# MLSALT7 REINFORCEMENT LEARNING COURSEWORK

Author: Yeziwei (Maggie) Wang

## Contents

## 1 Introduction

This coursework explores four different basic algorithms in reinforcement learning by applying them on grid-world models. There are three different grid model used: smallworld, gridworld and cliffworld. Value iteration and policy iteration are applied on grid world to find the optimal policy. Both algorithms are examples of dynamic programming, which is a model-based reinforcement learning method. On the other hand, SARSA and Q-learning are both model-free reinforcement learning algorithms. In the experiment, they both are applied on smallworld individually and then compared against each other on cliffworld. The following concepts from reinforcement learning are also discussed in the report: differences and similarities between on-policy and off-policy learning, Dynamic Programming, Monte Carlo methods and Temporal Different Learning.

## 2 Experiments

### 2.1 Value Iteration

---
**Algorithm 1:** Value Iteration

---
Initialise $V_0$ arbitrarily;

**repeat**

   $\left| \quad V_{k+1}(s) = \max_a \sum_{s',r} p(s',r|s,a)(r + \gamma V_k(s'));\right.$

**until** $\max_s |V_{k+1}(s) - V_k(s)| < \theta$;

---

As a type of a dynamic programming method, value iteration is a model-based algorithm. It evaluates the value functions and at the same time update the policy. Its update rule is the same as the policy evaluation except that it calculates the maximum value over all actions at a specific state. The action which gives the maximum value function is used to update the policy. This results in simple backup and fast convergence. The algorithm stops after the value function has converged like shown in pseudocode of Algorithm 1. Since the backup takes the maximum reward, the algorithm should converge to an optimal policy.

Figure 1 shows the value function and policy the algorithm converges to in gridworld. The algorithm successfully learns the correct policy from the start state to goal state avoiding the obstacle given the model of the environment.
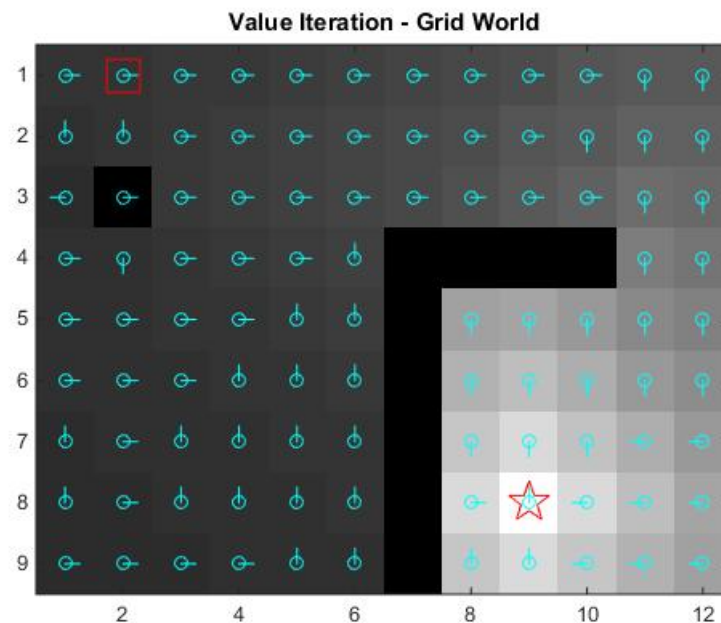


Figure 1: Gridworld: Value Iteration

In the algorithm implementation shown below, the policy is recorded every iteration while calculating the update for value function. The new policy over-write the previous one resulting in optimal policy to be saved at the end of the algorithm.

```
1   function [v, pi] = valueIteration(model, maxit)
2   % initialize the value function
3   v = zeros(model.stateCount, 1);
4
5   for i = 1:maxit,
6       % initialize the policy and the new value function
7       pi = ones(model.stateCount, 1);
8       v_ = zeros(model.stateCount, 1);
9       % perform the Bellman update for each state
10      for s = 1:model.stateCount,
```

```
11          [v_(s),pi(s)] = max(sum(squeeze(model.P(s,:,:))...
12                      .*(repmat(model.R(s,:),model.stateCount,1) ...
13                      + repmat(model.gamma*v,1,4)),1));
14      end
15      % exit early
16      if max(abs(v-v_)) < 0.00001,
17          fprintf('value function converge after %d iterations', i);
18          break;
19      end
20      v = v_; %update
21  end
```

## 2.2 Policy Iteration

---
**Algorithm 2:** Policy Iteration

---
Initialise $V$ and $\pi$ arbitrarily;

**repeat**

    Evaluate $V$ using $\pi$: $V_{k+1}(s) = \sum_a \pi(a,s) \sum_{s',r} p(s',r|s,a)(r + \gamma V_k(s'))$;

    Improve $\pi$ using $V$: $\pi'(s) = \arg\max_a \sum_{s',r} p(s',r|s,a)(r + V_\pi(s'))$

**until** $\max_s |V_{k+1}(s) - V_k(s)| < \theta$;

---

Policy iteration is also a model-based algorithm. It contains two steps in each iteration: 1) policy evaluation: updating the value function using current policy, 2) policy improvement: find the action that gives the maximum reward in the next step using updated value function. Alternating between these two steps can guarantee improvement over iterations and therefore the value function and policy will converge to optimal within finite iterations. This also provides fast convergence as each iteration is an improvement from previous one.

The final result of value function and optimal policy are shown in Figure 2. The optimal policy finds the path from start state to the goal state while avoiding obstacles in the world, where the black blocks are. The optimal policy is the same as what value iteration algorithm produced. So is the value function. The results from both value iteration and policy iteration shows great agreement between value function and policy, where the lighter the colour the closer it is to the target and all actions try to lead to where the target is.
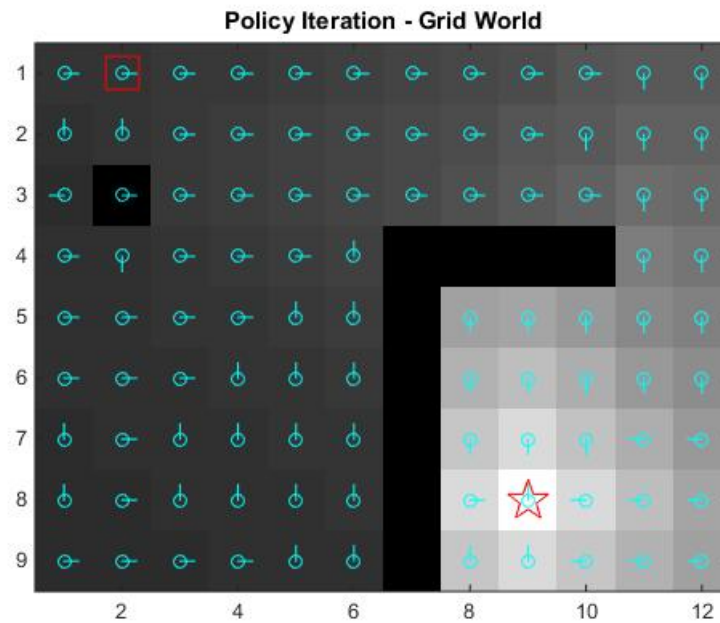
Figure 2: Gridworld: Policy Iteration

In the implementation of the algorithm, policy evaluation is conducted over all states before policy improvement is applied. This is to avoid updating policy before the evaluation is finished, so the policy improvement produces guaranteed improvement in value function and vice versa. The convergence condition is the same as value iteration.

```matlab
1  function [v, pi] = policyIteration(model, maxit)
2
3  % initialize the value function
4  v = zeros(model.stateCount, 1);
5  v_ = zeros(model.stateCount, 1);
6  pi = ones(model.stateCount, 1);
7
8  for i = 1:maxit,
9      v = v_;
10     for s = 1:model.stateCount  % policy evaluation
11         temp = sum(squeeze(model.P(s,:,:))...
12                 .*(repmat(model.R(s,:),model.stateCount,1)...
13                 + repmat(model.gamma*v,1,4)),1);
14         v_(s) = temp(pi(s));
15     end
16     for s = 1:model.stateCount  % policy imporovement
17         [¬,pi(s)] = max(sum(squeeze(model.P(s,:,:))...
18                 .*(repmat(model.R(s,:),model.stateCount,1)...
19                 + repmat(v_,1,4)),1));
20     end
21
22     % exit early
23     if max(abs(v_ - v))<0.0001
```

```
24          fprintf('Value function converged after %d itarations!', i)
25          break;
26      end
27
28  end
```

**Value Iteration and Policy Iteration**

|                  | Small world | Grid world | Cliff world |
|------------------|:-----------:|:----------:|:-----------:|
| Value Iteration  | 23          | 59         | 14          |
| Policy Iteration | 21          | 53         | 14          |

Table 1: Number of iterations till convergence

Both value iteration and policy iteration converge fast in three different world and all find optimal policy(Table 1). In both algorithm we have prior knowledge of the environment and at each update step the two algorithms update the policy with action that gives the maximum reward. Due to these reasons, the algorithm can quickly computes the optimal value function and policy. Both algorithms computes the same value function and optimal policy shown in Figure 3.



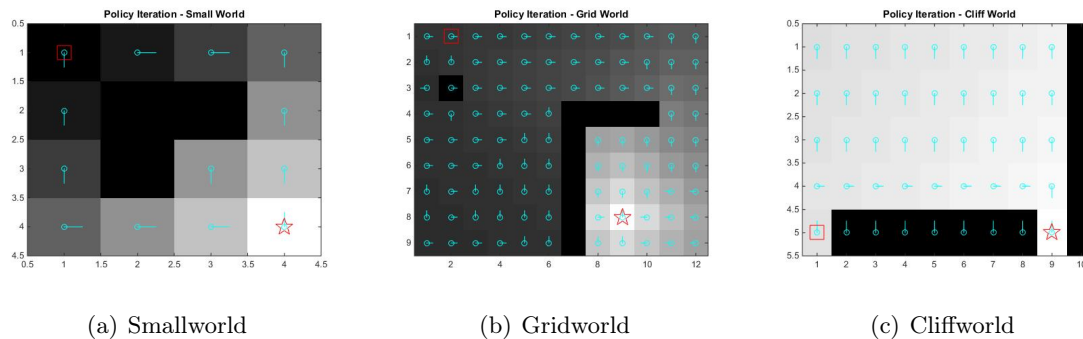(a) Smallworld　　　　(b) Gridworld　　　　(c) Cliffworld

Figure 3: Optimal value function and policy from value iteration and policy iteration

## 2.3  SARSA

---

**Algorithm 3:** SARSA algorithm

---

Initialise $Q$ arbitrarily, $Q(terminal, .) = 0$ ;

**repeat**

    Initialise s as the start state;

    Choose a $\epsilon$-greedily ;

    **repeat**

        Take action $a$, observe $r, s'$ ;

        Choose $a'$ $\epsilon$-greedily ;

        $Q(s, a) \leftarrow Q(s, a) + \alpha(r + \gamma Q(s', a') - Q(s, a))$ ;

        $s \leftarrow s', a \leftarrow a'$ ;

    **until** *s reaches the goal state*;

**until** $\max_s |V_{k+1}(s) - V_k(s)| < \theta$;

---

SARSA is an on-policy temporal-difference (TD) learning algorithm which has the advantages from both dynamic programming and Monte Carlo methods. Like dynamic programming, TD learning updates estimates based on other learned estimates, bootstrapping. Like Monte Carlo methods, it doesn't require knowledge of the environment and learns directly from raw experiences except that the update takes place at each time step:

$$V(s_t) \leftarrow V(s_t) + \alpha(r_{t+1} + \gamma V(s_{t+1}) - V(s_t))$$

where $\alpha > 0$ is the learning rate. SARSA learns the action-value function in the same manner. As an on-policy algorithm, it takes the action it evaluates. This is shown in the update step in Algorithm 3, where state and policy are both updated after the backup. The terminal state has zero value for Q-function. As SARSA learns from experience, for each episode the algorithm starts from the starting state and explores the world until it reaches the goal state. The Q-functon is updated after choosing next action $\epsilon$-greedily and it is updated using temporal difference in action-value function.
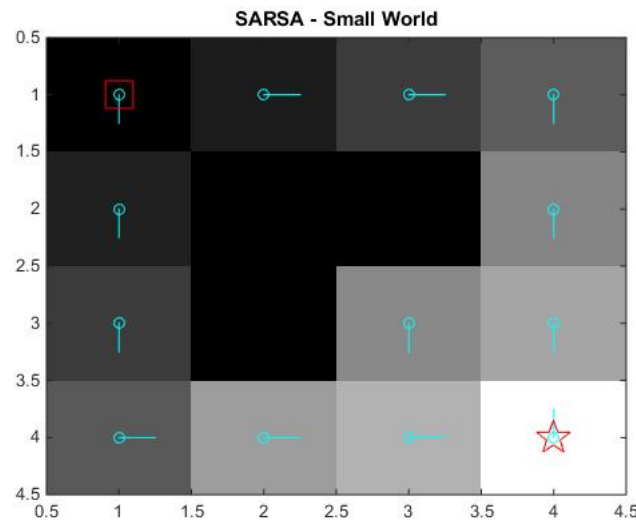
Figure 4: Smallworld: SARSA

The learning rate $\alpha$ is chosen as 0.3 and $\epsilon$ 0.1, so that the algorithm has a modest learning rate and keeps exploring throughout iterations. In smallworld the optimal policy is obtained. The same convergence criteria is used as in policy iteration. Although SARSA never reached the convergence criteria, it still produced optimal policy. This may be due to the exploration of the algorithm and the simplicity of the model itself. As the algorithm takes a random action every 10 steps, the value function is not able to converge to the one obtained from value iteration and policy iteration. The maximum number of steps for each episode is 30 and the maximum number of episode explored is 400. There are 16 states in smallworld, so a maximum number of 30 steps should be sufficient for the algorithm to reach the goal state. Once the goal stated is reached, current episode is terminated and the algorithm starts a new episode. Once the value function or when the maximum number of iteration is reached, the learning stops.

```matlab
1  function [v, pi, acc_R] = sarsa(model, maxit, maxeps, alpha, epsilon)
2  rand('seed',287)
3  % initialize the value function
4  Q = zeros(model.stateCount, 4);
5  v = zeros(model.stateCount, 1);
6  v_ = zeros(model.stateCount, 1);
7  pi = ones(model.stateCount, 1);
8  acc_R = zeros(maxeps,1);
9
10 for i = 1:maxeps,
11     % every time we reset the episode, start at the given startState
12     s = model.startState;
13     prob1 = rand;
14     if prob1 < epsilon  % epsilon greedily
15         a = datasample([1 2 3 4], 1); % radom action
16     else
17         [¬,a] = max(Q(s,:));  % optimal action
```

```matlab
18          end
19
20      for j = 1:maxit,
21          % PICK AN ACTION
22          p = 0;
23          r = rand;
24          for s_ = 1:model.stateCount,
25              p = p + model.P(s, s_, a);
26              if r ≤ p,
27                  break;
28              end
29          end
30          acc_R(i) = acc_R(i) + model.R(s,a); % accumulated reward
31          % pick next action epsilon greedily
32          prob2 = rand;
33          if prob2 < epsilon
34              a_ = datasample([1 2 3 4], 1);
35          else
36              [¬,a_] = max(Q(s_,:));
37          end
38
39          Q(s,a) = Q(s,a) + alpha*(model.R(s,a)+model.gamma*Q(s_,a_)-Q(s,a));
40          [v_(s),pi(s)] = max(Q(s,:));
41          s = s_; a = a_;
42          % once reach goal state, start a new episode
43          if s == model.goalState
44              fprintf('%d: reached the goal state! Start Again!\n', i)
45              break;
46          end
47      end
48      if max(abs(v_-v))<0.001
49          fprintf('value function converged after %d episodes!', i)
50          break;
51      end
52      v = v_;
53  end
```

## 2.4 Q-learning

Similar to SARSA, Q-learning is also a TD learning algorithm. Unlike SARSA, it is an off-policy algorithm, which means the policy it evaluates is different from the action it takes. The backup always uses the maximum Q-function without updating policy. As the algorithm is learning the optimal action-value function directly, it learns the optimal policy and suggests early convergence. As Q-learning is a model-free algorithm, its convergence is slower compared to model-based learning.

The optimal policy is obtained in the smallworld shown in Figure 5. The maxmum number

of steps in each episode is 30, as the grid only contains 16 states. If the algorithm couldn't find the goal state in 30 steps it should start again. And maximum munber of episode is 400. The $\epsilon$ and $\alpha$ values are set the same as in SARSA for the same reasons.

---

**Algorithm 4:** Q-Learning algorithm

---

Initialise $Q$ arbitrarily, $Q(terminal, .) = 0$;

**repeat**

    Initialise s ;

    **repeat**

        Choose a $\epsilon$-greedily ;

        Take action $a$, observe $r, s'$ ;

        $Q(s, a) \leftarrow Q(s, a) + \alpha(r + \gamma \max_{a'} Q(s', a') - Q(s, a))$ ;

        $s \leftarrow s'$ ;

    **until** *s is terminal*;
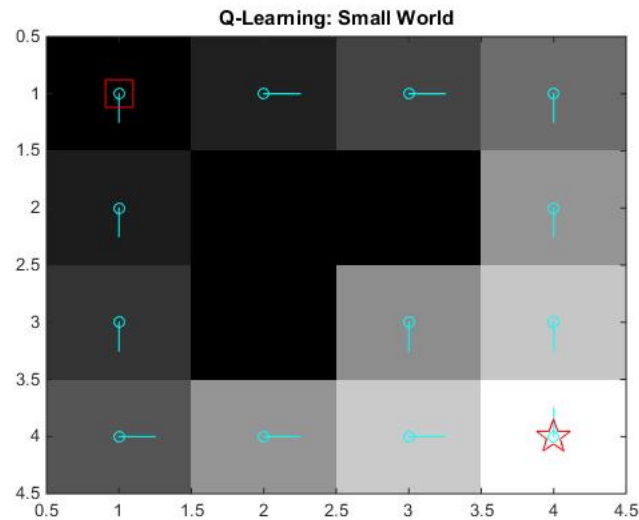
**until** $\max_s |V_{k+1}(s) - V_k(s)| < \theta$;

---



Figure 5: Smallworld: Q-learning

```matlab
function [v, pi, acc_R] = qLearning(model, maxit, maxeps, alpha, epsilon)
rand('seed',287)
% initialize the value function
Q = zeros(model.stateCount, 4);
v = zeros(model.stateCount, 1);
v_ = zeros(model.stateCount, 1);
pi = ones(model.stateCount, 1);
acc_R = zeros(maxeps, 1);

for i = 1:maxeps,
    s = model.startState;

    for j = 1:maxit,
```

```matlab
14          % PICK AN ACTION
15          prob1 = rand;
16          if prob1 < epsilon
17              a = datasample([1 2 3 4], 1);
18          else
19              [¬,a] = max(Q(s,:));
20          end
21
22          % look for next state s_ after action a
23          p = 0;
24          r = rand;
25          for s_ = 1:model.stateCount,
26              p = p + model.P(s, s_, a);
27              if r ≤ p,
28                  break;
29              end
30          end
31          acc_R(i) = acc_R(i) + model.R(s,a);
32          Q(s,a) = Q(s,a) + alpha*(model.R(s,a) + model.gamma*max(Q(s_,:)) - ...
                  Q(s,a));
33          [v_(s), pi(s)] = max(Q(s,:));
34          s = s_;
35
36          % start new episode after reaching goal state
37          if s == model.goalState
38              fprintf('%d: reached the goal state! Start Again!\n', i)
39              break;
40          end
41      end
42
43      if max(abs(v_-v))<0.00001
44          fprintf('value function converged after %d episodes!', i)
45          break;
46      end
47      v = v_;
48  end
```
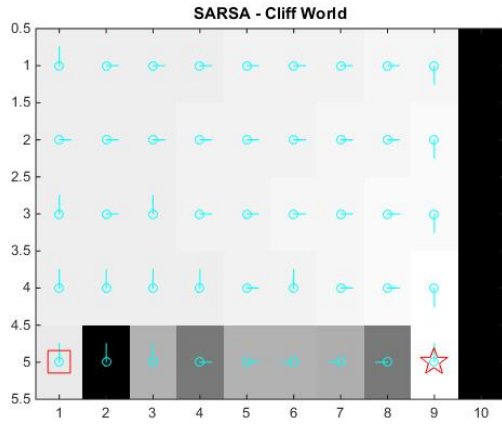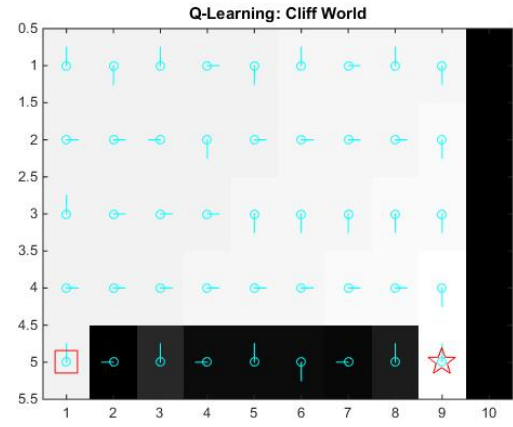
## 2.5 SARSA and Q-learning

SARSA and Q-learning are compared against each other in the cliffworld. The different performance is due to different update rule. Q-learning learns directly the optimal policy, which is the path along the edge of the cliff, by using maximum Q-function in its update rule. On the other hand, SARSA considers all actions in a particular state while updating Q-function, thus it learns a safer path further away from the cliff. This ensures that when algorithm chooses action $\epsilon$-greedily, it doesn't fall off the cliff.
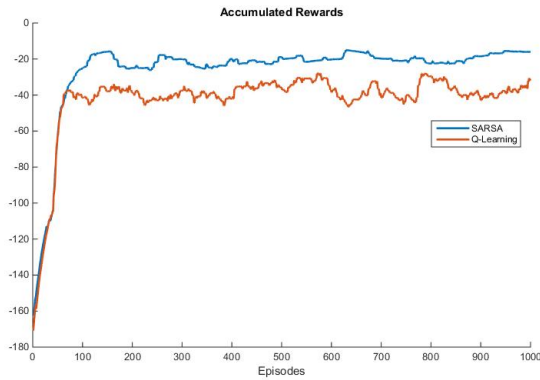
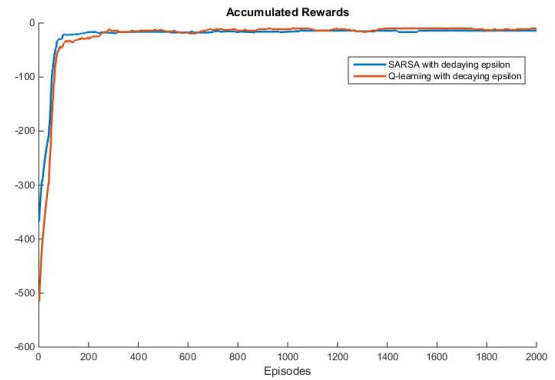<center>(a) SARSA             (b) Q-learning</center>

<center>Figure 6: Cliff World</center>

The accumulated rewards show the same results (Figure 7(a)). Although Q-learning learns the optimal policy, it's more risky. When the algorithm explores, it will fall off the cliff resulting in worse overall accumulated reward. On the contrary, SARSA learns a safer path further away from the cliff, so even if the agent explores, the accumulated reward would still be higher. The accumulated reward is calculated over 1000 iterations with a smoothing window size of 80 iterations.



<center>(a) $\epsilon = 0.1$             (b) Decaying $\epsilon$</center>

<center>Figure 7: Accumulated rewards over iterations</center>

If $\epsilon$ gradually reduces over time, the algorithm eventually stops exploring. In this case, the accumulated rewards from SARSA and Q-learning are going to merge as shown in Figure 7(b). However, this doesn't change the policy SARSA converges to, unless the algorithm visits all states infinite times.

```
1  %% --- accumulated reward --- %%
2  % smoothing accumulated reward cliff world
3  window_size = 80;
```

```matlab
4  for i=1:length(R3) % SARSA
5      if i ≤ window_size/2
6          smooth_R3(i) = mean(R3(1:i+window_size/2));
7      elseif i ≥ length(R3)-window_size/2
8          smooth_R3(i) = mean(R3(i-window_size/2:length(R3)));
9      else
10         smooth_R3(i) = mean(R3(i-window_size/2:i+window_size/2));
11     end
12 end
13
14 for i=1:length(Q_R3) % Q-learning
15     if i ≤ window_size/2
16         smooth_Q_R3(i) = mean(Q_R3(1:i+window_size/2));
17     elseif i ≥ length(Q_R3)-window_size/2
18         smooth_Q_R3(i) = mean(Q_R3(i-window_size/2:length(Q_R3)));
19     else
20         smooth_Q_R3(i) = mean(Q_R3(i-window_size/2:i+window_size/2));
21     end
22 end
23 % plot accumulated rewards
24 figure;hold on
25 plot(smooth_R3); plot(smooth_Q_R3)
26 title 'Accumulated Rewards'
27 xlabel 'Episodes'
```

**Explore combinations of $\epsilon$ and $\alpha$ for SARSA**

After trying different combinations of epsilon and alpha, the following observations are made.

First of all, to converge to a policy using SARSA the scale of $\epsilon$ and $\alpha$ should match. This means when using larger $\epsilon$, $\alpha$ should be increased accordingly. The level of exploration and learning rate should match in order to learn from experience. The more an agent explores, the learning rate should be higher in order to take into account all the experiences.

Secondly, the more the algorithm explores, the more it tends to converge to a safer path. This can be shown in figure 8. When $\epsilon = 0.3$ the algorithm has higher chance of falling off the cliff. In order to reach goal state, SARSA takes the path that's furtherest from the cliff. Also, as $\epsilon$ decreases, it starts to take paths that are closer to the cliff until the $\epsilon$ is small enough. When $\epsilon$ is very small, the algorithm hardly explores, it has less chance of falling off cliff, so it converges to the optimal path next to the cliff.

Finally, the more the algorithm explores, the more updates we have on all states. This will give better estimates for the final value function. Ideally, the cliff should be black in the plot. In practice, they tends to be different shades of grey. This is because they are not visited enough times to have their true value function estimated. When $\epsilon$ is larger, the algorithm explores more, therefore, the cliff has darker colour. The cliff next to the start state is black because it is very

likely to be visited. This means its value was updated more than other cliff positions. So the algorithm has more information to give a better estimate of its value function.



(a) $\epsilon = 0.3, \alpha = 0.5$

(b) $\epsilon = 0.1, \alpha = 0.3$

(c) $\epsilon = 0.1, \alpha = 0.1$

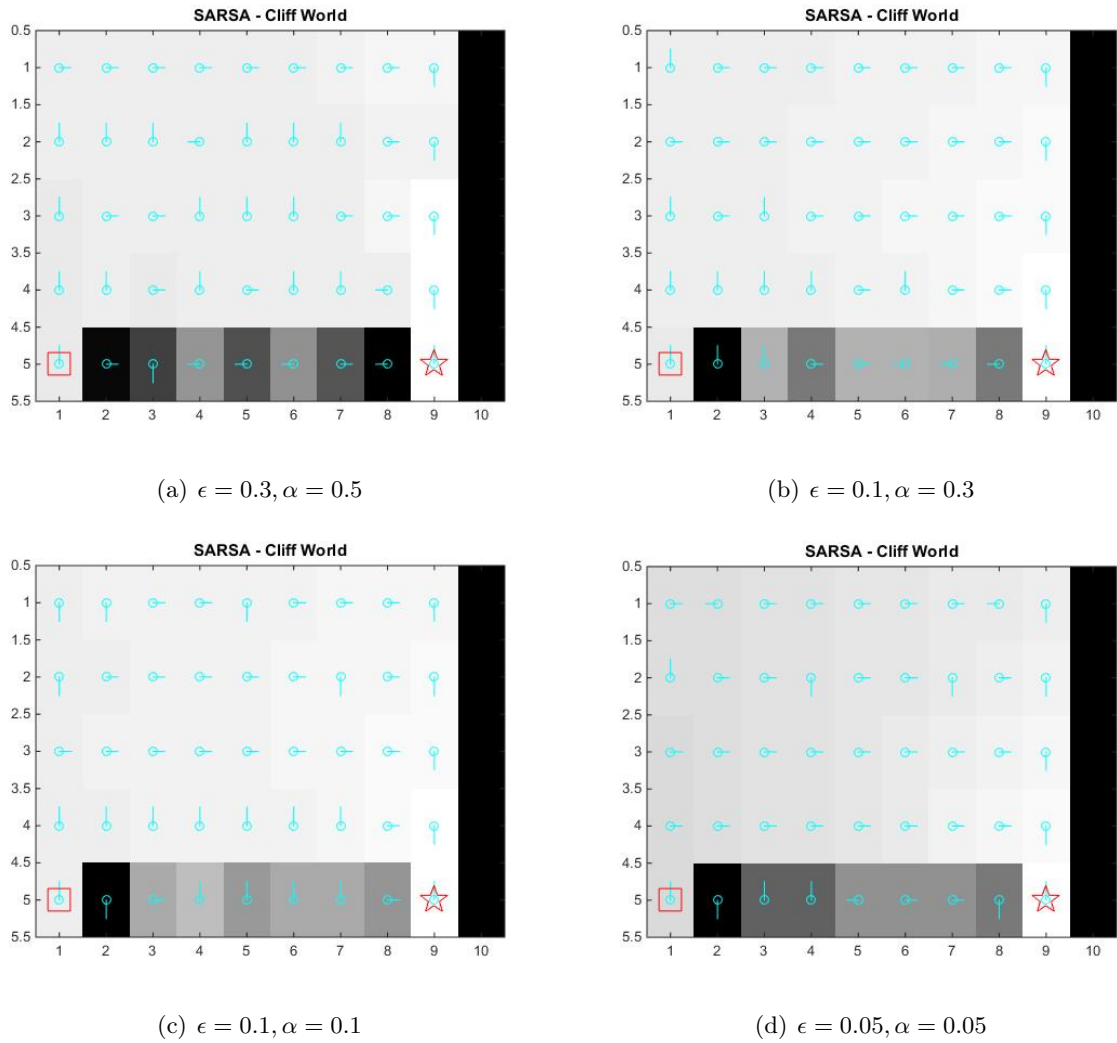(d) $\epsilon = 0.05, \alpha = 0.05$

Figure 8: Cliffworld:policy under different combinations of $\epsilon$ and $\alpha$
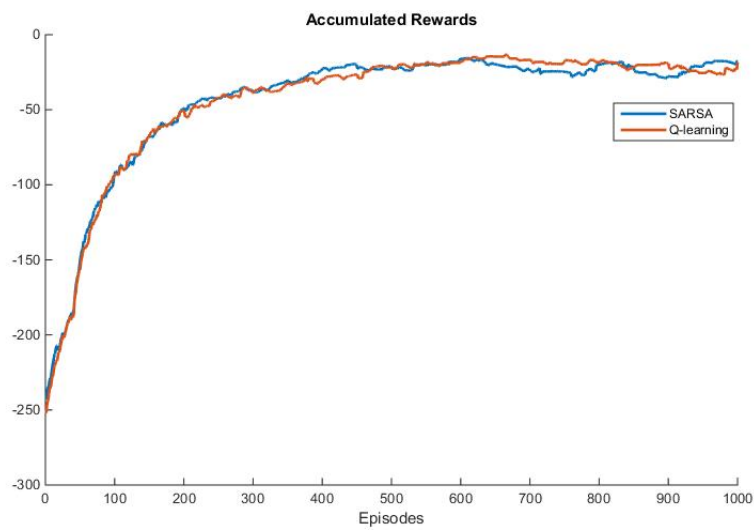


Figure 9: Accumulated reward when $\epsilon = 0.05, \alpha = 0.05$

Figure 9 shows the accumulated rewards from SARSA and Q-learning. In this particular case, SARSA was able to learn the optimal policy like Q-learning. As the epsilon is very small the accumulated reward from the two algorithms have the similar values.

## 3    Conclusion

There are three major types of reinforcement learning algorithm: dynamic programming, Monte Carlo methods and temporal difference learning. This report explored and implemented four basic algorithms in reinforcement learning, which shows the similarity and difference between the three types of methods mentioned.

Value iteration and policy iteration, examples of dynamic programming method, both learn from a known environment, therefore, they can both learn the optimal policy with fast convergence. Value iteration combines the two steps from policy iteration updating the value function and policy at the same time. This results in a simpler implementation. SARSA and Q-learning are TD learning algorithms, which takes the advantage from dynamic programming and Monte Carlo methods. SARSA is more aware of the risks involved in the actions, unlike Q-Learning which learns the optimal policy directly, it learns a safer non-optimal path. Different combinations of learning rate and epsilon value result in SARSA converging to different policy. This is due to the exploration nature of SARSA and on-policy methods. If SARSA has a decaying epsilon value and it visits all states an infinite number of times, it will eventually converge to the optimal policy and have the same accumulated rewards as Q-learning.

## References

[1] Sutton, R S and Barto A G (1998) *Reinforcement Learning: An Introduction.* The MIT Press