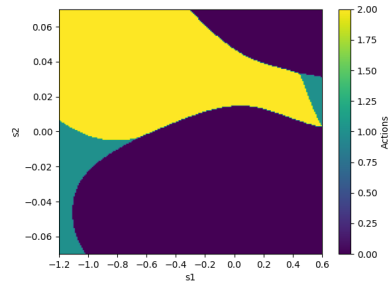


Figure 35: Temperature plot of the optimal policy

184 In the case of the optimal policy, there is a notable difference between the above
 185 figures and (29) & (30), namely that the optimal policy towards when s_1 is closer
 186 to 0.6 is to almost certainly *push right* in (30) whereas it is not necessarily the case
 187 in (35) depending on the speed s_2 . This might be due to the fact that including the
 188 constant term in the basis allows the function approximator to capture the mean
 189 or offset of the value function. This can be important if your value function has a
 190 non-zero mean or if there is a constant component that needs to be represented.

2.5 Question e :

2.5.1

The figure 36 shows the average total reward of the policy as a function of α , whereas figure 37 shows the average total reward of the policy as a function of λ :

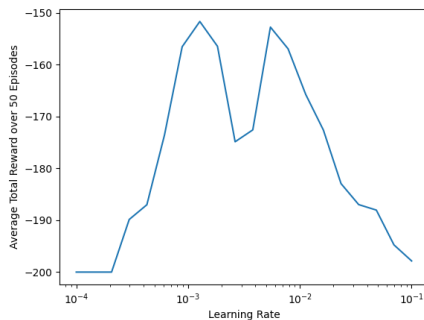


Figure 36: Average reward with respect to α

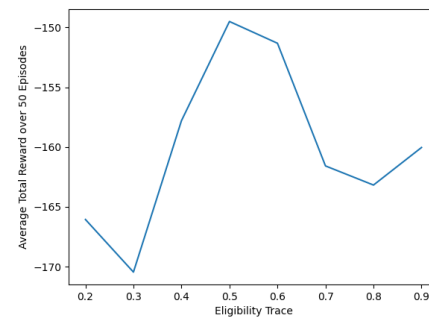


Figure 37: Average reward with respect to λ

Both the figures above have a quasi-bell shape, (37) more than (36).

When λ is close to 0, the algorithm considers only the immediate rewards and neglects the influence of future rewards. This leads to a myopic behavior, where the agent focuses on short-term gains.

As λ increases, the eligibility trace incorporates information about more distant future rewards, allowing the agent to make decisions that consider a longer time horizon.

If the learning rate is excessively low, it may cause the algorithm to converge slowly or become trapped in suboptimal solutions.

Conversely, an excessively high learning rate can induce oscillations or divergence, preventing the algorithm from converging to an optimal policy.

2.5.2

Here are some strategies for initializing Q-values:

Zero Initialization: Initialize all Q-values to zero. This is a straightforward and commonly used approach. However, it may lead to slow convergence or convergence to suboptimal solutions, especially if the true Q-values have a non-zero mean.

Random Initialization: Initialize Q-values with small random values. This introduces exploration and helps break symmetries. However, the impact of randomness may make it challenging to reproduce results.

Optimistic Initialization: Set Q-values to a high constant value. This encourages exploration in the early stages of learning, as the agent believes all actions have high potential rewards. Over time, Q-values are updated based on the observed rewards.

217 **Prior Knowledge Initialization:** If there is prior knowledge about the problem,
218 initialize Q-values based on that knowledge. This can provide a head start and
219 accelerate learning. However, it requires domain expertise.

220 In my opinion, a combination of optimistic initialization and exploration strategies
221 (such as epsilon-greedy) often works well. Optimistic initialization provides a bias
222 towards exploration in the early stages, and exploration strategies ensure ongoing
223 exploration throughout the learning process.

224 2.5.3

225 Our exploration strategy is a decreasing ϵ . At the beginning of the training, our ϵ
226 is large enough to allow exploration of all possible actions independently on the
227 q-value. ϵ decreases with respect to the number of episodes as the training allows
228 the algorithm to converge.

229 After implementing this strategy, our average episodic reward (38) increases faster
230 compared to a fixed value of epsilon (26)



Figure 38: Average reward with a decreasing epsilon wto. episodes