

# REINFORCEMENT LEARNING ASSIGNMENT № 2

Ruslan Burakov, s1569105

17/03/2015

## Note

All full code fragments are located in the end of the report in the appendix section. Naming convention. By one episode I mean period when agent is initialised randomly within maze and it tries to reach goal. So learning process consists of many episodes. Once the learning is over and the final policy is evaluated I call it the end of experiment. I repeat several experiments to determine mean convergence time and confidence intervals.

## 1 Function Approximation Discussion

According to the lectures Q function can be approximated by k basis function in the following expression:

$$Q(s, a) = w_1 f_1(s, a) + w_2 f_2(s, a) + \dots + w_k f_k(s, a) = w^T f(s, a) \quad (1)$$

$$s - \text{state}, a - \text{action} \quad (2)$$

For this question and following ones our basis functions can be partitioned into this product:

$$f_i(s, a) = f_i(s) f_i(a) \quad (3)$$

In the expression above each i'th basis function which depends only on action is

$$f_i(a) = I(a = a')$$

where  $I(a = a')$  is identity/indicator function of particular action. In the assignment we use 25 basis functions which depends only on the state  $f_i(s)$  and also we have 5 different actions which means we need 5 action basis functions  $I(a = a')$  and therefore by multiplying these domains we will get the following number of general basis functions:

$$|f_i(s, a)| = |f_i(s)| * |f_i(a)| = 25 * 5 = 125 \quad (4)$$

and due to this we get that in 1 we have 125 weights overall or 25 weight per 5 actions. This can be simplified even further during the learning step (here I omit indexes which means that these are vectors):

$$\delta \leftarrow r + \gamma \max_{a' \in A} w^T f(s', a') - w^T f(s, a) \quad (5)$$

$$w \leftarrow w + \eta \delta f(s, a) \quad (6)$$

$$(7)$$

by using the fact that  $f_i(a) = I(a = a')$  and noting that all actions which differs from action identity will make no contribution in the product  $w^T f(s, a)$ :

$$Q(s, a) = w_1 f_1(s) + w_2 f_2(s) + \dots + w_{25} f_{25}(s) = w_a^T f(s) \quad (8)$$

$$\delta \leftarrow r + \gamma \max_{a' \in A} w_{a'}^T f(s') - w_a^T f(s) \quad (9)$$

$$w_a \leftarrow w_a + \eta \delta f(s) \quad (10)$$

where each  $w_a$  is 25 vector of weights which corresponds to particular action  $a$  (11)

Relying on the fact that our task is discrete one I precompute values of each 25 basis functions for each cell of maze in order to save computational time. Here how it looks approximately in my code (particular cases for identity functions  $f_i(s) = I(s = s_i)$  or RBF are in the appendix) like this:

---

```

1  %...
2  %precompute basis functions in advance
3  % number of states
4  Sx = 10; % width of grid world
5  Sy = 10; % length of grid world
6  M = 2; % cover neighbouring area
7  %basis functions
8  F = zeros(Sx, Sy, 25);
9
10 %iterate through maze cells
11 for ix = 1:Sx
12     for iy = 1:Sy
13         %identity functions
14         %iterate through 25 basis functions
15         for kx = 1:5
16             for ky = 1:5
17                 %basis function index (can be from 1 to 25)
18                 fId = 5 * (ky - 1) + kx;
19                 %Here it is identity functions
20                 %It can be RBF as well
21                 F(ix, iy, fId) = M*(kx-1) < ix && ix <= M*kx &&...
22                               M*(ky-1) < iy && iy <= M*ky;
23             end
24         end
25     end
26 end
27
28 %EPISODES/LEARNING LOOP GOES BELOW
29 %...
```

---

These give sufficient speed up, especially when using RBF (almost 20 times). Also, in terms of computation speed there is basically no difference whether identity functions or RBF are used

as their precomputed before hand.

Additionally by analysing equations 8 it can be seen that it is possible to rewrite them as dot products between vectors and matrices. Code snippet which demonstrates dot products speed ups:

---

```
1      %...
2
3      %INITIALISATION
4      %number of actions
5      A = 5;
6      %weights matrix. For each actions keeps values of 25 basis functions
7      W = zeros(A,25);
8      %basis functions
9      F = zeros(Sx, Sy, 25);
10
11     %...
12
13     %INSIDE PARTICULAR EPISODE
14
15     %...
16
17     % choosing best action is the dot product
18     [V_s0, a0]=max(W * reshape(F(x0, y0, :), 25, 1));
19
20     %...
21
22     %learning
23     V_s1 = max(W * reshape(F(x1, y1, :), 25, 1));
24     %approximate Q(s0, a0) by dot product
25     Q_s0_a0 = W(a0, :) * reshape(F(x0, y0, :), 25, 1);
26     delta = r+gamma*V_s1 - Q_s0_a0;
27     %basic functions in state x0 y0
28     bfs = reshape(F(x0, y0, :), 1, 25);
29     %updating weights matrix
30     W(a0, :) = W(a0, :) + eta * delta * bfs;
31
32     %...
```

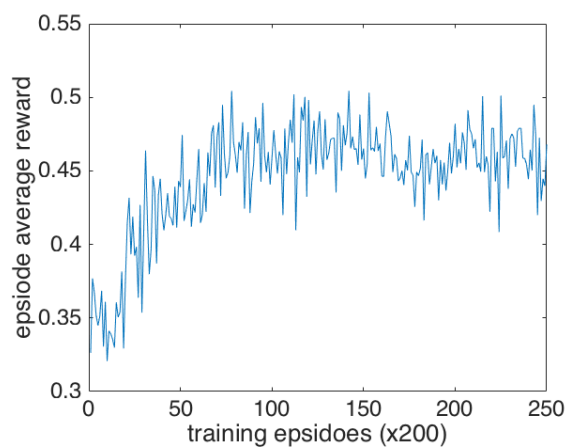
---

Dot products work much faster in contrast to usual for loops in MATLAB. That is why due to such implementation I am able to go through 150000 episodes with RBF functions on 25 by 25 maze in 2 minutes.

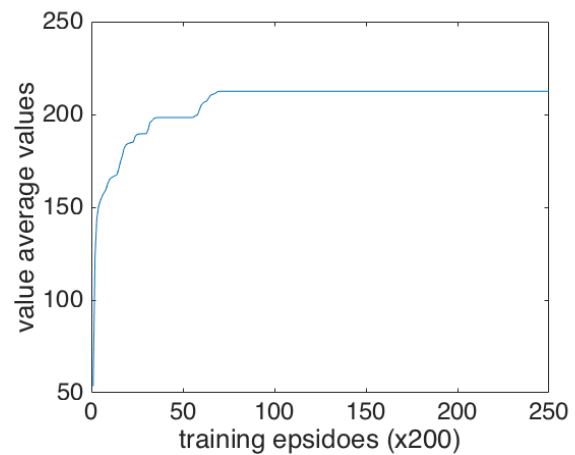
## 2 Identity functions representation in 5x5 maze

In this task we are using  $f_i(s) = I(s = s_i)$  as state basis functions. First of all, let's see results achieved by true Q function (table representation) for this task. I have used the following learning parameters:

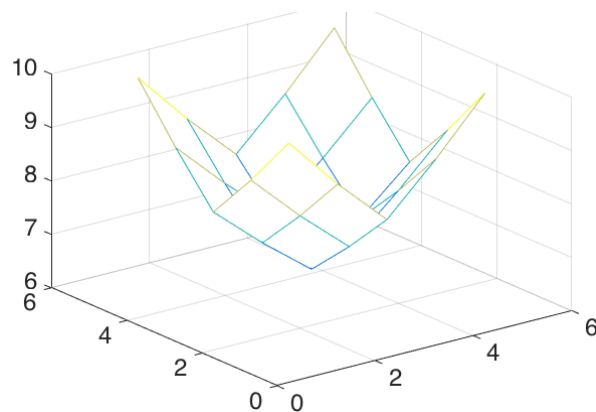
|                             |      |
|-----------------------------|------|
| Learning rate $\eta$        | 0.2  |
| Exploration rate $\epsilon$ | 0.1  |
| Discount rate $\gamma$      | 0.9; |



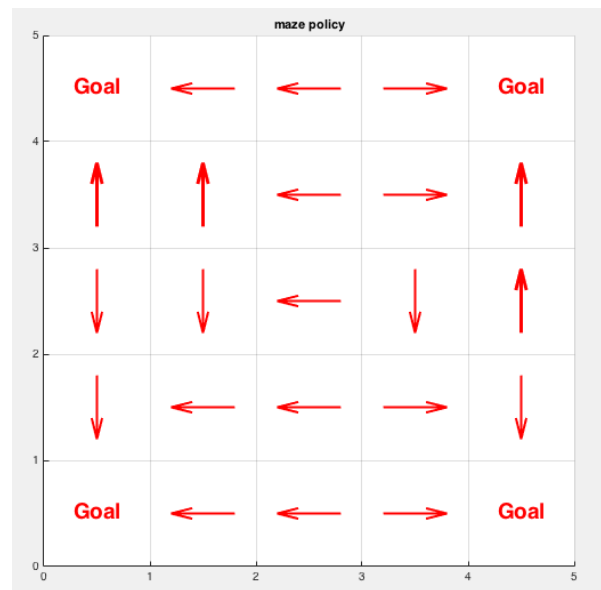
(a) Average reward



(b) Value function convergence



(c) Final value function on maze states



(d) Final policy

Figure 1: 5x5 maze. Table representation

In this problem our reward doesn't depend on how many steps agent did within one episode in order to reach reward (it is just 1 if agent has reached its goal). That is why reward in the end of one episode is not very representative value in terms of how well the learning process goes as it won't show how fast we reached our goal. Due to that I use average reward per episode, as the more time agent spent in the maze the smaller his average reward will be.

---

```

1    %...
2    %start episode
3    for u = 1:Sx*Sy
4        %learning and acting in maze
5        %...
6
7        % goto next episode once the goal is reached
8        if ismember([x0, y0], Goals, 'rows')
9            avg_r = r/u;
10           %display(avg_r, 'reached the goal!');
11           break;
12       end
13   end
14   %end of episode
15
16   %store average reward
17   avg_rs(t) = avg_r;
18   %...

```

---

All values below are averaged over 200 training episodes () to ease fluctuations due to random starting position and  $\epsilon$ -greedy learning.

The average reward is presented on Figure 1a. It reaches approximately 0.45 which corresponds to the  $1/0.45 = 2.222$  steps in the one episode on average. This is logical as in 5x5 maze with goals in the corners of the maze the average number of steps to reach goal should be between 2 and 3.

To determine dynamic of learning process more precisely I use value function  $V(sx, sy)$ . I sum up all entries of value function which corresponds to particular positions in the maze (we shouldn't have negatives ones because there are no rewards -> only 1 and 0) and rely on the fact that if entries of value function doesn't change then their sum doesn't change, as well. The situations when some entries declined by exactly the same amount other entries increased are extremely unlikely. Code snippet which is responsible for this:

---

```

1    %...
2    % after each episode:
3
4    %compute value function
5    VV = max(Q,[],3);%Size: (Sx,Sy)
6    %storing for future plotting
7    samples(t) = sum(reshape(VV, 1,numel(VV)));
8
9    %...
10
11   %after end of training
12   %averaging to remove most fluctuations

```

```

13     average_over = 200;
14     avg_samples = mean(reshape(samples, average_over, []));
15
16     %...

```

---

On Figure 1b we can see that values are almost not changed after approximately  $160 * 200 = 32000$  episodes. I determine convergence by computing standard deviation of the averaged value function (over 200 episodes) with window 50 and after that I find the first spike in its values from the end. Such procedure allows to avoid determining convergence too early when value function reaches local maximum (on Figure 1b the local maximum is reached near  $80 * 200 = 8000$ ). The disadvantage of this method that it doesn't allow to determine convergence online during learning that is why it is impossible to abort learning earlier. Additionally, especially for approximation with basis functions, value function reaches certain level and then starts to fluctuate around some particular value. That is why converges criteria must be adjusted in order to accommodate that (probably convergence criteria is not the best name for it) and its initial value is estimated through value function plots over episodes. Due to that plots must be shown together with these values. Ideally, convergence must be measured by comparing generated policy with optimal (which can handle several best actions) one for each certain amount of episodes. Code snippet responsible for this:

---

```

1     %...
2
3     %computing standard deviations over window 50
4     smp_stds = zeros(1,(size(avg_samples, 2) - 50));
5     for t=1:(size(avg_samples, 2) - 50)
6         smp_std = std(avg_samples(t:(t+50)));
7         smp_stds(t) = smp_std;%abs(avg_samples(t) - avg_samples(t-1))/↵
            avg_samples(t);
8     end
9
10    %find first spike bigger than 0.2 from the end
11    convg = find(smp_stds > 0.2, 1, 'last')
12    %convergence time
13    convg_t = convg * average_over
14    %if it is not empty then store in experiment array
15    %to compute confidence intervals
16    if ~isempty(convg_t)
17        convg_times = [convg_times, convg_t];
18    end
19
20    %...

```

---

Here, I set criteria for standard deviation to be bigger than 0.2 starting from the end in order to determine convergence.

After repeating experiment for 10 time I got the following results for convergence with 95% confidence interval (from here on it is always 95% confidence interval):

|                            |                           |
|----------------------------|---------------------------|
| Convergence criteria (STD) | 0.2                       |
| Number of experiments      | 10                        |
| Mean convergence time      | 15480 $\pm$ 2462 episodes |

Finally, on Figure 1c we can see value function  $V(s_x, s_y)$  in the end of learning (in the end of one of the experiments) for each of maze cells. It has symmetric convex shape with 4 peaks located exactly in the goals positions in the corners of the maze as it would be expected. On Figure 1d we can see the final generated policy (after the end of learning process for each maze state I choose actions absolutely greedily). It is optimal one because in this task there are many equally optimal policies given that sometimes distances to goals are the same and there is several best action. Overall, this indicates that learning has been successful.

Now let's compare above results with the exactly same task where Q functions was approximated by 25 basis identity state functions  $f_i(s) = I(s = s_i)$  I have used the same learning parameters:

|                             |      |
|-----------------------------|------|
| Learning rate $\eta$        | 0.2  |
| Exploration rate $\epsilon$ | 0.1  |
| Discount rate $\gamma$      | 0.9; |

I have received the following convergence results:

|                            |                           |
|----------------------------|---------------------------|
| Convergence criteria (STD) | 0.2                       |
| Number of experiments      | 10                        |
| Mean convergence time      | 16220 $\pm$ 1475 episodes |

Within its confidence intervals this results matches the convergence time obtained for tabular representation.

When basic function approximation is applied the approximated value functions can be computed from weights  $W(A, \text{number of basis functions})$  and basis functions  $F(Sx, Sy, \text{number of basis functions})$  using the following code snippet:

---

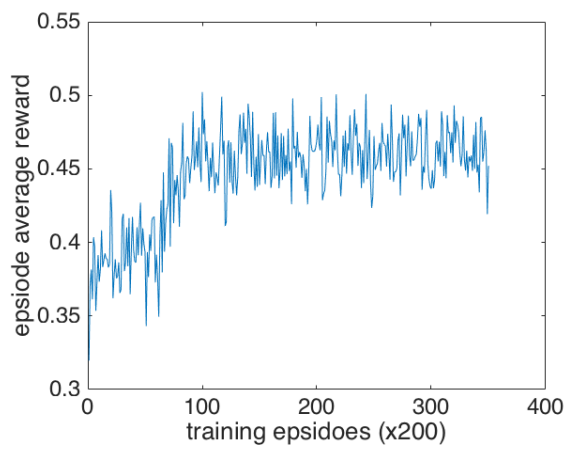
```

1  % reshape(F, Sx * Sy, 25) % Sx*Sy x 25
2  % W' % 25 x A
3  QQ = reshape(F, Sx * Sy, 25) * W'; % (Sx*Sy x 25) * (25 x A) = Sx*Sy x A
4  VV = max(QQ,[],2); %Sx*Sy
5  %approximated value function
6  VV = reshape(VV, Sx, Sy);

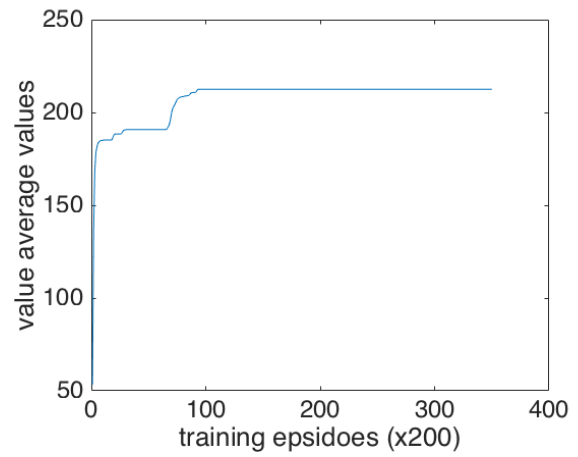
```

---

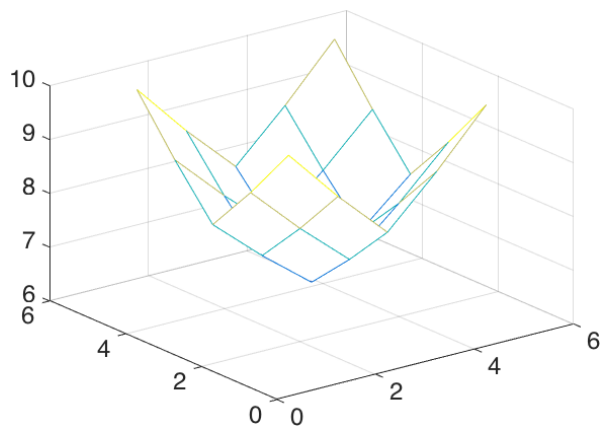
Next, Figures 2a, 2b, 2c are in good correspondence with its analogues from tabular representation on Figure 1. This is expected as when we are using 5x5 Identity functions for 5x5 maze



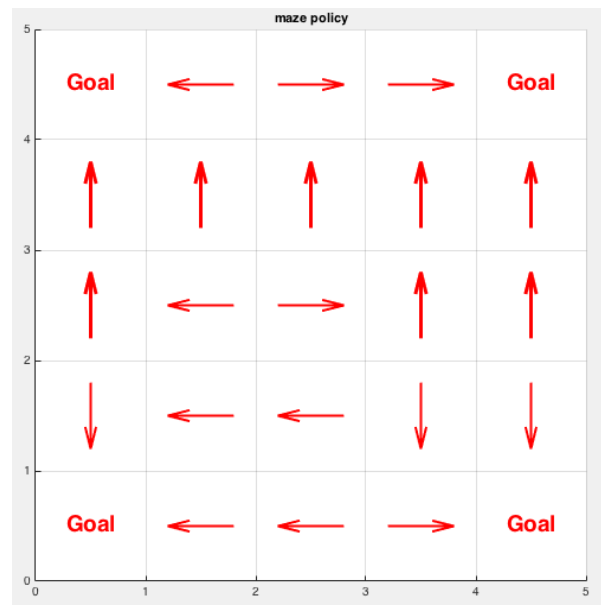
(a) Average reward



(b) Value function convergence



(c) Final value function on maze states



(d) Final policy

Figure 2: 5x5 maze. Identity Functions  $f_i(s) = I(s = s_i)$  representation



the tabular representation and basis functions representation should be essentially the same. This can be seen by comparing the number of parameters to be learn. In tabular representation we learn the Q values  $|Q(Sx, Sy, A)| = 5 * 5 * 5 = 125$ . In basis functions representation we learn weight matrix from Equation 8 ->  $|W(A, \text{number of basis functions})| = 5 * 25 = 125$ . The final policy depicted on Figure 2d differs from the one on Figure 1d but it is still optimal policy and the difference between them can be described by the fact that in this maze problem there are many equally optimal policies.

### 3 Identity functions representation in 10x10 maze

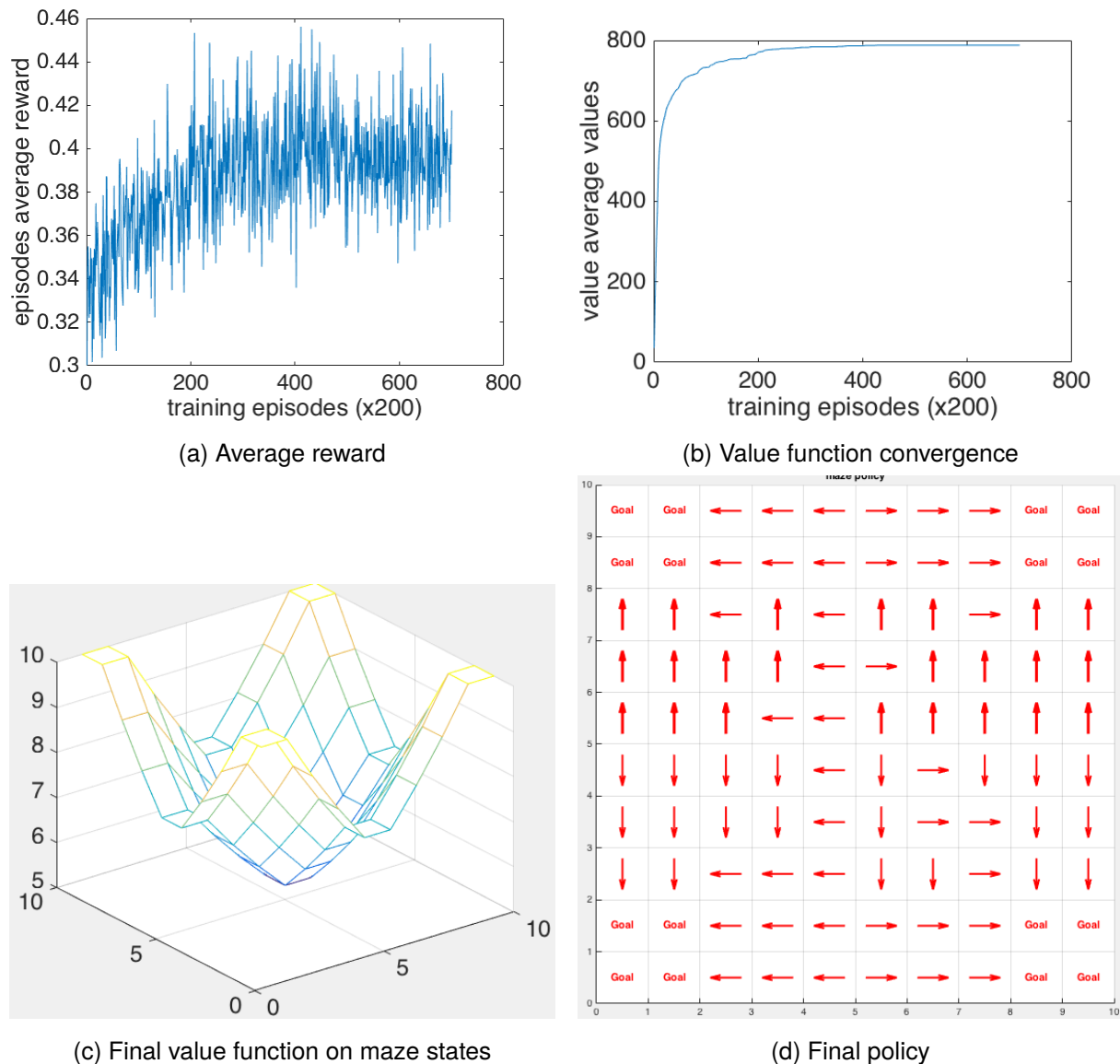


Figure 3: 10x10 maze. Table representation

In this question identity basis functions evenly cover patches 2x2 of maze cells. It has been already shown how state identity functions are computed in that case in the first code listing of question 1.

Here, I redefine goals as patches 2x2 placed in the corners of the maze because in this representation agent doesn't have enough information to reach final goal reliably once it is within goal patch 2x2. It happens because for all closest neighbours around goal states the same identity functions will be activated.

First, we compute true results in tabular representation of Q function for maze 10x10. For that I have used the same learning parameters as in previous question:

|                             |      |
|-----------------------------|------|
| Learning rate $\eta$        | 0.2  |
| Exploration rate $\epsilon$ | 0.1  |
| Discount rate $\gamma$      | 0.9; |

The converges results are:

|                            |                           |
|----------------------------|---------------------------|
| Convergence criteria (STD) | 0.2                       |
| Number of experiments      | 10                        |
| Mean convergence time      | 86240 $\pm$ 4727 episodes |

The convergence time is 5 times bigger than for the case of 5x5 maze which roughly matches the fact that there are 4 times more states to explore ( $10 * 10 = 100$  in contrast to  $5 * 5 = 25$  for the 5x5 maze).

On Figure 3a the average reward per episode has almost the same mean as on Figure 1a but it has significantly bigger variance in results. This can be described by the fact that now the maze has become bigger and average path to nearest the goal has become longer. However, the chances of appearing on the goal cell from the beginning has remained the same. It was  $4/25 = 0.16$  for 5x5 maze and now for 10x10 maze  $4 * 4/100 = 0.16$ . That is why starting point plays bigger role which leads to bigger variance in average reward per episode.

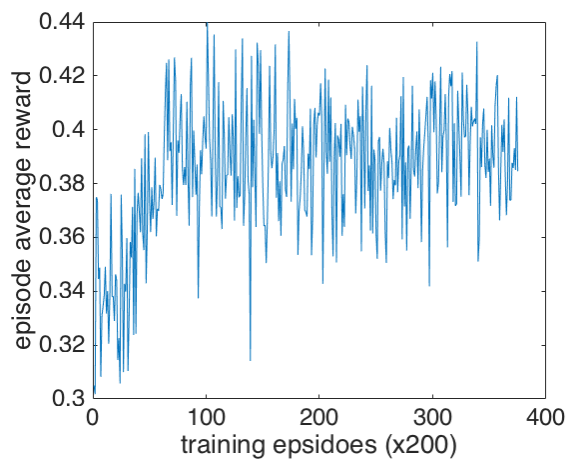
On Figure 3c we can see Value Function shape for 10x10 where goals are patches 2x2 in the corners. Once again it has 4 peaks exactly in corners locations. The policy shown on Figure 3d is also one of optimal ones.

Now, let's run experiments for identity functions approximation. It uses the same learning parameters as table representation above.

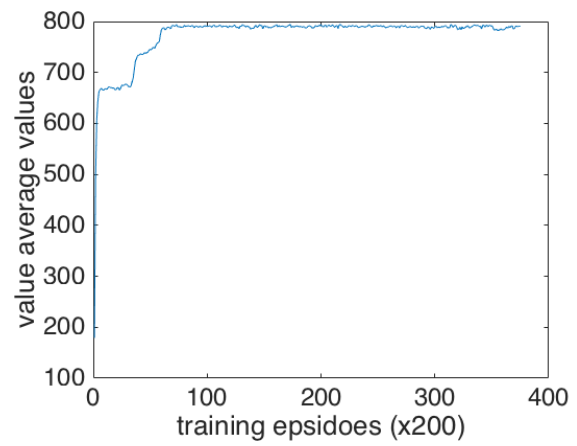
On Figure 4b we can see that values fluctuate stronger than in case of table representation on Figure 3b or identity basis functions representation for 5x5 maze on Figure 2b. This can be explained by the fact that now it is only approximation of real value function on Figure 3c as number of learning parameters is 4 times lower in that case ( $|Q(S_x, S_y, A)| = 10 * 10 * 5 = 500$  for tabular representation in comparison to the same  $|W(A, \text{number of basis functions})| = 5 * 25 = 125$  for function approximation). Due to that I relaxed convergence criteria and the convergence results are:

|                            |                           |
|----------------------------|---------------------------|
| Convergence criteria (STD) | 6.0                       |
| Number of experiments      | 10                        |
| Mean convergence time      | 13200 $\pm$ 1331 episodes |

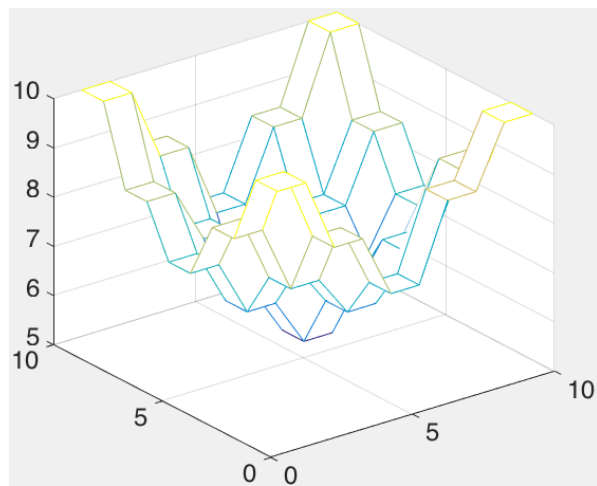
This time is essentially the same as for the case of identity basis functions approximation for 5x5 maze (16220  $\pm$  1475 episodes). Again, it can be explained by the fact that in 5x5 maze



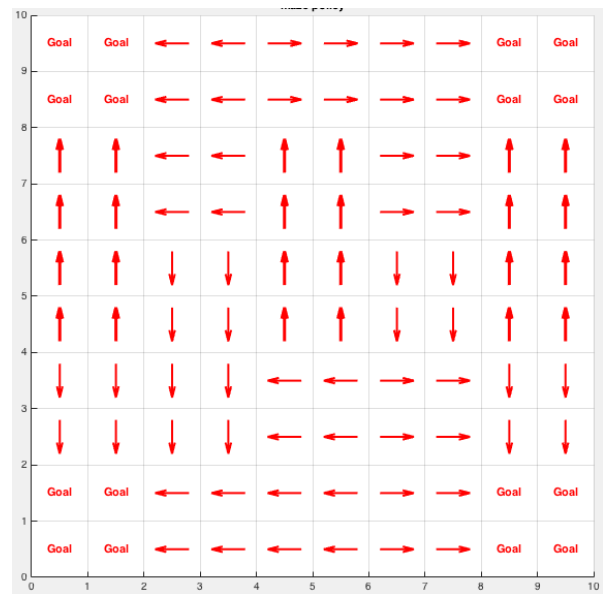
(a) Average reward



(b) Value function convergence



(c) Final value function on maze states



(d) Final policy

Figure 4: 10x10 maze. Identity basis functions representation

and 10x10 for identity basis functions approximation we have the same number of parameters to learn ( $|W(A, \text{number of basis functions})| = 5 * 25 = 125$ ). Here, we start to see the advantage of using basis functions approximation as convergence time is approximately 6 times faster than in plain tabular representation.

One of the most notable difference in comparison to table representation is on Figure 4c where the value functions is essentially 2 dimensional step function. This is logical as each 2x2 patch in the maze corresponds to the same 5 weights (for each action) in basis functions approximation. That is why once the maximum by action is taken then the 2x2 patches have the same value. Moreover, the policy generated on Figure 4d is optimal one and if you look carefully you will notice that each 2x2 patch shares the same actions. This is again due to the same reason that each 2x2 patch in the maze corresponds to the same 5 weights (for each action) in basis functions approximation.

## 4 Radial Basis Functions

For a Gaussian shape like basis functions I have used smoothed RBF kernel functions from the lectures as they proved to be more stable in contrast to usual Gaussians (it is easier to set their width parameter). Because I tried it myself and I saw in other people simulations when plain Gaussians are used then one goal becomes prevalent in comparison to all other. That is why many people just run their simulation for one goal in one of the corners at this stage. For smoothed RBF basis functions values must sum up to one for each state in the maze. They are precomputed in the following code snippet:

---

```

1  %...
2  %precompute basis functions in advance
3  % number of states
4  Sx = 10; % width of grid world
5  Sy = 10; % length of grid world
6  M = 2; % cover neighbouring area
7  %basis functions
8  F = zeros(Sx, Sy, 25);
9
10 %iterate through maze cells
11 for ix = 1:Sx
12     for iy = 1:Sy
13         %identity functions
14         %iterate through 25 basis functions
15         for kx = 1:5
16             for ky = 1:5
17                 %basis function index (can be from 1 to 25)
18                 fId = 5 * (ky - 1) + kx;
19                 %position of rbf in maze coordinates
20                 cx = kx*M - (M - 1)/2;
21                 cy = ky*M - (M - 1)/2;

```

```

22         E = exp(-(1/(2*SIGMA^2))*((ix-cx)^2+(iy-cy)^2));
23         F(ix, iy, fId) = E/Z;
24     end
25 end
26 Z = sum(F(ix, iy, :));
27 %smoothing
28 F(ix, iy, :) = F(ix, iy, :) / Z;
29 end
30 end
31
32 %EPISODES/LEARNING LOOP GOES BELOW
33 %...

```

---

Note here that basis functions are centered in the middle of MxM patches (by subtracting  $(M - 1)/2$ ). Also the Gaussian centers are computed in maze space that is why parameter SIGMA which corresponds to the width of RBF is also measured in maze space.

For this representation I still use wider patches for goals in the corners of the maze because as it was said in the previous question we may need to redefine goal and here we are still having approximation although much better one than simple identity functions. These wider goal patches allows algorithm to start convergence much faster especially on few initial steps in 25x25 maze as it finds goal with much higher probability (as there are more goals overall).

I use the following learning parameters unless it is said otherwise

|                             |      |
|-----------------------------|------|
| Learning rate $\eta$        | 0.2  |
| Exploration rate $\epsilon$ | 0.1  |
| Discount rate $\gamma$      | 0.9; |

For the 5x5 maze using RBF with width set to 0.5 I get the following convergence time result:

|                            |                           |
|----------------------------|---------------------------|
| Convergence criteria (STD) | 0.2                       |
| Number of experiments      | 10                        |
| Mean convergence time      | 11900 $\pm$ 1077 episodes |

The convergence time is slightly faster in comparison to results obtained in question 2 (we can compare them because convergence criteria is same as in question 2). This potentially can be described by the fact RBF gives additional positional information on early stages so inference can be done better. For example, let say we are located in position (1,1) in the 5x5 then values of my RBF functions will be:

| 0 | 1      | 2      | 3      | 4      | 5      |
|---|--------|--------|--------|--------|--------|
| 1 | 0.7753 | 0.1049 | 0.0003 | 0.0000 | 0.0000 |
| 2 | 0.1049 | 0.0142 | 0.0000 | 0.0000 | 0.0000 |
| 3 | 0.0003 | 0.0000 | 0.0000 | 0.0000 | 0.0000 |
| 4 | 0.0000 | 0.0000 | 0.0000 | 0.0000 | 0.0000 |
| 5 | 0.0000 | 0.0000 | 0.0000 | 0.0000 | 0.0000 |

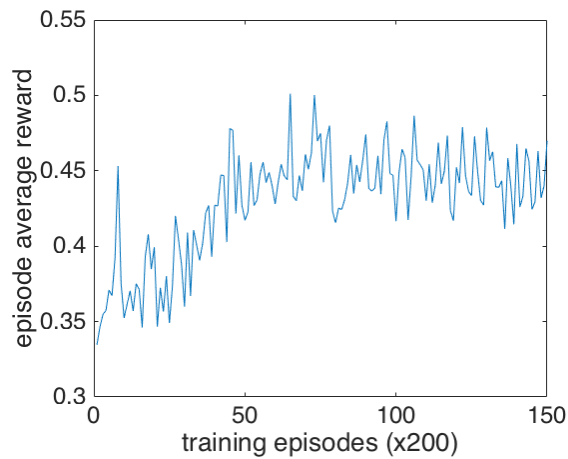
As we can see the highest values gives RBF functions located in (1,1). Next, much smaller ones are located in (1,2) and (2,1). After that we have a few tiny values and rest is zero.

Another example, if we are located in the middle of the maze (3,3):

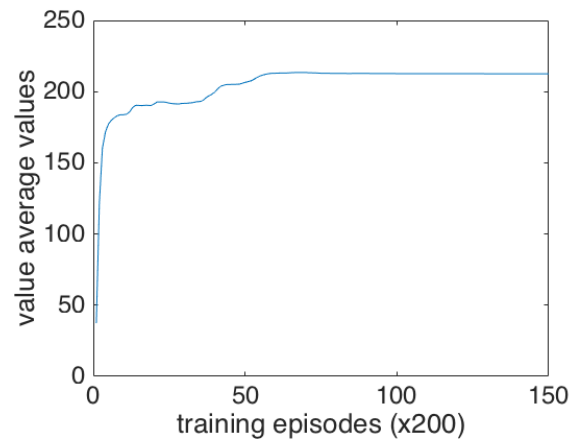
| 0 | 1      | 2      | 3      | 4      | 5      |
|---|--------|--------|--------|--------|--------|
| 1 | 0.0000 | 0.0000 | 0.0002 | 0.0000 | 0.0000 |
| 2 | 0.0000 | 0.0113 | 0.0837 | 0.0113 | 0.0000 |
| 3 | 0.0002 | 0.0837 | 0.6187 | 0.0837 | 0.0002 |
| 4 | 0.0000 | 0.0113 | 0.0837 | 0.0113 | 0.0000 |
| 5 | 0.0000 | 0.0000 | 0.0002 | 0.0000 | 0.0000 |

We have high value 0.6187 in the (3,3) and after that all other values goes down symmetri-  
cally from it.

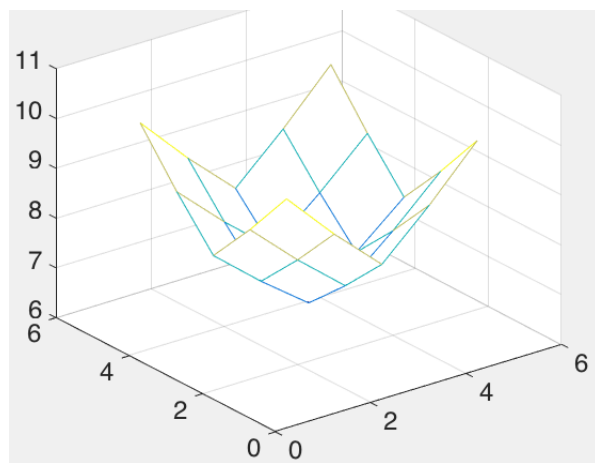
All the other results for 5x5 maze on Figure 5 are identical to Figures 1 and 2 in question 2.



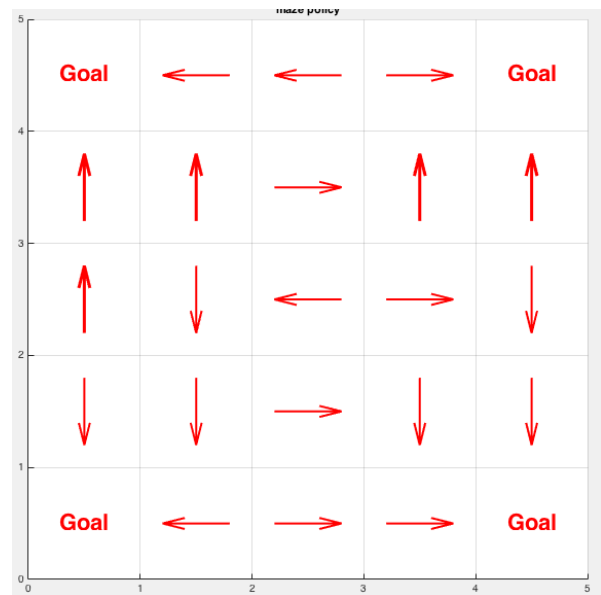
(a) Average reward



(b) Value function convergence



(c) Final value function on maze states



(d) Final policy

Figure 5: 5x5 maze. RBF representation with width 0.5

Essentially for RBF when the width (SIGMA) shrinks it becomes more and more peaked. In

that case for 10x10 and 25x25 mazes value functions should become a step function over 2x2 and 5x5 patches once width shrinks to certain value.

For the 10x10 maze experiment with width (SIGMA) equals to 1.0 the value function is never converged under convergence criteria of 6.0 for 100000 episodes due to relatively high fluctuation as it can be seen on Figure 6b. That is why I increased that criteria to 25.0 but in that case it is not quite right to compare convergence results with previous questions. However, as it can be seen on Figure 6b that after approximately  $200 * 200 = 40000$  episodes value function reach a certain level and it keeps fluctuating around it.

Here is the convergence results for width (SIGMA) 1.0 10x10 maze:

|                            |                           |
|----------------------------|---------------------------|
| Convergence criteria (STD) | 25                        |
| Number of experiments      | 10                        |
| Mean convergence time      | $45840 \pm 8467$ episodes |

As you can see here the confidence interval is really large which is in agreement with Figure 6b as fluctuations are quite high in comparison to attained level of value function. But this convergence time is almost 2 times faster than convergence time for table representation in 10x10 maze (It was  $86240 \pm 4727$  episodes)

As it can be seen on Figure 6c the resulting shape value function approximation is quite similar to value function in the table representation on Figure 3. The produced policy on Figure 6d is an optimal one as well.

After that I have changed width (SIGMA) to value of 0.5 which makes RBF functions really peaked and very similar to identity functions. That is why here I used 6.0 as convergence criteria from question 3. Here is the convergence results for width (SIGMA) 0.5 on 10x10 maze:

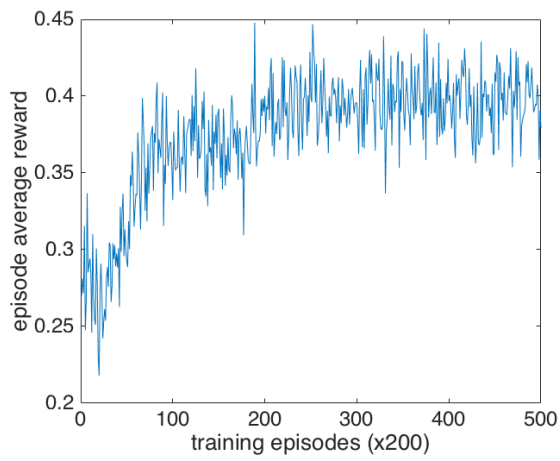
|                            |                           |
|----------------------------|---------------------------|
| Convergence criteria (STD) | 6.0                       |
| Number of experiments      | 10                        |
| Mean convergence time      | $15240 \pm 1172$ episodes |

This result within confidence interval is identical to convergence time achieved in 10x10 maze with identity functions in question 3.

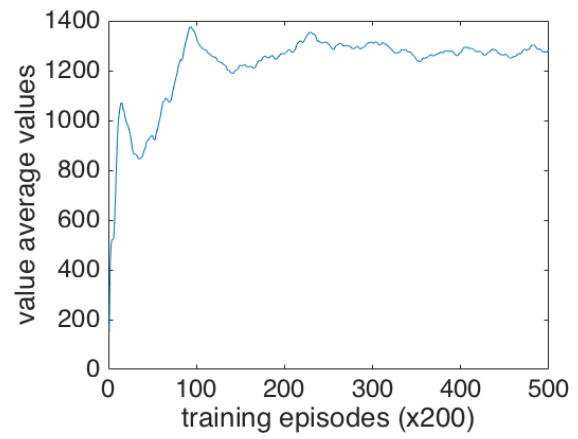
This is also confirmed by Figure 7. On Figure 7c we have value function which has almost step function shape similarly to Figure 4c. Also the policy shown on Figure 7d has actions only in groups by 2x2 patches as on Figure 4d.

For 25x25 case for table Q representation my simulation didn't converge even within 250000 learning episodes. The obtained result for that simulation are present on Figure 8. The resulting policy on Figure 8c almost optimal but there is few states for which action are chosen not optimally.

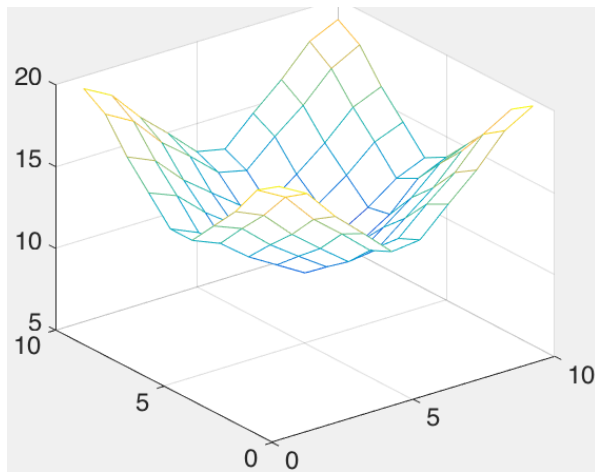
Now let's run simulation for RBF functions in 25x25 maze. It catches general pattern fairly quickly but it fluctuates a lot as it can be seen on 9b. The important distinction why it converges in contrast to Figure 8b is that on Figure 9b it stops growing at some point and resulting generated policy on Figure 8 is an optimal one. The received shape of value function on Figure 9a is similar to Figure 8a but it is more smooth one. By setting convergence criteria to the level of fluctuation on Figure 9b the following convergence time was computed:



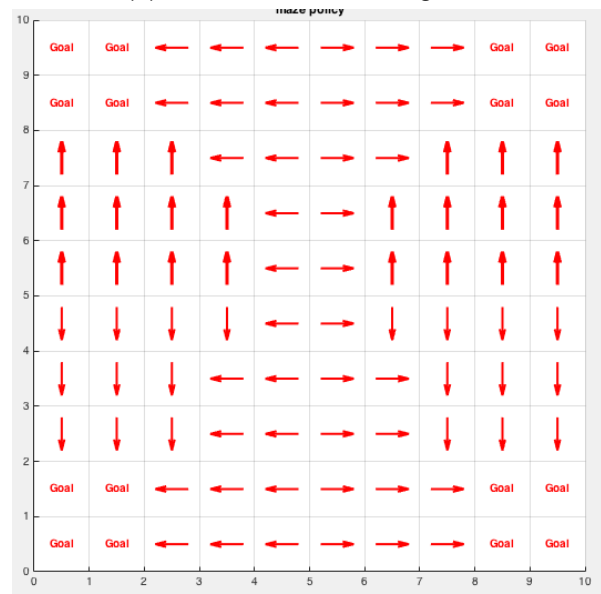
(a) Average reward



(b) Value function convergence



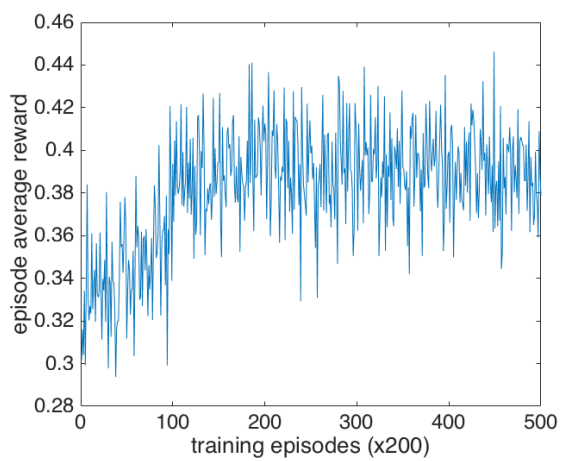
(c) Final value function on maze states



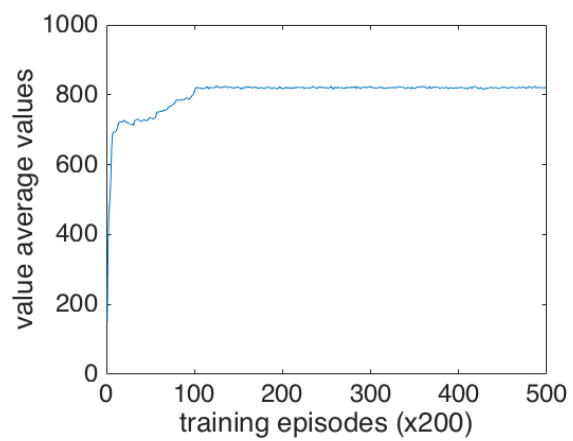
(d) Final policy

Figure 6: 10x10 maze. RBF representation with width 1.0

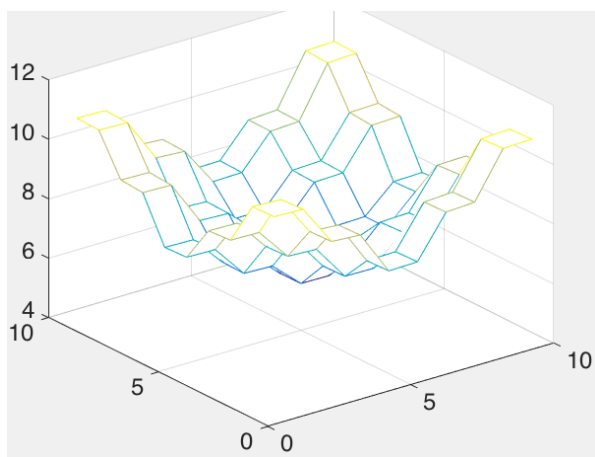




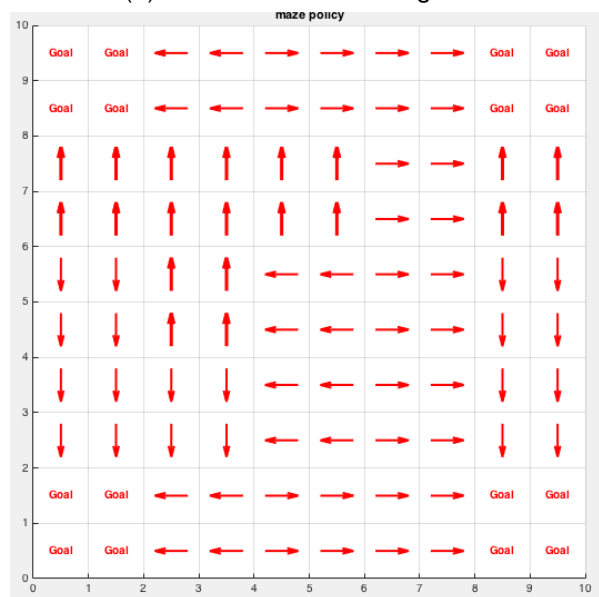
(a) Average reward



(b) Value function convergence

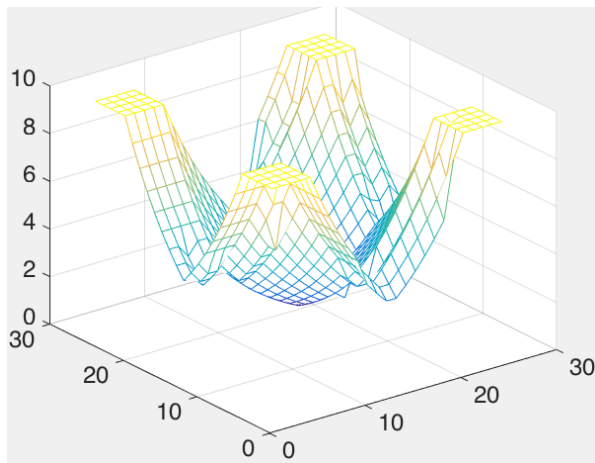


(c) Final value function on maze states

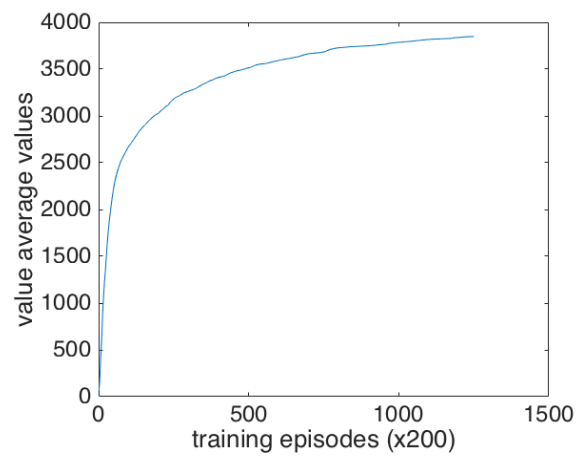


(d) Final policy

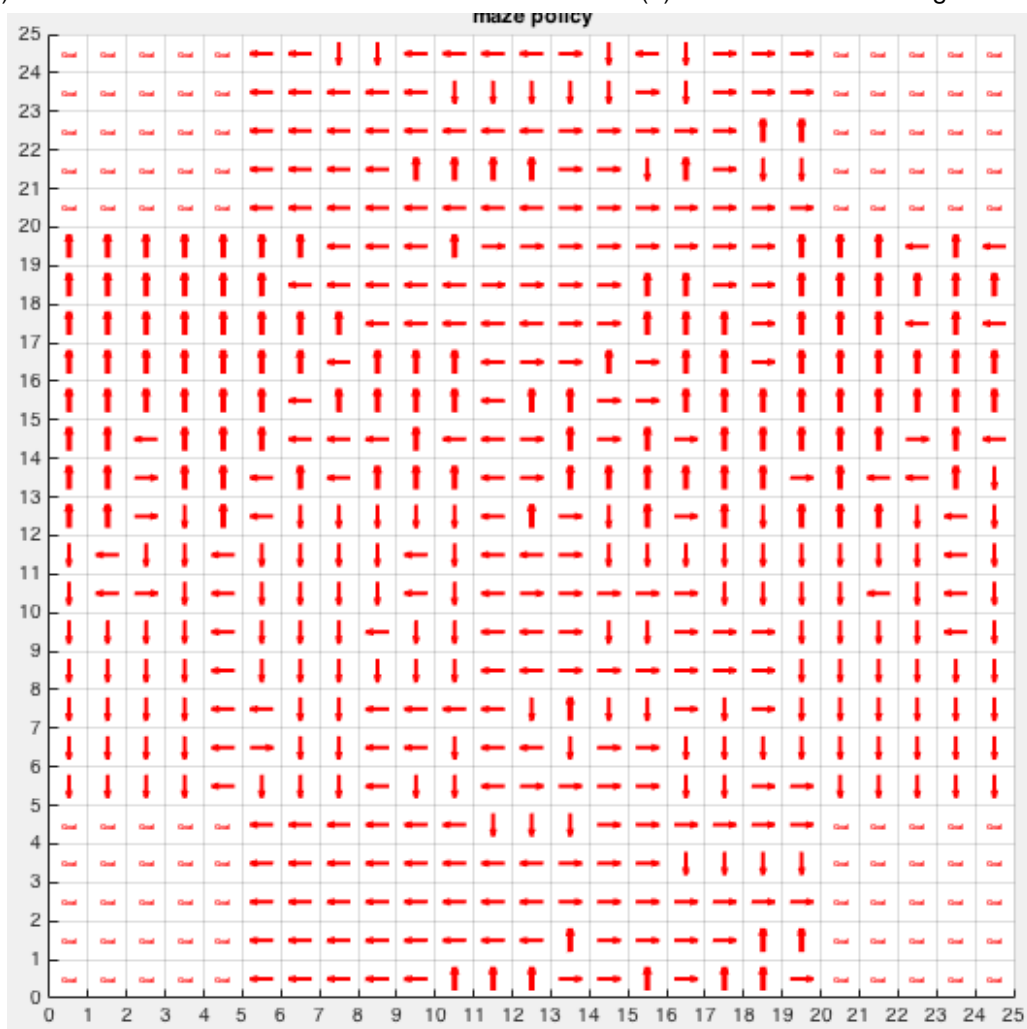
Figure 7: 10x10 maze. RBF representation with width (SIGMA) 0.5



(a) Final value function on maze states



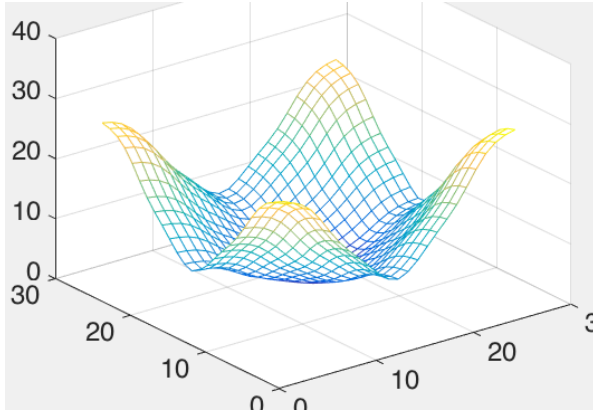
(b) Value function convergence



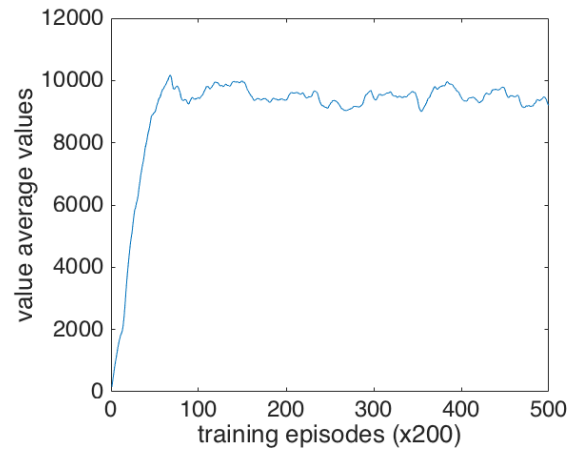
(c) Final policy. Not completely optimal one

Figure 8: 25x25 maze. Table representation

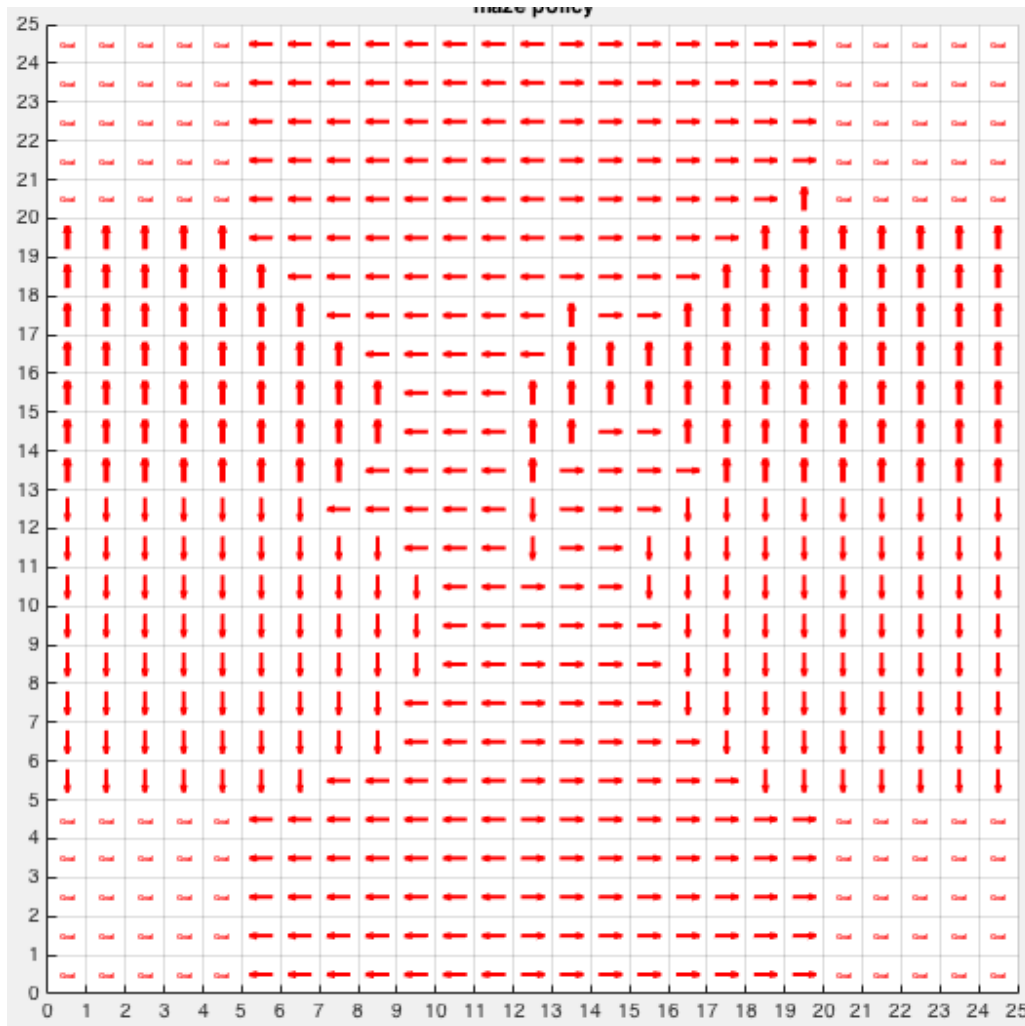
|                            |                           |
|----------------------------|---------------------------|
| Convergence criteria (STD) | 300                       |
| Number of experiments      | 10                        |
| Mean convergence time      | 15540 $\pm$ 2556 episodes |



(a) Final value function on maze states



(b) Value function convergence



(c) Final policy

Figure 9: 25x25 maze. RBF representation with width (SIGMA) 2.5

Now, let set smaller width in order to get more peaked distribution and we should get kind

of step function. So for 25x25 maze I set width (SIGMA) and the results I got are on Figure 10. As we can see on Figure 10a it looks like a step function and it is also confirmed by policy on Figure 10c as actions are the same for patches 5x5. Overall, we can conclude that setting small width (SIGMA) relative to patch size results in highly peaked basis functions which leads to almost identity functions approximation.

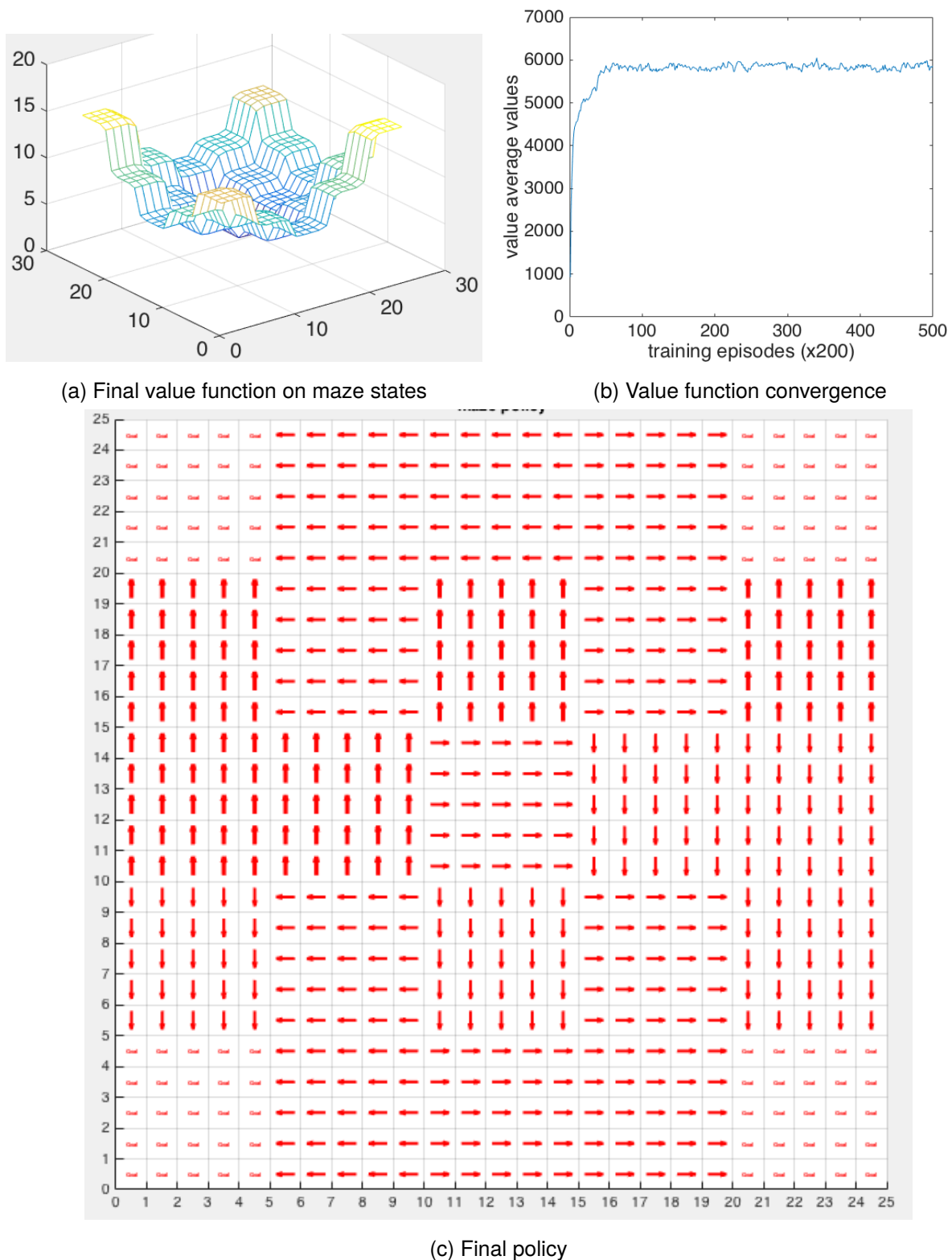
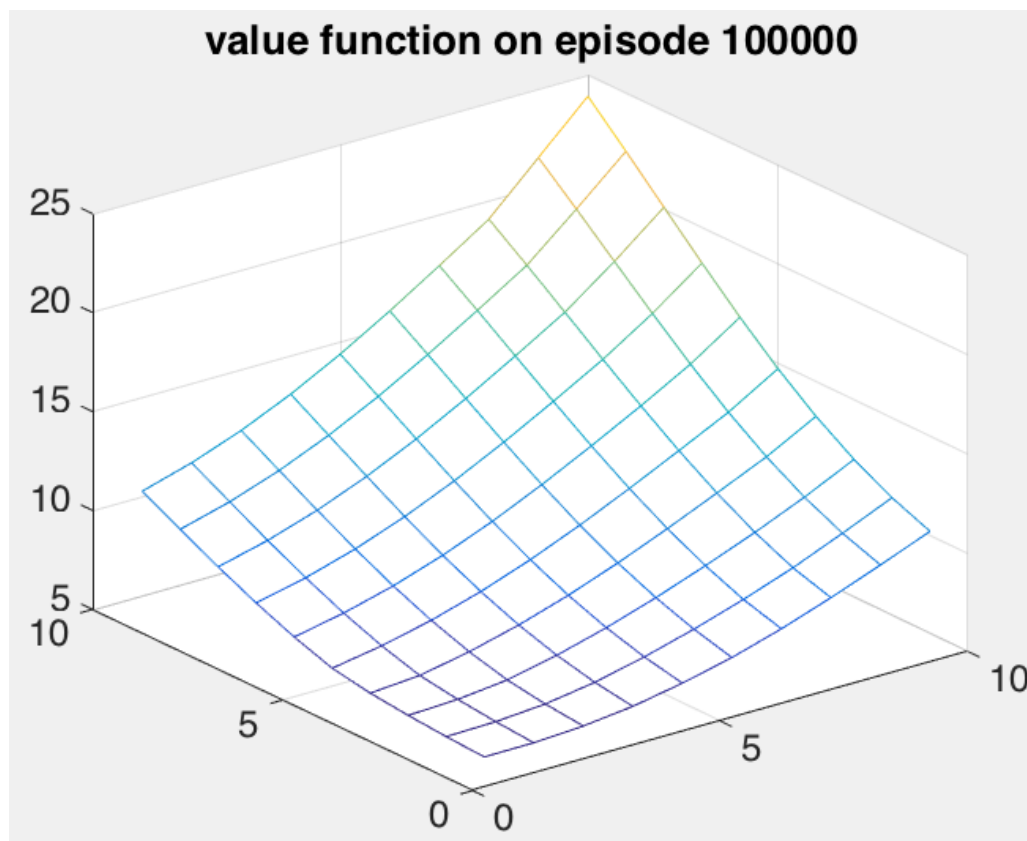


Figure 10: 25x25 maze. RBF representation with width (SIGMA) 1.0

Also for the case of very large width (SIGMA) 4.0 the produced value function has the following shape:



(a) Value function with RBF for width (SIGMA) 4.0

As width (SIGMA) 4.0 is very large in comparison to the size of the maze it means that basically RBF produce everywhere the same values as they are broad and therefore they are able to cover the whole maze. That is why very soon one solution dominates all other. Due to the same instability usual non-smoothed RBF are not suitable for this task as they produce result similar to Figure 11a for all almost any width (SIGMA).

For the final part I have tried different exploration rate (greedy  $\epsilon$  values) for RBF approximation on the 10x10 maze with the following parameters:

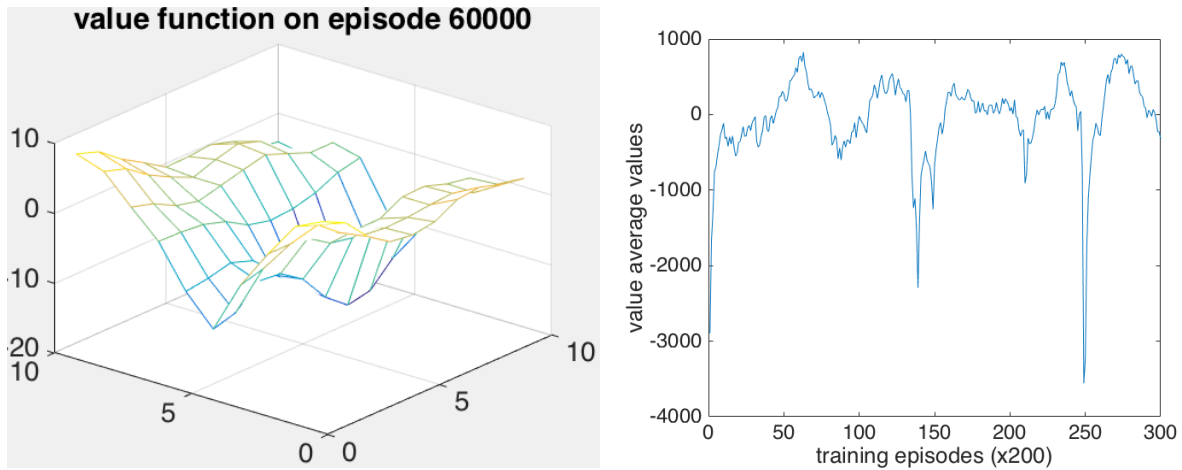
|                            |     |
|----------------------------|-----|
| Learning rate $\eta$       | 0.2 |
| Discount rate $\gamma$     | 0.9 |
| Width (SIGMA)              | 1.0 |
| Convergence criteria (STD) | 25  |
| Number of experiments      | 10  |

After running simulation for fixed  $\epsilon$  values 0.1, 0.2 and 0.4 I got these results

| Exploration rate $\epsilon$ | Convergence time |
|-----------------------------|------------------|
| 0.1                         | 45840 $\pm$ 8467 |
| 0.2                         | 23620 $\pm$ 4008 |
| 0.4                         | 14740 $\pm$ 3409 |

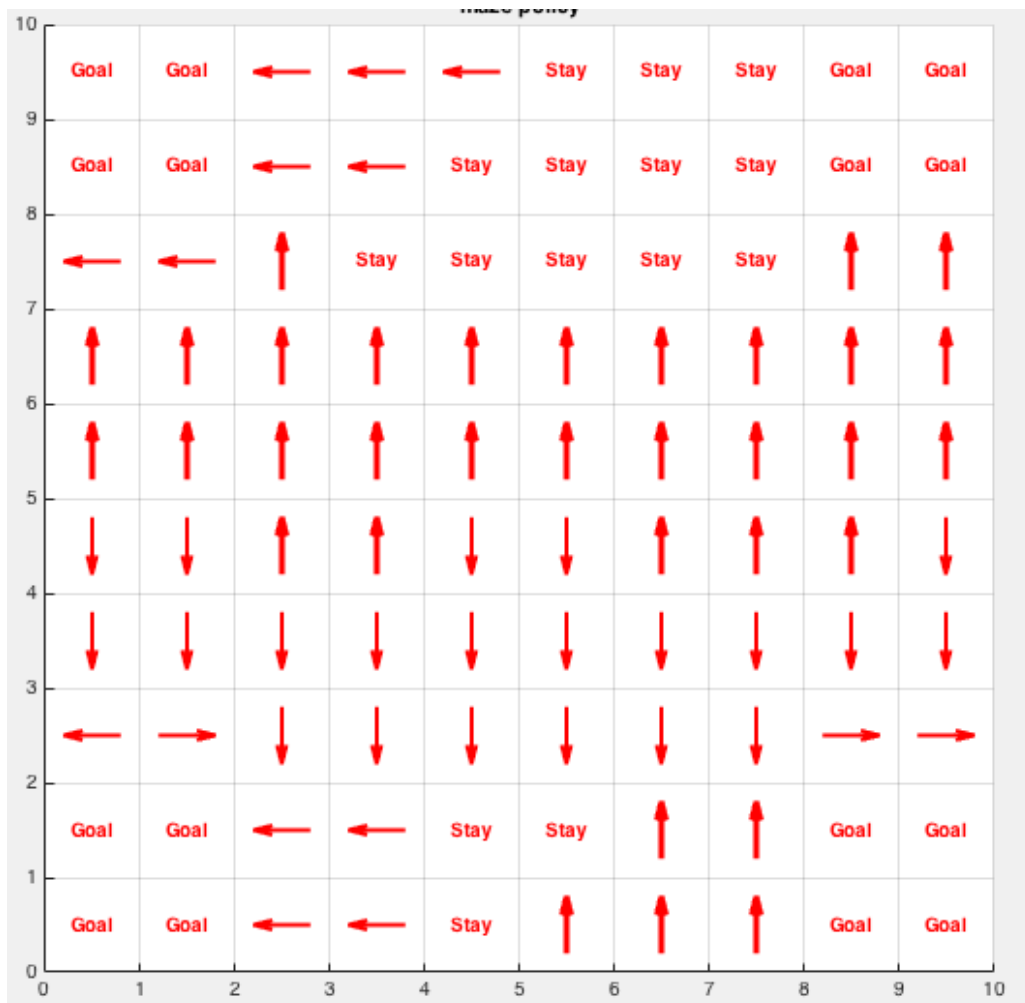
According to this table approximated value function convergence happens faster for bigger exploration rates.

## 5 Obstacles



(a) Final value function on maze states

(b) Value function convergence



(c) Final policy for maze with cliff

Figure 12: 10x10 maze with cliff. RBF representation with width (SIGMA) 1.0

The learning parameters for this task:

|                             |      |
|-----------------------------|------|
| Learning rate $\eta$        | 0.2  |
| Exploration rate $\epsilon$ | 0.1  |
| Discount rate $\gamma$      | 0.9; |

As my obstacles I have created a cliff (reward -10) through the middle of the maze which splits world in two parts. Now the agent must avoid that cliff as punishment significantly out weights reward. However, the agent appears only on one side of the cliff that is why potentially it always has an opportunity to avoid cliff and reach the goal. But, it seems that RBF doesn't allow agent to avoid cliff reliably enough that is why especially on first several episodes it spends a lot of time in the maze. The results are presented on Figure 12. Figure 12c showing policy is especially interesting as agent learns to stay in some states. This can be explained by the fact that agent doesn't receive any discount if it stays on one place in contrast to stepping into cliff which is very dangerous. The only opportunity to escape states where agent is learnt to stay is  $\epsilon$  greedy exploration. But still judging by Figure 12b task becomes very unstable and it doesn't seem to converge. However, if we run simulation for table representation of the task it will produce more symmetric value for the value functions although values will be much smaller than usual as discount for stepping into the cliff is too big. This is shown on Figure 13

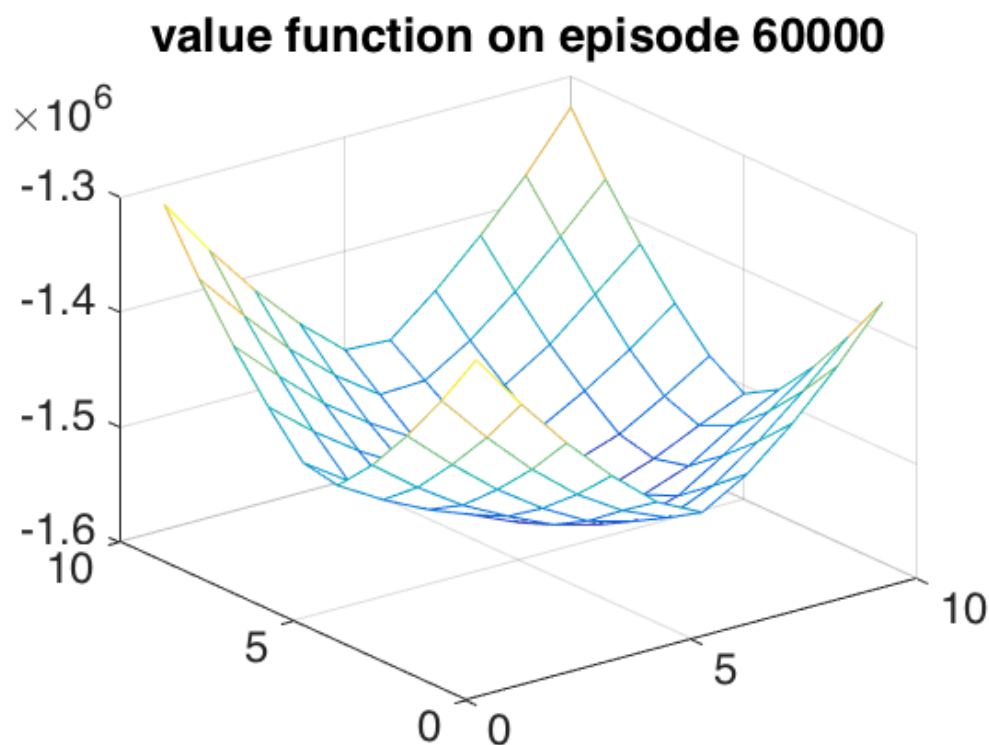


Figure 13: Final value function for 10x10 maze with obstacles using table representation of Q

Finally I tried to run simulation with the same conditions but in the case of 25x25 maze. Algorithm ran very slowly because agent basically learnt to stay on the same place as it is more safe as it doesn't have any discount for that. That is why I has changed the rules by

prohibiting stay action and I also reduced the punishment for stepping into cliff to -1. After that my value functions converged to Figure 14 which essentially has the same shape as in previous questions. It turned out that RBF representation is very sensible to the punishment reward. If punishment becomes lower than -1 then agent prefer to stay on the same place instead of acting and episodes take extremely long time to be completed.

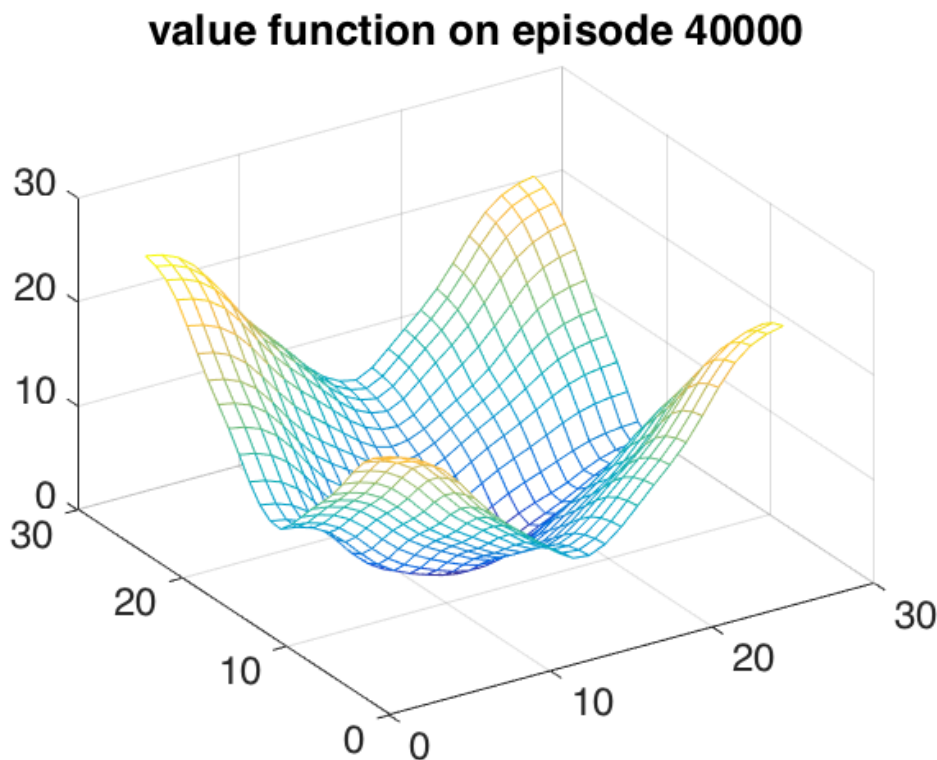


Figure 14: Value function for 25x25 maze with obstacles using RBF

## 6 Conclusion

Here we have seen examples how function approximation leads to faster convergence of the value function although sometime making it less stable. RBF basis functions predictably provided better idea about position of the agent in the world than identity functions. Also, and important observations that RBF with very small width (SIGMA) turns into essentially identity functions. On the other hand, very big value of width (SIGMA) leaves agent disoriented about its true position in the maze as for each state basis functions return almost the same value.



## 7 Code Appendix

main\_v1 runs simulation for tabular Q function:

---

```
1 function [samples, avg_rs, convg_times, actions_map] = main_v1()
2 % a simple illustration of Q-learning: a walker in a 2-d maze
3
4 % number of states
5 Sx = 10; % width of grid world
6 Sy = 10; % length of grid world
7 M = 2; % cover neighbour area
8
9 S = Sx*Sy; % number of states in grid world
10 %4 corners
11 GoalsOrig = [[1, 1];
12              [Sx, 1];
13              [1, Sy];
14              [Sx, Sy]];
15
16 %spread goals across patches MxM
17 Goals = zeros(4 * M * M, 2);
18 shift = M - 1;
19 k = 1;
20 for i = 1:4
21     goal = GoalsOrig(i, :);
22     for p = -shift:shift
23         for q = -shift:shift
24             cand = goal + [p, q];
25             if 1 <= cand(1) && cand(1) <= Sx &&...
26                 1 <= cand(2) && cand(2) <= Sy
27                 Goals(k, :) = cand;
28                 k = k + 1;
29             end
30         end
31     end
32 end
33
34 % number of actions
35 A = 5; % number of actions: E, S, W, N and 0
36 % total number of learning trials
37 T = 60000;
38
39 eta = 0.2;
40 gamma = 0.9;
41 epsilon = 0.1;
```

```

42
43 NUM_EXPS = 1; %number of experiments
44 convg_times = [];
45 actions_map = zeros(Sx, Sy);
46
47 for exp_id = 1:NUM_EXPS
48     samples = zeros(T, 1);
49     avg_rs = zeros(T, 1);
50     %initialisation
51     Q = 0.1*rand(Sx, Sy, A);
52     % run the algorithm for T trials
53     r = 0;
54     for t=1:T
55         % set the starting state
56         x0=randi(Sx);
57         y0=randi(Sy);
58         % start trial
59         avg_r = 0;
60         for u=1:S*S
61             % exploration (epsilon-greedy)
62             if (rand(1)<epsilon)
63                 a0=randi(A);
64             else
65                 [V_s0,a0]=max(Q(x0, y0,:)); % we only need the a0 here
66             end;
67
68             % reward is for reaching the goal and staying there
69             if ismember([x0, y0], Goals, 'rows') %&& (a0==5)
70                 r = r + 1;
71             else
72                 r = r + 0;
73             end
74
75             if x0 == int32(Sx/2) % && y0 == int32(Sy/2)
76                 r = r - 10;
77             end
78
79             %%you can add obstacles simply by adding punishment
80             %%e.g. somewhere in the middle
81             % if (abs(rem(s0,Sx)-Sy/2)<2)&&(abs(idivide(cast(s0,'int8'),cast(↵
82                 Sy,'int8'))-Sx/2)<2)
83             % r=-10;
84             % end;
85             %%The agent can still move through the obstacle, so it's more like↵
86             a
87             %%pothole.

```

```

87     % now moving left, right, up, down, or not
88     % and don't step outside the track
89     % E, S, W, N and 0
90     x1 = x0;
91     y1 = y0;
92     if (a0==1)
93         x1 = min(max(x0+1, 1), Sx); %east
94     elseif (a0==2)
95         y1 = min(max(y0-1, 1), Sy); %south
96     elseif (a0==3)
97         x1 = min(max(x0-1, 1), Sx); %west
98     elseif (a0==4)
99         y1 = min(max(y0+1, 1), Sy); %north
100    end
101
102    % now the learning step
103    V_s1 = max(Q(x1, y1,:));
104    Q(x0, y0, a0)=(1-eta)*Q(x0, y0, a0)+eta*(r+gamma*V_s1);
105
106    % goto next trial once the goal is reached
107    if ismember([x0, y0], Goals, 'rows')
108        avg_r = r/u;
109        %display(avg_r, 'reached the goal!');
110        break;
111    end
112
113    x0=x1;
114    y0=y1;
115 end
116
117 avg_rs(t) = avg_r;
118 VV = max(Q,[],3);
119 samples(t) = sum(reshape(VV, 1,numel(VV)));
120
121 % % you may prefer using a mesh plot as it is a 2D example.
122 if (rem(t, 1000) == 0)
123     display(t, 'current trial');
124 end
125
126 if (rem(t,1000)==0 && NUM_EXPS == 1)
127     figure(6);
128     clf;
129     zlim([0 1/(1-gamma)]);
130     %samples(
131     mesh(VV);
132     %%you may also like to have a look at (a=2, choose also other ↵
        actions):

```

```

133     %Q2 = reshape(Q(:,2),Sx,Sy);
134     %mesh(Q2);
135     title(['value function on episode ', num2str(t)]);
136     fig=gcf;
137     set(findall(fig, '-property', 'FontSize'), 'FontSize', 24)
138     drawnow
139
140     %t % don't know how to put current time in the figure
141     end
142
143 end
144 %end of epoch
145 display(exp_id, 'end of experiment');
146
147 average_over = 200;
148 avg_samples = mean(reshape(samples, average_over, []));
149
150 if NUM_EXPS == 1
151     figure(7);
152     clf;
153     title('value average values');
154     plot(avg_samples);
155     ylabel('value average values');
156     xlabel(['training episodes (x', num2str(average_over), ')']);
157     fig=gcf;
158     set(findall(fig, '-property', 'FontSize'), 'FontSize', 24)
159
160
161     figure(9);
162     clf;
163     title('episode average reward');
164     plot(mean(reshape(avg_rs, average_over, [])));
165     ylabel('episodes average reward');
166     xlabel(['training episodes (x', num2str(average_over), ')']);
167     fig=gcf;
168     set(findall(fig, '-property', 'FontSize'), 'FontSize', 24)
169 end
170
171 %     rel_diff = zeros(1, numel(avg_samples) - 1);
172 %     for t=2:size(avg_samples, 2)
173 %         rel_diff(t) = abs(avg_samples(t) - avg_samples(t-1))/avg_samples(t-1);
174 %     end
175 %
176 %     convg = find(rel_diff > 0.00001, 1, 'last')
177 %     convg_t = convg * average_over
178 %     if ~isempty(convg_t)

```

```

179 %     convg_times = [convg_times, convg_t];
180 %     end
181
182     smp_stds = zeros(1,(size(avg_samples, 2) - 50));
183     for t=1:(size(avg_samples, 2) - 50)
184         smp_std = std(avg_samples(t:(t+50)));
185         smp_stds(t) = smp_std;%abs(avg_samples(t) - avg_samples(t-1))/↵
            avg_samples(t);
186     end
187
188     convg = find(smp_stds > 0.2, 1, 'last')
189     convg_t = convg * average_over
190     if ~isempty(convg_t)
191         convg_times = [convg_times, convg_t];
192     end
193
194     if NUM_EXPS == 1
195         figure(10);
196         clf;
197         title('averaged value function standard deviation over 50');
198         plot(smp_stds);
199         ylabel('standard deviation');
200         xlabel(['training episodes (x', num2str(average_over), ')']);
201         fig=gcf;
202         set(findall(fig,'-property','FontSize'),'FontSize',24)
203     end
204
205     if NUM_EXPS == 1
206         %plot policy
207         actions_map = zeros(Sx, Sy);
208         %initial state in the left bottom corner with
209         %pasenger and must go to G
210
211         for xx=1:Sx
212             for yy=1:Sy
213                 %perform greedy action
214                 [V_s,aa]=max(Q(xx, yy,:));
215                 actions_map(xx, yy) = aa;
216             end
217         end
218
219         %lbls_actions = grid_to_lbls(actions_map);
220         %lbls_actions = actions2str(lbls_actions);
221         figure(31);
222         clf;
223         draw_maze(actions_map, 'maze policy', Sx, Sy, Goals);
224     end

```

```

225     end
226 %end of experiments loop
227 SE = std(convg_times)/sqrt(length(convg_times));
228 display(mean(convg_times), 'mean convergence time');
229 display(SE, 'standard error');
230 display(SE * 1.96, 'confidence interval');
231 end

```

---

main\_v2 runs simulation for identity basis functions:

---

```

1 function [samples, avg_rs, convg_times, actions_map] = main_v2()
2 % a simple illustration of Q-learning: a walker in a 2-d maze
3
4 % number of states
5 Sx = 10; % width of grid world
6 Sy = 10; % length of grid world
7 M = 2; % cover neighbour area
8
9 S = Sx*Sy; % number of states in grid world
10 %4 corners
11 GoalsOrig = [[1, 1];
12              [Sx, 1];
13              [1, Sy];
14              [Sx, Sy]];
15
16 %spread goals across patches MxM
17 Goals = zeros(4 * M * M, 2);
18 shift = M - 1;
19 k = 1;
20 for i = 1:4
21     goal = GoalsOrig(i, :);
22     for p = -shift:shift
23         for q = -shift:shift
24             cand = goal + [p, q];
25             if 1 <= cand(1) && cand(1) <= Sx &&...
26                 1 <= cand(2) && cand(2) <= Sy
27                 Goals(k, :) = cand;
28                 k = k + 1;
29             end
30         end
31     end
32 end
33
34 % number of actions
35 A = 5; % number of actions: E, S, W, N and 0

```

```

36 % total number of learning trials
37 T = 75000;
38
39 eta = 0.2;
40 gamma = 0.9;
41 epsilon = 0.1;
42
43 NUM_EXPS = 10; %number of experiments
44 convg_times = [];
45 actions_map = zeros(Sx, Sy);
46
47
48 F = zeros(Sx, Sy, 25);
49 %precompute basis functions in advance
50 for ix = 1:Sx
51     for iy = 1:Sy
52         %identity functions
53         for kx = 1:5
54             for ky = 1:5
55                 fId = 5 * (ky - 1) + kx;
56                 F(ix, iy, fId) = M*(kx-1) < ix && ix <= M*kx &&...
57                                     M*(ky-1) < iy && iy <= M*ky;
58             end
59         end
60     end
61 end
62
63 for exp_id = 1:NUM_EXPS
64     samples = zeros(T, 1);
65     avg_rs = zeros(T, 1);
66     %initialisation
67     %weights
68     W = zeros(A,25);%0.1*rand(A, 25);
69     % run the algorithm for T trials
70     for t=1:T
71         % set the starting state
72         x0=randi(Sx);
73         y0=randi(Sy);
74         % start trial
75         avg_r = 0;
76         for u=1:S*S
77             % exploration (epsilon-greedy)
78             if (rand(1)<epsilon)
79                 a0=randi(A);
80             else
81                 [V_s0, a0]=max(W * reshape(F(x0, y0, :), 25, 1)); % we only ↵
                                     need the a0 here

```

```

82     end;
83
84     % reward is for reaching the goal and staying there
85     if ismember([x0, y0], Goals, 'rows') %&& (a0==5)
86         r=1;
87     else
88         r=0;
89     end
90
91     %%you can add obstacles simply by adding punishment
92     %%e.g. somewhere in the middle
93     % if (abs(rem(s0,Sx)-Sy/2)<2)&&(abs(idivide(cast(s0,'int8'),cast(↵
94         Sy,'int8'))-Sx/2)<2)
95     % r=-10;
96     % end;
97     %%The agent can still move through the obstacle, so it's more like↵
98     a
99     %%pothole.
100
101     % now moving left, right, up, down, or not
102     % and don't step outside the track
103     % E, S, W, N and 0
104     x1 = x0;
105     y1 = y0;
106     if (a0==1)
107         x1 = min(max(x0+1, 1), Sx); %east
108     elseif (a0==2)
109         y1 = min(max(y0-1, 1), Sy); %south
110     elseif (a0==3)
111         x1 = min(max(x0-1, 1), Sx); %west
112     elseif (a0==4)
113         y1 = min(max(y0+1, 1), Sy); %north
114     end
115
116     % now the learning step
117     %approximate value function V(s1)
118     V_s1 = max(W * reshape(F(x1, y1, :), 25, 1));
119     %approximate Q(s0, a0)
120     Q_s0_a0 = W(a0, :) * reshape(F(x0, y0, :), 25, 1);
121     delta = r+gamma*V_s1 - Q_s0_a0;
122     %basic functions in state x0 y0
123     bfs = reshape(F(x0, y0, :), 1, 25);
124     W(a0, :) = W(a0, :) + eta * delta * bfs;
125
126     % goto next trial once the goal is reached
127     if ismember([x0, y0], Goals, 'rows')
128         avg_r = r/u;

```



```

127         %display(avg_r, 'reached the goal!');
128         break;
129     end
130
131     x0=x1;
132     y0=y1;
133 end
134
135     avg_rs(t) = avg_r;
136     % reshape(F, Sx * Sy, 25) % Sx*Sy x 25
137     % W' % 25 x A
138     QQ = reshape(F, Sx * Sy, 25) * W'; % Sx*Sy x A
139     VV = max(QQ,[],2); %Sx*Sy
140     VV = reshape(VV, Sx, Sy);
141
142
143     %VV = max(Q,[],3);
144     samples(t) = sum(reshape(VV, 1,numel(VV)));
145
146     % % you may prefer using a mesh plot as it is a 2D example.
147     if (rem(t, 1000) == 0)
148         display(t, 'current trial');
149     end
150
151     if (rem(t,1000)==0 && NUM_EXPS == 1)
152         figure(6);
153         clf;
154         zlim([0 1/(1-gamma)]);
155         %samples(
156         mesh(VV);
157         %%you may also like to have a look at (a=2, choose also other ↵
158         %Q2 = reshape(Q(:,2),Sx,Sy);
159         %mesh(Q2);
160         title(['value function on episode ', num2str(t)]);
161         fig=gcf;
162         set(findall(fig,'-property','FontSize'),'FontSize',24)
163         drawnow
164
165         %t % don't know how to put current time in the figure
166     end
167
168 end
169 %end of epoch
170 display(exp_id, 'end of experiment');
171
172     average_over = 200;

```

```

173     avg_samples = mean(reshape(samples, average_over, []));
174
175     if NUM_EXPS == 1
176         figure(7);
177         clf;
178         title('value average values');
179         plot(avg_samples);
180         ylabel('value average values');
181         xlabel(['training epsidoes (x', num2str(average_over), ')']);
182         fig=gcf;
183         set(findall(fig, '-property', 'FontSize'), 'FontSize', 24)
184
185         figure(9);
186         clf;
187         title('episode average reward');
188         plot(mean(reshape(avg_rs, average_over, [])));
189         ylabel('episode average reward');
190         xlabel(['training epsidoes (x', num2str(average_over), ')']);
191         fig=gcf;
192         set(findall(fig, '-property', 'FontSize'), 'FontSize', 24)
193     end
194
195     smp_stds = zeros(1, (size(avg_samples, 2) - 50));
196     for t=1:(size(avg_samples, 2) - 50)
197         smp_std = std(avg_samples(t:(t+50)));
198         smp_stds(t) = smp_std; %abs(avg_samples(t) - avg_samples(t-1))/↔
            avg_samples(t);
199     end
200
201     convg = find(smp_stds > 6, 1, 'last')
202     convg_t = convg * average_over
203     if ~isempty(convg_t)
204         convg_times = [convg_times, convg_t];
205     end
206
207     if NUM_EXPS == 1
208         figure(10);
209         clf;
210         title('averaged value function standard deviation over 50');
211         plot(smp_stds);
212         ylabel('standard deviation');
213         xlabel(['training episodes (x', num2str(average_over), ')']);
214         fig=gcf;
215         set(findall(fig, '-property', 'FontSize'), 'FontSize', 24)
216     end
217
218     if NUM_EXPS == 1

```

```

219     %plot policy
220     actions_map = zeros(Sx, Sy);
221     %initial state in the left bottom corner with
222     %passenger and must go to G
223     QQ = reshape(F, Sx * Sy, 25) * W'; % Sx*Sy x A
224     Q = reshape(QQ, Sx, Sy, A);
225     for xx=1:Sx
226         for yy=1:Sy
227             %perform greedy action
228             [V_s,aa]=max(Q(xx, yy,:));
229             actions_map(xx, yy) = aa;
230         end
231     end
232
233     %lbls_actions = grid_to_lbls(actions_map);
234     %lbls_actions = actions2str(lbls_actions);
235     figure(31);
236     clf;
237     draw_maze(actions_map, 'maze policy', Sx, Sy, Goals);
238 end
239 end
240 %end of experiments loop
241 SE = std(convg_times)/sqrt(length(convg_times));
242 display(mean(convg_times), 'mean convergence time');
243 display(SE, 'standard error');
244 display(SE * 1.96, 'confidence interval');
245 end

```

---

main\_v3 runs simulation for RBF functions:

---

```

1 function [samples, avg_rs, convg_times, actions_map] = main_v3()
2     % a simple illustration of Q-learning: a walker in a 2-d maze
3
4     % number of states
5     Sx = 10; % width of grid world
6     Sy = 10; % length of grid world
7     M = 2; % cover neighbouring area
8
9     S = Sx*Sy; % number of states in grid world
10    %4 corners
11    GoalsOrig = [[1, 1];
12                 [Sx, 1];
13                 [1, Sy];
14                 [Sx, Sy]];
15    % GoalsOrig = [[1, 1];

```

```

16 %           [Sx, Sy]];
17
18 %spread goals across patches MxM
19 Goals = zeros(size(GoalsOrig, 1) * M * M, 2);
20 shift = M - 1;
21 k = 1;
22 for i = 1:size(GoalsOrig, 1)
23     goal = GoalsOrig(i, :);
24     for p = -shift:shift
25         for q = -shift:shift
26             cand = goal + [p, q];
27             if 1 <= cand(1) && cand(1) <= Sx &&...
28                 1 <= cand(2) && cand(2) <= Sy
29                 Goals(k, :) = cand;
30                 k = k + 1;
31             end
32         end
33     end
34 end
35
36 %Goals = GoalsOrig;
37
38 % number of actions
39 A = 5;    % number of actions: E, S, W, N and 0
40 % total number of learning trials
41 T = 60000;
42
43 eta = 0.2;
44 gamma = 0.9;
45 epsilon = 0.1;%0.2;%0.4
46
47 NUM_EXPS = 1; %number of experiments
48 convg_times = [];
49 actions_map = zeros(Sx, Sy);
50
51 %1.0 for 10 is smooth like true function, 0.5 is almost like identity ↔
    function
52 %2.5 for 25 is smooth like true function (epsilon must be 0.1), 1.0 for ↔
    25 is almost like identity function
53 %0.5
54 SIGMA = 1.0;%1.0;
55
56 F = zeros(Sx, Sy, 25);
57
58 %precompute basis functions in advance
59
60

```

```

61     for ix = 1:Sx
62         for iy = 1:Sy
63             %identity functions
64             for kx = 1:5
65                 for ky = 1:5
66                     fId = 5 * (ky - 1) + kx;
67                     %position of rbf in maze coordinates
68                     cx = kx*M - (M - 1)/2;
69                     cy = ky*M - (M - 1)/2;
70                     E = exp(-(1/(2*SIGMA^2))*((ix-cx)^2+(iy-cy)^2));
71                     Z = 1;%2*pi*SIGMA^2;
72                     F(ix, iy, fId) = E/Z;
73                 end
74             end
75             Z = sum(F(ix, iy, :));
76             F(ix, iy, :) = F(ix, iy, :) / Z;
77         end
78     end
79
80     for exp_id = 1:NUM_EXPS
81         samples = zeros(T, 1);
82         avg_rs = zeros(T, 1);
83         %initialisation
84         %weights
85         W = zeros(A, 25);%0.1*rand(A, 25);
86         % run the algorithm for T trials
87         for t=1:T
88             % set the starting state
89             x0=randi(Sx);
90             y0=randi(Sy);
91             % start trial
92             avg_r = 0;
93             r = 0;
94             for u=1:S*S
95                 % exploration (epsilon-greedy)
96                 if (rand(1)<epsilon)
97                     a0=randi(A);
98                 else
99                     [V_s0, a0]=max(W * reshape(F(x0, y0, :), 25, 1)); % we only ↔
100                     need the a0 here
101                 end;
102
103                 % reward is for reaching the goal and staying there
104                 if ismember([x0, y0], Goals, 'rows') %&& (a0==5)
105                     r = r + 1;
106                 else
107                     r = r + 0;

```

```

107     end
108
109     %%you can add obstacles simply by adding punishment
110     %%e.g. somewhere in the middle
111     % if (abs(rem(s0,Sx)-Sy/2)<2)&&(abs(idivide(cast(s0,'int8'),cast(↵
        Sy,'int8'))-Sx/2)<2)
112     % r=-10;
113     % end;
114     %%The agent can still move through the obstacle, so it's more like↵
        a
115     %%pothole.
116     if x0 == int32(Sx/2) % && y0 == int32(Sy/2)
117         r = r - 10;
118     end
119     % now moving left, right, up, down, or not
120     % and don't step outside the track
121     % E, S, W, N and 0
122     x1 = x0;
123     y1 = y0;
124     if (a0==1)
125         x1 = min(max(x0+1, 1), Sx); %east
126     elseif (a0==2)
127         y1 = min(max(y0-1, 1), Sy); %south
128     elseif (a0==3)
129         x1 = min(max(x0-1, 1), Sx); %west
130     elseif (a0==4)
131         y1 = min(max(y0+1, 1), Sy); %north
132     end
133
134     % now the learning step
135     %approximate value function V(s1)
136     V_s1 = max(W * reshape(F(x1, y1, :), 25, 1));
137     %approximate Q(s0, a0)
138     Q_s0_a0 = W(a0, :) * reshape(F(x0, y0, :), 25, 1);
139     delta = r+gamma*V_s1 - Q_s0_a0;
140     %basic functions in state x0 y0
141     bfs = reshape(F(x0, y0, :), 1, 25);
142     W(a0, :) = W(a0, :) + eta * delta * bfs;
143
144     % goto next episode once the goal is reached
145     if ismember([x0, y0], Goals, 'rows')
146         avg_r = r/u;
147         %display(avg_r, 'reached the goal!');
148         break;
149     end
150
151     x0=x1;

```

```

152     y0=y1;
153 end
154
155     avg_rs(t) = avg_r;
156     % reshape(F, Sx * Sy, 25) % Sx*Sy x 25
157     % W' % 25 x A
158     QQ = reshape(F, Sx * Sy, 25) * W'; % (Sx*Sy x 25) * (25 x A) = Sx*Sy↔
        x A
159     VV = max(QQ,[],2); %Sx*Sy
160     VV = reshape(VV, Sx, Sy);
161
162
163     %VV = max(Q,[],3);
164     samples(t) = sum(reshape(VV, 1,numel(VV)));
165
166     % % you may prefer using a mesh plot as it is a 2D example.
167     if (rem(t, 1000) == 0)
168         display(t, 'current trial');
169     end
170
171     if (rem(t,1000)==0 && NUM_EXPS == 1)
172         figure(6);
173         clf;
174         zlim([0 1/(1-gamma)]);
175         %samples(
176         mesh(VV);
177         %%you may also like to have a look at (a=2, choose also other ↔
            actions):
178         %Q2 = reshape(Q(:,2),Sx,Sy);
179         %mesh(Q2);
180         title(['value function on episode ', num2str(t)]);
181         fig=gcf;
182         set(findall(fig,'-property','FontSize'),'FontSize',24)
183         drawnow
184
185         %t % don't know how to put current time in the figure
186     end
187
188 end
189 %end of epoch
190 display(exp_id, 'end of experiment');
191
192 average_over = 200;
193 avg_samples = mean(reshape(samples, average_over, []));
194
195 if NUM_EXPS == 1
196     figure(7);

```

```

197     clf;
198     title('value average values');
199     plot(avg_samples);
200     ylabel('value average values');
201     xlabel(['training episodes (x', num2str(average_over), ')']);
202     fig=gcf;
203     set(findall(fig, '-property', 'FontSize'), 'FontSize', 24)
204
205     figure(9);
206     clf;
207     title('episode average reward');
208     plot(mean(reshape(avg_rs, average_over, [])));
209     ylabel('episode average reward');
210     xlabel(['training episodes (x', num2str(average_over), ')']);
211     fig=gcf;
212     set(findall(fig, '-property', 'FontSize'), 'FontSize', 24)
213 end
214
215 smp_stds = zeros(1, (size(avg_samples, 2) - 50));
216 for t=1:(size(avg_samples, 2) - 50)
217     smp_std = std(avg_samples(t:(t+50)));
218     smp_stds(t) = smp_std; %abs(avg_samples(t) - avg_samples(t-1))/↵
        avg_samples(t);
219 end
220
221 convg = find(smp_stds > 25, 1, 'last')
222 convg_t = convg * average_over
223 if ~isempty(convg_t)
224     convg_times = [convg_times, convg_t];
225 end
226
227 if NUM_EXPS == 1
228     figure(10);
229     clf;
230     title('averaged value function standard deviation over 50');
231     plot(smp_stds);
232     ylabel('standard deviation');
233     xlabel(['training episodes (x', num2str(average_over), ')']);
234     fig=gcf;
235     set(findall(fig, '-property', 'FontSize'), 'FontSize', 24)
236 end
237
238 if NUM_EXPS == 1
239     %plot policy
240     actions_map = zeros(Sx, Sy);
241     %initial state in the left bottom corner with
242     %pasenger and must go to G

```



```

243     QQ = reshape(F, Sx * Sy, 25) * W'; % Sx*Sy x A
244     Q = reshape(QQ, Sx, Sy, A);
245     for xx=1:Sx
246         for yy=1:Sy
247             %perform greedy action
248             [V_s,aa]=max(Q(xx, yy,:));
249             actions_map(xx, yy) = aa;
250         end
251     end
252
253     %lbls_actions = grid_to_lbls(actions_map);
254     %lbls_actions = actions2str(lbls_actions);
255     figure(31);
256     clf;
257     draw_maze(actions_map, 'maze policy', Sx, Sy, Goals);
258 end
259 end
260 %end of experiments loop
261 SE = std(convg_times)/sqrt(length(convg_times));
262 display(mean(convg_times), 'mean convergence time');
263 display(SE, 'standard error');
264 display(SE * 1.96, 'confidence interval');
265 end

```

---

supporting utility to draw policy:

---

```

1 function [] = draw_maze(lbls, name, Sx, Sy, Goals)
2     %lbls - labels - cell array of strings to display in the maze grid
3     %size X*Y
4     %name - plot title
5     %draw location names with light green color
6
7     hold on;
8     grid on;
9     h = findobj(gca,'Type','Text');
10    delete(h);
11
12    title(name);
13    xlim([0 Sx]);
14    ylim([0 Sy]);
15    set(gca,'XTick',0:Sx);
16    set(gca,'YTick',0:Sy);
17    %to have nice grid cells as squares
18    axis square;
19

```

```

20 % grid domains
21 xg = 1:Sx;
22 yg = 1:Sy;
23 %label coordinates, subtract 0.5 for centering
24 [xlbl, ylbl] = meshgrid(xg - 0.5, yg - 0.5);
25 drawArrow = @(x,y,varargin) quiver( x(1),y(1),x(2)-x(1),y(2)-y(1),0, ↵
    varargin{:} );
26 for ix = 1:Sx
27     for iy = 1:Sy
28         xx = xlbl(ix, iy);
29         yy = ylbl(ix, iy);
30         a = lbls(ix, iy);
31         if ismember([ix, iy], Goals, 'rows')
32             text(xx, yy, 'Goal',...
33                 'Color', 'r',...
34                 'FontWeight', 'bold', 'FontUnits', 'normalized', 'FontSize', ↵
35                     1/( 5* Sy),...
36                     'HorizontalAlignment', 'center',...
37                     'VerticalAlignment', 'middle');
38         elseif a == 1
39             %right
40             drawArrow([xx, xx], [yy - 0.3, yy + 0.3], 'MaxHeadSize',10,'↵
41                 'Color','r','LineWidth',3);
42         elseif a == 2
43             %down
44             drawArrow([xx + 0.3, xx - 0.3], [yy,yy], 'MaxHeadSize',12,'Color↵
45                 ', 'r', 'LineWidth', 2);
46         elseif a == 3
47             %left
48             drawArrow([xx, xx], [yy + 0.3, yy - 0.3], 'MaxHeadSize',12,'Color↵
49                 ', 'r', 'LineWidth', 2);
50         elseif a == 4
51             %up
52             drawArrow([xx - 0.3, xx + 0.3], [yy,yy], 'MaxHeadSize',12,'Color↵
53                 ', 'r', 'LineWidth', 2);
54         elseif a == 5
55             text(xx, yy, 'Stay',...
56                 'Color', 'r',...
57                 'FontWeight', 'bold', 'FontUnits', 'normalized', 'FontSize', ↵
58                     1/( 5* Sy),...
59                     'HorizontalAlignment', 'center',...
60                     'VerticalAlignment', 'middle');
61         end
62     end
63 end

```

```
60 %      text(xlbl(:), ylbl(:), lbls(:), 'color', 'r', ...
61 %          'HorizontalAlignment', 'center', 'VerticalAlignment', 'middle', ...
62 %          'FontWeight', 'bold', 'FontUnits', 'normalized', 'FontSize', ←
        1/(5*Y));
63
64     hold off;
65     drawnow;
66 end
```

---