# Deep SARSA and Deep Q-Learning

Yehia Ragab

## I. INTRODUCTION

The link to the GitHub Repository is here.

Reinforcement Learning is a sub-field of Machine Learning where an agent learns to make optimal decisions by interacting with a defined environment, and observing the feedback, in the form of rewards and punishments, and using that feedback to improve future actions that maximise rewards and minimise punishments. Deep Reinforcement Learning represents a significant breakthrough in the field of Reinforcement Learning, as it leverages Neural Networks to enhance traditional Reinforcement Learning. This report focuses on Deep SARSA (State-Action-Reward-State-Action) and Deep Q-Learning. We will implement and test a Deep Q-Learning and a Deep SARSA agent and compare their algorithmic implementations and their performances in a controlled environment.

In the realm of Reinforcement Learning, Chess represents a huge challenge as it is a very complex game with many rules and permutations that it is impossible to solve through brute force. However, Chess endgames with 7 pieces or less have already been solved with optimal play, and the tablebase, called Szygy Endgames tablebase, is publicly available. For this assignment, we were tasked with solving a chess endgame in a 4x4 board involving 3 pieces, the 2 kings and a queen. The goal of this assignment is to analyse the differences in performance and time to convergence of Deep SARSA and Deep Q-Learning.

In this report, we will explain what Deep Q-Learning and Deep SARSA are, explain our methodology for our implementation, and present and analyse our results.

## II. METHODOLOGY

### A. Deep SARSA and Deep Q-Learning

Deep SARSA and Deep Q-Learning are deep reinforcement learning algorithms that are very similar in concept. We will first discuss the similarities, then we will explain the differences.

*1) Similarities:* Both algorithms involve initialising the weights and biases of the neural network with random numbers between 0 and 1. The dimensions of the weight of the neural network are as follows: $(i, h)$, $(h, o)$, where $i$ is the dimension of the input, in our case, the features $X$, $h$ is the number of hidden nodes, and $o$ the dimensions of the output, in our case the number of possible actions for the agent king and queen.

The environment is initialised for each episode, which is defined by making the agent play the game once. The initial state of the environment is defined as $s$. We run forward propagation on our neural network with $s$ as our input. The output of the neural network is an array of **Q-values**, which are estimations of future rewards for each corresponding action. Using our chosen policy, which in our case is $\epsilon$-**greedy**, we pick the first action $a$. $\epsilon$ defines the percentage of times the algorithm will pick the best possible action based on the **Q-values**, compared to the percentage of times the algorithm will pick a random action based on all possible actions. The agent performs that action, and observes $R$, the reward, and $s'$, the new state. We define the **target Q-values** as a copy of the **Q-values** that were output by the neural network. If the $s'$ state is an end-state, meaning either the agent won, or drew, we set the $a$-**th** value of target **Q-values** to the reward $R$ observed. Otherwise, we run forward propagation on $s'$ and select the optimal action $a'$ based on a policy, which is different for Q-Learning and SARSA. We get the predicted Q-Value of $a'$, defined as **Q-value**[$a'$] and set the $a$-**th** value of **target Q-values** to $R + \gamma$ **Q-value**[$a'$]. We then compute the **error**, defined as **target Q-values - Q-values**. We then run backward propagation and update the weights and biases of the neural network. We repeat the process until the game reaches an endstate, with $a'$ becoming $a$ and $s'$ becoming $s$.

*2) Exploration vs Exploitation:* The difference between Q-Learning and SARSA, as explained before, is in the policy used to pick the best possible action for all states other than the initial state $s$. For SARSA, we use an $\epsilon$-**greedy** policy to choose the action we will use when computing the error, which means that we don't always use the action that maximises the **Q-values** that were output by the neural network for backpropagation. For Q-Learning, after making the initial action $a$, we always pick the optimal action corresponding to the highest **Q-value** for backpropagation, meaning we use a greedy policy. Deep SARSA is said to be an on-policy method, meaning it uses the chosen initial policy to dictate the future action for each state in an episode, while Q-Learning is said to be an off-policy algorithm, meaning that after making the initial action $a$, the algorithm picks the best possible action based on a greedy policy.

The main advantage of Q-learning is that it has the potential to converge more quickly than SARSA, as it explores less than SARSA. However, the disadvantage is that it may converge to a local optimum instead of a global optimum.

## III. IMPLEMENTATION DETAILS

### A. Neural Network

Our neural network used the suggested dimensions, that is to say there is 1 hidden layer with 200 hidden nodes.

## B. Forward propagation

For our forward propagation function, we used ReLU activation. We computed the dot product of $X$, our features, and $W1$, the weights of the first layer, and added the biases $b1$. We then applied the ReLU activation function **h_activation**, then computed the **Q-values** by getting the dot product of **h_activation** and $W2$, then added $b2$.

## C. Backward propagation

For our forward propagation function, we used $X$, the **error** obtained, and **h_activation**. We then computed the outer product of **h_activation** and **delta2**, (the error) , which yield **dW2**. We updated the **W2** and **b2**., then applied the inverse of the ReLU function, which involves taking the heaviside of **h_activation** and multiplying it by the dot product of **W2** and **delta2**, which yields **delta1**. We then compute the outer product of **X** and **delta1** to get **dW1**. We then update the .

## D. Rotating the board

In chess, there is always at least 1 square between both kings, both horizontally and vertically. With our chess board being 4x4, that means that if we divide the board into 4 2x2 squares, the kings will always be on different quarters. Furthermore, all our pieces are able to move in all directions in equal measure, so the orientation of the board doesn't change anything. The optimal move in a board with the state $s$ is the mirror move of the optimal move of a board with the state $s_{mirror}$. The same logic applies to a rotated or transposed board. For this reason, we implemented as function that ensures the agent's king is always in the top left quarter of the board by rotating the board until this requirement is satisfied. Since the enemy king can't be in the same quarter, the king will be in one of the remaining quarters. We can ensure that enemy king will be on one of the quarters on the right (excluding the coordinate (3,2)) by transposing the board if it isn't already the case. This makes it so that the boards are all normalised and the number of possible states drastically decreases. This rotating function is applied after every move is made 8.

## E. Features

After normalising the board, there are 4 possible spots for the agent's king, and 7 possible spots for the enemy king, which means we can reduce the number of features from the original 58 to 37, with the first 4 numbers indicating the position of the agent's king, the next 16 numbers indicating the position of the agent's queen, and the next 7 numbers indicating the position of the enemy king. We removed the last 10 numbers which indicated whether the enemy king was in check and the allowed moves the enemy kings can make.

## F. Optimizer

For our optimizer, which is used to update the weights and biases of the neural network, we used a modified version of the Adam optimizer 9 that was in the lab, with $beta1 = 0.9$, $beta2 = 0.999$, and $\epsilon = 1 \times 10^{-8}$.

## G. Code organisation

The Agent is defined in a class called Agent. The class is initialised with all the necessary parameters 10. This class works for both the Q-Learning implementation and the SARSA implementation. The class has the following functions:

- **get_action**: This function takes the state, the policy and allowed_a as inputs and returns an action that depends on the policy and all the possible actions. The functions returns the chosen action, the Q-values, of the state, and h_activation 11.
- **forward_prop**: This function takes the state as input and returns the Q-values and h_activation 11.
- **backward_prop**: This function rakes in the state, the error, and h_activation as input and updates the weights 11.
- **episode**: This function runs an episode of the training loop and returns the reward and number of steps of the episode. It also prints the average reward and average number of steps over the previous 100 episodes which happens every 100 episodes 12.
- **run_training_loop**: This function takes in the number of episodes one would like to run the training loop for and runs the training loop. It also returns the average reward and average number of steps of the training loop.
- **testing_episode**: This function runs 1 testing episode, which is an episode where epsilon is set to 0 and there is no backpropagation.
- **test**: This function takes in the number of episodes one would like to test, and prints the average Reward and number of steps of the trained agent.

## H. Parameters

*1) $\epsilon_f$ and $\beta$:* $\epsilon_f$ is the value of $\epsilon$ that is chosen for the $\epsilon$-greedy policy. It is set to decay for every episode that is run so that over time, the agent reduces the amount of exploration it does in favour of exploitation. The formula for $\epsilon_f$ is:

$$\epsilon_f = \frac{\epsilon}{1 + (\beta * n)} \tag{1}$$

where $\epsilon$ is the initial value chosen, $\beta$ being the decaying rate, and $n$ being the number of the episode. Increasing $\beta$ would reduce the amount of exploration the agent does over time, as $\epsilon_f$ would decay quicker. Although it might converge quicker, it also make the agent more prone to being stuck in a local optimum. Decreasing $\beta$ would increase the amount of exploration compared to smaller $\beta$ over time which decreases the likelihood of being stuck in a local optimum, but it also increases the amount of episodes needed to find the global optimum, which is the goal of the agent.

*2) $\gamma$:* $\gamma$, also known as the discount factor, plays a crucial role in the Bellman equation, which is used to calculate the error. The Bellman equation used in the code if the environment isn't in an endstate is as follows:

$$TargetQvalue[action] = R + \gamma \times Qvalues_{next}[action_{next}] \tag{2}$$

A small $\gamma$ makes the agent more concerned with immediate rewards compared to future rewards, while a big $\gamma$ makes the agent more concerned with future rewards compared to immediate rewards. In many cases, a bigger $\gamma$ tends to increase variance and make the learning process unstable.

### I. Plots

All plots were created using matplotlib. For each plot, 3 different agents with the same parameters were trained separately, then final rewards were the average of the reward at each episode, and the exponential moving average function from pandas was used on the averaged rewards. The alpha chosen was 0.0025, which is the value that struck the best balance between reducing noise and showing the trend. All agents were trained on 10,000 episodes, and in every plot, the blue line represents the Q-Learning agent and the green line represents the SARSA agent.

## IV. RESULTS AND DISCUSSION

### A. Initial Implementation

The initial implementation for this assignment was Deep Q-Learning and Deep SARSA with the suggested parameters which are:

- $hidden\_layers = 1$
- $hidden\_nodes = 200$
- $\epsilon_0 = 0.2$
- $\beta = 0.00005$
- $\gamma = 0.85$
- $learning\_rate = 0.0035$

As we can see from the plots 1, the performance and convergence of both agents is very similar. The main difference here is that the SARSA agent makes fewer moves per game on average at the start.

### B. Effects of changing $\beta$ and $\gamma$

We tested the agent with 3 different values of $\beta$ and 3 different values of $\gamma$. The values of $\gamma$ we tested were 0.75, 0.85 (base), and 0.95. The values of $\beta$ we tested were 0.00005 (base), 0.0001, and 0.00025.

*1) $\beta$:* The features implemented decrease the number of episodes needed for convergence, which makes the default $\beta$ value a bottleneck in our implementation. For this reason, we decided to only test higher $\beta$ values. For both plots 2 and 3, the agent converges quicker than the agents with the base parameters. The average number of moves on average is lower for SARSA in both cases.

*2) $\gamma$:* The average reward for the Q-Learning agent is higher than the average reward for the SARSA agent for the first 5,000 episodes when we lower $\gamma$ to 0.75, as seen in 4. This can be because with a lower $\gamma$, the SARSA agent has fewer opportunities to update the Q-Values for the optimal action, meaning it converges slower.

When we increase $\gamma$ to 0.95, the inverse happens. The average reward for the SARSA agent is higher than the average reward for the SARSA agent for the first 5,000 episodes, as

seen in 5. The average number of moves is also lower for the SARSA agent.

However, in both cases, the difference isn't considerable when comparing them to the agents trained with the base parameters and when comparing both agents with each other.

### C. Additional Experiments

*1) Original Features vs Rotated board:* The plots 6 show that the altered features perform significantly better than the original features. Each episode also involves less moves. Additionally, the altered features run in half the time of the original features due to the considerably smaller input vector (27 vs 58). These results were expected because the agent has fewer permutations to learn from, and is therefore able to find the optimal solution faster.

*2) Finding a better agent:* We decided to experiment further to try and find a better performing agent. To do so, we tried to increase $\beta$ to 0.0005, the learning rate to 0.005, and decreased $\epsilon$ to 0.175. We were able to get the best performing agent both in terms of convergence speed and average reward 7.

## V. CONCLUSION

In this environment, the performance of Deep Q-Learning and Deep SARSA is very similar. This might be due to the simplicity of the task at hand. We were able to train agents that converge in a small number of episodes and reach the optimal reward. We looked at the effect of changing certain parameters and changing the shape of the environment and input features.
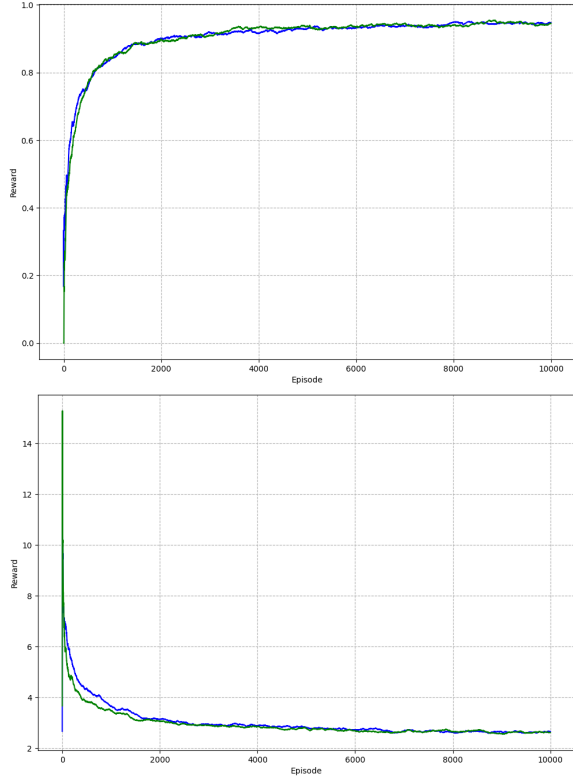
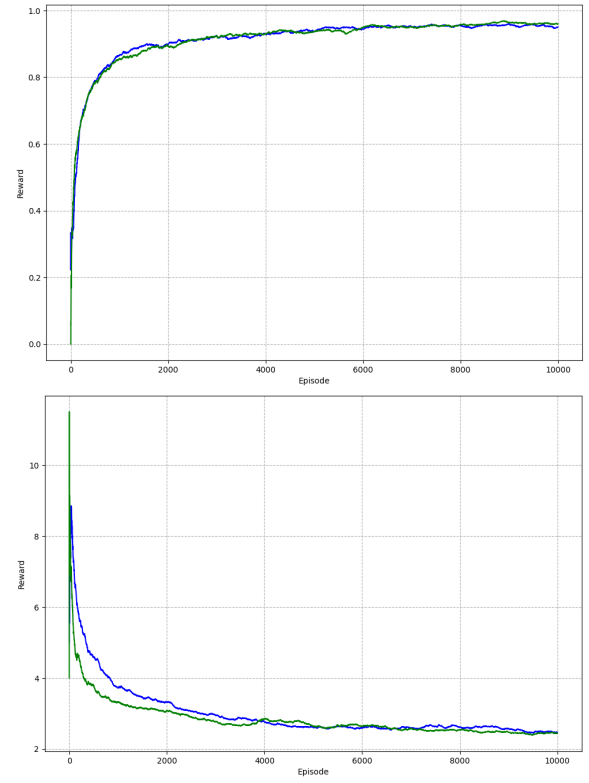Fig. 1. Exponential Moving Average Reward and Number of steps for the agents with the base parameters.



Fig. 2. Exponential Moving Average Reward and Number of steps for the agents with the base parameters with $\beta = 0.0001$.
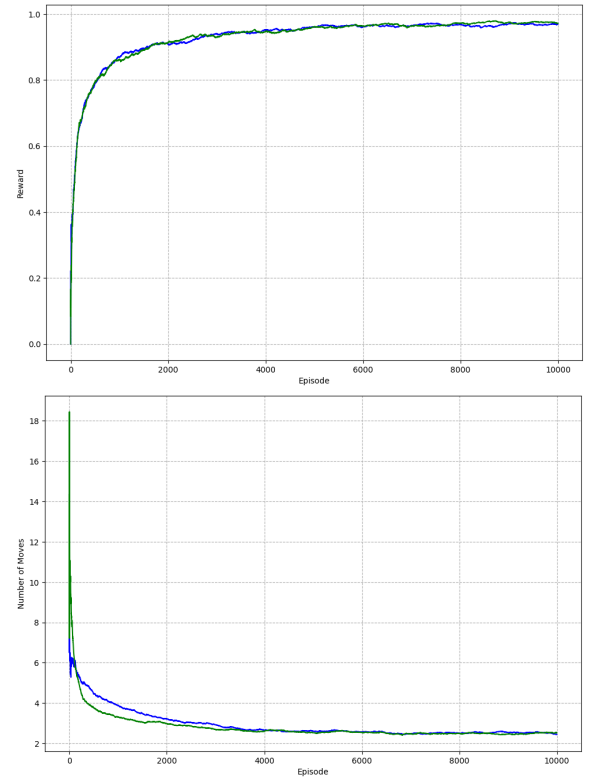


Fig. 3. Exponential Moving Average Reward and Number of steps for the agents with the base parameters with $\beta = 0.00025$.
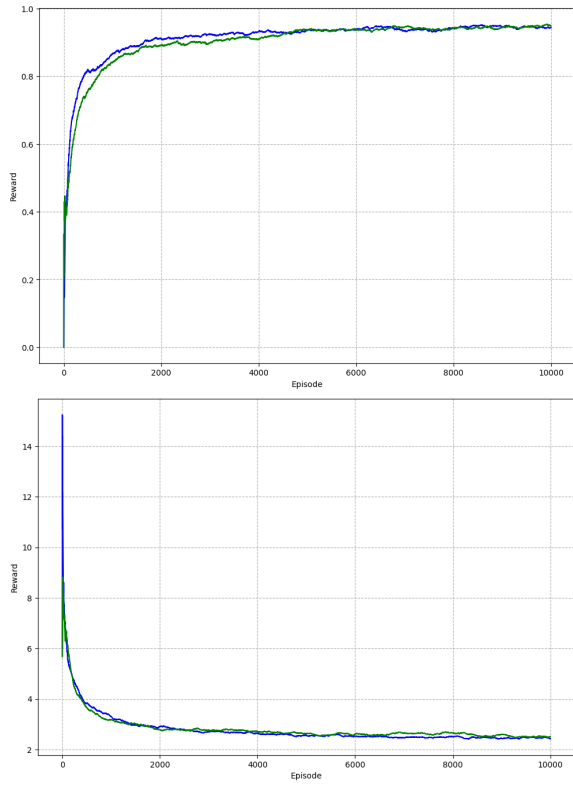
Fig. 4. Exponential Moving Average Reward and Number of steps for the agents with the base parameters with $\gamma = 0.75$.
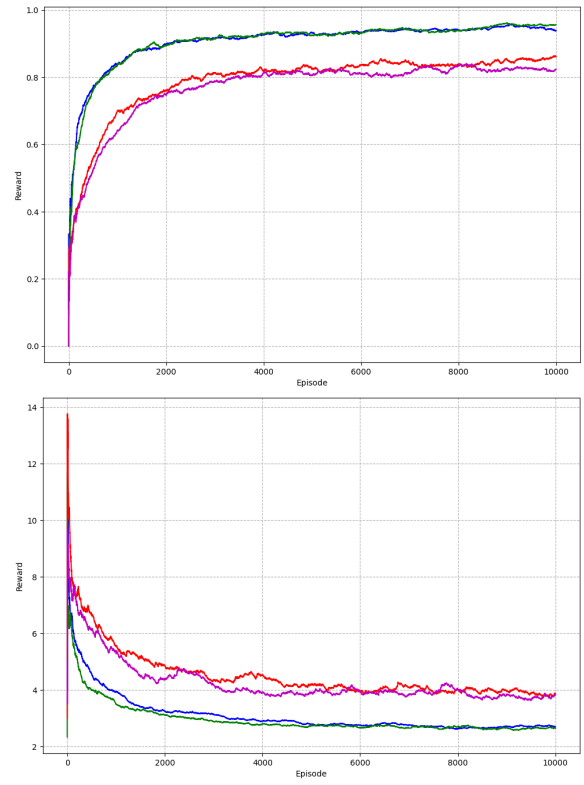


Fig. 6. Comparison between original features and altered features with rotated board. Red represents the Q-Learning Agent with the original features and magenta represents the SARSA agent with the original features.
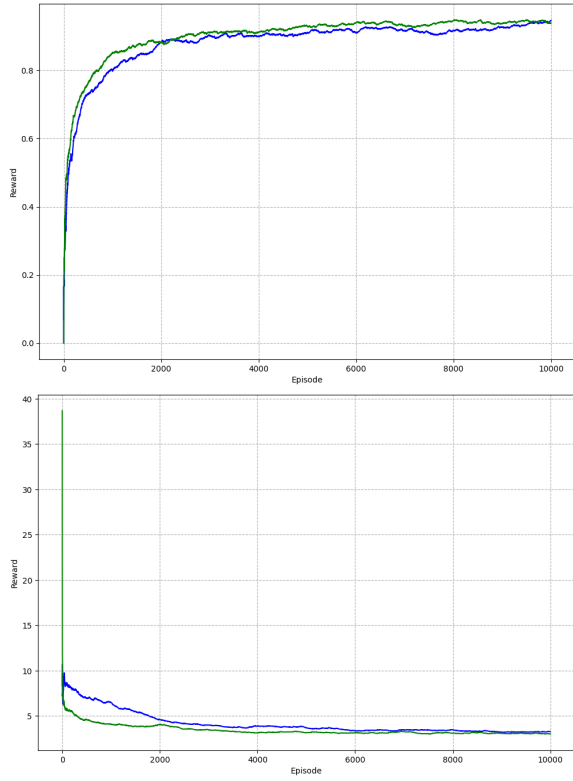


Fig. 5. Exponential Moving Average Reward and Number of steps for the agents with the base parameters with $\gamma = 0.95$.
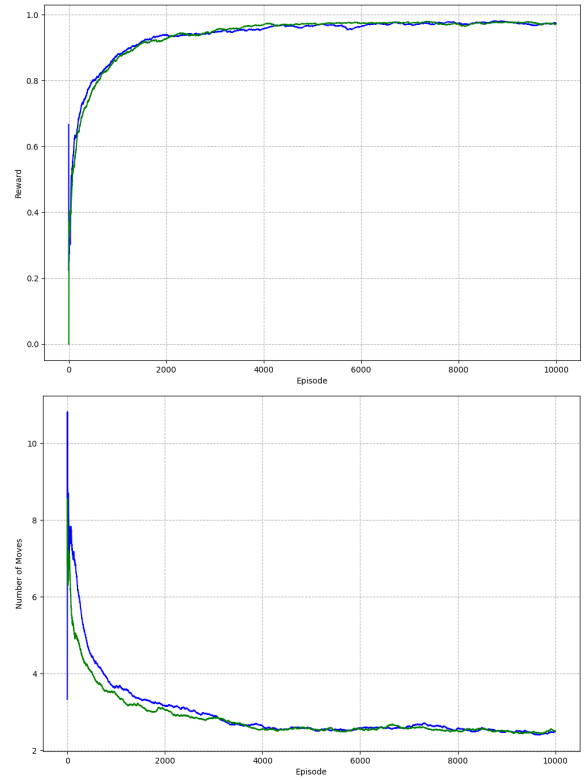


Fig. 7. The best agent found.

```python
215    def rotate_board(self, Board, p_k1, p_k2, p_q1, k2=False):
216 >       """ ...
241
242        if not k2:
243            if p_k1[1] >= 2:
244                # k1 isn't in c1
245                if p_k1[0] < 2:
246                    # k1 is in c2
247                    Board = np.rot90(Board, 1)
248                else:
249                    # k1 is in c4
250                    Board = np.rot90(Board, 2)
251            elif p_k1[0] >= 2:
252                # k1 is in c3
253                Board = np.rot90(Board, 3)
254
255            p_k2 = np.array(np.where(Board==3)).reshape(2,)
256
257            if p_k2[1] < 2 or (p_k2[1] == 2 and p_k2[0] == 3):
258                # k2 isn't in h1
259                Board = Board.T
260                p_k2 = np.array(np.where(Board==3)).reshape(2,)
261
262            p_k1 = np.array(np.where(Board==1)).reshape(2,)
263            p_q1 = np.array(np.where(Board==2)).reshape(2,)
264
265            return Board, p_k1, p_k2, p_q1
```

Fig. 8. The function that rotates the board.

```python
1   class Adam:
2
3       def __init__(self, params, beta1=0.9, beta2=0.999, epsilon=1e-8):
4
5           self.params = params
6           self.beta1 = beta1
7           self.beta2 = beta2
8           self.epsilon = epsilon
9
10          # Initialize the momentums.
11          self.mt = [np.zeros_like(p) for p in params]
12          self.vt = [np.zeros_like(p) for p in params]
13
14          # Initialize the counter.
15          self.counter = 0
16
17
18      def Compute(self, grads):
19
20          self.counter += 1
21
22          # Update the momentums.
23          self.mt = [self.beta1 * mt + (1 - self.beta1) * g for mt, g in zip(self.mt, grads)]
24
25          # Update the variances.
26          self.vt = [self.beta2 * vt + (1 - self.beta2) * g**2 for vt, g in zip(self.vt, grads)]
27
28          # Bias-corrected momentums.
29          mt_hat = [mt / (1 - self.beta1**self.counter) for mt in self.mt]
30          vt_hat = [vt / (1 - self.beta2**self.counter) for vt in self.vt]
31
32          # Compute the Adam updates.
33          new_gradients = [mth / (np.sqrt(vth) + self.epsilon) for mth, vth in zip(mt_hat, vt_hat)]
34
35          return new_gradients
```

Fig. 9. The implementation of Adam.

```python
1   class Agent:
2
3       def __init__(self, input_dim, hidden_dim, output_dim, learning_rate, epsilon, beta, gamma, env, algorithm):
4
5           self.W1 = np.random.randn(input_dim, hidden_dim) / 10
6           self.b1 = np.zeros((hidden_dim))
7
8           self.W2 = np.random.randn(hidden_dim, output_dim) / 10
9           self.b2 = np.zeros((output_dim))
10
11          self.learning_rate = learning_rate
12          self.optimizer = Adam([self.W1, self.W2, self.b1, self.b2])
13
14          self.rewards = []
15          self.n_of_moves = []
16          self.n = 0
17          self.epsilon = epsilon
18          self.beta = beta
19          self.gamma = gamma
20          self.env = env
21          self.algorithm = algorithm
22          self.alpha = 0.99
23          self.error = []
24
```

Fig. 10. Initialising an Agent.

```python
25    def get_action(self, X, epsilon_f, allowed_a):
26
27        qvalues, h_activation = self.forward_prop(X)
28        a, _ = np.where(allowed_a==1)
29
30        if np.random.rand() < epsilon_f:
31            action = np.random.permutation(a)[0]
32        else:
33            allowed_qvalues = qvalues[a]
34            max_index = np.argmax(allowed_qvalues)
35            action = a[max_index]
36
37        return action, qvalues, h_activation
38
39    def forward_prop(self, X):
40        h = np.dot(X, self.W1) + self.b1
41        # ReLU activation
42        h_activation = np.maximum(0, h)
43        qvalues = np.dot(h_activation, self.W2) + self.b2
44
45        return qvalues, h_activation
46
47    def backward_prop(self, X, delta2, h_activation):
48
49        dW2 = np.outer(h_activation, delta2)
50
51        delta1 = np.heaviside(h_activation, 0) * np.dot(self.W2, delta2)
52
53        dW1 = np.outer(X, delta1)
54
55        new_gradients = self.optimizer.Compute([dW1, dW2, delta1, delta2])
56
57        self.W1 += self.learning_rate * new_gradients[0]
58        self.W2 += self.learning_rate * new_gradients[1]
59        self.b1 += self.learning_rate * new_gradients[2]
60        self.b2 += self.learning_rate * new_gradients[3]
```

Fig. 11. get_action, froward_prop, and backward_prop.

```python
62    def episode(self):
63
64        epsilon_f = self.epsilon / (1 + self.beta * self.n)
65        policy = epsilon_f if self.algorithm=='sarsa' else 0
66        self.n += 1
67        Done = 0
68        i = 1
69        _, X, allowed_a = self.env.Initialise_game()
70
71        while Done==0:
72
73            a_agent, qvalues, h_activation = self.get_action(X, epsilon_f, allowed_a)
74
75            _,X_next,allowed_a_next,R,Done=self.env.OneStep(a_agent)
76
77            target = np.copy(qvalues)
78
79            if Done==1:
80                target[a_agent] = R
81            else:
82                a_next, qvalues_next, _ = self.get_action(X_next, policy, allowed_a_next)
83                target[a_agent] = R + self.gamma * qvalues_next[a_next]
84
85                X = np.copy(X_next)
86                allowed_a = np.copy(allowed_a_next)
87                i += 1
88
89            error = target - qvalues
90            self.error.append(error[a_agent])
91            self.backward_prop(X, error, h_activation)
92
93
94        self.rewards.append(R)
95        self.n_of_moves.append(i)
```

Fig. 12. Training episode.