



Problem Set III

Karan Chawla
karangc2@illinois.edu

Course Instructor: **Prof. G. Chowdhary**

April 28, 2017

1

¹This assignment was done with Wyatt McAllistar

1 QUESTION 1

Consider a 5 by 5 grid world. We would like to find a path from one end of the grid to the other end, that is from location (1,1) to location (5,5). A reward of 1 is obtained when the agent reaches (5,5), a reward of -0.01 is obtained for every transition. There are 5 actions: stay, go left, go up, go right, go down. Each action results in a stochastic transition, with 90% probability of the intended transition, and 10% probability of a random transition.

What would be the state transition matrix, what would be the reward function? Assuming the state transition and reward function is known, implement the following algorithms to find this path:

1. Value iteration (Algorithm 2 from Geramifard et al.).
2. TBVI (Algorithm 3 from Geramifard et al.).

If the state transition and the reward matrices are not known, implement the following reinforcement learning algorithm:

1. SARSA (Algorithm 6 from Geramifard et al.).

For all algorithms, implement the following value function approximations:

1. Tabular.
2. Radial Basis Function network, limit the maximum number of basis to 20.
3. (Bonus) Gaussian Process (5 points extra for each algorithm you implement GP for).

An implementation of Q learning with Tabular, Radial Basis, and GP function approximation has been provided. The main file is *gridworldQLGPMmain.m*. It has been commented for your review. You need not use this implementation if you do not wish to.

For implementing the SARSA algorithm, you could change the Q learning update law with the SARSA one in the given code.

For implementing the MDP algorithms, you have you change the code to use the transition and reward models in the computation of the policy. Else, you can write code from scratch, which in this case could be easier since an episodic treatment that is done in the RL code is not needed for solving MDPs when the reward and transition models are known.

Note that the transitions in this grid world are stochastic, as opposed to the deterministic transitions in the grid world in Problem Set 2.

Present your answer as plots and a discussion. For the dynamic programming problems, you will evaluate the performance of your algorithm over a series of 100 Monte-Carlo runs. The stochasticity comes from the random transitions.

The RL algorithms are set to run in 5 executions of 200 episodes each with 100 evaluations (samples).

The files given to you plot the results with their mean and standard deviations. This is the kind of plot you should be presenting.

SOLUTION

For value iteration, the policy and reward converged for both tabular and radial basis function implementations, and are shown in Table 1, and Table 2 below.

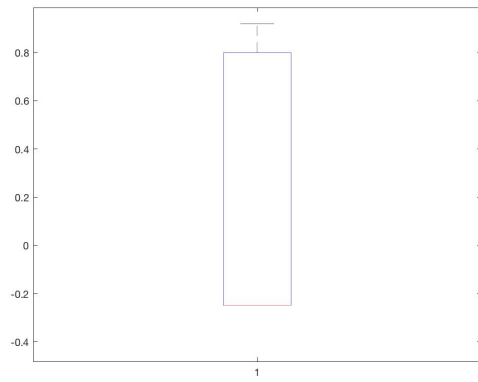
Table 1: Value Iteration Tabular

Value					Policy				
0.7372	0.8129	0.8979	0.9912	1.0000	2	2	2	2	1
0.6729	0.7409	0.8169	0.9001	0.9912	2	2	2	2	3
0.6108	0.6717	0.7409	0.8169	0.8979	3	2	2	3	3
0.5539	0.6085	0.6717	0.7409	0.8129	3	2	3	3	3
0.5029	0.5539	0.6108	0.6729	0.7372	2	2	2	3	3

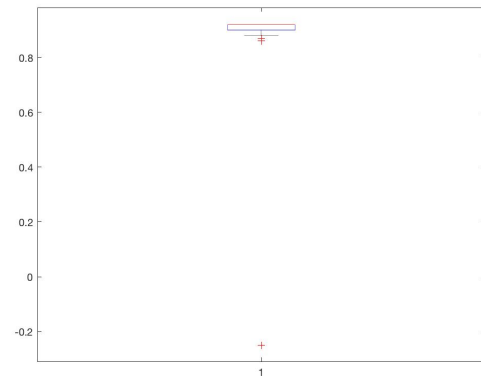
Table 2: Value Iteration RBF

Value					Policy				
0.5675	0.5823	0.6363	0.8334	1.0000	2	2	2	2	1
0.5772	0.6014	0.6379	0.7443	0.8333	2	2	2	3	3
0.5729	0.6002	0.6255	0.6373	0.6361	2	2	3	3	3
0.5491	0.5719	0.5996	0.6011	0.5836	3	3	3	3	3
0.5313	0.5490	0.5727	0.5785	0.5707	3	2	3	3	3

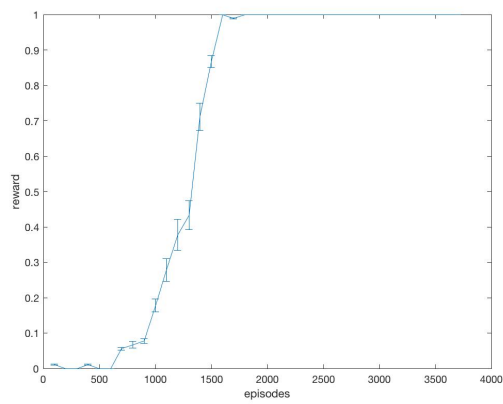
The figures on the following page show the six learning algorithms. It can see, there is a large difference between on policy and off policy learning for this trivial example. Because the state space is really small, and because there are a range of optimal paths, the on-policy algorithms such as SARSA and TBVI do not converge well. In the case of the RBF implementation, TBVI will converge, but SARSA never does not. Q learning always converges as it is an off-policy method.



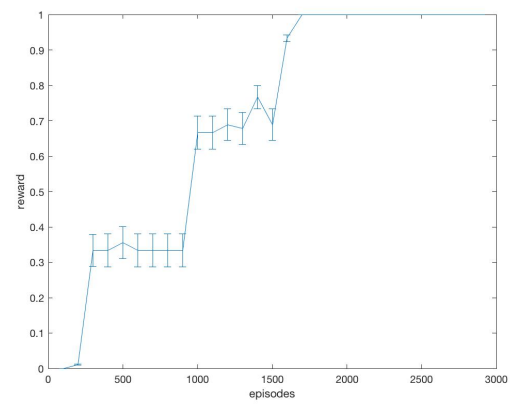
(a) TBVI Tabular



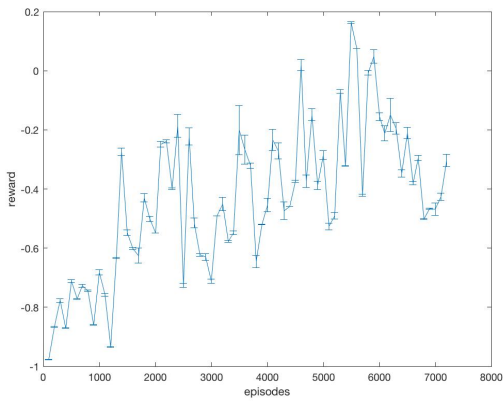
(b) TBVI RBF



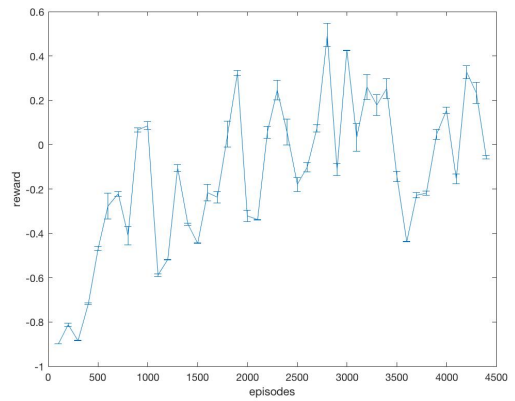
(c) Q Tabular



(d) Q RBF



(e) SARSA Tabular



(f) SARSA RBF

Figure 1.1: Learning Algorithms

The figures above shown the six learning algorithms.

MATLAB CODE FOR VALUE ITERATION

```

1 clear all; close all; clc;
2
3 %Initialize Variables
4 basis      = 1;                % basis
5 N_grid     = 5;                % grid size
6 N_state    = N_grid*N_grid;    % state size
7 N_act      = 5;                % action size
8 noise      = 0.1;              % transition uncertainty
9 gamma      = 0.9;              % discount factor
10 eta       = 0.05;              % tolerance
11 alpha0     = 0.5;              % initial Learning Rate
12 alpha_exp  = 0.5;              % learning rate
13
14 % create value function table and policy
15 v          = zeros(N_state, 1); % value function
16 old_v      = zeros(N_state, 1); % previous value function
17 policy     = zeros(N_state,1);  % policy
18
19 %Create state transition matrix
20 P          = zeros(N_state,N_state,N_act); %(A, Sold, Snew)
21
22 params.N_grid = N_grid;
23 params.N_state = N_state;
24 params.rbf_c = [1 1; 1 5; 2 4; 3 3; 4 2; 5 1; 5 5];
25 params.s     = size(params.rbf_c,1)+1; % RBF standard deviation
26 params.mu    = ones(params.s,1)*1;    % RBF average
27 params.bw    = 1;                     % RBF bias parameter
28
29 for s=1:N_state
30     %states: stay, right, up, left, down
31     new_s(1) = s;                % stay
32     new_s(2) = s+1*(mod(s,N_grid)~=0); % right
33     new_s(3) = s+N_grid*(s+N_grid<=N_state); % up
34     new_s(4) = s-1*(mod(s,N_grid)~=1); % left
35     new_s(5) = s-N_grid*(s-N_grid>=1); % down
36
37     % for all actions
38     for a=1:N_act
39         % compute transition probability for desired
40         P(s,new_s(a),a) = (1-noise);
41         % for all actions
42         for a_rand=1:N_act
43             % compute uncertain transition probability
44             P(s,new_s(a_rand),a) = P(s,new_s(a_rand),a) + noise/(N_act-1);
45         end
46     end
47 end
48
49 % choose basis
50
51
52 % tabular
53 if basis==0

```

```

54 % initialize delta
55 delta=eta;
56 % while delta is not within the tolerance
57 while delta>=eta
58     % initialize delta
59     delta=0;
60     % create a temporary set of value and policy vectors
61     va = zeros(N_state, N_act);
62     v_ = zeros(N_state,1);
63     %for all initial states
64     for s=1:N_state
65         % for all actions
66         for a=1:N_act
67             % for all final states
68             for s_prime=1:N_state
69                 % compute the new value function
70                 va(s,a)=va(s,a)+P(s,s_prime,a)*(reward(s_prime)+gamma*v(s_prime,:));
71             end
72         end
73         % store the new maximum value and corresponding action
74         [v_(s,:), action] = max(va(s,:));
75         policy(s,:) = action;
76         old_v(s,:) = v(s,:);
77         % update the value function and recompute delta
78         v(s,:) = v_(s,:);
79         delta = max(delta,abs(v(s,:)-old_v(s,:)));
80     end
81 end
82
83 elseif basis==1 % RBF
84     % initialize theta, phi, and delta
85     theta = zeros(params.s,1);
86     phi = VI_RBF(params);
87     delta=eta;
88     i=1;
89     while delta>=eta
90         delta=0;
91         for s=1:N_state
92             v=phi(s,:)*theta;
93             %find action that maximizes value
94             va=zeros(1,N_act);
95             for a=1:N_act
96                 va(a)=P(s,:,a)*(reward(1:N_state)'+gamma*phi*theta);
97             end
98             v_=max(va);
99             % update RBF
100             alpha = alpha0/i^alpha_exp;
101             theta_old = theta;
102             theta = theta + alpha*(v_-v)*phi(s,:)';
103             theta=theta/norm(theta);
104             % update delta
105             delta = max(delta,max(theta-theta_old));
106         end
107         i=i+1;

```

```

108     end
109
110     % update value function using RBF
111     v = phi*theta;
112     Va=zeros(N_state,N_act);
113     %compute maximum value
114     for a=1:N_act
115         Va(:,a)=P(:, :,a)*v;
116     end
117     %update policy;
118     [~, policy]=max(Va,[],2);
119 end
120
121 % display output
122 V_new=flipud(reshape(v,[N_grid N_grid]))'/max(v);
123 policy=flipud(reshape(policy,[N_grid N_grid]))');
124 display(V_new)
125 display(policy)

```

MATLAB CODE FOR TBVI

```

1 %% Gridworld Simulation
2 clear all;
3 close all;
4 clc;
5 %% Gridworld Parameters
6 basis          = 1;           % 0 for tabular, 1 for RBF
7 N_grid         = 5;           % size of grid
8 N_state        = N_grid*N_grid; % size of state
9 N_act          = 5;           % number of actions
10 N_eps          = 20;          % max length of trajectory
11 N_time         = 1000;        % length of time
12 N_converge     = 6;           % number of convergences
13 N_exec         = 5;           % number of executions
14 N_dim          = 2;           % state dimension
15
16 eta            = 0.01;        % convergence tol
17 s_init         = [1;1];       % initial state
18 a_init         = 2;           % initial action
19 s_goal         = [N_grid;N_grid]; % goal state
20
21 rew_goal       = 1;           % goal reward
22 rew_trans      = -0.01;       % transition reward
23 noise          = 0.1;         % transition uncertainty
24
25 params.N_grid  = N_grid;      % size of grid
26 params.s_goal  = s_goal;      % goal state
27 params.rew_goal = rew_goal;    % goal reward
28 params.rew_trans = rew_trans; % transition reward
29 params.N_dim   = N_dim;       % state dimension
30 params.N_act   = N_act;       % number of actions
31 params.noise    = noise;      % transition uncertainty

```

```

32
33 %% Learning Parameters
34 gamma          = 0.9;           % discount factor
35 params.gamma    = gamma;        % discount factor
36
37 alpha_init      = 0.5;           % initial learning rate
38 alpha_dec       = 0.5;           % decay rate
39 mu              = 1;
40 eps_init        = 0.8;           % initial exploration rate
41 eps_dec         = 0.1;           % exploration rate decay
42
43 max_points      = 25;           % max centres allowed for RBF
44 tol             = 1e-4;         % tolerance
45
46 if basis==0
47     params.N_s = N_state;        % Number state-features
48     params.N_sa = params.N_s*params.N_act; % Number state-action-features
49     params.state_action_slicing_on=1;
50     gpr = onlineGP_RL(0,0,0,0,params);
51     params.basis=0;
52 elseif basis==1
53     params.c = [5 5; 1 5; 5 1; 1 1]';
54     params.N_s = size(params.c,2)+1; % Number state-features
55     rbf_mu = ones(params.N_s,1)*mu; % RBF mu
56     params.mu=ones(params.N_s,1)*1; % RBF mu
57     params.bw=1;                 % RBF bias
58     params.N_sa = params.N_s*params.N_act; % Number state-action-features
59     params.state_action_slicing_on=1;
60     gpr = onlineGP_RL(0,0,0,0,params);
61     params.basis=1;
62 end
63
64 %% Algorithm Execution
65 %Create state transition matrix
66 P=zeros(N_act,N_state,N_state);
67 for s=1:N_state
68     %states: stay, right, up, left, down
69     new_s(1) = s; % stay
70     new_s(2) = s+1*(mod(s,N_grid)~=0); % right
71     new_s(3) = s+N_grid*(s+N_grid<=N_state); % up
72     new_s(4) = s-1*(mod(s,N_grid)~=1); % left
73     new_s(5) = s-N_grid*(s-N_grid>=1); % down
74
75     % for all actions
76     for a=1:N_act
77         % compute probability of desired transition
78         P(s,new_s(a),a) = (1-noise);
79         % for all actions
80         for a_rand=1:N_act
81             % compute probability of uncertain transition
82             P(s,new_s(a_rand),a) = P(s,new_s(a_rand),a) + noise/N_act;
83         end
84     end
85 end

```



```

86
87 % initialize
88 theta = zeros(params.N_sa,1);
89 ctr = 0;
90 eval_ctr = 0;
91 n_conv=0;
92 i=0;
93
94 % while the time is less than the required time
95 % while the number of convergences is less than that required
96 while i<=N_time && n_conv<N_converge
97     % initialize
98     s_old = s_init;
99     break_cmd = 0;
100     delta = 0;
101     k=1;
102     % while less than the number of episodes
103     % while not commanded to break
104     while k<=N_eps && ~break_cmd
105         ctr = ctr+1;
106
107         % set the exploration probability
108         p_eps = eps_init/(ctr)^eps_dec;
109
110         % check whether to explore
111         r = sample_discrete([p_eps 1-p_eps]);
112
113         % explpre
114         if r==1
115             p = 1/N_act.*ones(1,N_act);
116             action = sample_discrete(p);
117         % exploit
118         else
119             [Q_opt,action] = Q_greedy_act(theta,s_old,params,gpr);
120         end
121
122         % compute next state and reward
123         s_old_lin = sub2ind([N_grid N_grid],s_old(1),s_old(2));
124         s_new=zeros(2,N_exec);
125         rew=zeros(1,N_exec);
126         % for all executions
127         for j=1:N_exec
128             % calculate next state
129             s_new(1,j) = sample_discrete(P(s_old_lin,:,action));
130             [s_new(1,j),s_new(2,j)]=ind2sub([N_grid N_grid],s_new(1,j));
131             % calculate reward
132             [rew(j),break_cmd] = reward2(s_new(:,j),params);
133         end
134         % recompute learning rate
135         alpha = alpha_init/ctr^alpha_dec;
136         % recompute feature vector
137         phi_old = Q_feature(s_old,action,params);
138         % calculate the value
139         val_old = Q_value(theta,s_old,action,params);

```

```

140     v_new = 0;
141     % for all executions
142     for j=1:N_exec
143         % compute optimal action
144         [Q_opt,a_op] = Q_greedy_act(theta,s_new(:,j),params,gpr);
145         % compute new value
146         v_new=v_new+rew(j)+gamma*Q_value(theta,s_new(:,j),a_op,params);
147     end
148     v_new = v_new/N_exec;
149     % compute error
150     err = (v_new - val_old);
151
152     % update basis vector
153     theta_old=theta;
154     theta = theta + alpha*(err.*phi_old);
155     delta = max(delta,abs(max(theta-theta_old)));
156
157     % reset state
158     s_old = s_new(:,1);
159     k=k+1;
160 end
161 % check for break condition
162 if (delta<eta && break_cmd)
163     n_conv = n_conv+1;
164 end
165 i=i+1;
166 end
167
168 % update policy
169 policy=zeros(N_grid,N_grid);
170 for i=1:N_grid
171     for j=1:N_grid
172         [~,policy(i,j)] = Q_greedy_act(theta,[i;j],params,gpr);
173     end
174     policy(N_grid,N_grid)=1;
175 end
176 policy = flipud(policy');
177
178 %% Monte Carlo runs
179 rew_eval = zeros(1,100);
180 for eval_ctr = 1:100
181
182     s_old = s_init;
183     for ctr=1:N_state
184
185         [~,action] = Q_greedy_act(theta,s_old,params,gpr);
186         s_next = gridworld_trans(s_old,action,params);
187         [rew,break_cmd] = reward2(s_next,params);
188         rew_eval(eval_ctr) = rew_eval(eval_ctr) + rew;
189
190         if break_cmd
191             break;
192         end
193

```

```

194         s_old = s_next;
195     end
196 end
197 boxplot(rew_eval);
198 rew_mean=mean(rew_eval);
199 rew_var=var(rew_eval);
200 display(rew_mean);
201 display(rew_var);

```

MATLAB CODE FOR SARSA

```

1 %% Gridworld Simulation
2 clear all;
3 close all;
4 clc;
5 %% Gridworld Parameters
6 basis          = 1;           % 0 for tabular, 1 for RBF
7 N_grid         = 5;           % grid size
8 N_state        = N_grid*N_grid; % state size
9 N_act          = 5;           % number of actions
10
11 N_dim          = 2;           % state dimension
12 s_init         = [1;1];      % initial state
13 a_init         = 2;           % initial action
14 s_goal         = [N_grid;N_grid]; % goal state
15
16 N_obstacle     = 0;           % number of obstacles
17 obs_list       = [];         % coordinates of obstacles
18
19 rew_goal       = 1;           % reward for goal
20 rew_trans      = -0.01;       % reward for transition
21 noise          = 0.1;         % transition uncertainty
22
23 params.N_grid  = N_grid;      % grid size
24 params.s_goal  = s_goal;      % goal state
25 params.rew_goal = rew_goal;   % reward for goal
26 params.rew_trans = rew_trans; % reward for transition
27 params.N_dim   = N_dim;       % state dimension
28 params.N_act   = N_act;       % number of actions
29 params.noise    = noise;      % transition uncertainty
30
31 %% Learning Parameters
32 N_length       = 100;         % length of episode
33 N_eps          = 200;         % number of episodes
34 N_exec         = 3;           % number of executions
35 N_freq         = 100;         % frequency of evaluation
36 N_eval         = 30;          % evaluation iterations
37 N_budget       = 25;          % max points in stack
38
39 gamma          = 0.9;         % discount factor
40 params.gamma   = gamma;       % discount factor
41

```

```

42 params.N_budget      = N_budget;           % max points in stack
43 data_method          = 2;                   % 1 For cyclic and 2 for SVD
44 params.epsilon_data_select=0.2;
45 stack_index          = 0;                   % initial stack index
46 points_in_stack      = 0;                   % initial stack length
47
48 alpha_init           = 0.5;                 % initial learning rate
49 alpha_dec             = 0.5;                 % learning rate decay
50 eps_init              = 0.8;                 % initial exploration rate
51 eps_dec               = 0.1;                 % exploration rate decay
52 N_pts                 = 25;                 % max centres allowed for RBF
53 tol                   = 1e-4;               % tolerance
54
55 if basis==0
56     params.N_s = N_state;                    % state-features
57     params.N_sa = params.N_s*params.N_act; % state-action-features
58     params.basis=basis;
59     params.state_action_slicing_on=1;
60     params.basis=0;
61     gpr = onlineGP_RL(0,0,0,0,params);
62 elseif basis==1
63     params.c = [5 5; 1 5; 5 1; 1 1]';
64     params.N_s = size(params.c,2)+1;         % number of state-features
65     params.mu=ones(params.N_s,1)*1;         % RBF mu
66     params.bw=1;                             % RBF bias
67     params.N_sa = params.N_s*params.N_act; % state-action features
68     params.state_action_slicing_on = 1;
69     params.basis = 1;
70     gpr = onlineGP_RL(0,0,0,0,params);
71 end
72
73 %% Algorithm Execution
74 rew_exec = zeros(N_exec,1);
75 eval_ctr = zeros(N_exec,1);
76
77 % for all executions
78 for i =1:N_exec
79     % reset Q function
80     theta = zeros(params.N_sa,1);
81
82     step_ctr = 1;
83     % for all episodes
84     for j = 1:N_eps
85         %fprintf('Episode: %d/%d, Execution: %d/%d \n',j,N_eps,i,N_exec);
86
87         % reset to initial state
88         s_old = s_init;
89
90         % set exploration probability
91         p_eps = eps_init/(step_ctr)^eps_dec;
92         % check if exploring
93         r = sample_discrete([p_eps 1-p_eps]);
94         % explore
95         if r==1

```

```

96         p = 1/N_act.*ones(1,N_act);
97         action = sample_discrete(p);
98     % exploit
99     else
100         [Q_op,action] = Q_greedy_act(theta,s_old,params,gpr);
101     end
102     % for lenfth of evaluation
103     for k = 1: N_length
104         % check if it is time to evaluate
105         if(mod(step_ctr,N_freq) == 0)
106             % evaluate reward
107             eval_ctr(i) = eval_ctr(i) + 1;
108             rew_eval = zeros(1,N_eval);
109             % for number of evals
110             for eval_count = 1:N_eval
111                 % reset state
112                 s_prv = s_init;
113                 % for legnth of evaluation
114                 for step_count = 1:N_length
115                     % calculate optimal action
116                     [Q_op,action]=Q_greedy_act(theta,s_prv,params,gpr);
117                     s_next = gridworld_trans(s_prv,action,params);
118                     % calculate reward
119                     [rew, break_cmd] = reward2(s_next,params);
120                     rew_eval(eval_count) = rew_eval(eval_count)+rew;
121                     % check break condition
122                     if break_cmd
123                         break;
124                     end
125                     % update state
126                     s_prv = s_next;
127                 end
128             end
129             % update reward
130             rew_exec(i,eval_ctr(i)) = mean(rew_eval);
131         end
132         step_ctr = step_ctr + 1;
133
134         % get nextsState
135         s_new = gridworld_trans(s_old,action,params);
136
137         % calculate reward
138         [rew,break_cmd] = reward2(s_new,params);
139
140         % set exploration rate
141         p_eps = eps_init/(step_ctr)^eps_dec;
142         % check if exploring
143         r = sample_discrete([p_eps 1-p_eps]);
144         % explore
145         if r==1
146             p = 1/N_act.*ones(1,N_act);
147             a_new = sample_discrete(p);
148         % exploit
149         else

```

```

150         % calculation optimal action
151         [Q_op,a_new] = Q_greedy_act(theta,s_old,params,gpr);
152     end
153
154     % calculate learning rate
155     alpha = alpha_init/(step_ctr)^alpha_dec;
156     % calculate feature vector
157     phi_old = Q_feature(s_old,action,params);
158     % calculate value
159     v_old = Q_value(theta,s_old,action,params);
160     % update value
161     v_new = Q_value(theta,s_new,a_new,params);
162     % compute error
163     err = (rew + gamma*v_new - v_old);
164     % update RBF parameter
165     theta = theta + alpha*(err.*phi_old);
166
167     % reset state
168     s_old = s_new;
169     action = a_new;
170
171     % check break condition
172     if break_cmd
173         break;
174     end
175 end
176 end
177 end
178
179 %% Post Process
180 % find minimum number of evaluations
181 min_eval = min(eval_ctr);
182 rew_exec = rew_exec(:,1:min_eval);
183 rew_total = zeros(1,min_eval);
184 std_total = zeros(1,min_eval);
185
186 % update reward
187 for m =1:min_eval
188     rew_total(m) = mean(rew_exec(:,m));
189     std = var(rew_exec(:,m));
190     std_total(m) = 0.1*std;
191 end
192
193 %% Plots
194 t = 1:min_eval;
195 t = t.*N_freq;
196 errorbar(t,rew_total,std_total);
197 xlabel('episodes')
198 ylabel('reward')

```
