

# RESEARCHING THE ABILITY OF REINFORCEMENT LEARNING AGENTS TO ADAPT TO CHANGES IN THEIR TRAINING ENVIRONMENT

Tomos Ody (ODY14527411)

Department of Computer Science, University of Lincoln

*Abstract*—The quick brown fox jumped over the lazy dog Lorem ipsum dolor sit amet, consectetur adipiscing elit. Ut purus elit, vestibulum ut, placerat ac, adipiscing vitae, felis. Curabitur dictum gravida mauris. Nam arcu libero, nonummy eget, consectetur id, vulputate a, magna. Donec vehicula augue eu neque. Pellentesque habitant morbi tristique senectus et netus et malesuada fames ac turpis egestas. Mauris ut leo. Cras viverra metus rhoncus sem. Nulla et lectus vestibulum urna fringilla ultrices. Phasellus eu tellus sit amet tortor gravida placerat. Integer sapien est, iaculis in, pretium quis, viverra ac, nunc. Praesent eget sem vel leo ultrices bibendum. Aenean faucibus. Morbi dolor nulla, malesuada eu, pulvinar at, mollis ac, nulla. Curabitur auctor semper nulla. Donec varius orci eget risus. Duis nibh mi, congue eu, accumsan eleifend, sagittis quis, diam. Duis eget orci sit amet orci dignissim rutrum.

## I. INTRODUCTION

### A. Rationale

Much of the Reinforcement Learning (RL) literature is focused around specific best case applications to specific problem domains. This is an improvement over past trends described in Kaelbling's survey of the subject from 1996 Kaelbling, Littman, and Moore 1996 where research was being done on poorly defined problem domains. More recently RL research has become more focused into complex and specific problem domains. However there is little research comparing the various RL methods to one another in a comprehensive manner. Often with brief discussions of the various methods without evidence validating the chosen RL methodology. More often than not the choice of method comes down to examination of the problem domain to be explored.

The focus of this project is in comparing the performance of the more popular Q-Learning (QL) control policy (Christopher J. C. H. Watkins and Dayan 1992) against the lesser known State-Action-Reward-State-Action (SARSA) control policy (Rummery and Niranjan 1994) in as comprehensive a manner as possible. With the aim of understanding how applicable these methods are to various problem domains and situations. In theory SARSA should always outperform QL since it always evaluates it's chosen action (Sutton and A G. Barto 2018). And is the more computationally expensive method as a result.

Q-Learning and SARSA are very closely related temporal difference (TD) problems and so many optimizations of one

are applicable to the other. This makes them simple to compare against one another as well as understand the differences in reinforcement learning environments. However QL overshadows SARSA in the literature quite dramatically, to the point SARSA is considerably harder to research than QL in the academic literature, however I have been unable to find an obvious reason why.

### B. Aims and Objectives

The aim of this project is to understand the differences in performance of QL and SARSA and by extension offline and online control policies for TD problems. With focus being drawn to how the environments used for training affect the performance of the various control policies.

Further research will be conducted to understand the scope of TD problem domains as well as explore potential optimisations of QL and SARSA. This will allow the analysis to be more complete and present additional opportunities for comparison between the control policies.

The reason for the prevalence of QL over SARSA in the academic literature should be explored. Experimental data will be considered alongside the literature to inform the evaluation of these control policies for TD problems.

### C. Hypothesis

### D. Background

1) *Reinforcement Learning*: RL is distinct from most Artificial Intelligence (AI) disciplines in that instead of relying on a large pool of ground truth examples for the learning to be done upon. Instead RL relies on the design of the learning environment to teach an AI how to perform a task. In most cases this is done through reward functions with the RL agent inferring from these rewards how the task should be performed (Kaelbling, Littman, and Moore 1996).

This however relies on the designed learning environment being simple enough for an agent to learn from it's interactions with the environment alone. This leads to the main challenge in RL, exploration vs exploitation. This describes the balance between the agent exploring the environment to update all actions with their corresponding rewards and exploiting the already learned actions, taking the optimal action. This is important as it not only decides how the agent learns the environment

but also its ability to interact with the environment effectively. Since if the agent always takes random actions to explore the environment fully it will never learn the optimal strategy. Likewise if the agent only follows the optimal strategy it will not learn better strategies if they exist. This can drastically affect the time it takes to train the agent since being too close to one extreme or the other causes the agent to either never learn the optimal strategy or learn an incorrect optimal strategy by never encountering the optimal strategy (Sutton and A G. Barto 2018).

2) *Q-Learning*: Q-Learning was initially proposed by Watkins in their Phd thesis Christopher John Cornish Hellaby Watkins 1989 building on MDPs to solve more complex problems without the difficulty of building an equally complex model to solve it. Q-Learning was formalised 3 years later in Christopher J. C. H. Watkins and Dayan 1992. ———  
===== SORT! (watkins page 55) Harvard “It provides agents with the capability of learning to act optimally in Markovian domains by experiencing the consequences of actions, without requiring them to build maps of the domains.”

Q-Learning simplifies the environment into a Q-table with the dimensions

*[number of states, number of actions]*

Each point in the table corresponds to a Q-value for an action at the given state. This table allows for the agent to learn from the experience of each state and action. By following the maximum Q-value the optimal path should be followed given that all Q-values are representative of the environment. However, the discrete nature of the Q-table means that QL is limited to discrete problems or the environment must be discretized leading to a loss in state and action accuracy (Gaskett, Wettergreen, and Zelinsky 1999).

3) *SARSA*:

4) *Temporal Difference Learning*:

#### E. Report Structure

#### CONTENTS

<b>I</b>	<b>Introduction</b>	1
I-A	Rationale . . . . .	1
I-B	Aims and Objectives . . . . .	1
I-C	Hypothesis . . . . .	1
I-D	Background . . . . .	1
	I-D1 Reinforcement Learning . . .	1
	I-D2 Q-Learning . . . . .	2
	I-D3 SARSA . . . . .	2
	I-D4 Temporal Difference Learning	2
I-E	Report Structure . . . . .	2
<b>II</b>	<b>Literature Review</b>	2
II-A	Academic basis . . . . .	2
II-B	Scope of reinforcement learning . . .	3
II-C	Reinforcement learning techniques . .	3
<b>III</b>	<b>Methodology</b>	4
<b>IV</b>	<b>Design and Development</b>	4

<b>V</b>	<b>Experiments and Evaluation</b>	4
<b>VI</b>	<b>Discussions and Reflective Analysis</b>	4
<b>VII</b>	<b>Conclusion</b>	4
	<b>References</b>	4
	<b>Appendix</b>	4
	<b>Acknowledgements</b>	7

## II. LITERATURE REVIEW

### A. Academic basis

The standard RL model consists of an agent interacting with an environment via perception and action. At each step of the training process the agent receives some state information. Based on this information the agent will then take an action based on it's observations of the environment. Depending on the action taken a scalar reinforcement signal is sent to the agent as a reflection of success. The agent aims to maximise the long-run sum of the reinforcement signal. This process leads to the agent choosing the rewarded actions leading to the development of the desired behaviour in the agent. As outlined in Kaelbling, Littman, and Moore 1996, a survey of RL field from 1996. This basic outline is corroborated by Busoniu, Babuska, and De Schutter 2008, a survey of multiagent RL from 2008 and Kober, Bagnell, and Peters 2013, a survey of RL in robotics from 2013. This basic conception of RL models is consistent across all literature and appears in papers up to the current day.

This learning method does not fit neatly into the two main machine learning paradigms of supervised and unsupervised learning. Instead of being directly trained by correct examples or inferring patterns from the data, a RL agent learns from it's own interactions with the environment quantified by a scalar reward described in the book Reinforcement Learning: An Introduction Sutton and A G. Barto 2018, cited in many RL works. This agent driven learning style leads to a trade-off between exploration and exploitation (Sutton and A G. Barto 2018; Kaelbling, Littman, and Moore 1996; Busoniu, Babuska, and De Schutter 2008; Kober, Bagnell, and Peters 2013). Due to the need for the agent to explore the environment as well as perform actions to maximise the positive reward (exploitation). The trade-off stems from the incompatibility of these two behaviours which are both required for an effective agent. Since the agent must explore the environment to understand the positive and negative options possible yet avoid negative actions to succeed (Kaelbling, Littman, and Moore 1996). For most RL problems Markov Decision Processes (MDP) are applied as the methodology for modelling decisions to calculate action rewards accounting for the exploration-exploitation trade-off (Kober, Bagnell, and Peters 2013). MDPs work by considering rewards temporally, meaning that a series of actions can be attributed to the delayed reward function (Kaelbling, Littman, and Moore 1996; Busoniu, Babuska, and De Schutter 2008; Kober, Bagnell, and Peters 2013). This is in contrast to more simple greedy strategies which reward

the agent directly for each individual action. However these methods can fall victim to unlucky reward sampling early in the training process (Kaelbling, Littman, and Moore 1996). A useful framing problem for understanding the balancing of exploration and exploitation is the K-Armed Bandit problem. This is used in most survey papers (Kaelbling, Littman, and Moore 1996; Busoniu, Babuska, and De Schutter 2008; Kober, Bagnell, and Peters 2013) to explore the performance of RL methods in the scope of exploration vs exploitation. The more recent papers Busoniu, Babuska, and De Schutter 2008; Kober, Bagnell, and Peters 2013 analysing these methods to balance exploration and exploitation tend towards the more dynamic methods like MDPs as opposed to the more simplistic greedy strategies mentioned in Kaelbling, Littman, and Moore 1996. The other simplistic methods suggested by *ibid.* are not mentioned in the more recent surveys of RL Busoniu, Babuska, and De Schutter 2008; Kober, Bagnell, and Peters 2013, only in Sutton and A G. Barto 2018 are these methods mentioned as introductory methods to RL. The literature seems to be in agreement on many of the details of RL methods, however many methods mentioned in the older Kaelbling, Littman, and Moore 1996 seem to be less widely used in the surveys written in recent years (Busoniu, Babuska, and De Schutter 2008; Kober, Bagnell, and Peters 2013).

Most of the methods discussed in RL are very old, for example MDPs were originally proposed by Bellman 1958 in 1958 which is cited in each of the RL surveys Kaelbling, Littman, and Moore 1996; Busoniu, Babuska, and De Schutter 2008; Kober, Bagnell, and Peters 2013 as well as the seminal Sutton and A G. Barto 2018. The vast majority of improvements made to RL methods are minor improvements to MDPs.

### *B. Scope of reinforcement learning*

Reinforcement learning is a very adaptable machine learning approach due to the wide range of approaches and methodologies that can be used to form a solution. However Kaelbling, Littman, and Moore 1996 concludes that RL performs well on small problems but poorly on larger problems due to the efficiency of generalising the problems. This point is reaffirmed by Wirth et al. 2017, adding that an efficient generalisation of a complex problem requires expert domain knowledge. This suggests that whilst versatile, RL is only generally suited to simple problems.

One of the primary areas of RL research is robotics since the domain of controlling a robot fits neatly into a RL problem as discussed in Kober, Bagnell, and Peters 2013, whereas other machine learning disciplines are very difficult to utilise for robotic control systems. This view of RL in robotics is shared by Smart and Pack Kaelbling 2002. Many other hurdles are introduced by using reinforcement learning in a real context such as observation inaccuracies and the expense of hardware Kober, Bagnell, and Peters 2013, however this does not pertain to this project. Robotic RL agents instead present a strong base of real agent-environment interactions (*ibid.*).

Up till this point single-agent systems have been discussed, however Multi-Agent Reinforcement Learning (MARL) is another large body of work with multiple agents training in

the same environment. This domain is similar to single-agent RL but multiple agents open up a new set of challenges and advantages. As described in Busoniu, Babuska, and De Schutter 2008 the main challenge faced is an exacerbation of the exploration-exploitation trade-off due to inter-agent interactions adding a new level of complexity to the issue. Multiple agents can also be turned into an advantage in MARL by allowing agents to share their experience with other agents. Additionally more advanced reward functions can encourage collaborative actions between agents. One clear advantage over other RL domains is that MARL can benefit from parallel computing increasing the computational efficiency.

Most RL environments could be described as games, unsurprisingly the most active domain of RL is playing video games. Video games present a host of environments easily adaptable to RL applications, however the low complexity tolerance of RL still takes effect. This leads to either very simplistic games being used to train complete agents to perform well such as in Bellemare, Veness, and Bowling 2012 or more high level applications like Amato and Shani 2010 where a RL agent controls the strategies of an artificial intelligence agent designed to play the game. These applications both use MDPs with Q-Learning (Christopher J. C. H. Watkins and Dayan 1992) to consider the higher level temporal planning required for playing video games.

A criticism is raised in Shoham, Powers, and Grenager 2003, a critical survey of MARL using MDPs, as to the definition of many RL solutions in that they are often performed from an ill defined domain basis in how learning best occurs. It outlines a four stage method of properly defining a RL problem, describing how too many RL solutions work off previous bases without exploring their applicability to the problem being addressed. This demonstrated by the contrast between Amato and Shani 2010 where little academically supported domain knowledge is presented versus Ng et al. 2006 where the domain knowledge is comprehensively explored before designing the RL solution. This suggests that certain areas of RL problems are ill guided.

### *C. Reinforcement learning techniques*

Whilst not encompassing all of RL, MDPs represent the vast majority of modern machine learning solutions. This is due to the simplicity with which a complex problem can be represented, taking into account time and previous actions as outlined in Andrew G. Barto and Mahadevan 2003. As a result of this focus there is a wealth of different MDP strategies in the literature to draw from. An MDP is a framework for modelling actions with controllable yet undetermined outcomes to be optimised Bellman 1958. Utilising an MDP a complex problem can be abstracted down to a comparably computationally simple problem, this however requires much domain specific knowledge to be properly exercised Shoham, Powers, and Grenager 2003.

MDPs are a mathematical framework to quantify and rationalise decisions made in a Stochastic Game (SG) Bellman 1958. SGs are step by step games which quantify an environment probabilistically Shapley 1953, this hugely simplifies

reward and action computation in RL problems Busoniu, Babuska, and De Schutter 2008.

In both Bellemare, Veness, and Bowling 2012; Amato and Shani 2010 Q-Learning is used, it is an approach to MDPs which allows a RL agent to explore an environment to experience the various reward states associated with each action. The Q-Learning task will use these immediate rewards to plan temporally how to maximise the reward overall as described in Christopher J. C. H. Watkins and Dayan 1992.

Much of the research pertaining to RL problems seems much too condensed, with single authors appearing in multiple papers such as Barto in Andrew G. Barto and Mahadevan 2003 as well as the book Sutton and A G. Barto 2018. Kaelbling also appears in both Kaelbling, Littman, and Moore 1996; Smart and Pack Kaelbling 2002. Other concentrations like this appear all across the RL literature which makes me sceptical of the academic basis of many of these methodologies.

### III. METHODOLOGY

#### IV. DESIGN AND DEVELOPMENT

#### V. EXPERIMENTS AND EVALUATION

#### VI. DISCUSSIONS AND REFLECTIVE ANALYSIS

#### VII. CONCLUSION

#### REFERENCES

- Amato, Christopher and Guy Shani (2010). “High-level Reinforcement Learning in Strategy Games”. In: *Proceedings of the 9th International Conference on Autonomous Agents and Multiagent Systems*. Vol. 1. Toronto, Canada, pp. 75–82. ISBN: 978-0-9826571-1-9.
- Barto, Andrew G. and Sridhar Mahadevan (2003). In: *Discrete Event Dynamic Systems* 13.4, pp. 341–379. DOI: 10.1023/a:1025696116075.
- Bellemare, Marc G., Joel Veness, and Michael H. Bowling (2012). “Investigating Contingency Awareness Using Atari 2600 Games”. In: *AAAI*.
- Bellman, Richard (1958). “Dynamic programming and stochastic control processes”. In: *Information and Control* 1.3, pp. 228–239. DOI: 10.1016/s0019-9958(58)80003-0.
- Busoniu, L., R. Babuska, and B. De Schutter (2008). “A Comprehensive Survey of Multiagent Reinforcement Learning”. In: *IEEE Transactions on Systems, Man, and Cybernetics Part C (Applications and Reviews)* 38.2, pp. 156–172. DOI: 10.1109/tsmcc.2007.913919.
- Gaskett, Chris, David Wettergreen, and Alexander Zelinsky (1999). “Q-Learning in Continuous State and Action Spaces”. In: *Proceedings of the 12th Australian Joint Conference on Artificial Intelligence: Advanced Topics in Artificial Intelligence*. AI ’99. London, UK, UK: Springer-Verlag, pp. 417–428. ISBN: 3-540-66822-5. URL: <http://dblp.acm.org/citation.cfm?id=646077.678231>.
- Kaelbling, L. P., M. L. Littman, and A. W. Moore (1996). “Reinforcement Learning: A Survey”. In: *Journal of Artificial Intelligence Research* 4, pp. 237–285. DOI: 10.1613/jair.3013.
- Kober, Jens, J. Andrew Bagnell, and Jan Peters (2013). “Reinforcement learning in robotics: A survey”. In: *The International Journal of Robotics Research* 32.11, pp. 1238–1274. DOI: 10.1177/0278364913495721.

Ng, Andrew Y. et al. (2006). “Autonomous Inverted Helicopter Flight via Reinforcement Learning”. In: *Experimental Robotics IX*. Ed. by Marcelo H. Ang and Oussama Khatib. Berlin, Heidelberg: Springer Berlin Heidelberg, pp. 363–372. DOI: 10.1007/11552246\_35.

Rummery, Gavin Adrian and M. Niranjan (1994). “On-line Q-learning using connectionist systems”. In:

Shapley, L. S. (1953). “Stochastic Games”. In: *Proceedings of the National Academy of Sciences* 39.10, pp. 1095–1100. DOI: 10.1073/pnas.39.10.1095.

Shoham, Yoav, Rob Powers, and Trond Grenager (June 2003). “Multi-Agent Reinforcement Learning: a critical survey”. In:

Smart, W.D. and L. Pack Kaelbling (2002). “Effective reinforcement learning for mobile robots”. In: *Proceedings 2002 IEEE International Conference on Robotics and Automation*. DOI: 10.1109/robot.2002.1014237.

Sutton, R S. and A G. Barto (2018). *Reinforcement Learning: An Introduction*. Second edition. MIT Press.

Watkins, Christopher J. C. H. and Peter Dayan (1992). “Q-learning”. In: *Machine Learning* 8.3, pp. 279–292. DOI: 10.1007/BF00992698.

Watkins, Christopher John Cornish Hellaby (May 1989). “Learning from Delayed Rewards”. PhD thesis. Cambridge, UK: King’s College. URL: [http://www.cs.rhul.ac.uk/~chrisw/new\\_thesis.pdf](http://www.cs.rhul.ac.uk/~chrisw/new_thesis.pdf).

Wirth, Christian et al. (2017). “A Survey of Preference-Based Reinforcement Learning Methods”. In: *Journal of Machine Learning Research* 18, 136:1–136:46.

#### APPENDIX

```

1 import sys
2 import gym
3 import numpy as np
4 import time
5 from timeit import default_timer as timer
6 import matplotlib.pyplot as plt

# Initialize environment and q table
def init_env(initialisation = 'greedy'):
    env = gym.make('CartPole-v0').env
    n_states, n_actions = env.observation_space.n, env.action_space.n

# Initialise q-table with supplied type
if initialisation == 'zeros':
    Q = np.ones((n_states, n_actions))
elif initialisation == 'random':
    Q = np.random.uniform((n_states, n_actions))
elif initialisation == 'zeros':
    Q = np.zeros((n_states, n_actions))

return Q, env, n_states, n_actions

# e-Greedy algorithm for action selection from the q table by state with flag
# to force greedy method for testing
def e_greedy(Q, epsilon, n_actions, s, greedy=False):
    if greedy or np.random.rand() < epsilon:
        a = np.argmax(Q[s])
    else:
        a = np.random.randint(0, n_actions)

    return a

# Standard Q-Learning control method
def q_lrn(alpha, gamma, epsilon, episodes, max_steps, n_tests, \
    initialisation = 'greedy', render=True, test=True):
    # Initialise environment to get q table and action and state information
    Q, env, n_states, n_actions = init_env(initialisation)

    # Create list to store the timestep rewards

```

```

40 timestep_reward = []
41 timestep_max = []
42 timestep_min = []
43 resolution = 10
44 res = 0
45 timestep_reward_resolution = np.zeros(resolution)
46
47 # Iterate through episodes
48 for episode in range(episodes):
49     # Reset environment for new episode and get initial state
50     s = env.reset()
51
52     # Create values for recording rewards , steps and task completion
53     total_reward = 0
54     steps = 0
55     done = False
56
57     # Loop the task until task is completed or max steps are reached
58     while True:
59         steps += 1
60
61         a = e_greedy(Q, epsilon, n_actions, s)
62
63         # Get next state from the chosen action and record reward
64         s_, reward, done, info = env.step(a)
65         total_reward += reward
66
67         # Select next action greedily based on next state
68         a_ = e_greedy(Q, epsilon, n_actions, s_, greedy=True)
69
70         # If the task is completed do not use next q-value to update Q
71         if done:
72             Q[s, a] += alpha * (reward - Q[s, a])
73         # If not done use the Bellman equation to update Q
74         else:
75             Q[s, a] += alpha * (reward + (gamma * Q[s_, a_] - Q[s, a]))
76
77         # Set next state to current state (Q-Learning) control policy
78         s = s_
79
80         # If the task is completed break and record the reward for graphing
81         if done:
82             timestep_reward_resolution[res] = total_reward
83
84             # If the resolution count is reached reset counter
85             if res == resolution - 1:
86                 res = 0
87             # If not iterate the resolution counter
88             else:
89                 res += 1
90
91         # Average the reward for the resolution and record for plotting
92         if episode % resolution == 0:
93             timestep_reward.append(np.average(\
94                 timestep_reward_resolution))
95             timestep_max.append(np.max(timestep_reward_resolution))
96             timestep_min.append(np.min(timestep_reward_resolution))
97
98         break
99
100     # If the max steps are reached (agent stuck in action loop)
101     if steps == max_steps:
102         timestep_reward_resolution[res] = total_reward
103
104         if res == resolution - 1:
105             res = 0
106         else:
107             res += 1
108
109         if episode % resolution == 0:
110             timestep_reward.append(np.average(\
111                 timestep_reward_resolution))
112
113         break
114
115     # Check if testing flag is passed (on by default)
116     if test:
117         # Call testing class and return metrics for display
118         avg_rwd, avg_stp, failed = test_qtable(Q, env, n_tests, n_actions, \
119             max_steps, render)
120
121     # Return the time stepped rewards, averages for rewards and steps and the
122     # the failure percentage
123     return timestep_reward, timestep_max, timestep_min, avg_rwd, avg_stp, failed
124
125 # Adapted SARSA control method

```

```

126 def sarsa(alpha, gamma, epsilon, episodes, max_steps, n_tests, \
127     initialisation = 'random', render=True, test=True):
128     # Initialise environment to get q table and action and state information
129     Q, env, n_states, n_actions = init_env(initialisation)
130
131     # Create list to store the timestep rewards
132     timestep_reward = []
133     resolution = 10
134     res = 0
135     timestep_reward_resolution = np.zeros(resolution)
136
137     # Iterate through episodes
138     for episode in range(episodes):
139         # Reset environment for new episode and get initial state
140         s = env.reset()
141
142         # Create values for recording rewards , steps and task completion
143         total_reward = 0
144         steps = 0
145         done = False
146
147         # Get initial action using e_greedy method
148         a = e_greedy(Q, epsilon, n_actions, s)
149
150         # Loop the task until task is completed or max steps are reached
151         while True:
152             steps += 1
153
154             # Get next state from the chosen action and record reward
155             s_, reward, done, info = env.step(a)
156             total_reward += reward
157
158             # Select next action based on next state using
159             a_ = e_greedy(Q, epsilon, n_actions, s_)
160
161             # If the task is completed do not use next q-value to update Q
162             if done:
163                 Q[s, a] += alpha * (reward - Q[s, a])
164             # If not done use the Bellman equation to update Q
165             else:
166                 Q[s, a] += alpha * (reward + (gamma * Q[s_, a_] - Q[s, a]))
167
168             # Set next state and action to current state and action
169             s, a = s_, a_
170
171             # If the task is completed break and record the reward for graphing
172             if done:
173                 timestep_reward_resolution[res] = total_reward
174
175                 # If the resolution count is reached reset counter
176                 if res == resolution - 1:
177                     res = 0
178                 # If not iterate the resolution counter
179                 else:
180                     res += 1
181
182             # Average the reward for the resolution and record for plotting
183             if episode % resolution == 0:
184                 timestep_reward.append(np.average(\
185                     timestep_reward_resolution))
186
187             break
188
189             # If the max steps are reached (agent stuck in action loop)
190             if steps == max_steps:
191                 timestep_reward_resolution[res] = total_reward
192
193                 if res == resolution - 1:
194                     res = 0
195                 else:
196                     res += 1
197
198                 if episode % resolution == 0:
199                     timestep_reward.append(np.average(\
200                         timestep_reward_resolution))
201
202             break
203
204         # Check if testing flag is passed (on by default)
205         if test:
206             # Call testing class and return metrics for display
207             avg_rwd, avg_stp, failed = test_qtable(Q, env, n_tests, n_actions, \
208                 max_steps, render)
209
210         # Return the time stepped rewards, averages for rewards and steps and the
211         # the failure percentage

```

```

212     return timestep_reward, avg_rwd, avg_stp, failed
213
214 # Test function to test the Q-table
215 def test_qtable(Q, env, n_tests, n_actions, max_steps, render, delay=0.1):
216
217     # Create array to store total rewards and steps for each test
218     step_rewards = np.zeros([n_tests, n_tests])
219     # Set failed counter to zero
220     failed = 0
221
222     # Iterate through each test
223     for test in range(n_tests):
224
225         # Reset the environment and get the initial state
226         s = env.reset()
227         # Set done flag to false
228         done = False
229
230         # Set step and reward counters to zero
231         steps = 0
232         total_reward = 0
233
234         # Set e-greedy parameters greedy flag sets greedy method to be used
235         epsilon = 0.9
236         greedy = True
237
238         # Loop until test conditions are met iterating the steps counter
239         while True:
240             steps += 1
241
242             # Get action by e-greedy method
243             a = e_greedy(Q, epsilon, n_actions, s, greedy)
244
245             # If render flag has been passed render the environment
246             if render:
247                 time.sleep(delay)
248                 env.render()
249                 print(f'Action {a} chosen for state {s}')
250
251             # Get state by applying the action to the environment and add reward
252             s, reward, done, info = env.step(a)
253             total_reward += reward
254
255             # If task is completed break
256             if done:
257                 #print(f'Episode reward: {total_reward}, Steps: {steps}')
258                 # If render flag is set render and pause
259                 if render:
260                     time.sleep(1)
261                 break
262
263             # If step limit is reached end and record the failure
264             if steps == max_steps:
265                 # If render flag is set render and wait for user input
266                 if render:
267                     print(f'Action {a} chosen for state {s}, maximum steps {a}')
268                     env.render()
269                     input()
270
271                 #total_reward = 0
272                 #steps = 0
273                 failed += 1
274                 break
275
276             # Record total rewards and steps
277             step_rewards[0, test] = total_reward
278             step_rewards[1, test] = steps
279
280         # Get averages of the steps and rewards and failure percentage for tests
281         avg_rwd = np.average(step_rewards[0])
282         avg_stp = np.average(step_rewards[1])
283         fail_percent = (failed / n_tests) * 100
284
285         # Print average test values for all tests
286         print(f'Average reward {avg_rwd}, steps {avg_stp}, failure {fail_percent} %')
287         print('-----')
288
289     return avg_rwd, avg_stp, failed
290
291 # Plotting function to plot timestep rewards to show how the average agent
292 # reward increases over the training period by the specified resolution
293 def plot(data, maxim, minim, mthd):
294     plt.title(mthd)
295     plt.xlabel(' Episodes ')
296     plt.ylabel(' Rewards ')
297     plt.plot(data)

```

```

298     plt.plot(maxim)
299     plt.plot(minim)
300
301     plt.show()
302
303     # Start timer, set run length, set length of hyperparameters, and define
304     # corresponding lists of hyperparameters to be used
305     start = timer()
306
307     runs = 10
308
309     initialisation = 'ones' # ones, zeros or random
310
311     epsilon = 0.9
312
313     max_steps = 500
314     n_tests = 1000
315
316     eps = 1
317     gammas = 1
318     alphas = 1
319
320     arr_eps = [50000, 2500, 5000, 10000, 50000]
321     arr_gamma = [0.999, 0.99, 0.9, 0.9999]
322     arr_alpha = [0.05, 0.02, 0.03, 0.04, 0.01, 0.06, 0.07, 0.08, 0.09, 0.1]
323
324     # Iterate through each level of hyperparameters, episodes, gammas and alphas
325     for e in range(eps):
326         for g in range(gammas):
327             for a in range(alphas):
328
329                 # Set hyper parameters
330                 episodes = arr_eps[e]
331                 gamma = arr_gamma[g]
332                 alpha = arr_alpha[a]
333
334                 # Check if runs os greater then 3 to a void indexing errors
335                 if runs >= 3:
336                     # If so create aggregate array to average the total runs value
337                     aggr = np.zeros([runs, runs, runs])
338
339                 # Iterate through runs
340                 for i in range(runs):
341                     # Check of first argument is passed
342                     if len(sys.argv) > 1:
343                         # Check method flag for s (SARSA) or q (Q-Learning)
344                         if sys.argv[1] == 's':
345                             # Check second argument is passed
346                             if len(sys.argv) > 2:
347                                 # Check if render flag is set or not
348                                 if sys.argv[2] == 'r':
349                                     # Start SARSA function with render flag set
350                                     data, maxim, minim, avg_rwd, avg_stp, failed = sarsa(alpha, \
351                                         gamma, epsilon, episodes, max_steps, \
352                                         n_tests, initialisation, render=True)
353                                 elif sys.argv[2] == 'n':
354                                     data, maxim, minim, avg_rwd, avg_stp, failed = sarsa(alpha, \
355                                         gamma, epsilon, episodes, max_steps, \
356                                         n_tests, initialisation, render=False)
357                             else:
358                                 print(sys.argv[2], 'is not an option use "r" or "n"')
359                                 sys.stdout.flush()
360
361                         else:
362                             print('Wrong render option "r" or "n" or "\'')
363                             sys.stdout.flush()
364
365                 # Same as above for Q-Learning function
366                 elif sys.argv[1] == 'q':
367                     if len(sys.argv) > 2:
368                         if sys.argv[2] == 'r':
369                             data, maxim, minim, avg_rwd, avg_stp, failed = q_lrn(alpha, \
370                                 gamma, epsilon, episodes, max_steps, \
371                                 n_tests, initialisation, render=True)
372                         elif sys.argv[2] == 'n':
373                             data, maxim, minim, avg_rwd, avg_stp, failed = q_lrn(alpha, \
374                                 gamma, epsilon, episodes, max_steps, \
375                                 n_tests, initialisation, render=False)
376                     else:
377                         print(sys.argv[2], 'is not an option use "r" or "n"')
378                         sys.stdout.flush()
379
380                 else:
381                     print('Wrong render option "r" or "n" or "\'')
382                     sys.stdout.flush()
383
384             else:

```

```

384         print(sys.argv[1], 'is not valid, use "q" or "t"')
385
386     # Check of third argument is passed
387     if len(sys.argv) > 3:
388         # Check plot flag and plot data if so
389         if sys.argv[3] == 'q':
390             plot(data, maxim, minim, sys.argv[1])
391         elif sys.argv[3] == 't':
392             pass
393         else:
394             print(sys.argv[3], 'is not valid, use "q" or "t"')
395             # to plot rewards
396     else:
397         print('Please give argument "q" or "t" at position 3')
398
399     else:
400         print('Please a method "q" or "t" at position 3')
401
402     # If the runs threshold is met record values from testing
403     if runs >= 3:
404         aggr[0, i], aggr[1, i], aggr[2, i] = avg_rwd, avg_stp, \
405             failed
406
407     # Print the method used for clarity
408     print(sys.argv[1])
409
410     # If runs threshold is met print the averages and standard
411     # deviation of the rewards, steps and failures
412     if runs >= 3:
413         print('Total average reward', np.average(aggr[0]), \
414             np.std(aggr[0]), 'Steps', np.average(aggr[1]), \
415             np.std(aggr[1]), 'Failures', np.average(aggr[2]), \
416             np.std(aggr[2]))
417
418     # Print hyper parameters for testing
419     print('Episodes', episodes, 'gamma', gamma, 'Alpha', alpha)
420
421 # End timer and print time
422 end = timer()
423 print('Time', end-start)

```

---

## ACKNOWLEDGEMENTS