

Mastering Simplified Chess Endgames with Reinforcement Learning

Source code available at https://github.com/timZurichLoaf/Reinforcement_learning_chess_endgame

Zheng Luo

Faculty of Business, Economics and Informatics

University of Zurich

Zurich, Switzerland

zheng.luo@uzh.ch | UZH ID: 21-738-901

Abstract—When we reflect on how an infant learns to play, walk and not touch a hot radiator, there is no explicit teacher. It reveals the very nature of learning process, to learn by interacting with the environment. Reinforcement learning is a computational approach mimicking the interactive learning. This study will review the two mainstream Temporal Difference (TD) learning methods, Q-learning and State-action-reward-state-action (SARSA), along with experience replay technique from a theoretical perspective. A simplified chess endgame will be then used as a sample task to examine the TD learning methods and the dominant role of goal, or the administration of rewards, in shaping the behavior of an agent. Deep reinforcement learning algorithm and the effect of hyper-parameters will also be analyzed in detail.

I. INTRODUCTION

A TD agent relies on the estimation of rewards of state-action pairs (Q values) to decide its moves. As it explores the environment by trial-and-error, the estimation is updated based on its observation. Conventional TD methods face the sequential dependency. The experience replay technique first studied by [Lin92] is a successful attempt to cope with it. In the simplified 'over-killing' endgame on a 4-by-4 chessboard, a reinforcement learning agent learns to checkmate. Despite the small chessboard, all conventional chess rules specified in [Sch03] apply. SARSA and Q-learning are used in training a neural network to approximate Q-values. As an agent plays, moves per game and rewards received in each episode are recorded to evaluate the models. The effect of hyper-parameters is explored theoretically and empirically with SARSA in light of the trade-off between exploration and exploitation. Finally, the administration of reward is modified to train an impatient agent eager to end a game as soon as possible, by punishing long games and neutralizing the outcome.

II. METHODS

A. Theoretical Review of TD methods, Q-learning and SARSAⁱ

Q-learning and SARSA are two TD learning methods, which combine the idea of Monte Carlo and that of dynamic programming (DP). TD resembles Monte Carlo method in a

way that an agent learns directly from raw experience without much a posteriori knowledge or a given model of the environment's dynamics. Similar to DP, TD updates estimations based on existing learning outcomes, meaning that neither of them waits until the final reward is available.

B. Comparison between Q-learning and SARSAⁱ

Both as TD methods, the primary difference between Q-learning and SARSA lies in the approaches to estimate Q values, namely on-policy and off-policy. SARSA as an on-policy method, estimates $q_\pi(s, a)$ of current state-action pair (s, a) by its immediate reward and the subsequent $q_\pi(s', a')$ following a given policy π . Intuitively, a SARSA agent guesses the reward of current state-action pair based on its guess of the next. Q-learning, on the other hand, as an off-policy method, estimates $q_\pi(s, a)$ by considering the optimal outcome of all possible actions $\max_a q_\pi(s', a)$, given a following state s' . Thus, Q-learning directly approximates the optimal action-value function q^* , regardless of the policy. Besides, whenever a rewarding move has been discovered, the information propagates fast to its preceding state-action pairs.

Later in this study, an over-killing chess endgame experiment shows that a Q-learning agent accumulates fewer rewards and takes more moves per game than a SARSA agent, as a rewarding move may drive it to explore some costly intermediate moves. However, in the demonstration game, the Q-learning agent adopts an optimal strategy, while SARSA sticks to a cautious time-consuming approach.

C. Introduction of Experience Replayⁱⁱ

Experience replay was first studied by [Lin92]. When the agent executes each action, a quadruple $(S_t, A_t, R_{t+1}, S_{t+1})$ is stored in the replay memory of a fixed size to accumulate experiences over time. At each time step, Q-learning updates a mini-batch by sampling uniformly from the replay memory

ⁱAnswer to Task 1: Describe Q-learning, SARSA and explain their differences.

ⁱⁱAnswer to Task 2: Describe the experience replay technique and cite relevant papers.

and then proceeds according to a given policy. The experience replay complements the conventional Q-learning by making good use of its off-policy nature, which is independent of connected trajectories. Uncorrelated successive updates contribute to a reduction of the variance. Nevertheless, [Lin92] points out that experience replay is applicable to a static environment, where the rewards almost don't change. Otherwise, past experience may become irrelevant or even harmful.

[Mni+13; Mni+15] extend the idea of [Lin92] by introducing a second network and holding the second network's weights constant for a certain number C of updates. The output of this duplicate network with temporarily constant weights serves as the Q-learning target during C updates and leads to the following expression,

$$w_{t+1} = w_t + \alpha [R_{t+1} + \gamma \max_a Q^n(S_{t+1}, a, \tilde{w}_t) - Q^n(S_t, A_t, w_t)] \nabla Q^n(S_t, A_t, w_t) \quad (1)$$

Thus experience replay brings Q-learning closer to simpler supervised learning and speeds up the learning process.

D. Implementation of SARSA and Q-learning for 'over-killing' endgames

1) *Environment setup*: In the 'over-killing' endgames on a 4-by-4 chessboard, two white pieces, a king and a queen, are controlled by the agent and a black king by a random agent. Each episode of the game is initialized in a state where the opponent king is not threatened, as demonstrated in Figure 1.

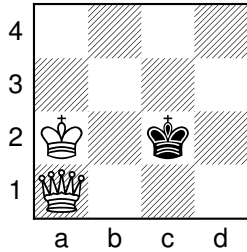


Fig. 1. An example of the initial state of an episode, where the opponent king (black) is not threatened.

2) *Artificial Neural Network*: An artificial neural network (ANN) with 1 hidden layer of 200 units is used to approximate the true Q-value functions. A default learning rate η is set to 0.0035 to avoid overshooting in gradient descent of back-propagation. To combat vanishing gradient, Rectified Linear Unit (ReLU) is applied to both the hidden layer and the output layer. The default architecture of the ANN is demonstrated in Table I.

TABLE I
DEFAULT ARCHITECTURE OF THE ANN

Layer	Input		Hidden		Output	
Feature	#dimension	#neurons	Activation	#dimension	Activation	
Value	58	200	ReLU	32	ReLU	

3) *SARSA*: In the course of learning, an ϵ -greedy policy is adopted with $\epsilon = 0.2$ to encourage early exploration. A decay factor $\beta = 5 \cdot 10^{-5}$ reduces the actual ϵ_n used in the n^{th} episode to $\frac{\epsilon}{1+\beta \cdot n}$ to facilitate exploitation of accumulated knowledge in the late stage.

Given a current state-action pair (s, a) , a SARSA agent anticipates the next move a' according to the specified policy and trains the ANN with the observed reward r and the Q-value of the next move $q(a')$. If s is a terminal stage of the episode, then the agent considers only r . A reward $r = 1$ is given in case of checkmate, otherwise, $r = 0$.

4) *Q-learning*: Unlike SARSA, in non-terminal scenarios, a Q-learning agent speculates on the seemingly most rewarding next move a^* among all feasible moves and updates the neural network with the observed reward r and Q-value of this anticipated best move $q(a^*)$. However, a^* may not be the move to take. The Q-learning agent takes the next move a' according to the same ϵ -greedy policy as SARSA and proceeds to an immediate subsequent state-action pair (s', a') , which is the core logic of the off-policy learning method.

To curb the exploding weights of the neural network, Sigmoid activation is applied to the output layer. ReLU activation remains in use for the hidden layer.

5) *Change of the administration of reward*: With the default game setup and SARSA learning method, a negative unit reward $r = -1$ is imposed on each move taken, but no reward is granted at the terminal state, regardless of the result.

III. RESULTS

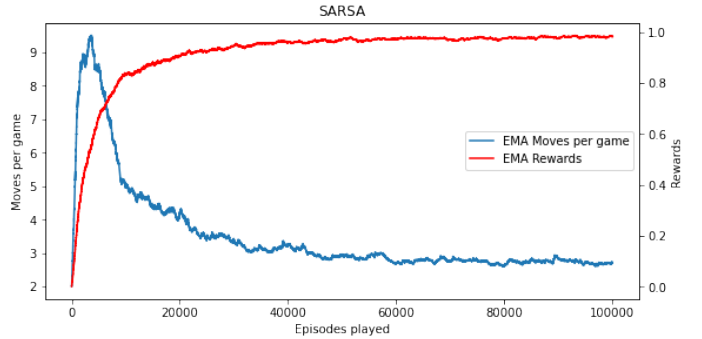


Fig. 2. Exponential Moving Average (EMA) moves per game and rewards received by a SARSA agent over 100k episodes.

A. Learning with SARSAⁱⁱⁱ

As a SARSA agent learns to checkmate without much knowledge of the game, except for the basic rules, it proceeds randomly at the beginning and may find itself in limbo performing certain inefficient moves. Figure 2 shows a spike of exponential moving average (EMA) moves per game in the early phase of the training. As the training progresses, with better knowledge, the agent checkmates more often and

ⁱⁱⁱSolution to Task 3: Implement SARSA and produce plots that show the reward per game and the number of moves per game vs training time with exponential moving average.

faster, which leads to a steep hike in rewards and a drop in moves per game before either reaches a plateau, where playing more episodes does not improve the performance significantly. Having trained for 100k episodes, the agent accumulates an reward of 0.92779 by finishing a game in 3.54345 moves on average.

B. Hyper-parameter analysis with SARSA^{iv}

With a SARSA agent, hyper-parameters are scrutinized in a grid-search fashion for $\beta \in \{5 \cdot 10^{-7}, 3 \cdot 10^{-5}, 5 \cdot 10^{-5}, 7 \cdot 10^{-5}\}$ and $\gamma \in \{0.2, 0.85, 0.99\}$. As reinforcement learning relies on the agent-environment interaction, meaning no separation between training and testing, the average moves per game and rewards over 100k episodes are used as the performance measures. The discount factor γ dominates how much the agent deducts from a future reward while considering an immediate move. The smaller γ is, the more myopic an agent is. β regulates the exploration behavior of an agent. Given the exploration rate ϵ , in the n^{th} episode, the actual exploration rate is given by $\epsilon_n = \frac{\epsilon}{\beta \cdot n}$. The larger β is, the less often an agent explores as the training progresses.

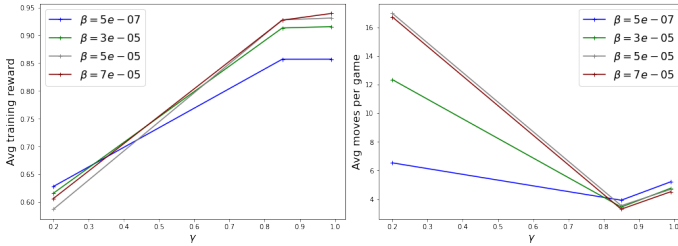


Fig. 3. Average moves per game and rewards received by a SARSA agent after 100k episodes with different combinations of discount factor γ and decay factor β .

As Figure 3 shows, the results of the experiment agree with the intuition that a rather small discount factor γ encourages an agent to focus only on the immediate effect of a move rather than the subsequent states where it might be better off, it therefore fails to explore efficient strategies. An increasing γ leads to more rewards and fewer moves per game in general. However, games last longer when γ goes from 0.85 to 0.99. One explanation is that by discounting future rewards very slightly, an agent might take unnecessary moves to where it is not significantly better than its current state.

The exploration decay factor β has a say in the trade-off between exploration and exploitation. Large β works better by allowing the agent to deviate from its sub-optimal strategy when the agent fails to learn decent strategies with a small discount factor γ .

With $\gamma = 0.99$ and $\beta = 7 \cdot 10^{-5}$, the SARSA agent accumulates the highest average reward of 0.93943 and finishes a game within 4.50880 moves on average over 100k episodes.

^{iv}Solution to Task 4: Change the discount factor γ and the speed β of the decaying trend of ϵ and analyse the results.

C. Q-learning^v

As illustrated in Figure 4, a Q-learning agent has bumpier curves compared with SARSA. The Q-learning agent achieves lower average rewards of 0.86206 and higher average moves per game of 11.23885 after 100k episodes. However, in spite of its sub-ideal performance measures, Q-learning might have explored a more efficient approach to tackle the game, which will be discussed later with a demonstration game. Besides, as neither of the curves reaches a plateau, one might argue that further training could yield a better outcome.

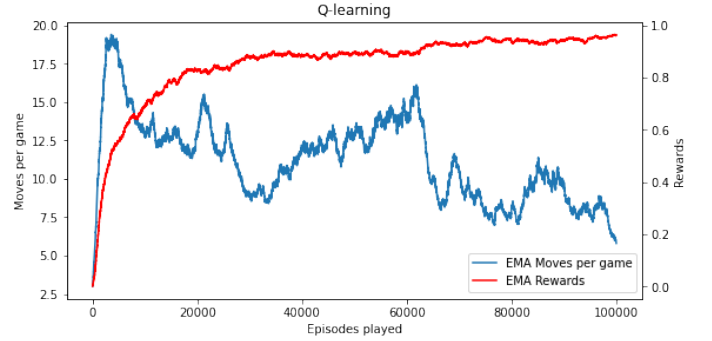


Fig. 4. EMA moves per game and rewards received by a Q-learning agent over 100k episodes.

D. An impatient agent^{vi}

By changing the administration of reward, the SARSA learning agent is now encouraged to end a game quickly, since the negative reward for each move penalizes long games. The absence of reward at a terminal state makes a checkmate and a draw indifferent for the agent, thus neutralizes the motivation to win the game. Figure 5 shows that the SARSA agent learns to end a game within 2 moves on average.

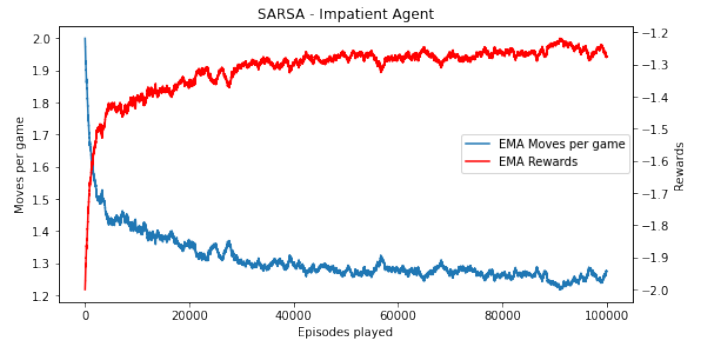


Fig. 5. EMA moves per game and rewards received by an impatient SARSA agent aiming at finishing the game ASAP over 100k episodes.

^vSolution to Task 5: Implement Q Learning and compare the new results with the previous results.

^{vi}Solution to Task 6: Change the administration of reward and interpret the results

E. A demonstration game

Figure 1 visualizes the initial setup of a demonstration game, with **Qa1Ka2 Kc2**.

A SARSA agent checkmates within 6 moves. The first 3 moves in Figure 6 are analyzed. The first move **Qd4 Kc1** is optimal, but the following **Qb2 Kd1** is not. The agent doesn't proceed aggressively to keep the opponent king in check. Instead, it backs off with **Qb3 Kd2** to avoid a stalemate.

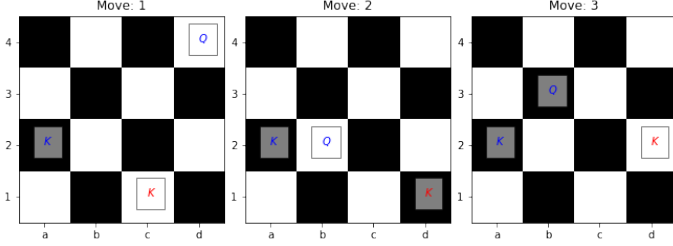


Fig. 6. Sub-optimal moves of a SARSA agent

With the same initial setup and identical first move, Figure 7 shows that a Q-learning agent pursues a more aggressive par-human approach by taking **Kb3 Kb2** to expand the squares covered by its king and leave only one possible option for the opponent. With **Qb2#**, it checkmates in 3 moves.

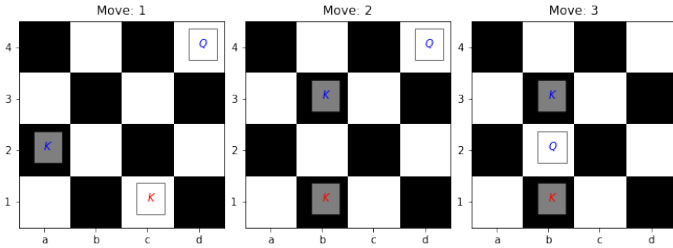


Fig. 7. Par-human optimal moves of a Q-learning agent

As demonstrated in Figure 8, an impatient SARSA agent doesn't care about the outcome but only wants to end the game as soon as possible. **Qb2 Kd3** leaves the opponent king in a stalemate. With **Kb3**, the game ends in a draw.

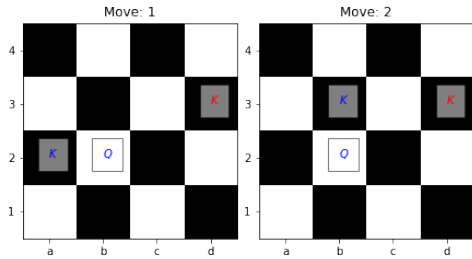


Fig. 8. Impatient agent only wants to end the game ASAP

IV. CONCLUSION

ANN is a powerful tool to approximate Q-values, so deep reinforcement learning can attempt tasks with high-dimensional state-action space. In the simplified chess task,

a SARSA agent generates relatively smooth EMA curves of moves per game and rewards, while a Q-learning agent explores more effective strategies powered by its off-policy update rule. To facilitate comparison, the number of episodes is fixed at 100k. Training the Q-learning agent for more episodes is expected to generate a better outcome.

By design, reinforcement learning aims at solving goal-directed tasks, where policy and reward structure shape the behavior of an agent.

REFERENCES

- [Lin92] Long-Ji Lin. "Self-improving reactive agents based on reinforcement learning, planning and teaching". In: *Machine learning* 8.3 (1992), pp. 293–321.
- [Sch03] Eric Schiller. *Official Rules Of Chess*. Cardoza Publishing, 2003.
- [Mni+13] Volodymyr Mnih et al. "Playing Atari with Deep Reinforcement Learning". In: *CoRR* abs/1312.5602 (2013). arXiv: 1312.5602. URL: <http://arxiv.org/abs/1312.5602>.
- [Mni+15] Volodymyr Mnih et al. "Human-level control through deep reinforcement learning". In: *nature* 518.7540 (2015), pp. 529–533.

APPENDIX

Detailed implementation of TD algorithms is attached in the appendix, where several highlights are worth mentioning.

An objected-oriented Artificial Neural Network is implemented to facilitate easy access to intermediate results, such as values of hidden neurons, in forward-feed and back-propagation.

A set of functions, namely the Chess Visualization Toolkit, is developed to facilitate visual inspection of the endgames, translation from the algebraic notation to matrix indices and easy customization of the chessboard.

A. Neural Network

```

1 # OOD ANN implementation
2 class ANN: # ANN object w/ 1 hidden layer
3     ↳ w/ 200 units by default
4
5     def __init__(self, N_in = 58, N_a =
6         ↳ 32, random_seed = 9, act1 = 'Relu',
7         ↳ act2 = 'Relu'): # constructor w/ the
8         ↳ size of input, the size of output
9         ↳ and random_seed as optional
10        ↳ parameters
11            self.N_h=200                                ##
12        ↳ NUMBER OF HIDDEN NODES (A NETWORK
13        ↳ WITH ONE HIDDEN LAYER WITH SIZE 200)
14
15        ### Random seed
16        np.random.seed(random_seed)
17
18        ### Xavier initialization

```

```

11     self.W1 = np.random.randn(self.N_h
12     ↪ , N_in) * np.sqrt(1 / (N_in)) #
13     ↪ input layer, of shape (200, 58)
14     self.W2 = np.random.randn(N_a,
15     ↪ self.N_h) * np.sqrt(1 / (self.N_h))
16     ↪ # hidden layer, of shape (32, 200)
17
18     self.b1 = np.zeros((self.N_h,)) #
19     ↪ of size 200
20     self.b2 = np.zeros((N_a,)) # of
21     ↪ size 32
22
23     ### Parameterize the choices of
24     ↪ activation functions
25     self.act1 = 1 if act1 == 'Relu'
26     ↪ else 2
27     self.act2 = 1 if act2 == 'Relu'
28     ↪ else 2
29     self.act2 = 3 if act2 == 'Sigmoid'
30     ↪ else self.act2
31
32     ### Initiate neurons and
33     ↪ activations
34     self.z1 = np.zeros((self.N_h,))
35     self.a1 = np.zeros((self.N_h,))
36     self.z2 = np.zeros((N_a,))
37     self.a2 = np.zeros((N_a,))
38
39     def forwardfeed(self, X):
40     ↪ # Forwardfeed
41     ↪ # input -> hidden
42     self.z1 = self.W1 @ X + self.b1 #
43     ↪ of size 200
44
45     self.a1 = self.relu(self.z1) if
46     ↪ self.act1 == 1 else self.z1
47
48     ↪ # hidden -> output
49     self.z2 = self.W2 @ self.a1 + self
50     ↪ .b2 # of size 32
51
52     self.a2 = self.relu(self.z2) if
53     ↪ self.act2 == 1 else self.z2 # of
54     ↪ size 32
55     return self.a2
56
57     def relu(self, x): # rectified linear
58     ↪ unit activation to cope w/ vanishing
59     ↪ gradient
60     return (x > 0).astype(int) * x
61
62     def backpropagation(self, X, delta,
63     ↪ a_agent, eta):
64     ↪ # Backpropagation
65     ↪ # Gradients
66
67     dz1, dz2 = self.calc_gradient(
68     ↪ delta, a_agent)
69
70     ↪ # Descent (update)
71     self.W2[a_agent, :] = self.W2[
72     ↪ a_agent, :] + eta * dz2 * self.a1
73     self.b2[a_agent] = self.b2[a_agent
74     ↪ ] + eta * dz2
75
76     self.W1 = self.W1 + eta * np.outer
77     ↪ (dz1, X)
78     self.b1 = self.b1 + eta * dz1
79
80     return self.W1, self.W2
81
82     def calc_gradient(self, delta, a_agent
83     ↪ ):
84     ↪ # Gradients
85     dz2 = delta * self.heavy_side(self
86     ↪ .a2[a_agent]) if self.act2 == 1 else
87     ↪ delta # of size 1
88
89     dz1 = delta * self.W2[a_agent, :]
90     ↪ * self.heavy_side(self.a1) if self.
91     ↪ act1 == 1 else dz2 * self.W2[a_agent
92     ↪ , :] * self.a1 # of size 200
93     return dz1, dz2
94
95     def heavy_side(self, x): # pseudo
96     ↪ derivative of ReLu
97     return (x > 0).astype(int)
98
99     def load_weights(self, W1, W2, b1, b2)
100    ↪ : # load trained model
101    self.W1 = W1
102    self.W2 = W2
103    self.b1 = b1.flatten()
104    self.b2 = b2.flatten()

```

Listing 1. Object-oriented implementation of the Artificial Neural Network

B. SARSA

```

1 # TRAINING LOOP of SARSA
2
3 nn = ANN() # instantiate a Neural Network
4 ↪ object w/ default activation
5 ↪ functions ('Relu')
6
7 for n in range(N_episodes):
8
9     epsilon_f = epsilon_0 / (1 + beta * n)
10    ↪ # DECAING EPSILON
11    Done=0
12    ↪ # SET DONE TO ZERO (BEGINNING OF
13    ↪ THE EPISODE)

```

```

9     i = 1
10    ↪    ## COUNTER FOR NUMBER OF ACTIONS
11
12    S,X,allowed_a=env.Initialise_game()
13    ↪    ## INITIALISE GAME
14
15    # Choose action A (w/o taking it
16    ↪ actually)
17    ## pre-action meditation
18    a,_=np.where(allowed_a==1) # find
19    ↪ allowed actions
20
21    ## find Q-values
22    Qvalues = nn.forwardfeed(X) # generate
23    ↪ a vector of estimated Q values by
24    ↪ ANN
25    Qvalues_of_allowed_a = Qvalues[a] #
26    ↪ select only the Q values of allowed
27    ↪ actions
28
29    ## Given Q-values, pick an action in
30    ↪ an e-Greedy way
31    a_agent = a[EpsilonGreedy_Policy(
32    ↪ Qvalues_of_allowed_a, epsilon_f)][0]
33    ↪ # idx of the chosen action
34
35    q = np.copy(Qvalues[a_agent]) # copy
36    ↪ the Q value of the chosen action
37
38    while Done==0:
39    ↪    ## START THE EPISODE
40
41        # take the action a_agent &
42    ↪ observe
43        S_next,X_next,allowed_a_next,R,
44    ↪ Done=env.OneStep(a_agent)
45
46        ## THE EPISODE HAS ENDED, UPDATE
47    ↪ ...BE CAREFUL, THIS IS THE LAST STEP
48    ↪ OF THE EPISODE
49        if Done==1:
50
51            # Record the reward and moves
52    ↪ taken
53            R_save[n]=np.copy(R)
54            N_moves_save[n]=np.copy(i)
55
56            # Update Weights
57            ## Backpropagation
58            delta = (R - q) # Compute the
59    ↪ delta
60            nn.backpropagation(X, delta,
61    ↪ a_agent, eta) # backpropagation via
62    ↪ gradient descent
63
64            break
65
66    else:
67        # post-action reflection (
68    ↪ choose the next action w/o actually
69    ↪ taking it either)
70        a_next,_=np.where(
71    ↪ allowed_a_next==1) # find allowed
72    ↪ actions
73
74        # Find Qvalues
75        Qvalues_next = nn.forwardfeed(
76    ↪ X_next) # generate a vector of
77    ↪ estimated Q values by ANN
78
79        Qvalues_of_allowed_a_next =
80    ↪ Qvalues_next[a_next] # select only
81    ↪ the Q values of allowed actions
82
83        a_agent_next = a_next[
84    ↪ EpsilonGreedy_Policy(
85    ↪ Qvalues_of_allowed_a_next, epsilon_f
86    ↪ )][0] # idx of the chosen action
87
88        q_new = np.copy(Qvalues_next[
89    ↪ a_agent_next]) # copy the Q value of
90    ↪ the chosen action
91
92        # Update Weights
93        ## Backpropagation
94
95        delta = (R + gamma * q_new - q
96    ↪ ) # Compute the delta
97
98        nn.backpropagation(X, delta,
99    ↪ a_agent, eta) # backpropagation via
100    ↪ gradient descent
101
102        # NEXT STATE AND CO. BECOME ACTUAL
103    ↪ STATE...
104        S=np.copy(S_next)
105        X=np.copy(X_next)
106        allowed_a=np.copy(allowed_a_next)
107        a_agent = np.copy(a_agent_next)
108        q = np.copy(q_new)
109
110        i += 1 # UPDATE COUNTER FOR
111    ↪ NUMBER OF ACTIONS
112
113    print_progress(n, N_episodes) #
114    ↪ odometer to show training process
115
116    print_progress(None, N_episodes) #
117    ↪ indicate whether training is over

```



```

80 print('e-Greedy agent w/ 1QN SARSA,
    ↳ Average reward:', np.mean(R_save), '
    ↳ Number of steps: ', np.mean(
    ↳ N_moves_save)) # print performance
    ↳ measures

```

Listing 2. Training logic of SARSA

C. Q-learning

```

1 # TRAINING LOOP of Q-learning
2
3 nn = ANN(act2 = 'Sigmoid') # to curb the
    ↳ exploding weights, use sigmoid
    ↳ activation function in the output
    ↳ layer
4
5
6 for n in range(N_episodes):
7
8     epsilon_f = epsilon_0 / (1 + beta * n)
    ↳ ## DECAYING EPSILON
9     Done=0
    ↳ ## SET DONE TO ZERO (BEGINNING OF
    ↳ THE EPISODE)
10    i = 1
    ↳ ## COUNTER FOR NUMBER OF ACTIONS
11
12    S,X,allowed_a=env.Initialise_game()
    ↳ ## INITIALISE GAME
13
14
15    while Done==0:
    ↳ ## START THE EPISODE
16
17        # Choose action A (w/o taking it
    ↳ actually)
18        ## pre-action meditation
19        a,_=np.where(allowed_a==1) # find
    ↳ allowed actions
20
21        ## find Q-values
22        Qvalues = nn.forwardfeed(X)
23        Qvalues_of_allowed_a = Qvalues[a]
24
25        ## Given Q-values, pick an action
    ↳ in an e-Greedy way
26
27        a_agent = a[EpsilonGreedy_Policy(
    ↳ Qvalues_of_allowed_a, epsilon_f)][0]
    ↳ # idx of the chosen action
28
29        q = np.copy(Qvalues[a_agent])
30
31        # take the action a_agent &
    ↳ observe

```

```

32    S_next,X_next,allowed_a_next,R,
    ↳ Done=env.OneStep(a_agent)
33
34
35    ## THE EPISODE HAS ENDED, UPDATE
    ↳ ...BE CAREFUL, THIS IS THE LAST STEP
    ↳ OF THE EPISODE
    if Done==1:
36
37        # Record the reward and moves
    ↳ taken
38        R_save[n]=np.copy(R)
39        N_moves_save[n]=np.copy(i)
40
41        # Update Weights
42        ## Compute the delta
43        delta = (R - q)
44
45        ## Backpropagation
46        nn.backpropagation(X, delta,
    ↳ a_agent, eta)
47
48        break
49
50    else:
51        # post-action reflection (
    ↳ choose the next action w/o actually
    ↳ taking it either)
52        a_next,_=np.where(
    ↳ allowed_a_next==1) # find allowed
    ↳ actions
53
54        # Find Qvalues
55        Qvalues_next = nn.forwardfeed(
    ↳ X_next)
56
57        Qvalues_of_allowed_a_next =
    ↳ Qvalues_next[a_next]
58        a_agent_next = a_next[
    ↳ EpsilonGreedy_Policy(
    ↳ Qvalues_of_allowed_a_next, epsilon_f
    ↳ )][0] # idx of the chosen action
59
60        q_new = np.max(Qvalues_next)
61
62        # Update Weights
63        ## Compute the delta
64        delta = (R + gamma * q_new - q
    ↳ ) # no punishment
65
66        ## Backpropagation
67        nn.backpropagation(X, delta,
    ↳ a_agent, eta)
68
69
70

```

```

71     # NEXT STATE AND CO. BECOME ACTUAL
    ↪ STATE...
72     S=np.copy(S_next)
73     X=np.copy(X_next)
74     allowed_a=np.copy(allowed_a_next)
75
76
77     i += 1 # UPDATE COUNTER FOR
    ↪ NUMBER OF ACTIONS
78
79
80     print_progress(n, N_episodes)
81
82     print_progress(None, N_episodes)
83
84     print('e-Greedy agent w/ 1QN Q-learning,
    ↪ Average reward:', np.mean(R_save), '
    ↪ Number of steps: ', np.mean(
    ↪ N_moves_save))

```

Listing 3. Training logic of Q-learning

D. Hyper-parameter analysis

```

1  # TRAINING LOOP of SARSA w/ Hyper-
    ↪ parameter tuning
2
3  # Grid-search for hyperparameters
4  beta_list = [5e-07, 3e-05, 5e-05, 7e-05]
    ↪ # Decay factor of EPSILON
5
6  gamma_list = [0.2, 0.85, 0.99] #
    ↪ THE DISCOUNT FACTOR
7
8  performance = []
9
10 for beta in beta_list:
11
12     for gamma in gamma_list:
13
14         # SAVING VARIABLES
15         R_save = np.zeros([N_episodes, 1])
16         N_moves_save = np.zeros([
    ↪ N_episodes, 1])
17
18         # instantiate a new ANN object for
    ↪ each combination of hyper-
    ↪ parameters
19         nn = ANN() # instantiate a Neural
    ↪ Network object w/ default activation
    ↪ functions ('Relu')
20
21         for n in range(N_episodes):
22
23             '''
24             ** Game initialization **

```

```

25         Some code is hidden for
    ↪ simplicity
26         Please refer to detailed
    ↪ source code in the main training
    ↪ logic of SARSA
27         '''
28
29         while Done==0:
30             ## START THE EPISODE
31
32             ## THE EPISODE HAS ENDED,
    ↪ UPDATE...BE CAREFUL, THIS IS THE
    ↪ LAST STEP OF THE EPISODE
33             if Done==1:
34
35                 # Record the reward
    ↪ and moves taken
36                 R_save[n]=np.copy(R)
37                 N_moves_save[n]=np.
    ↪ copy(i)
38
39                 '''
40                 ** Episode terminates
    ↪ **
41
42                 Some code is hidden
    ↪ for simplicity
43                 Please refer to
    ↪ detailed source code in the main
    ↪ training logic of SARSA
44                 '''
45
46                 break
47
48             else:
49
50                 '''
51                 ** Game continues **
52
53                 Some code is hidden
    ↪ for simplicity
54                 Please refer to
    ↪ detailed source code in the main
    ↪ training logic of SARSA
55                 '''
56
57                 # Given a set of hyper-parameters,
    ↪ save performance metrics
58                 pd.DataFrame(R_save).to_csv(f'
    ↪ R_save_w_beta_{beta}_gamma_{gamma}
    ↪ _100k.csv')
59
60                 pd.DataFrame(N_moves_save).to_csv(
    ↪ f'N_moves_save_w_beta_{beta}_gamma_{

```



```
→ gamma}_100k.csv')
```

```
performance.append([beta, gamma,
→ np.mean(R_save), np.mean(
→ N_moves_save)])
```

Listing 4. Hyper-parameter analysis with SARSA

E. Change of the administration of reward

```
1 # TRAINING LOOP of impatient SARSA
```

```
2 nn = ANN(act2 = None) # use linear
3 → activation function (no activation)
4 → of the output layer to approximate
5 → negative Q-values
```

```
6 '''
7 ** ONLY different implementation from
8 → SARSA is exhibited in this snippet
9 → of code **
```

```
10 Some code is hidden for simplicity
11 Please refer to detailed source code in
12 → the main training logic of SARSA
13 '''
```

```
14 for n in range(N_episodes):
```

```
15 R_accu = 0 # initiate accumulative
16 → Reward of an episode
```

```
17 '''
18 ** Game initialization **
```

```
19 Some code is hidden for simplicity
20 Please refer to detailed source code
21 → in the main training logic of SARSA
22 '''
```

```
23 while Done==0:
24 → ## START THE EPISODE
```

```
25 # take the action A & observe
26 S_next, X_next, allowed_a_next, R,
27 → Done=env.OneStep(a_agent)
```

```
28 ## THE EPISODE HAS ENDED, UPDATE
29 → ...BE CAREFUL, THIS IS THE LAST STEP
30 → OF THE EPISODE
31 if Done==1:
```

```
32 # Record the reward and moves
33 → taken
34 R_save[n] = R_accu - 1 #
35 → record the accumulative reward of
36 → the episode
```

```
37 # Update Weights
38 ## Backpropagation
39 delta = (0 - q) # Neutralize
40 → the impact of result
```

```
41 '''
42 ** Episode terminates **
```

```
43 Some code is hidden for
44 → simplicity
45 Please refer to detailed
46 → source code in the main training
47 → logic of SARSA
```

```
48 '''
49 break
50 else:
51 '''
52 ** Game continues **
```

```
53 Some code is hidden for
54 → simplicity
55 Please refer to detailed
56 → source code in the main training
57 → logic of SARSA
```

```
58 '''
59 # Update Weights
60 ## Backpropagation
61 R_accu -= 1 # immediate reward
62 → -1 for each step the agent takes (
63 → penalize long games)
```

```
64 delta = (-1 + gamma * q_new -
65 → q) # Compute the delta w/ penalty -1
66 → for each step the agent takes
```

```
67 nn.backpropagation(X, delta,
68 → a_agent, eta) # backpropagation via
69 → gradient descent
```

Listing 5. Impatient agent with SARSA

F. Chess Visualization

```
1 # chess_vis.py
2 # define a series of functions to
3 → facilitate visualization
4 import numpy as np
5 import matplotlib.pyplot as plt
6 from degree_freedom_queen import *
7 from degree_freedom_king1 import *
8 from degree_freedom_king2 import *
```

```

8 from generate_game import *
9 from Chess_env import *
10
11 # define a function to translate piece
12     ↳ location from board to the
13     ↳ environment index system: a0 -> [3,
14     ↳ 0]
15 def trans(pos):
16     y = ord(pos[0]) - 97 # char c -> c* in
17     ↳ ASCII -> c* - 97
18     x = 4 - int(pos[1]) # char n -> int n
19     ↳ - 1
20     return (x, y)
21
22 # define a function to initiate the game
23     ↳ lie the demo case
24 def demo_initialise_game(env, k1, q1, k2):
25     # START THE GAME BY SETTING PIECES
26
27     size_board = 4; # 4-by-4 board
28
29     env.Board = np.zeros([size_board,
30     ↳ size_board], dtype=int) # initialize
31     ↳ the board given size w/ a numpy
32     ↳ zero matrix
33
34     k1_x, k1_y = k1 # read the x,y
35     ↳ coordinates of my king
36     env.p_k1 = np.array([k1_x, k1_y]) # 1
37     ↳ for my king, pass its initial
38     ↳ position to the env.
39     env.Board[k1_x, k1_y] = 1 # update env
40     ↳ .board w/ king's position
41
42     q1_x, q1_y = q1 # read the x,y
43     ↳ coordinates of my queen
44     env.p_q1 = np.array([q1_x, q1_y]) # 2
45     ↳ for my queen, pass its initial
46     ↳ position to the env.
47     env.Board[q1_x, q1_y] = 2 # update env
48     ↳ .board w/ queen's position
49
50     k2_x, k2_y = k2 # read the x,y
51     ↳ coordinates of opponent king
52     env.p_k2 = np.array([k2_x, k2_y]) # 3
53     ↳ for opponent king, pass its initial
54     ↳ position to the env.
55     env.Board[k2_x, k2_y] = 3 # update env
56     ↳ .board w/ opponent king's position
57
58     # Allowed actions for the agent's king
59     env.dfkl_constrain, env.a_k1, env.dfkl
60     ↳ = degree_freedom_king1(env.p_k1,
61     ↳ env.p_k2, env.p_q1, env.Board)
62
63     # Allowed actions for the agent's
64     ↳ queen
65     env.dfql_constrain, env.a_q1, env.dfql
66     ↳ = degree_freedom_queen(env.p_k1,
67     ↳ env.p_k2, env.p_q1, env.Board)
68
69     # Allowed actions for the enemy's king
70     env.dfk2_constrain, env.a_k2, env.
71     ↳ check = degree_freedom_king2(env.
72     ↳ dfk1, env.p_k2, env.dfql, env.Board,
73     ↳ env.p_k1)
74
75     # ALLOWED ACTIONS FOR THE AGENT, ONE-
76     ↳ HOT ENCODED
77     allowed_a=np.concatenate([env.a_q1,env
78     ↳ .a_k1],0)
79
80     # FEATURES (INPUT TO NN) AT THIS
81     ↳ POSITION
82     X=env.Features()
83
84     return env.Board, X, allowed_a
85
86 # define a function to print the chess
87     ↳ board w/ given position of pieces
88 def print_board(S, num = None, figsize =
89     ↳ (5, 5), ax = None):
90
91     # instantiate a subplots obj w/
92     ↳ default figsize = (5, 5)
93     if ax == None:
94         fig, ax = plt.subplots(
95     ↳ figsize = figsize)
96
97     # create a blank board
98     blank_board = np.add.outer(range
99     ↳ (4), range(4))%2 # create grids
100     ax.imshow(blank_board, cmap="
101     ↳ binary_r") # use 2D raster to
102     ↳ generate the board
103     ax.set_xticks([0.0, 1.0, 2.0,
104     ↳ 3.0]) # formulate x ticks
105     ax.set_xticklabels(['a', 'b', 'c',
106     ↳ 'd']) # formulate x ticks
107     ax.set_yticks([0.0, 1.0, 2.0,
108     ↳ 3.0]) # formulate y ticks
109     ax.set_yticklabels([4, 3, 2, 1]) #
110     ↳ formulate y ticks
111
112     # put pieces on the board
113     board_idx = np.array([[ (x, y) for
114     ↳ y in range(4)] for x in range(4)]) #
115     ↳ generate a positional index table

```

```

74     pieces = ['K', 'Q', 'K'] # piece
    ↪ names
75     colors = ['blue', 'blue', 'red'] #
    ↪ piece colors (blue -> player, red
    ↪ -> opponent)
76     for i in range(3): # assign each
    ↪ piece
77         y_loc, x_loc = board_idx[S
    ↪ == i + 1][0]
78         ax.text(x_loc, y_loc,
    ↪ pieces[i], style='italic', color =
    ↪ colors[i],
79                 bbox={'facecolor':
    ↪ 'white', 'alpha': 0.5, 'pad': 10})
80
81     # format title
82     if num == None:
83         ax.set_title(f'Overkilling
    ↪ Endgame')
84     else:
85         ax.set_title(f'Move: {num}
    ↪ ')
86
87     # show figure
88     if ax == None:
89         fig.show()

```

Listing 6. Chess Visualization Toolkit