# REINFORCEMENT LEARNING ASSIGNMENT № 1

Ruslan Burakov, s1569105                                        11/02/2015

## Note

All the code fragments are located in the end of the report in the appendix section. In order to reduce number of pages I don't attach unit tests.
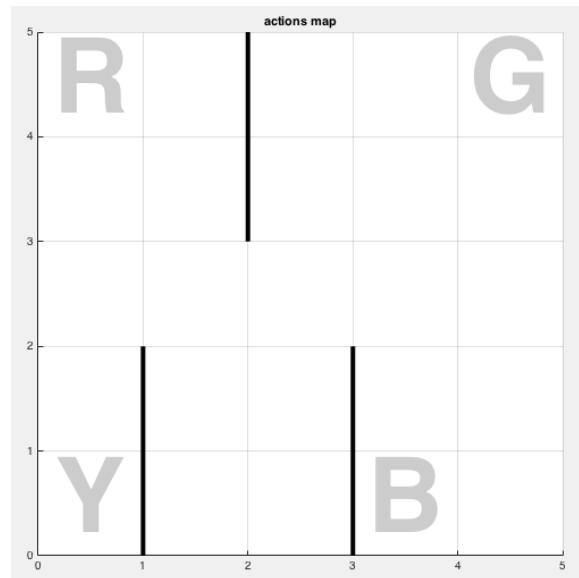
## 1    Theoretical estimation



Figure 1: Taxi problem maze

### 1.1    Typical time horizon of the problem

The time horizon is the time needed for value functions in order to converge:

$$V = E(r_0 + \gamma r_1 + \gamma^2 r_2 + \dots) = E(\sum_{t=0}^{T} \gamma^t r_t) \tag{1}$$

Where $E$ - means expectation, $V$ - value function, $r_t$ - immediate reward on the time-step t, $\gamma$ - discount factor, and T - is either $\infty$ or expected duration of the episode. Here in the taxi problem we deal with the episodic problem as we here a terminal state (the passenger is dropped out in the target destination). The duration of the episode for the optimal policy is somewhere between 1 step which describes the best case when we start with the passenger within taxi in the exact drop-off location (for example in R location on Figure 1) and the worst scenario of $9*2 = 18$ when we start in Y location and the passenger in G, and drop-off location is Y again. The probability

of the best case is small $4\frac{1}{25}\frac{1}{5}\frac{1}{4} \approx 0.008$ (using the number of possible states in the maze and the fact that there are only 4 best possible cases: R, G, Y, B) and the worst case probability is small as well $\frac{1}{25}\frac{1}{5}\frac{1}{4} \approx 0.002$. That is why I expect average time horizon to be close to 10 steps. Using formula for sum of geometric progression in the lecture it was shown that time horizon can be approximately connected to discount factor $\gamma$ by $time\ horizon = \frac{1}{1-\gamma}$. Therefore, the $\gamma$ is $1 - \frac{1}{T} = 1 - \frac{1}{10} = 0.9$

## 1.2 Mean (immediate) reward for a good policy

In the good policy agent doesn't bump into the walls or try to pick/drop off passenger in the wrong location so cumulative reward for one episode must be positive. There are two positive rewards in this problem: successful pick up and successful drop off. For successful pick on average we get $\frac{1}{5} * 0 + \frac{4}{5} * 1 = 0.8$ where 0 corresponds to the case when the passenger appears within the taxi from the beginning. The drop off reward is calculated using the time since last pick up $r_{dropoff} = \frac{10}{(time\ since\ pick)}$ and is determined by $4x5$ possible combinations of passenger and drop off locations. The knowledge of the number of optimal actions to do since pick up (for example to drop off, from G to R it will take 9 actions Figure1) allows to calculate average drop off reward $\frac{1}{16}(4*10 + 4*\frac{10}{9} + 2*\frac{10}{5} + 4*\frac{10}{8} + 2*\frac{10}{6}) = 3.549$ (omitting the case when the passenger is already in the taxi). Summing reward for the pick and average reward for drop off gives **4.549** as our estimation of average cumulative (means sum of all received immediate rewards during one episode) reward in one episode.

## 1.3 maximal value of a state-action pair

The maximal value state action pair is the any of best cases discussed above. It is the taxi in the drop off position and the passenger within it and its action is drop off. In my implementation once the **correct** drop off location is reached the passenger is out of the taxi and its location becomes drop off location and this is reflected in the states. The reason why it is important is that Q-learning converges to the right policy but not necessary right values [SB98]. Here I consider case for G drop off location but it will be the same for any other drop off location. Q(taxi in G, passenger within taxi, drop off in G, action - drop off) = Q(G, IN, G, drop) and Q(taxi in G, passenger in G, drop off in G, action - pick up) = Q(G, G, G, pick). So in my case I when I try to drop off my Q(G, IN, G, drop) updated as following:

$$Q_{t+1}(G, IN, G, drop) = Q_t(G, IN, G, drop) + \alpha(r_{drop} + \gamma max_a Q_t(G, G, G, a) - Q_t(G, IN, G, drop)) \tag{2}$$

The maximum value of the next state (although in reality there is nothing after terminal state but in the tutorial examples provided to us it was solved in the same way) there will be Q(G, G, G, pick) and:

$$Q_{t+1}(G, IN, G, drop) = Q_t(G, IN, G, drop) + \alpha(r_{drop} + \gamma Q_t(G, G, G, pick) - Q_t(G, IN, G, drop)) \tag{3}$$

The update formula for Q(G, G, G, pick) (remembering that reward for pick up is 1):

$$Q_{t+1}(G,G,G,pick) = Q_t(G,G,G,pick) + \alpha(1 + \gamma Q_t(G,IN,G,drop) - Q_t(G,G,G,pick)) \quad (4)$$

Assuming that our Q values will eventually converge (t goes towards very big value) $Q_{t+1} = Q_t$ we will receive the following system of linear equations:

$$r_{drop} + \gamma Q(G,G,G,pick) - Q(G,IN,G,drop) = 01 + \gamma Q(G,IN,G,drop) - Q(G,G,G,pick) = 0 \quad (5)$$

Solving this system we obtain:

$$Q(G,IN,G,drop) = \frac{r_{drop} + \gamma}{1 - \gamma^2} \quad (6)$$

estimating average reward $r$ for successful drop off similarly to previous section we get:

$$r_{drop} = (10 + 10/6 + 10/9 + 10/9) * \frac{1}{4} = 3.47 \quad (7)$$

After that substituting all the values into equation 6 we receive:

$$Q(G,IN,G,drop) = \mathbf{23.15} \quad (8)$$

## 1.4 number of trials

Well each of $5*5*5*4 = 500$ states must be visited at least several times. Let's set the number of times visited to the determined before time horizon of 10 steps. So roughly it is 5000 trials but of course it is worth investigating it more during computation
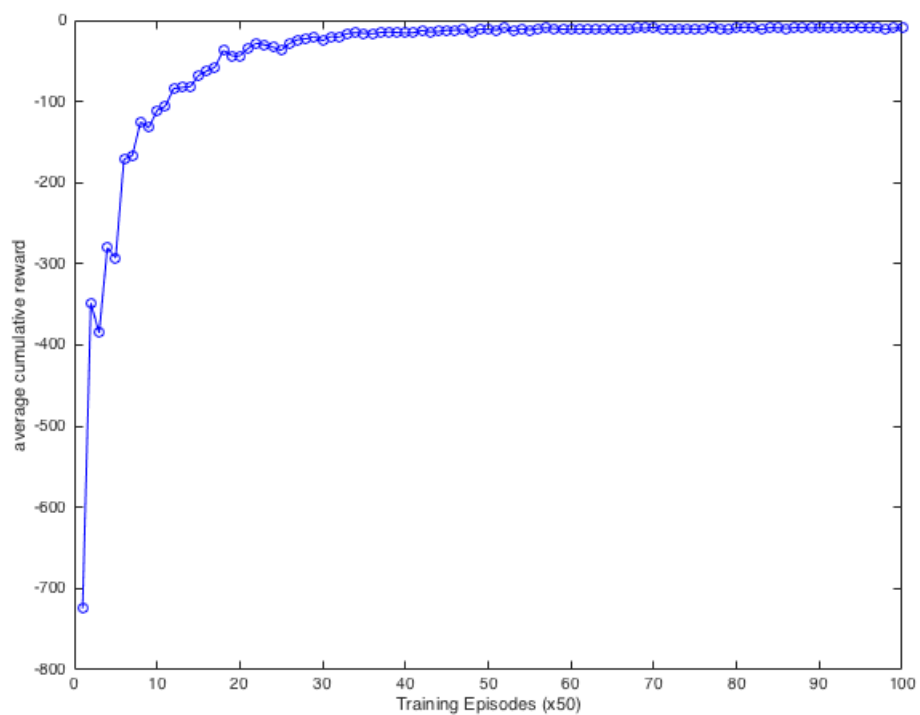
## 1.5 Stochastic vs Deterministic

We are using mixed exploring strategy (epsilon greedy, sometimes acting deterministically but with probability epsilon we are acting randomly). In the maze the results of our actions are deterministic given chosen action and current state. Given state representation optimal sequence of actions looks determined because all rewards are deterministic. So it points that it is deterministic problem but I note that from agent point of view the drop off reward is not always deterministic. Because the time since the last pick up is not a part of our state representation the agent can only guess and predict average reward for its drop off action although it is not driven from any distribution.
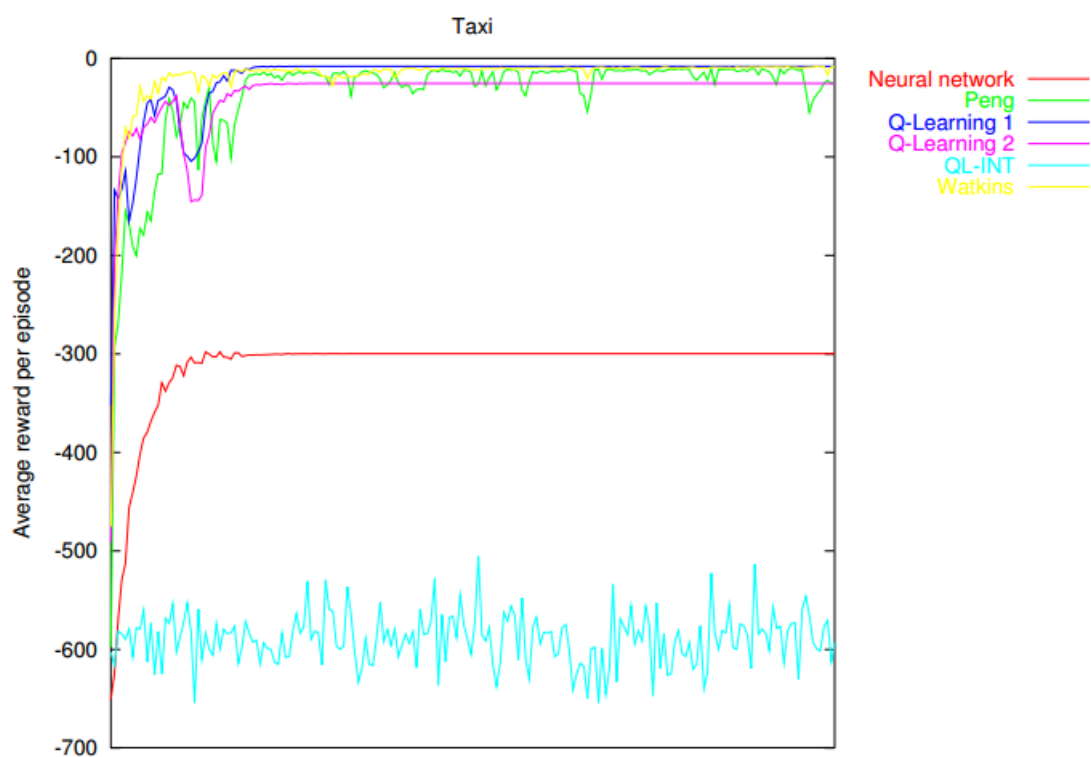
# 2 Solve the problem using Q-learning

## 2.1 Test correctness by comparing with known results

I wasn't sure whether my implementation (The Q learning code is in Q_learning_episode.m.) was right or wrong. That is why I have implemented reward function (code in trans_fun_paper.m) from the original benchmark [Dut+05] from which this homework was created in order to see the similarity between their results and my. I used parameters from there for QL1:

(a) Our reproduction of the paper [Dut+05] Q-learning 1 result



(b) Original plot from the paper [Dut+05] for taxi problem

Figure 2

| $\epsilon$ | $\epsilon$ decay | $\gamma$ | $\alpha$ | $\alpha$ decay |
|------|------|------|------|------|
| 0.5 | 0.999 | 0.9 | 0.5 | 0.99999 |

The result obtained after 5000 episodes can be seen on Figure 2. Here cumulative reward is averaged per 50 episodes.

As it can be seen the results for Q-learning 1 are converging to optimal policy with average cumulative reward around -5 in episode after $40 * 50 = 2000$ time steps. In the paper [Dut+05] the different reward function is used and -5 is the expected average cumulative reward for it. So it means our simulation converges to optimal policy.

How Q values converge can be judged by the maximum absolute difference in Q values between two episode of the problem. This measure indicates the order of fluctuation in Q values during the experiment. I am using this metric instead of variance because it shows the performance of the most unstable state so it plays role of upper-bound. After that by plotting the average of this metric (over 50 episodes, for example), it is possible to judge how fast the Q values are converging. On the other hand, the policy may converge to optimal earlier.
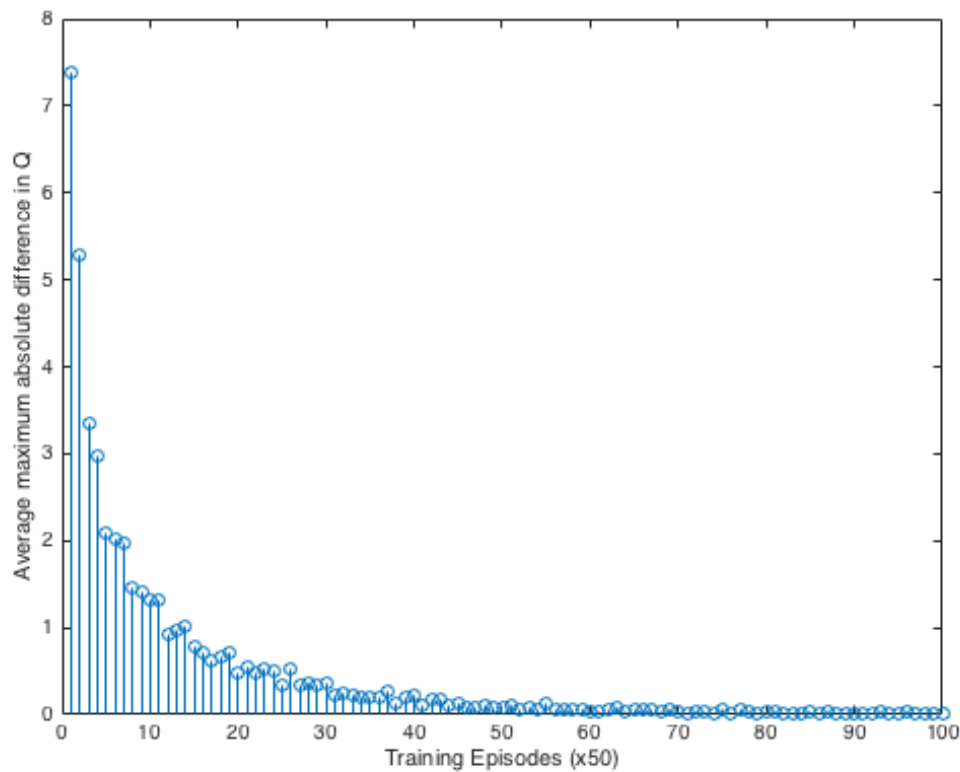


Figure 3: Q maximum absolute difference for reward function taken from paper [Dut+05]

Therefore, Figure 3 points out that after $45 * 50 = 2250$ time steps the Q-values magnitude of fluctuations goes below 0.4.

The match between paper results and my gave me additional confidence that q-learning implementation is correct.

Finally, I project final Q-values on the partially fixed state where passenger is within taxi and the drop off location is G on Figure 4. As it may be expected, the value in location G is the biggest
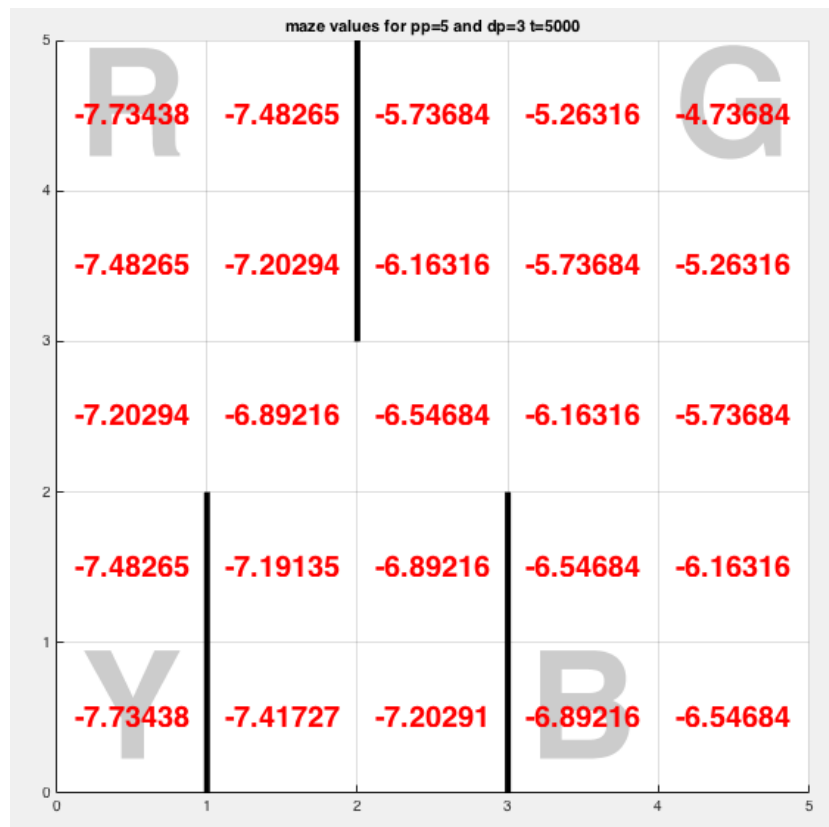
Figure 4: Values for case when passenger is within taxi and drop off location is in G for reward function [Dut+05]

one and all other values decreases very symmetrically from it depending on the distance.

## 2.2  Q learning performance with homework reward function

I tried Q-learning with parameters from the previous answer (with the exploration parameters such $\epsilon$ and $\alpha$ with corresponding exponential decays taken similarly from the paper [Dut+05] from previous section):

| $\epsilon$ | $\epsilon$ decay | $\gamma$ | $\alpha$ | $\alpha$ decay |
|---|---|---|---|---|
| 0.5 | 0.999 | 0.9 | 0.5 | 0.9999 |

On Figure 5 (note here it is averaged over 100 episodes) it can be seen that our rough estimation of cumulative reward in one episode was correct. The average cumulative reward value is close to 3.5 and the predicted maximum value was 4.549.

In order to draw the maximum value case as 2-maze I project my final Q-values on the partially fixed state where passenger is within taxi and the drop off location is G. The values obtained are on Figure 6. You can see that in the location G which corresponds to maximum value state the value is equal to 22.95. This is quite close to our estimation from question 1 which is 23.15.
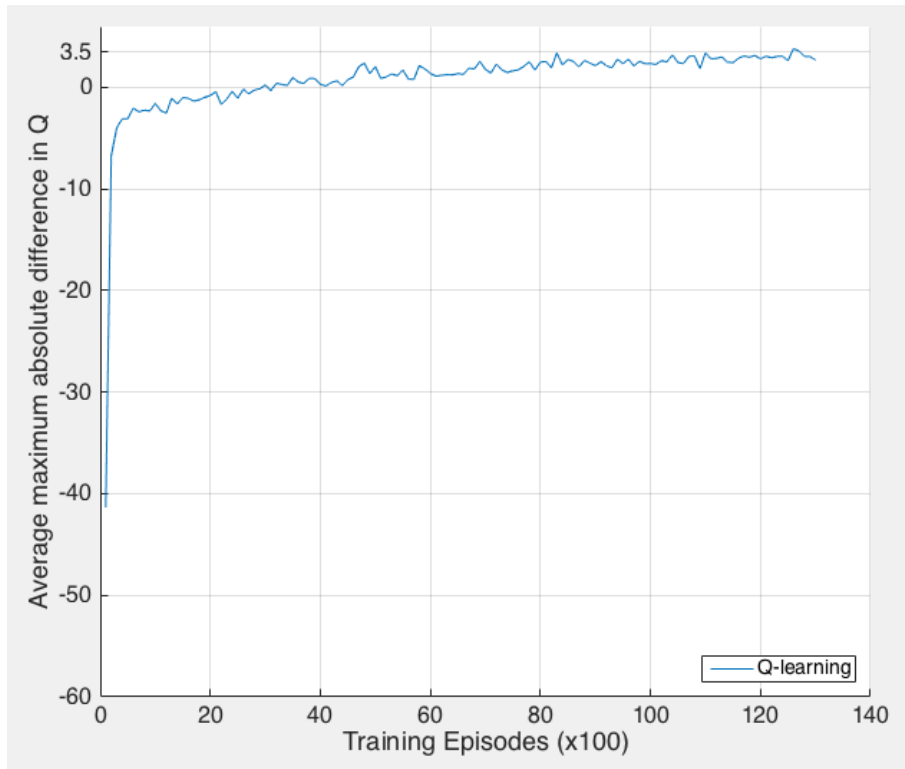
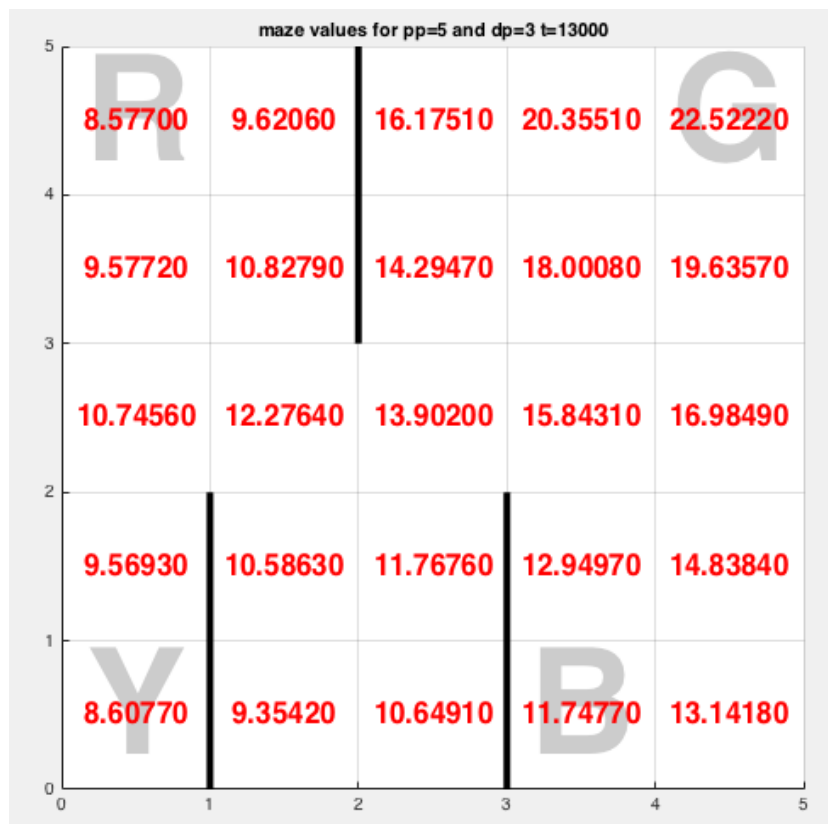Figure 5: Q learning average cumulative reward for homework reward function



Figure 6: Values for case when passenger is within taxi and drop off location is in G for howework reward function

## 2.3  Represent your solution using an example

I use the same example as in previous question where passenger is within taxi and the drop off location is G on Figure 6. It can be seen that in contrast to results obtained with reward function from the paper [Dut+05] the values are much less symmetric depending on the distance to the G location. The possible explanation of this is that fluctuation of values in simulation with homework reward function are much bigger. It explained further in next question. The policy generated for this example is on Figure 7. The generated policy is optimal in the sense that it will try to reach drop off location as fast as possible.
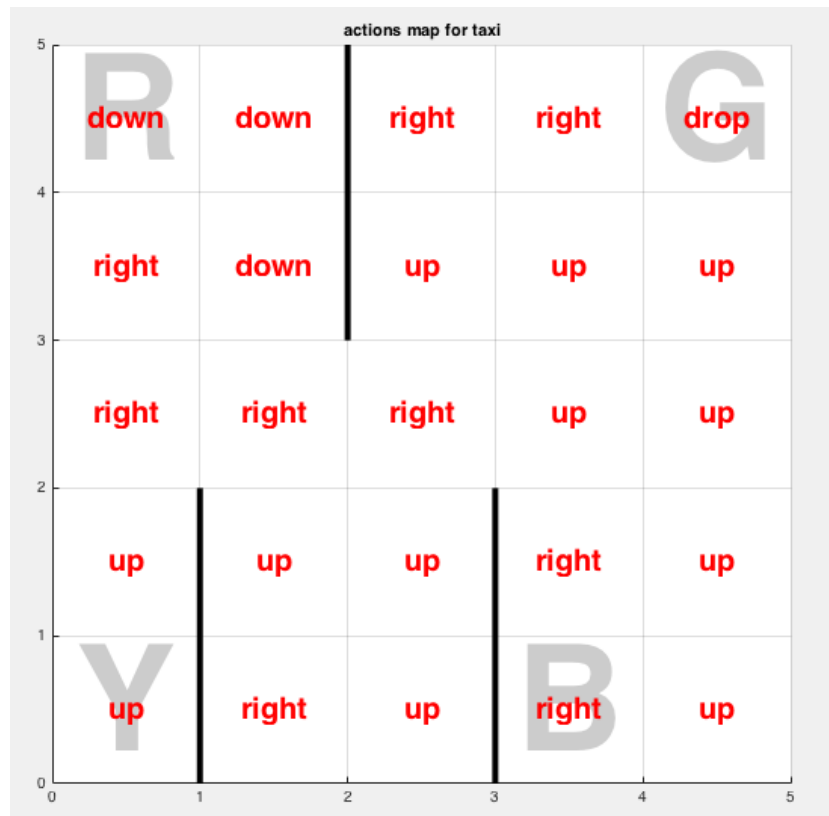


Figure 7: Best actions for case when passenger is within taxi and drop off location is in G

## 2.4  The time course of the values and the reward

On Figure 5 we can see that average cumulative reward converged (stabilized) on $100 * 100 = 10000$ time step. Approximately at the same time step maximum absolute difference on Figure 8 reaches the value of 0.7. However, not that the maximum fluctuation here is more than 5 times bigger than on Figure 3. This can be described by the fact that in the paper [Dut+05] the expected cumulative reward per episode is linear function of distance between pick up and drop off state, as they punish by -1 for each movement and give reward 0 for successful drop off. In our case the reward for successful drop off is 10/(time since pick) which is hyperbolic function and it changes much more suddenly. So depending on the random initial state in the episode the cumulative reward may vary significantly even with the optimal policy. In order to account such randomness the $\alpha$ parameter must be lower so more of the old values must be preserved and
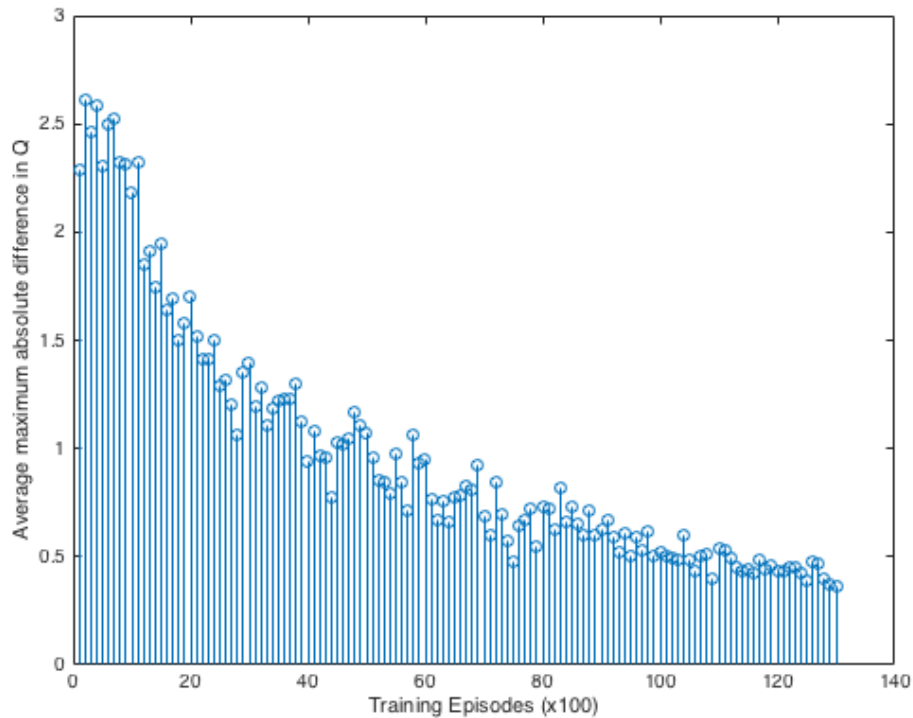
Figure 8: Q maximum absolute difference for homework reward function

the averaging of the Q-values will happen much smoother. But in that case the required number of time steps may increase drastically.

## 2.5 the convergence time depending on exploration strategy

Previously I was using exponential decays for exploration parameter $\epsilon$ of $\epsilon$-greedy strategy. That is why I have decided to fix value of $\epsilon$ (exponential decay parameter equals to 1.0) and see what convergence time I will achieve. All other parameters are taken from previous question and fixed:

| $\gamma$ | $\alpha$ | $\alpha$ decay |
|---|---|---|
| 0.9 | 0.5 | 0.9999 |

Judging by results obtained on Figure 9 the convergence time is roughly the same on the $3 * 100 = 300$ time step. Simulation with $\epsilon$ 0.5 converges to the cumulative reward close to -5 which is way below expected optimal one of 4.53. But this is expected as it acts randomly in half of the cases. The best result which is very close to the expected optimal cumulative reward is achieved by strategy with $\epsilon$ 0.01. This implies that in our problem even so little of exploration is enough to act nearly optimally. It can be described by the fact that we don't need to learn so much about all possible states. The most important information is that once we picked person we must reach drop off location with maximum speed.

After that I used greedy search approach and given that the most effective $\epsilon$ so far is 0.01 I have decided to find out what fixed value of $\alpha$ (interpolation parameter) will give the best results
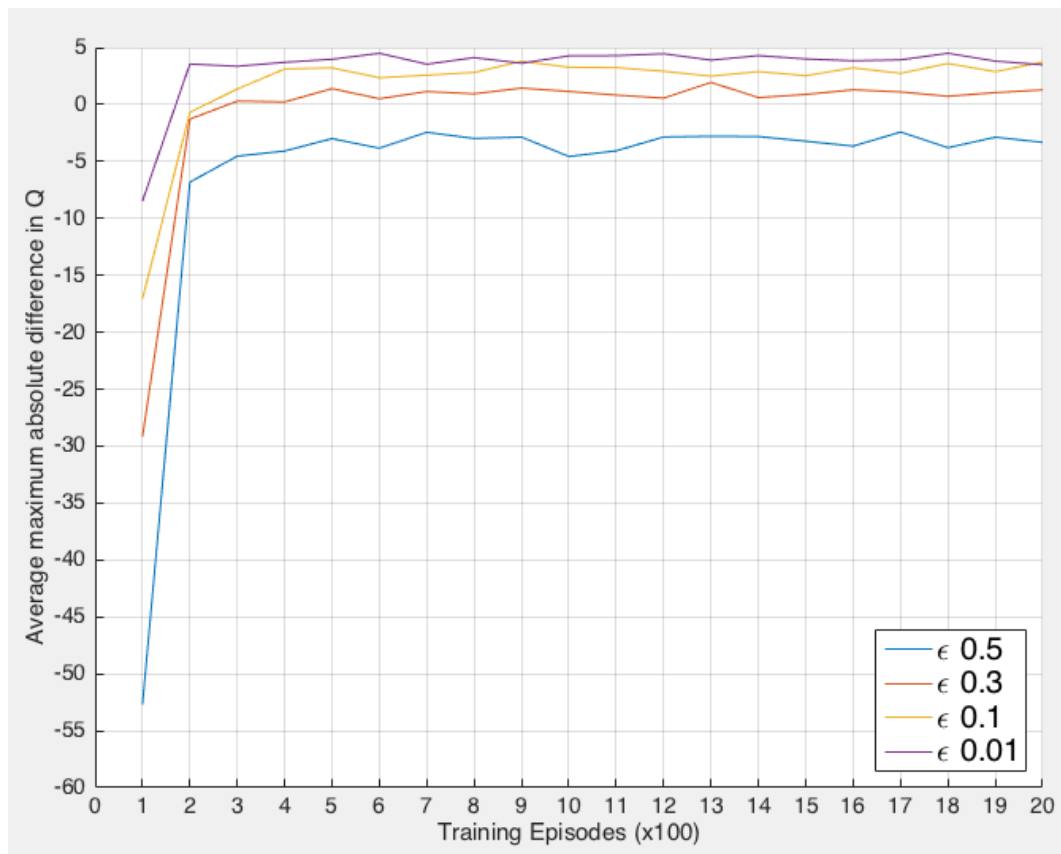
Figure 9: Average cumulative reward in one episode for different $\epsilon$

| $\gamma$ | $\epsilon$ | $\epsilon$ decay | $\alpha$ decay |
|---|---|---|---|
| 0.9 | 0.01 | 1.0 | 1.0 |

On Figure 10 the difference between difference $\alpha$ is almost negligible. Although $\alpha$ 0.9 converges slightly slower than others. This is expected as primal exploration parameter is $\epsilon$ and here with $\epsilon$ 0.01 we mostly update states greedily.

Finally, based on above results, I tried to get better average cumulative reward (closer to theoretical predicted one 4.549) by exploring for longer period of time.

| $\epsilon$ | $\epsilon$ decay | $\gamma$ | $\alpha$ | $\alpha$ decay |
|---|---|---|---|---|
| 0.5 | 0.9998 | 0.9 | 0.5 | 1.0 |

As we can see on Figure 11 the average cumulative reward is much closer to the theoretically predicted one.
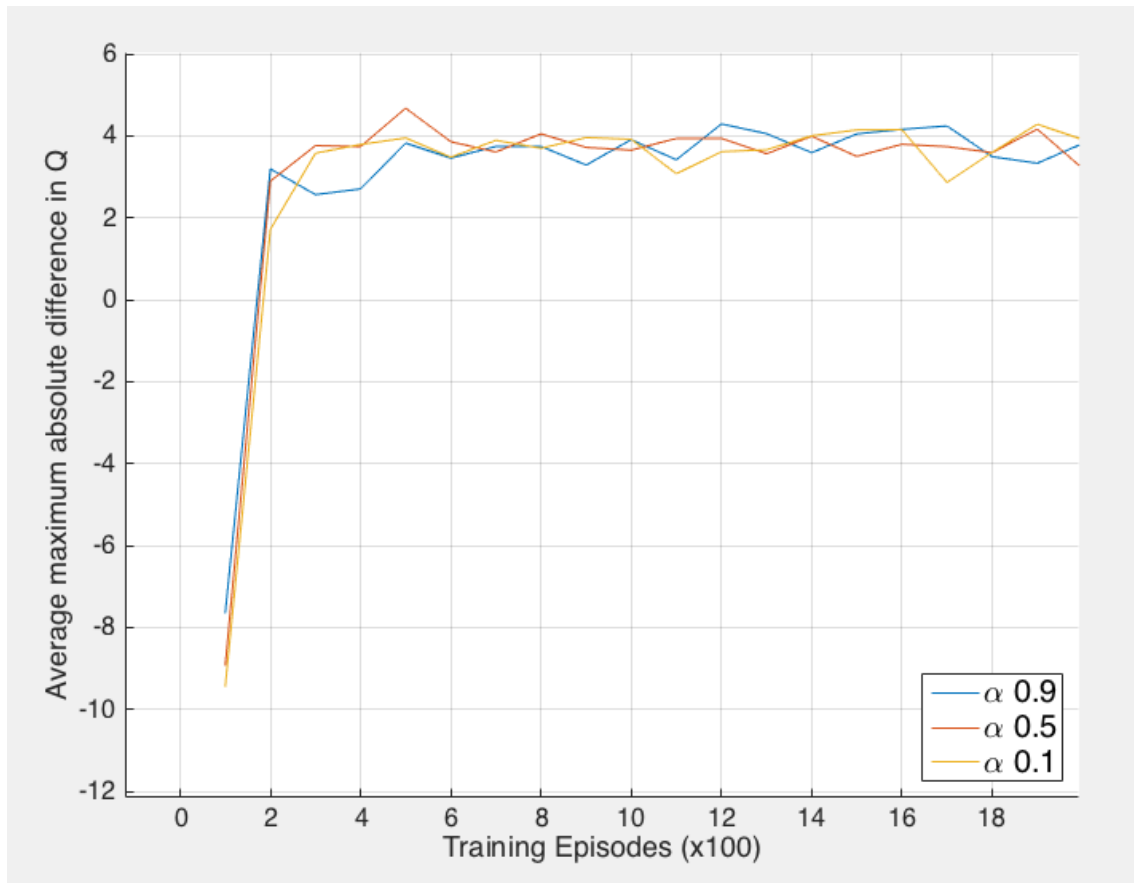
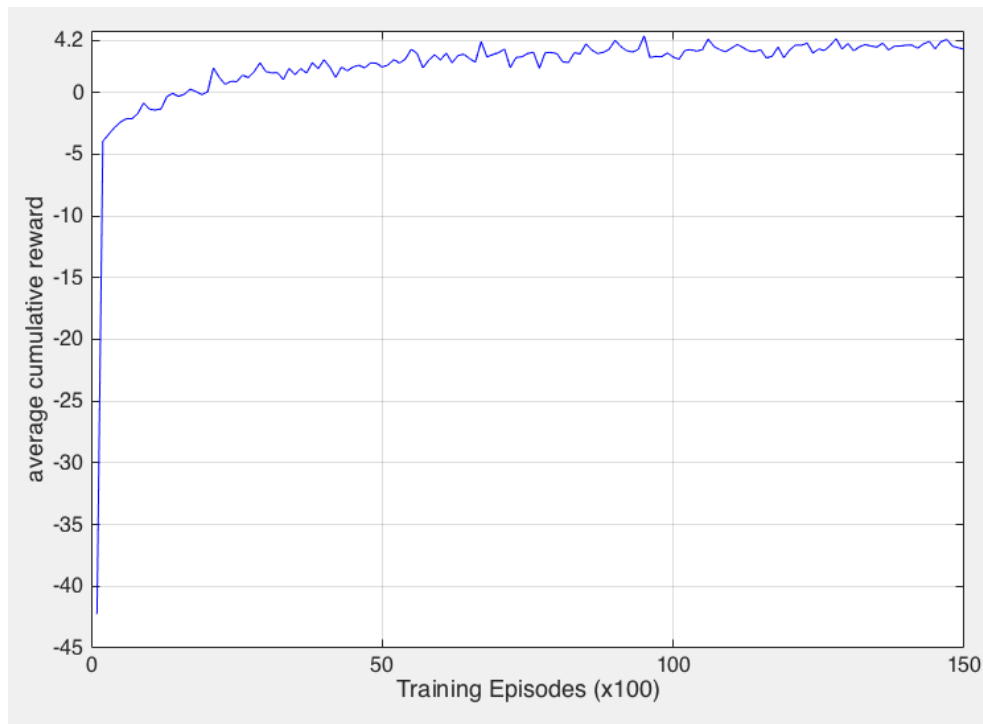Figure 10: Average cumulative reward in one episode for different $\alpha$



Figure 11: Average cumulative reward in one episode for Q-learning with longer exploring phase

## 3 SARSA

The SARSA code is in SARSA_episode.m. I am using initial parameters from section 2.2 for SARSA, so they can be comparable.

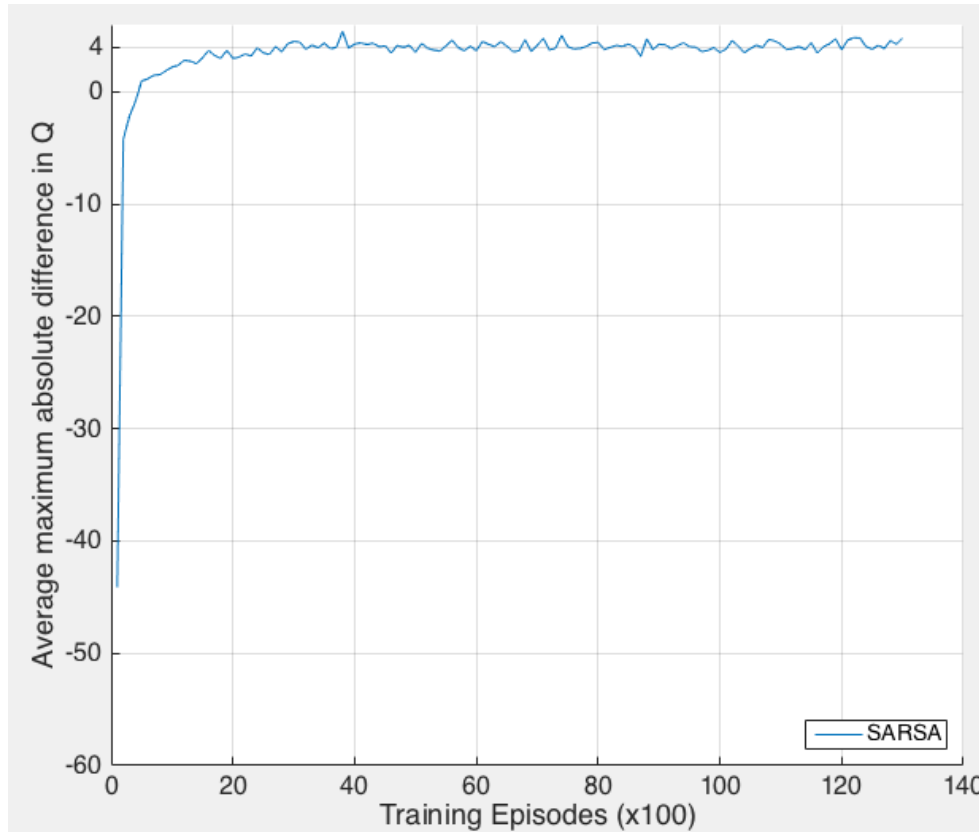| $\epsilon$ | $\epsilon$ decay | $\gamma$ | $\alpha$ | $\alpha$ decay |
|------|--------|-----|-----|--------|
| 0.5 | 0.9998 | 0.9 | 0.5 | 0.9999 |



Figure 12: Average cumulative reward in one episode using SARSA

Figure 12 indicates that SARSA converges to approximately 4.1 cumulative reward which is not far from of the expected 4.549 for the taxi problem.

Values for case when passenger is within taxi and drop off location is in G are on Figure 13. In my case that epsilon is decaying parameter so it becomes small relatively quickly. For instance, $0.5 * (0.999)^3 910 = 0.01$. So on time step 3910 exploration rate is 0.01 and behaviour of SARSA and Q-learning becomes very similar as they both start to act very greedily. That is why on Figure 13 the value G converges to almost the same expected value of 23.15 as in case of Q-learning.

The same reasoning can be applied to describe similarity of Q maximum absolute difference generated by SARSA on Figure 14 to the Q maximum absolute difference dynamic obtained using Q-learning on Figure 8.

The policy generated by SARSA for the case when passenger is within taxi and drop off location is in G on Figure 15. The notable distinction is that in 3 cells $(\{2,5\}, \{3,1\}, \{3,5\})$ the action chosen differs from what we would consider optimal. It means that SARSA didn't explore
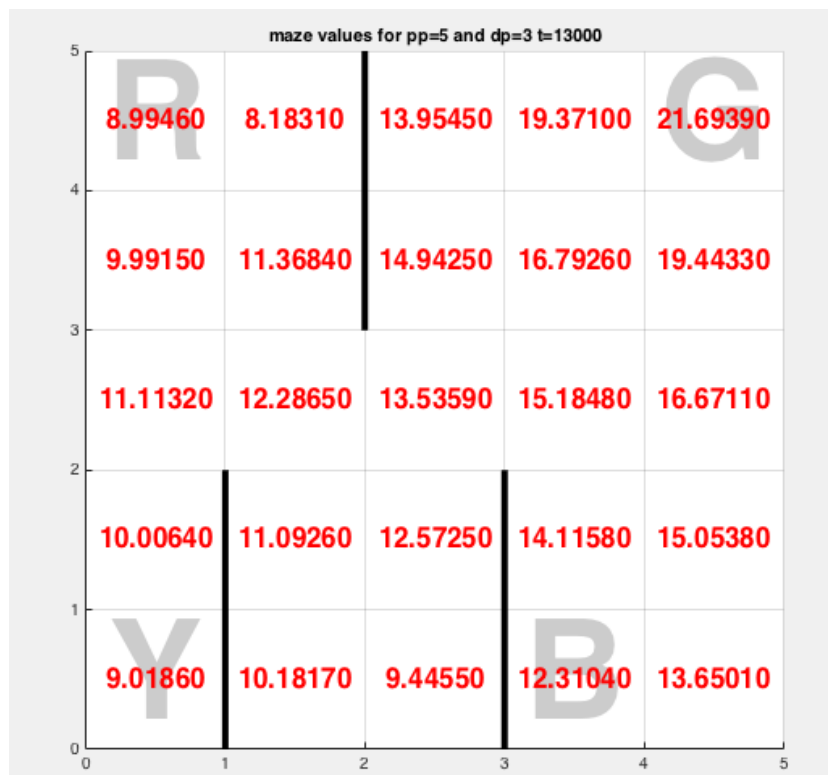
Figure 13: Final Values using SARSA for the case when passenger is within taxi and drop off location is in G
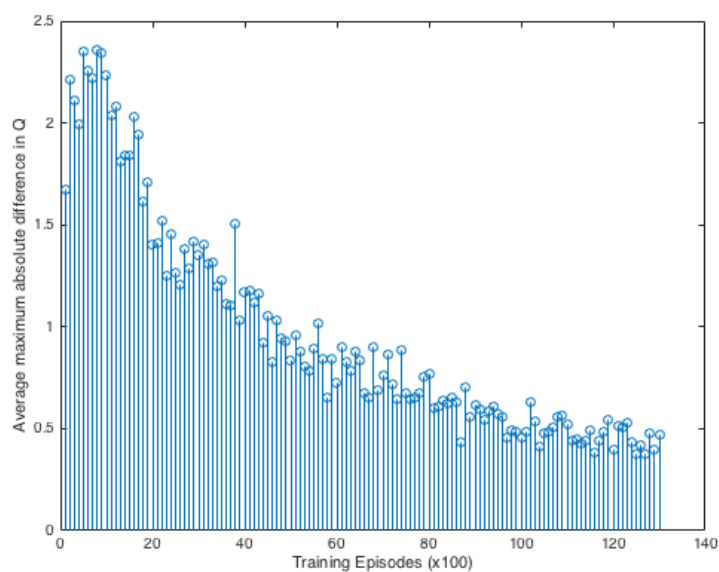


Figure 14: Q maximum absolute difference using SARSA

Figure 15: Generated policy using SARSA for the case when passenger is within taxi and drop off location is in G

these states enough. This can be describe by the fact that the probability of visiting these states is very small. Agent can appear in them only when exploring or when it was generated with the passenger within taxi in these particular cells and its drop off location is G, so it is $\frac{1}{25}\frac{1}{5}\frac{1}{4} \approx 0.002$. Due to this non-optimality doesn't noticeably affect average cumulative reward on Figure 12. Let's try another parameters with longer exploration (parameters the same as in the end of section 2.5).

| $\epsilon$ | $\epsilon$ decay | $\gamma$ | $\alpha$ | $\alpha$ decay |
|---|---|---|---|---|
| 0.5 | 0.9998 | 0.9 | 0.5 | 1.0 |

The results on Figure 16 for average cumulative reward are much more closer to theoretical value of 4.549. And this time the generated policy on Figure 17 is be optimal one.

Finally, I have run tests to find best fixed $\epsilon$ value using parameters as in previous question

| $\gamma$ | $\alpha$ | $\alpha$ decay |
|---|---|---|
| 0.9 | 0.5 | 0.9999 |

The results obtained on Figure 18 are very similar to those of Q-learning from previous question. The discussion why it can be like that in next question.

Figure 16: Average cumulative reward in one episode using SARSA with longer exploration phase
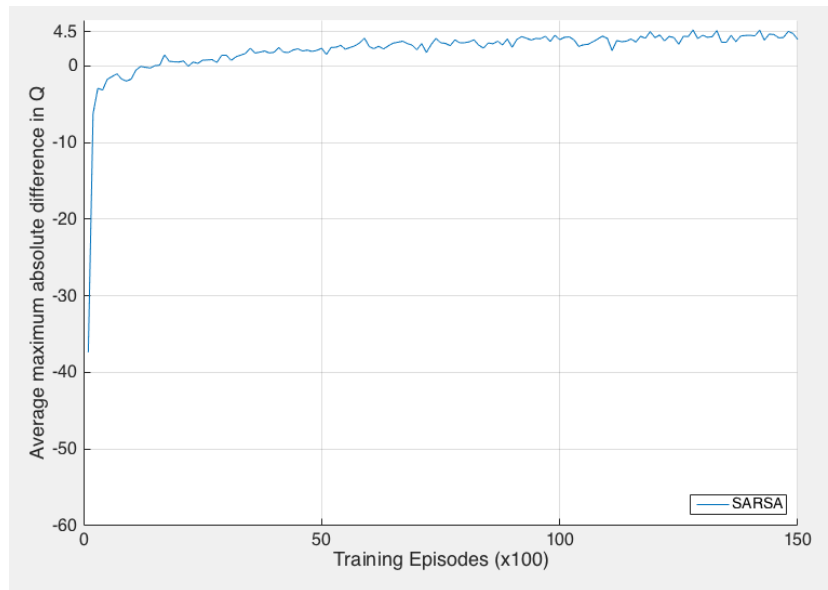


Figure 17: Generated policy using SARSA with longer exploration phase for the case when passenger is within taxi and drop off location is in G

Figure 18: Average cumulative reward in one episode for different $\epsilon$ using SARSA

# 4   Comparison of Q-learning and SARSA



Figure 19: Average cumulative reward comparison for Q-learning and SARSA in the case with longer exploration phase

In this task the difference between Q-learning and SARSA is almost negligible in contrast to cliff example from [SB98]. Comparing results for the cases with longer exploration phases from previous sections on Figure 19 we can see how close obtained results are. Such performance can be explained by the fact that when $\epsilon$ goes to 0 both algorithms start to act similarly. Another reason is that taxi problem doesn't have adequate cliff analogue. We can unsuccessfully pick or drop passenger in almost any state, so SARSA couldn't learn to avoid some dangerous states knowing that it may act randomly like cliff example. The walls of the maze can't act as cliff as well, because the penalty for them is -1 the average reward is bigger than that. And once again in the maze almost in any state there are walls around the agent. Although probably in this problem SARSA requires longer exploration phase as Figure 15 shows SARSA wasn't able to converge to completely optimal policy with short exploration phase in contrast to Q-Learning on Figure 7.

# 5 Improvements

Accordingly to the Berkley A.I. lecture 11 Reinforcement Learning II [Kle12] we can compress the information in our task by using features instead of direct states and approximate Q-values on:

## Approximate Q-Learning

$$Q(s,a) = w_1 f_1(s,a) + w_2 f_2(s,a) + \ldots + w_n f_n(s,a)$$

Q-learning with linear Q-functions:

$$\text{transition} = (s, a, r, s')$$

$$\text{difference} = \left[ r + \gamma \max_{a'} Q(s', a') \right] - Q(s, a)$$

$$Q(s,a) \leftarrow Q(s,a) + \alpha \, [\text{difference}] \qquad \text{Exact Q's}$$

$$w_i \leftarrow w_i + \alpha \, [\text{difference}] \, f_i(s,a) \qquad \text{Approximate Q's}$$

So by choosing let's 3 features: distance in the maze to the passenger location, distance to the drop off position, whether passenger reached location or not (either 1 or 0) we can achieve much better convergence as algorithm will be able to generalize problem. For example, by picking up passenger in the position R, the agent will learn that picking up the passenger in general is good thing and it will adjust weight w corresponding for the distance to passenger location (during learning this weight should become negative as the lower distance to the passenger the bigger Q value must become). The magnitude of weight for the distance to the drop off position must be much smaller than magnitude of weight corresponding to distance to passenger. This should be like because there is no point of trying reaching drop off position without passenger (I assume that when passenger is within taxi the distance to the passenger is zero). And the weight corresponding to passenger reached desired location should become something huge and positive as it is our final goal. On the other hand, in this scheme it is assumed that we can determine distance in the maze, when in our case agent learns the position of locations: R, B, Y and G all by itself.

# 6 Conclusion

In this taxi problem we used reward function which is different from the one described in [Dut+05]. It behaves non-linearly depending on the time since last pick up and as it was discussed earlier in section 2.4. As it was described in section 2.4 this non-linear reward function results in situation when it is very hard for our Q function to reach the stage when the values remains static. On the other hand, it is minor distinguish as we are mostly interested in the generated policy and in taxi task with homework reward function it converges to the optimal one.

One more. in section 4 we found the possible reasons why Q-learning and SARSA act so similarly in the taxi problem

# 7 Appendix

Main launch code where simulation starts and all learning parameters such as $\alpha$ (note that in the code I use eta instead for shortness), $\gamma$ and $\epsilon$ are specified:

```matlab
1  function [Q, rewards_per_episode, max_Q_diffs] = rl_taxi_sim(trans_f, ↩
       learning)
2      %reinforcement learning taxi problem simulation. Can be found here:
3      % Dutech, A., Edmunds, T., Kok, J., Lagoudakis, M., Littman, M., ↩
           Riedmiller, M., Whiteson, S.
4      % (2005) Reinforcement learning benchmarks and bake-offs II. Workshop ↩
           at Advances in Neural
5      % Information Processing Systems conference
6
7
8      % define maze global variables
9      define_maze_global();
10     global X Y PSNG_POS DROP_POS A MAX_EPISODE_TIME DROP_LOCS
11
12     %local paramaters
13     % total number of learning trials
14     T = 5000;
15     eps = 0.5;
16     eps_decay = 0.999;
17     gamma = 0.9;
18     eta = 0.5;
19     eta_decay = 0.99999;
20     over_episode = 50;
21     %initialisation
22     Q = 0.1*rand(X, Y, PSNG_POS, DROP_POS, A);
23
24     rewards_per_episode = [];
25     max_Q_diffs = [];
26
27     for t=1:T
28         [Q_new, total_episode_reward] = learning(Q, eta, gamma, eps, ↩
               trans_f);
29
30         rewards_per_episode = [rewards_per_episode, total_episode_reward];
31
32         Q_diff = abs(Q_new - Q);
33         max_diff = max(Q_diff(:));
34         max_Q_diffs = [max_Q_diffs, max_diff];
35
36         %paramaters decay
37         eta = eta * eta_decay;
```

```
38          eps = eps * eps_decay;
39
40          %print and values to control converges
41          if (rem(t, over_episode)==0)
42              display(eps, 'current eps');
43              display(eta, 'current eta');
44
45              figure(1);
46              avg_rewards = mean(reshape(rewards_per_episode, over_episode, ↩
                    []));
47              plot(avg_rewards, 'b');
48              ylabel('average cumulative reward');
49              xlabel(['Training Episodes (x', num2str(over_episode),')']);
50
51              %maze values for a case when
52              %the passenger within taxi and drop off point is in G or in my
53              %representation state with pp=5 and dp=3
54              MV = max(Q,[],5);
55              MV = MV(:, :, 5, 3);
56              MV = reshape(MV, X, Y);
57              %round to 4 digits on the left
58              MV = round(MV, 4);
59              lbls_MV = grid_to_lbls(MV, '%2.5f');
60              figure(2);
61              draw_maze(lbls_MV, strcat('maze values for pp=5 and dp=3 t=', ↩
                    num2str(t)));
62
63              figure(3);
64              max_diffs = mean(reshape(max_Q_diffs, over_episode, []));
65              stem(max_diffs);
66              xlabel(['Training Episodes (x', num2str(over_episode),')']);
67              ylabel('Average maximum absolute difference in Q');
68          end
69          Q = Q_new;
70      end
71  end
```

Defining supporting global variable which are the same for all simulations, such drop off locations, walls positions, etc. (except probably - sometimes I change it):

```
1  function [] = define_maze_global()
2      global X Y PSNG_POS DROP_POS DROP_LOCS A WALLS LOC_NAMES ↩
           MAX_EPISODE_TIME
3      X = 5;
4      Y = 5;
5      PSNG_POS = 5;
```

```
6      DROP_POS = 4;
7      % number of actions
8      A = 6;
9      walls = [[[1,1], [2, 1]];...
10             [[1,2], [2, 2]];...
11             [[2,4], [3, 4]];...
12             [[2,5], [3, 5]];...
13             [[3,1], [4, 1]];...
14             [[3,2], [4, 2]]];
15     mirrored_walls = zeros(size(walls));
16     mirrored_walls(:, [3, 1, 4, 2]) = walls(:, [1, 3, 2, 4]);
17     WALLS = [walls; mirrored_walls];
18     LOC_NAMES = {'Y', 'R', 'G', 'B'};
19     DROP_LOCS = [[1,1]; [1,5]; [5,5]; [4,1]];
20     MAX_EPISODE_TIME = 1500;
21 end
```

State representation of taxi problem:

```
1  function state = create_rnd_state()
2      %x  - x position in the maze
3      %y  - y position in the maze
4      %pp - passenger position
5      %dp - drop position
6      global X Y PSNG_POS DROP_POS
7      state = struct('x', randi(X),...
8                     'y', randi(Y),...
9                     'pp',randi(PSNG_POS),...
10                    'dp',randi(DROP_POS));
11 end
```

```
1  function copy_s = copy_state(s)
2      %s - state
3      copy_s = struct('x', s.x,...
4                      'y', s.y,...
5                      'pp',s.pp,...
6                      'dp',s.dp);
7  end
```

Support drawing utilities:

```
1  function lbls = grid_to_lbls( grid_values, varargin)
2      %transpose because the first axe is X, and the second is Y
3      %it is important for draw_maze
```

```matlab
 4
 5        flat_V = reshape(grid_values',[],1);
 6        if length(varargin) == 0
 7            fmt = '%d';
 8        else
 9            fmt = char(varargin(1));
10        end
11        lbls = strtrim(cellstr(num2str(flat_V, fmt)));
12    end
```

```matlab
 1    function [] = draw_maze(lbls, name)
 2        %lbls - labels - cell array of strings to display in the maze grid
 3        %size X*Y
 4        %name - plot title
 5        global WALLS X Y DROP_LOCS LOC_NAMES
 6        %draw location names with light green color
 7
 8        hold on;
 9        grid on;
10        h = findobj(gca,'Type','Text');
11        delete(h);
12        for i = 1:size(DROP_LOCS,1)
13            %substract shifts for centering, must be in sync with font size
14            dlx = DROP_LOCS(i,1) - 0.5;
15            dly = DROP_LOCS(i,2) - 0.5;
16            loc_name = LOC_NAMES(i);
17            %
18            text(dlx, dly, loc_name,...
19                'Color', [0.8 0.8 0.8],...
20                'FontWeight', 'bold', 'FontUnits', 'normalized', 'FontSize', ↩
                    1/Y,...
21                'HorizontalAlignment', 'center',...
22                'VerticalAlignment', 'middle');
23        end
24        title(name);
25        xlim([0 X]);
26        ylim([0 Y]);
27        set(gca,'XTick',0:X);
28        set(gca,'YTick',0:Y);
29        %to have nice grid cells as squares
30        axis square;
31
32        % grid domains
33        xg = 1:X;
34        yg = 1:Y;
```

```matlab
35        %label coordinates, substract 0.5 for centering
36        [xlbl, ylbl] = meshgrid(xg - 0.5, yg - 0.5);
37        text(xlbl(:), ylbl(:), lbls(:),'color','r',...
38            'HorizontalAlignment','center','VerticalAlignment','middle',...
39            'FontWeight', 'bold', 'FontUnits', 'normalized', 'FontSize', 1/(5*↩
                Y));
40        %plot maze walls
41
42        for i = 1:size(WALLS, 1)
43            wall = WALLS(i, :);
44            %centering
45            wall = wall - 0.5;
46            x1 = wall(1);
47            y1 = wall(2);
48            x2 = wall(3);
49            y2 = wall(4);
50            %vertical
51            if y1 == y2
52                wx = (x2 - x1)*0.5 + x1;
53                plot([wx, wx], [y1 - 0.5, y1 + 0.5],'-k', 'LineWidth',4);
54            end
55            %horizontal
56            if x1 == x2
57                wy = (y2 - y1)*0.5 + y1;
58                plot([x1 - 0.5, x1 + 0.5], [wy, wy],'-k', 'LineWidth',4);
59            end
60        end
61        hold off;
62        drawnow;
63  end
```

```matlab
1   function [ lbls_M ] = actions2str(lbls_M_raw)
2       %lbls_M_raw - array of string cells, with numeric actions
3       lbls_M = lbls_M_raw(1:numel(lbls_M_raw));
4       map = {'up', 'down', 'left', 'right', 'pick', 'drop'};
5       lbls_M(strcmp('0', lbls_M)) = {''};
6       lbls_M(strcmp('1', lbls_M)) = {'up'};
7       lbls_M(strcmp('2', lbls_M)) = {'down'};
8       lbls_M(strcmp('3', lbls_M)) = {'left'};
9       lbls_M(strcmp('4', lbls_M)) = {'right'};
10      lbls_M(strcmp('5', lbls_M)) = {'pick'};
11      lbls_M(strcmp('6', lbls_M)) = {'drop'};
12      for a1=1:6
13        for a2=1:6
14          action_combo = a2 + a1 * 6;
```

```matlab
15          lbls_M(strcmp(num2str(action_combo), lbls_M)) = {[map{a1}, '/', ↵
                map{a2}]};
16        end
17      end
18  end
```

```matlab
1  function [] = plot_policy(Q)
2    define_maze_global()
3    global X Y
4    actions_map = zeros(X, Y);
5    %initial state in the left bottom corner with
6    %pasenger and must go to G
7    s = create_rnd_state();
8    s.pp = 5;
9    s.dp = 3;
10
11   for x=1:X
12     for y=1:Y
13       %perform action
14        s.x = x;
15        s.y = y;
16       [v, a] = greedy(s, Q);
17       actions_map(s.x, s.y) = a;
18     end
19   end
20
21   lbls_actions = grid_to_lbls(actions_map);
22   lbls_actions = actions2str(lbls_actions);
23   figure(4);
24   draw_maze(lbls_actions, 'actions map for taxi');
25  end
```

Choosing action:

```matlab
1  function [v, a] = greedy(s, Q)
2    %s - state, v - value of state s, a - next action
3    [v, a] = max(Q(s.x, s.y, s.pp, s.dp, :));
4  end
```

```matlab
1  function [v, a] = eps_greedy(s, Q, eps)
2    %s - state, v - value of state s, a - next action, eps - epsilon
3    global A
4    [v, a] = greedy(s, Q);
```

```matlab
5       if (rand(1)<eps)
6           a = randi(A);
7       end;
8   end
```

Transfer-reward function:

```matlab
1   function [ns, r] = trans_fun(s, a, drive_time)
2       %s - state, a - action, ns - next state, r - reward
3       %drive_time - time since passenger pick up
4       global DROP_LOCS WALLS X Y
5
6       % number of actions
7
8       ns = copy_state(s);
9       %default - no reward
10      r = 0;
11      % a == 1, go north, up
12      if a == 1
13          ns.y = ns.y + 1;
14      end;
15      % a == 2, go south, down
16      if a == 2
17          ns.y = ns.y - 1;
18      end;
19      % a == 3, go west, left
20      if a == 3
21          ns.x = ns.x - 1;
22      end;
23      % a == 4, go east, right
24      if a == 4
25          ns.x = ns.x + 1;
26      end;
27      %check for collisions
28      hit_walls = any(ismember(WALLS,[s.x, s.y, ns.x, ns.y],'rows'));
29      within_maze = 1 <= ns.x && ns.x <= X && 1 <= ns.y && ns.y <= Y;
30      out_of_maze = ~within_maze;
31      if hit_walls || out_of_maze
32          %restore previous position
33          ns.x = s.x;
34          ns.y = s.y;
35          r = -1;
36      end;
37
38      % a == 5, pick up passenger
39      if a == 5
```

```
40          %default failed
41          r = -1;
42          %check for already pick up
43          if s.pp == 5
44              return;
45          end;
46          %check for the same location of taxi and passenger
47          psng_loc = DROP_LOCS(s.pp, :);
48          if all([s.x, s.y] == psng_loc)
49              ns.pp = 5;
50              r = 1;
51          end;
52      end;
53      % a == 6, drop off passenger
54      if a == 6
55          %default failed
56          r = -1;
57          %check for passenger presence
58          if s.pp ~= 5
59              return;
60          end;
61          %check for the same location of taxi and passenger
62          drop_loc = DROP_LOCS(s.dp, :);
63          if all([s.x, s.y] == drop_loc)
64              ns.pp = s.dp;
65              r = 10/drive_time;
66          end;
67      end;
68  end
```

Transfer-reward function from paper [Dut+05]

```
1   function [ns, r] = trans_fun_paper(s, a, drive_time)
2       %s - state, a - action, ns - next state, r - reward
3       %drive_time - time since passenger pick up
4       global DROP_LOCS WALLS X Y
5
6       % number of actions
7
8       ns = copy_state(s);
9       %default - no reward
10      r = -1;
11      % a == 1, go north, up
12      if a == 1
13          ns.y = ns.y + 1;
14      end;
```

```matlab
15      % a == 2, go south, down
16      if a == 2
17          ns.y = ns.y - 1;
18      end;
19      % a == 3, go west, left
20      if a == 3
21          ns.x = ns.x - 1;
22      end;
23      % a == 4, go east, right
24      if a == 4
25          ns.x = ns.x + 1;
26      end;
27      %check for collisions
28      hit_walls = any(ismember(WALLS,[s.x, s.y, ns.x, ns.y],'rows'));
29      within_maze = 1 <= ns.x && ns.x <= X && 1 <= ns.y && ns.y <= Y;
30      out_of_maze = ~within_maze;
31      if hit_walls || out_of_maze
32          %restore previous position
33          ns.x = s.x;
34          ns.y = s.y;
35          r = -1;
36      end;
37
38      % a == 5, pick up passenger
39      if a == 5
40          %default failed
41          r = -10;
42          %check for already pick up
43          if s.pp == 5
44              return;
45          end;
46          %check for the same location of taxi and passenger
47          psng_loc = DROP_LOCS(s.pp, :);
48          if all([s.x, s.y] == psng_loc)
49              ns.pp = 5;
50              r = -1;
51          end;
52      end;
53      % a == 6, drop off passenger
54      if a == 6
55          %default failed
56          r = -10;
57          %check for passenger presence
58          if s.pp ~= 5
59              return;
60          end;
61          %check for the same location of taxi and passenger
```

```
62        drop_loc = DROP_LOCS(s.dp, :);
63        if all([s.x, s.y] == drop_loc)
64            ns.pp = s.dp;
65            r = 0;
66        end;
67     end;
68 end
```

Q-learning code:

```
1 function [Q_final, total_episode_reward] = Q_learning_episode(Q, eta, ←
    gamma, eps, trans_f)
2   %v - value, a - action, s - state, r - reward, q - Q-value
3   %n in front of them means next
4   global DROP_LOCS MAX_EPISODE_TIME
5   s = create_rnd_state();
6   %trial
7   drive_time = 0;
8   total_episode_reward = 0;
9
10  for step=1:MAX_EPISODE_TIME
11      %passenger is within taxi increase time
12      if s.pp == 5
13          drive_time = drive_time + 1;
14      end
15
16      %perform action
17      [v, a] = eps_greedy(s, Q, eps);
18      [ns, r] = trans_f(s, a, drive_time);
19      total_episode_reward = total_episode_reward + r;
20
21      %learning
22      nv = max(Q(ns.x, ns.y, ns.pp, ns.dp, :));
23      q = Q(s.x, s.y, s.pp, s.dp, a);
24      Q(s.x, s.y, s.pp, s.dp, a) = (1-eta)*q + eta*(r+gamma*nv);
25
26      %check for terminal conditions
27      drop_loc = DROP_LOCS(s.dp, :);
28      %when we drop off passanger in the drop of location,
29      %previously we had passanger within taxi
30      if s.pp == 5 && a == 6 && all([ns.x, ns.y] == drop_loc)
31          break;
32      end
33      %update old state
34      s = ns;
35  end
```

```
36    Q_final = Q;
37 end
```

SARSA code:

```
1  function [Q_final, total_episode_reward] = SARSA_episode(Q, eta, gamma, ←
      eps, trans_f)
2      %v - value, a - action, s - state, r - reward, q - Q-value
3      %n in front of them means next
4      global DROP_LOCS MAX_EPISODE_TIME
5      s = create_rnd_state();
6      %trial
7      drive_time = 0;
8      total_episode_reward = 0;
9      %initial action
10     [v, a] = eps_greedy(s, Q, eps);
11     for step=1:MAX_EPISODE_TIME
12         %passenger is within taxi increase time
13         if s.pp == 5
14             drive_time = drive_time + 1;
15         end
16
17         %perform action
18         [ns, r] = trans_f(s, a, drive_time);
19         total_episode_reward = total_episode_reward + r;
20
21         %estimate next state
22         [nv, na] = eps_greedy(ns, Q, eps);
23         %learning
24         q = Q(s.x, s.y, s.pp, s.dp, a);
25         Q(s.x, s.y, s.pp, s.dp, a) = (1-eta)*q + eta*(r+gamma*nv);
26
27         %check for terminal conditions
28         drop_loc = DROP_LOCS(s.dp, :);
29         %when we drop off passanger in the drop of location,
30         %previously we had passanger within taxi
31         if s.pp == 5 && a == 6 && all([ns.x, ns.y] == drop_loc)
32             break;
33         end
34         %update old state
35         s = ns;
36         %new action
37         a = na;
38     end
39     Q_final = Q;
40 end
```

# References

[SB98]    Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press, 1998.

[Dut+05]  Alain Dutech et al. "Reinforcement learning benchmarks and bake-offs II". In: *Advances in Neural Information Processing Systems (NIPS)* 17 (2005).

[Kle12]   Dan Klein. *UC Berkeley CS188 Intro to AI*. `http://ai.berkeley.edu/lecture_videos.html`. 2012.