

Bits, Brain & Behaviour - Coursework 2

Question 1: Monte-Carlo Control in Easy21

(a) For Monte-Carlo (MC) Control, on-policy ϵ -greedy first-visit (FV) MC Iterative Learning Algorithm is used in MATLAB. FV MC is used instead of every-visit (EV) MC because, according to Singh and Sutton (1996) *Reinforcement Learning with Replacing Eligibility Traces*, for large numbers of trials the FV estimate has a lower mean-squared error (MSE) than EV estimate (Lower variance and also unbiased estimator). This method also uses running means instead of batch averaging, without the need to store lists of returns for each state-action pair (s,a) thus requires less memory. Algorithm also avoids the unlikely assumption of exploring starts, since starting an episode with player sum > 10 is not realistic. Algorithm used obtains the Return (R) from the first appearance of a state-action pair (s,a) in the generated trace, which is basically final reward (r_{final}) since $\gamma = 1$, and updates the state-action value function estimate $\hat{Q}(s, a)$ using Equation 1 (online-averaging).

$$\hat{Q}(s, a) \leftarrow \hat{Q}(s, a) + \alpha[r_{final} - \hat{Q}(s, a)] , \text{ where } \alpha = 1/N(s, a) \quad (1)$$

The step-size α is the reciprocal of cumulative number of visits to (s,a), $N(s,a)$. Since process is non-stationary, it is useful to gradually decrease the sensitivity to recent return values. Step size chosen also satisfies the Robbins-Monroe condition. The ϵ -greedy policy used (Equation 2) satisfies GLIE criterion, as it asymptotically converges to the deterministic optimal policy, ensuring infinite exploration. The constant N_0 , empirically chosen to be 100 (Fig.1), determines the ratio of exploration and exploitation done by the algorithm (See Part b). A high value of N_0 allows high exploration (low exploitation) since when ϵ is closer to 1 each action becomes equiprobable, but algorithm takes longer to converge. Algorithm was tested on 1,000,000 episodes due to computational power limitations. Also, if $\hat{Q}(s, a_1) = \hat{Q}(s, a_2)$, then $\pi(s, a) = 0.5$ for both.

$$\pi(s, a) = \begin{cases} 1 - \frac{\epsilon}{2} & , \text{ if } a^* = \arg \max_a Q(s, a) \\ \frac{\epsilon}{2} & , \text{ if } a \neq a^* \end{cases} , \text{ where } \epsilon = \frac{N_0}{N_0 + N(s)} . \quad (2)$$

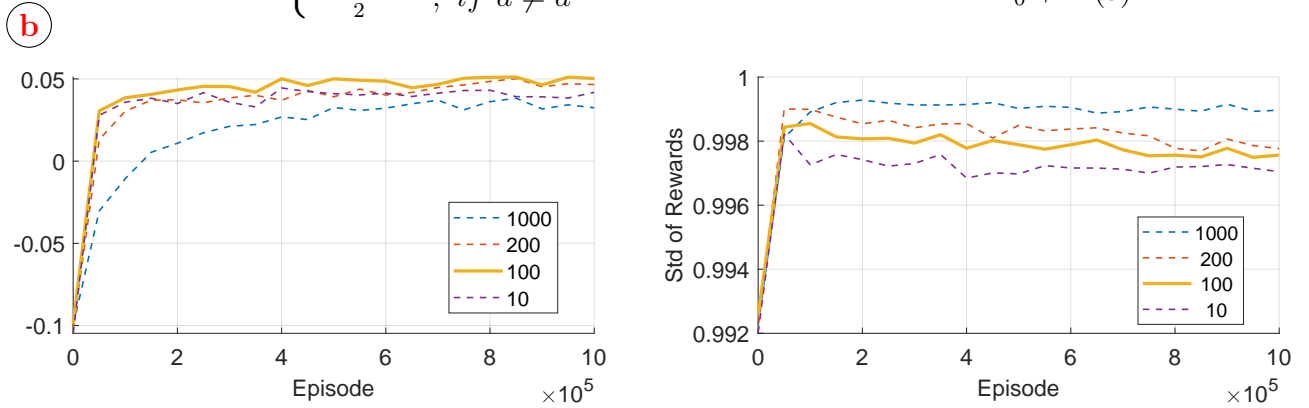


Figure 1: Learning Curve for MC Control showing mean and standard deviation of rewards vs episodes. Plot shows the learning curves for different values of N_0 (Solid line for chosen N_0 value).

Learning Curves were obtained by using validation sets of simulations, since plots of rewards vs episodes contain no useful information. Ideally after each training step, training temporarily stops, current policy is applied to several numbers of episodes of the game (validation set) and performance of system at this stage is estimated. Due to computational power limitations, and using 1,000,000 episodes, the validation step was performed after every 20,000 episodes for a validation set size of 100,000 episodes. Size of validation set was chosen such that that any variation in mean rewards would be smaller than changes in mean value itself, providing clearer trends. As expected, this approach makes the training longer, however it is a clean and methodologically appropriate way of obtaining such learning curves. Also does not have the issue of determining appropriate window size in methods involving sliding windows. Fig.1 shows that for $N_0 = 100$ the value of mean rewards converges to highest value and std of

rewards is low after the end of process. This is because lower N_0 values converge faster with low exploration, thus lower mean values, and higher N_0 values do not fully converge.

(c) Fig.2 shows that for player sum values close to 21, the value function is maximum. It also shows that the expected returns of high dealer values are much lower than those for lower dealer values, regardless of player's sum. Local max at player sum of 11 is due to the fact that player cannot go bust at this value. Moreover, it should be noted that plot of value function for MC Control is considerably noisy.

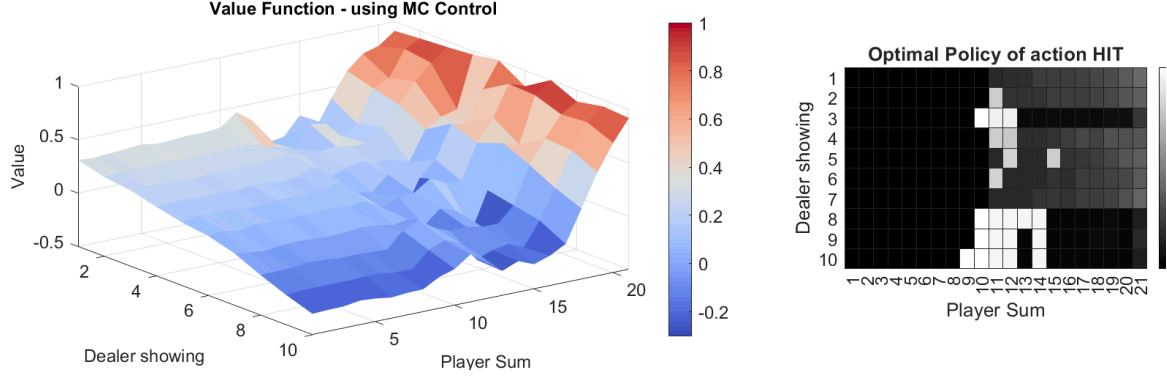


Figure 2: (Left) Optimal value function $V^*(s) = \max_a Q^*(s, a)$ for MC Control. (Right) Optimal policy for the action "hit", using MC Control after 1,000,000 episodes ($N_0 = 100$).

Question 2: TD Learning (SARSA) in Easy21

(a) For this question, on-policy learning TD control (SARSA) algorithm was implemented in MATLAB for 1,000,000 episodes. This method is based on bootstrapping instead of sampling of complete episodic traces like MC Control. Specifically the implemented algorithm performs bootstrap only one time-step ahead. In this method, the update rule for the state-action value function estimate $\hat{Q}(s, a)$ is given by Equation 3. The update rule depends on the current state (S), the action chosen from S (A), the reward (R) from choosing A , the next state S' and the next action (A') chosen in the new state.

$$\hat{Q}(S, A) \leftarrow \hat{Q}(S, A) + \alpha [R + \hat{Q}(S', A') - \hat{Q}(S, A)] , \text{ where } \alpha = 1/N(S, A) \quad (3)$$

Similarly as for MC Control, the step-size α is gradually decreased with visits to (S, A) , thus satisfying the Robbins-Monroe condition. Algorithm also uses an ϵ -greedy policy derived from \hat{Q} values, for choosing the next action, given by Equation 2, satisfying the GLIE criterion. Both ϵ and α thus ensure convergence to optimal solutions. By performing similar analysis for the most appropriate value of N_0 to be used, as was done in Q1 part b), this value was found to be 100. This value causes mean reward to converge to highest value, with relatively low standard deviation. Graphs are shown in Appendix A, fig.8.

(b) The learning curve showing mean and standard deviation of rewards was obtained using validation sets of simulations, which is described in Q.1 b). Again, 1,000,000 episodes were used and the validation step was performed after every 20,000 episodes for a validation set size of 100,000 episodes. To show the effect of step-size α on the curves, the learning curves were plotted (Fig.3) using the time-varying scalar step-size $\alpha_t = 1/(N(s_t, a_t))$ described in part a), and a range of constant values of α . Figure shows that the mean value of rewards converges to the highest value when step-size is time-varying. All constant α 's do not ensure convergence to the optimal solution. This is because a constant value of α close to 1, means high learning rate but in the expense of convergence to wrong optimal solutions (instability). Graphs of large α 's also have high-amplitude ripples. For constant α close to 0, the rate of convergence is much slower but with smaller ripples, so these do not converge in the 1,000,000 episodes used.

(c) Value function plot (Fig.4) was obtained using the time-varying scalar step-size, and $N_0 = 100$. Plot shows the same basic features as that for MC Control, however, in this case the plot of value function is less noisy (smoother).

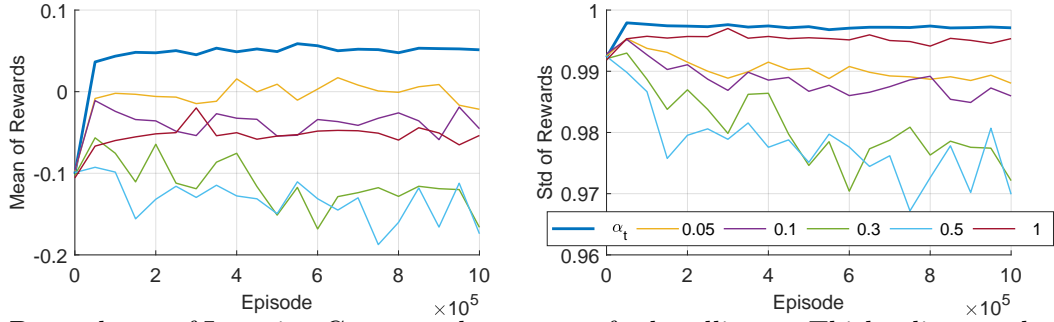


Figure 3: Dependency of Learning Curve on the strategy for handling α . Thicker line on plots belongs to the time-varying step-size.

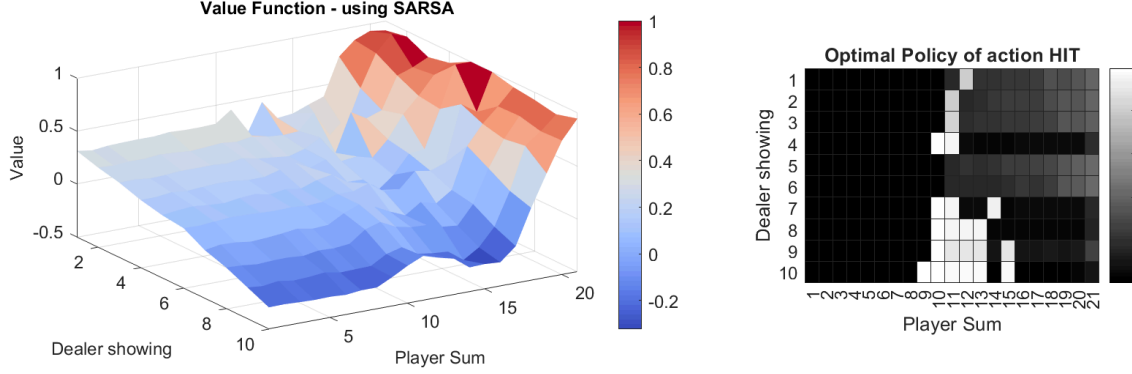


Figure 4: (Left) Optimal value function $V^*(s) = \max_a Q^*(s, a)$ for SARSA. (Right) Optimal policy for the action "hit", using SARSA after 1,000,000 episodes ($N_0 = 100$).

Question 3: Q Learning in Easy21

(a) Q Learning algorithm implemented in MATLAB is an off-policy TD Control algorithm, where the next action (a') is chosen using the behaviour policy, while updating action value estimate $\hat{Q}(s, a)$ in direction of the alternative (better) action dictated by the target policy (deterministic). This is the difference to SARSA. Update rule is shown in equation 4.

$$\hat{Q}(S, A) \leftarrow \hat{Q}(S, A) + \alpha [R + \max_{a'} \hat{Q}(s', a') - \hat{Q}(S, A)] , \text{ where } \alpha = 1/N(S, A) \quad (4)$$

Target Policy is greedy with respect to $\hat{Q}(s, a)$, while behaviour policy is ϵ -greedy with respect to $\hat{Q}(s, a)$. The step-size α used is again gradually decreased with visits to (S,A), satisfying the Robbins-Monroe condition, and the ϵ -greedy behaviour policy is given by Equation 2, satisfying the GLIE criterion, ensuring convergence. By performing similar analysis for the most appropriate value of N_0 to be used, as was done in Q1 part b), this value was found again to be 100. This value causes mean reward to converge to highest value, with relatively low standard deviation. Graphs are shown in Appendix A, fig.9.

(b) The learning curve showing mean and standard deviation of rewards was obtained using validation sets of simulations, which is described in Q.1 b). Again, 1,000,000 episodes were used and the validation step was performed after every 20,000 episodes for a validation set size of 100,000 episodes. To show the effect of ϵ -greediness on the curves, the learning curves were plotted (Fig.5) using the time-varying ϵ_t described in equation 2, and a range of constant values of ϵ . Generally, the higher the value of ϵ , the more the exploration done and the less the exploitation. Constant ϵ does not satisfy the GLIE criterion. The time-varying ϵ_t allows high exploration initially and gradually increases exploitation of agent. After ∞ episodes, exploration will reduce to zero. Fig.5 shows that as ϵ gets extremely close to 1 or 0, the mean value does not converge to optimal value. However, small values of ϵ converge to the optimal value in addition to ϵ_t , but only ϵ_t has the lowest variance. This is because $N(s)$ in ϵ_t equation has not yet increased to infinity and is simply a small value that is gradually decreasing.

(c) Value function plot (Fig.6) was obtained using the time-varying ϵ and $N_0=100$. Plot shows the same basic features as that for MC Control and SARSA, however, in this case the plot shows only smooth changes in the value functions of neighbouring states (smoothest plot).

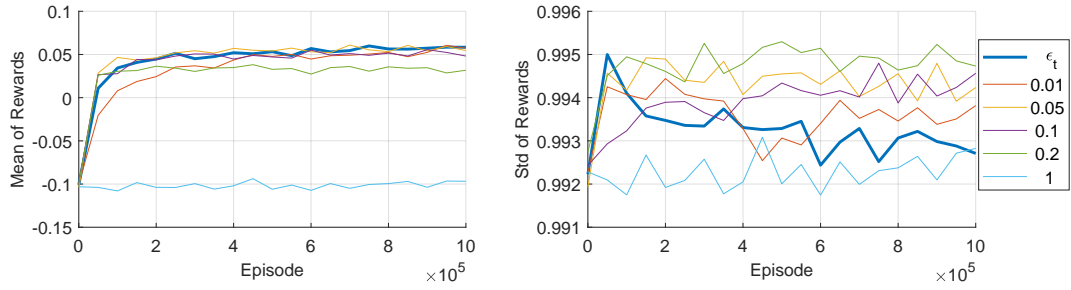


Figure 5: Dependency of Learning Curve on the strategy for handling ϵ -greediness. Thicker line on plots belongs to the time-varying ϵ_t .

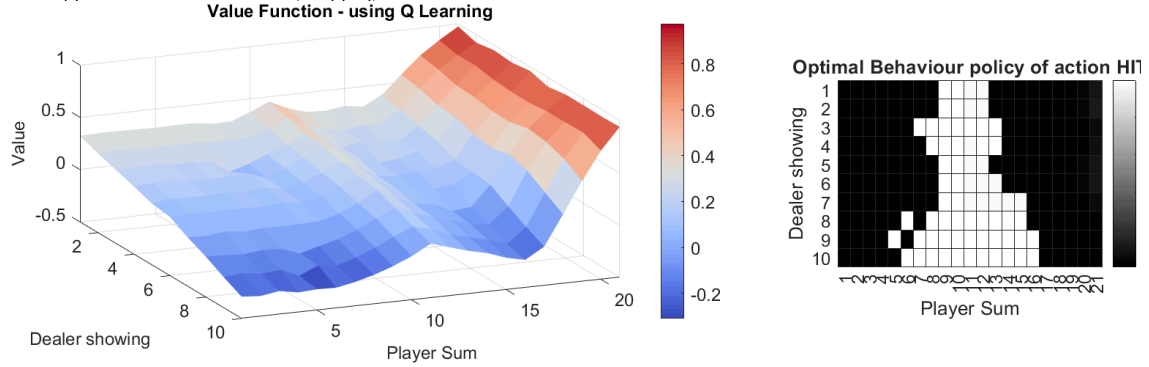


Figure 6: (Left) Optimal value function $V^*(s) = \max_a Q^*(s, a)$ for Q-learning. (Right) Optimal Behavioural policy for the action "hit", using Q-learning after 1,000,000 episodes ($N_0 = 100$).

Question 4: Compare the algorithms

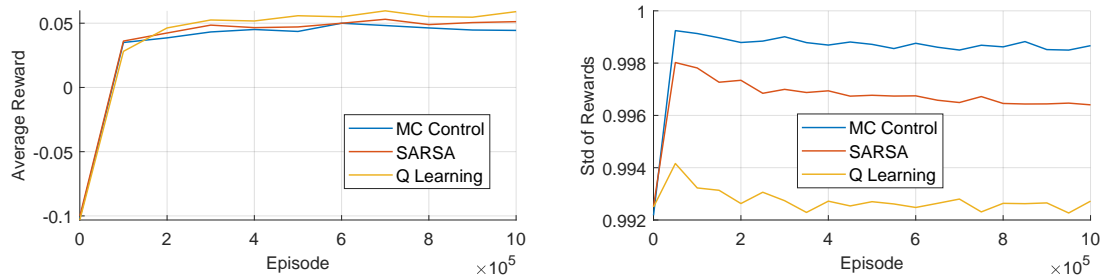


Figure 7: Learning curves of the three implementations. For mean plot, resolution is increased to identify differences. (validation steps=20,000 episodes, validation set size= 100,000 episodes)

The curves of mean rewards and standard deviation (std) of rewards (Fig.7) are equivalent to the changes that occur in the expected value and variance of value function estimate, respectively. Fig.7 shows that Q-Learning and SARSA (TD methods) outperform MC Control as they converge to higher mean reward value, while having a lower std of rewards. This can be explained by the Bias-Variance tradeoff of TD and MC methods. Both methods follow idea of generalised policy iteration (GPI). MC Control gives an unbiased estimate of $Q(s)$ during policy evaluation, since update of $\hat{Q}(s)$ (Eq.1) is based on sample returns, but introduces some bias during policy improvement. TD Control methods give biased estimates in both steps of GPI, since updates of $\hat{Q}(s)$ (Eq.3 and 4) are based on other $\hat{Q}(s)$ estimates and are sensitive to initial estimate values. MC Control estimates have high variance since sample returns are the sum of many random variables (i.e. rewards of each step), while TD estimates' update rule contains only a fixed number of three random variables, yielding proportionally less variance than MC estimates. The better performance of SARSA and Q-learning compared to MC Control, is due to the fact that TD methods exploit the Markov Property of process, have lower convergence time and probably have a lower MSE ($MSE = Bias^2 + Variance$), i.e. a better balance of bias and variance of estimates. Both curves also show that Q-Learning outperforms SARSA in mean and std. The only difference between update rules of the two methods, is the $\max()$ function in Q-Learning, which always considers the best possible outcome by directly learning the optimal policy, while choosing actions from the exploring ϵ -greedy behavior policy. SARSA, takes a more conservative approach avoiding optimal paths with high risk of large cost. Since we have a low-cost environment, Q-learning yields higher average rewards, with low std.

Appendix A - Learning curves of SARSA and Q-Learning for different N_0 values.

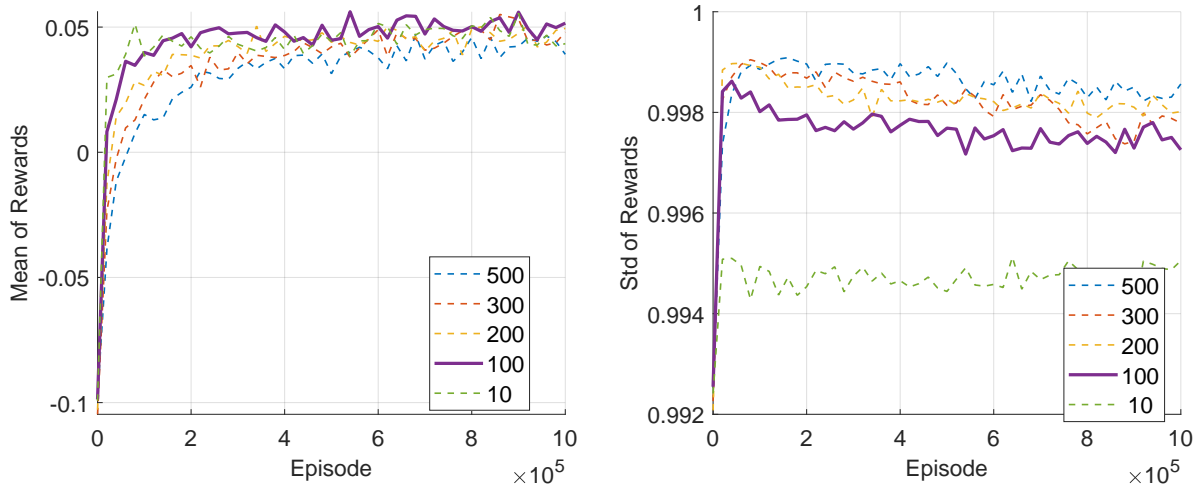


Figure 8: Learning curves of SARSA using different values of N_0 . Number of episodes=1,000,000. Thick Line represents the chosen N_0 value.

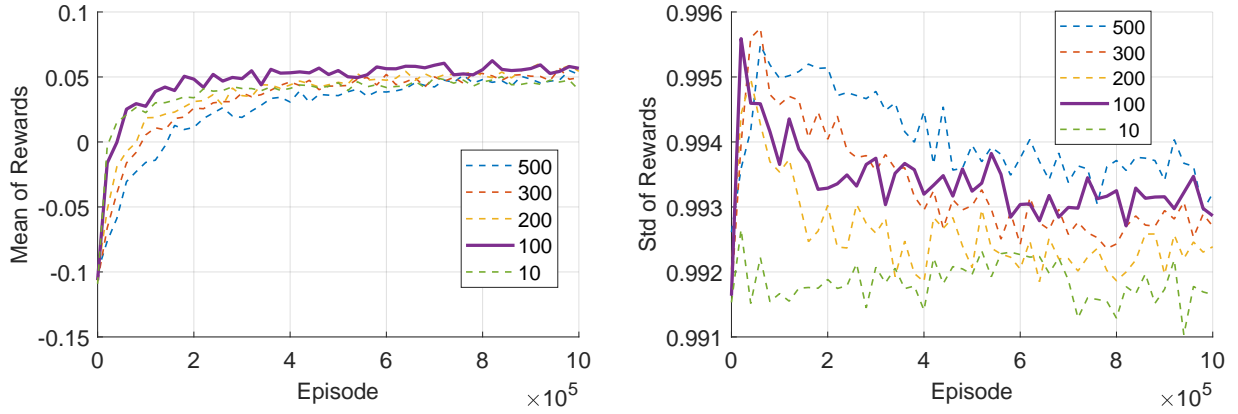


Figure 9: Learning curves of Q-learning using different values of N_0 . Number of episodes=1,000,000. Thick Line represents the chosen N_0 value.

Appendix B - Matlab code used for implementation of Easy21

```
1 classdef Easy21_v3 < handle
2     properties
3         max_length;
4         player_1stCard_val;
5         dealer_1stCard_val;
6         player_sum;
7         dealer_sum;
8
9         state;
10        player_goes_bust;
11        dealer_goes_bust;
12
13        player_card_val;
14        player_card_col;
15
16        dealer_card_val;
17        dealer_card_col;
18
19        ret;
20        terminal;
21        t;
22    end
23
24    methods
25        function init(obj,max_len)
26            obj.max_length = max_len;
27        end
28
29        function state = reset(obj)
30            obj.player_1stCard_val = randi([1 10]);
31            obj.dealer_1stCard_val = randi([1 10]);
32
33            obj.player_sum = obj.player_1stCard_val;
34            obj.dealer_sum = obj.dealer_1stCard_val;
35
36            obj.state = [obj.dealer_1stCard_val,obj.player_1stCard_val];
37
38            obj.player_goes_bust=false;
39            obj.dealer_goes_bust=false;
40
41            obj.ret=0;
42            obj.terminal = false;
43            obj.t=0;
44
45            state = obj.state;
46        end
47
48        function [state,Terminal,reward,Dealer_sum,Player_sum]=step(obj, ...
49            action)
50            %Action: 1 = hit , 0=stick
51            %colour: 1=black , -1=red
52
53            r=0; %Reward
54
55            previous_player_sum = obj.player_sum;
56
57            if action==1 %When player hits
58                obj.player_card_val = randi([1 10]);
```

```

59         x=rand; %Choosing random colour
60         if x<(1/3)
61             colour=-1;
62         else
63             colour=1;
64         end
65         obj.player_card.col=colour;
66
67         obj.player_sum = obj.player_sum + (obj.player_card.val * ...
68             obj.player_card.col);
69         obj.player_goes_bust = (obj.player_sum>21 || ...
70             obj.player_sum<1);
71
72         if obj.player_goes_bust == 1
73             r = -1;
74             obj.terminal = true;
75         end
76     end
77
78     if (action==0)&&(obj.terminal==0)%when player sticks
79         while (obj.dealer_sum>0)&&(obj.dealer_sum<17)
80             obj.dealer_card.val = randi([1 10]);
81             z=rand; %Choosing random colour
82             if z<(1/3)
83                 col=-1;
84             else
85                 col=1;
86             end
87             obj.dealer_card.col=col;
88
89             obj.dealer_sum = obj.dealer_sum + ...
90                 (obj.dealer_card.val * obj.dealer_card.col);
91             obj.dealer_goes_bust = (obj.dealer_sum>21 || ...
92                 obj.dealer_sum<1);
93
94             if obj.dealer_goes_bust == 1
95                 r = 1;
96                 obj.terminal = true;
97             end
98         end
99
100         obj.terminal = true; %game ends when dealer sticks
101         if ~obj.dealer_goes_bust
102             if obj.player_sum>obj.dealer_sum
103                 r=1;
104             elseif obj.player_sum<obj.dealer_sum
105                 r=-1;
106             elseif obj.player_sum==obj.dealer_sum
107                 r=0;
108             end
109         end
110     end
111
112     obj.t = obj.t + 1;
113     obj.ret = obj.ret + r;
114
115     if ((obj.terminal==0) && (obj.t==obj.max_length))
116         obj.terminal=true;
117         if obj.player_sum>obj.dealer_sum

```



```

116         r=1;
117         elseif obj.player_sum<obj.dealer_sum
118             r=-1;
119         elseif obj.player_sum==obj.dealer_sum
120             r=0;
121         end
122     end
123
124     if (obj.terminal==1) %Game has finished
125         Terminal = obj.terminal;
126         reward=r;
127         obj.state(2) = previous_player_sum;
128         state = obj.state;
129     end
130
131     if (obj.terminal==0)
132         obj.state(2) = obj.player_sum;
133         state = obj.state;
134         reward=r;
135         Terminal = obj.terminal;
136     end
137
138     Dealer_sum = obj.dealer_sum;
139     Player_sum = obj.player_sum;
140 end
141
142 end
143 end

```

```

1 clc
2 close all
3 clear all
4
5 %% Defining the Game
6
7 %Calling the object into this file
8 import Easy21.v3
9 Game = Easy21.v3;
10
11 max_length = 1000;
12 init(Game,max_length);
13 state(1,:) = reset(Game);
14 Terminal=0;
15
16 index=1;
17
18 while Terminal==0
19     x=rand;
20     if x<0.75
21         action=1;
22     else
23         action=0;
24     end
25     actions(index)=action;
26     [next_state,Terminal,reward,s_D,s_P]=step(Game, action);
27     Rew(index)=reward;
28     index=index+1;
29     if ~Terminal
30         state(index,:) = next_state;
31     end
32 end

```


Appendix C - Matlab code used for MC Control of Easy21

```
1 function [policy,action] = MC_get_action(Q,policy,state,No,count_state)
2
3     [ ~ , action_greedy] = max(Q(state(1),state(2),:));
4     if (Q(state(1),state(2),1)==Q(state(1),state(2),2))
5         action = randi([1 2])-1;
6         policy(state(1),state(2),1)=0.5;
7         policy(state(1),state(2),2)=0.5;
8     else
9         action_greedy=action_greedy-1;
10        epsilon = No/(No+(count_state(state(1),state(2))));
11        %Choosing the action
12        p=rand;
13        if p<(1-epsilon/2) %This was deduced from initial form since ...
14            number of possible actions =2
15            action=action_greedy;
16        else
17            action=1-action_greedy;
18        end
19
20        if action_greedy==0
21            policy(state(1),state(2),1)=(1-epsilon/2);
22            policy(state(1),state(2),2)=epsilon/2;
23        elseif action_greedy==1
24            policy(state(1),state(2),2)=(1-epsilon/2);
25            policy(state(1),state(2),1)=epsilon/2;
26        end
27    end
```

```
1 function [Q,Ret_episode,upd_policy,Ret,count_s,count_s_a]= ...
2     MC_FV_episode(obj,No,init_s,Q,Ret,policy,count_s,count_s_a)
3
4     states=[];
5     actions =[];
6     rewards=[];
7     history=[];
8     s=init_s;
9     states(1,:)=s;
10    index=1;
11
12    %Generating the episode
13    while true
14        %This is basically N(s)
15        count_s(states(index,1),states(index,2)) = ...
16            count_s(states(index,1),states(index,2)) + 1; %This increments ...
17            the visits to state s(1) and s(2)
18        [upd_policy,action] = ...
19            MC_get_action(Q,policy,states(index,:),No,count_s);
20        state_action = [states(index,:),action];
21
22        %The following are used to check if it is the first visit
23        %of state-action pair (s,a)
24        if index==1
25            M=0;
26        else
27            M=ismember(state_action,history,'rows');
28        end
29    end
```

```

27         if ~M
28             history=[history;states(index,:),action]; %Stores first ...
                occurence of (s,a)
29             count_s_a(states(index,1),states(index,2),action+1) = ...
30                 count_s_a(states(index,1),states(index,2),action+1) + 1;
31         end
32
33         actions=[actions;action];
34         state_action = [states(index,:),action];
35         [next_state,Terminal,reward,~,~]=step(obj, action);
36
37         rewards=[rewards;reward];
38
39         if Terminal
40             break
41         else
42             states=[states;next_state];
43         end
44         index=index+1;
45     end
46
47     Return = sum(rewards);
48     Ret_episode = Return;
49
50     %The rest has to be done for the history of states
51
52     for k=1:size(history,1)
53
54         alpha=(1/count_s_a(states(k,1),history(k,2),history(k,3)+1));
55         Ret(history(k,1),history(k,2),history(k,3)+1)=Return;
56
57         Q(history(k,1),history(k,2),history(k,3)+1) = ...
58             Q(history(k,1),history(k,2),history(k,3)+1)+ alpha*...
59             (Ret(history(k,1),history(k,2),history(k,3)+1) - ...
60                 Q(history(k,1),history(k,2),history(k,3)+1));
61     end
end

```

```

1  clc
2  close all
3  clear all
4
5  % Defining the Game
6
7  %Calling the object into this file
8  import Easy21.v3
9  Game = Easy21.v3;
10
11  No_vec=[1000 200 100 10]; %Testing the performance of MC Control with ...
        different No
12  %No_vec = 100; %Using the desired No value
13  for i=1:length(No_vec)
14
15      No=No_vec(i);
16
17      %Initialisation for MC
18
19      Q_MC = zeros(10,21,2); %Q(s,a)
20      Ret = zeros(10,21,2);
21      count_s_a = zeros(10,21,2); %N(s,a)
22      count_s = zeros(10,21); %N(s)

```

```

23
24     max_length = 1000;
25     n_episodes=1000000;
26     run_mean=0;
27     pol=0.5*ones(10,21,2); %Initial policy
28     Rewards=zeros(1,n_episodes);
29
30     index_test=50000;
31     size_test_set=150000;
32     Mean_vect=[];
33     Std_vect=[];
34     xaxis=[];
35
36     %Generating episodes and calculating Q function
37     for episode=1:n_episodes
38         if mod((episode/n_episodes),0.05)==0
39             episode/n_episodes
40         end
41         init(Game,max_length);
42         init_s = reset(Game);
43         [Q_MC,Return_episode,pol,Ret,count_s,count_s_a]=...
44             MC_FV_episode(Game,No,init_s,Q_MC,Ret,pol,count_s,count_s_a);
45         Rewards(episode)=Return_episode;
46
47         %This is when validation testing is done
48         if mod(episode,index_test)==0 || episode==1
49             [mean_test,std_test]=test(pol,size_test_set);
50             Mean_vect = [Mean_vect;mean_test];
51             Std_vect = [Std_vect;std_test];
52             xaxis=[xaxis;episode];
53         end
54     end
55
56     figure(1)
57     subplot(1,2,1)
58     hold on
59     grid on
60     plot(xaxis,Mean_vect)
61     ylabel('Mean of Rewards')
62     xlabel('Episode')
63
64     subplot(1,2,2)
65     hold on
66     grid on
67     plot(xaxis,Std_vect)
68     ylabel('Std of Rewards')
69     xlabel('Episode')
70 end
71
72 legend(num2str([1000; 200; 100 ;10]))
73
74 %%
75 h=gcf;
76 h.PaperPositionMode='auto';
77 set(h,'PaperOrientation','landscape');
78 set(findall(gcf,'-property','FontSize'),'FontSize',14)
79 print(gcf, '-dpdf', 'Learning Curve .pdf','-fillpage')
80
81 %%
82 %Plot of policy of hitting
83 figure(2);

```

```

84 subplot(1,2,2)
85 h = heatmap(pol(:, :, 2));
86 xlabel('Player Sum')
87 xlim([1 21])
88 ylabel('Dealer showing')
89 ylim([1 10])
90 title('Optimal Policy of action HIT')
91 colormap(gca, 'gray')
92
93 %Histogram of Rewards
94 figure;
95 histogram(Rewards)
96
97 %% Storing the learning curves mean
98 avg.Rewards_MC=Mean_vect;
99 save('avg.Rewards_MC.mat', 'avg.Rewards_MC')
100
101 %% Storing the learning curves std
102 std.Rewards_MC=Std_vect;
103 save('std.Rewards_MC.mat', 'std.Rewards_MC')
104 %% Optimal value plot
105
106 [V_MC, index] = max(Q_MC, [], 3);
107 index=index-1;
108
109 figure(2);
110 subplot(1,2,1)
111 surf(V_MC, 'EdgeColor', 'none')
112 ylabel('Player Sum')
113 ylim([1 21])
114 xlabel('Dealer showing')
115 xlim([1 10])
116 zlabel('Value')
117 title('Value Function - using MC Control')
118 colormap(gca, coolwarm)
119 hcb1=colorbar;
120
121 %%
122 h=gcf;
123 h.PaperPositionMode='auto';
124 set(h, 'PaperOrientation', 'landscape');
125 set(findall(gcf, '-property', 'FontSize'), 'FontSize', 12)
126 print(gcf, '-dpdf', 'Value Function.pdf', '-fillpage')

```

Appendix D - Matlab codes for SARSA Learning of Easy21

```
1 function [policy,action] = SARSA_get_action(Q,policy,state,No,count_state)
2
3     [ ~ , action_greedy] = max(Q(state(1),state(2),:));
4     if (Q(state(1),state(2),1)==Q(state(1),state(2),2))
5         action = randi([1 2])-1;
6         policy(state(1),state(2),1)=0.5;
7         policy(state(1),state(2),2)=0.5;
8     else
9         action_greedy=action_greedy-1;
10        epsilon = No/(No+(count_state(state(1),state(2))));
11
12        %Choosing the action
13        p=rand;
14        if p<(1-epsilon/2) %This was deduced from initial form since ...
15            number of possible actions =2
16            action=action_greedy;
17        else
18            action=1-action_greedy;
19        end
20
21        if action_greedy==0
22            policy(state(1),state(2),1)=(1-epsilon/2);
23            policy(state(1),state(2),2)=epsilon/2;
24        elseif action_greedy==1
25            policy(state(1),state(2),2)=(1-epsilon/2);
26            policy(state(1),state(2),1)=epsilon/2;
27        end
28    end
end
```

```
1 function [Q,Return_episode,upd_pol,count_s,count_s_a]=...
2     SARSA_episode(obj,No,initial_s,Q,pol,count_s,count_s_a,alpha_choice)
3
4     s=initial_s;
5     Terminal=false;
6     [upd_pol,action] = SARSA_get_action(Q,pol,s,No,count_s);
7
8     %Generating the episode
9
10    while ~Terminal
11
12        %This is basically N(s) - used for step size (alpha)
13        count_s(s(1),s(2)) = count_s(s(1),s(2)) + 1; %This increments the ...
14        visits to state s(1) and s(2)
15        count_s_a(s(1),s(2),action+1) = count_s_a(s(1),s(2),action+1) + 1;
16
17        [next_state,Terminal,reward,~,~]=step(obj, action);
18
19        %Obtain next action - ONLY IF NOT terminal
20
21        if ~Terminal
22            [upd_pol,next_action] = ...
23                SARSA_get_action(Q,upd_pol,next_state,No,count_s);
24            d=Q(next_state(1),next_state(2),next_action+1);
25        else
26            d=0;
27        end
28    end
```

```

27
28     if alpha_choice==0 %Then we use time-varying step-size
29         alpha=(1/count_s_a(s(1),s(2),action+1));
30     else
31         alpha=alpha_choice;
32     end
33
34     %Updating Q-value
35     Q(s(1),s(2),action+1) = Q(s(1),s(2),action+1) + alpha*(reward + d ...
36         - Q(s(1),s(2),action+1));
37     if ~Terminal
38         s=next_state;
39         action = next_action;
40     else
41         Return_episode = reward;
42     end
43 end

```

```

1  clc
2  close all
3  clear all
4
5  %Defining the Game
6
7  %Calling the object into this file
8  import Easy21.v3
9  Game = Easy21.v3;
10
11 %No_vec=[500 300 200 100 10]; %Testing the performance of MC Control with ...
12     different No
13 No_vec = 100; %Using the desired No value
14 %alpha = [0 0.01 0.05 0.1 0.3 0.5 1];
15 alpha=0;
16 for i=1:length(alpha)
17
18     No=No_vec;
19     %Initialisation for SARSA
20
21     Q_Sar = zeros(10,21,2); %Q(s,a)
22     count_s_a = zeros(10,21,2); %N(s,a)
23     count_s = zeros(10,21); %N(s)
24
25     max_length = 1000;
26     n_episodes=1000000;
27     pol=0.5*ones(10,21,2);
28     Rewards=zeros(1,n_episodes);
29
30     index_test=50000;
31     size_test_set=100000;
32     Mean_vect=[];
33     Std_vect=[];
34     xaxis=[];
35
36     for episode=1:n_episodes
37         if mod((episode/n_episodes),0.05)==0
38             episode/n_episodes
39         end
40         init(Game,max_length);
41         init_s = reset(Game);

```

```

42     [Q_Sar,Return_episode,pol,count_s,count_s_action]=...
43     SARSA_episode(Game,No,init_s,Q_Sar,pol,count_s,count_s_a,alpha(i));
44     Rewards(episode)=Return_episode;
45
46     %This is when validation testing is done
47     if mod(episode,index_test)==0 || episode==1
48         [mean_test,std_test]=test(pol,size_test_set);
49         Mean_vect = [Mean_vect;mean_test];
50         Std_vect = [Std_vect;std_test];
51         xaxis=[xaxis;episode];
52     end
53 end
54
55 figure(1)
56 subplot(1,2,1)
57 hold on
58 grid on
59 plot(xaxis,Mean_vect)
60 ylabel('Mean of Rewards')
61 xlabel('Episode')
62
63 subplot(1,2,2)
64 hold on
65 grid on
66 plot(xaxis,Std_vect)
67 ylabel('Std of Rewards')
68 xlabel('Episode')
69 end
70 %legend(num2str([0;0.01;0.05;0.1;0.3;0.5;1]))
71
72 %%
73 h=gcf;
74 h.PaperPositionMode='auto';
75 set(h,'PaperOrientation','landscape');
76 set(findall(gcf,'-property','FontSize'),'FontSize',14)
77 print(gcf, '-dpdf', 'Learning Curve .pdf','-fillpage')
78
79 %%
80 %Plot of policy of hitting
81 figure(2);
82 subplot(1,2,2)
83 h = heatmap(pol(:,:,2));
84 xlabel('Player Sum')
85 xlim([1 21])
86 ylabel('Dealer showing')
87 ylim([1 10])
88 title('Optimal Policy of action HIT')
89 colormap(gca,'gray')
90
91 %Histogram of Rewards
92 figure;
93 histogram(Rewards)
94
95 %% Storing Learning Curves mean
96 avg_Rewards_SARSA=Mean_vect;
97 save('avg_Rewards_SARSA.mat','avg_Rewards_SARSA')
98 %% Storing Learning Curves std
99 std_Rewards_SARSA=Std_vect;
100 save('std_Rewards_SARSA.mat','std_Rewards_SARSA')
101
102 %% Optimal value plot

```



```

103
104 [V_Sarsa, index] = max(Q_Sar,[],3);
105
106 figure(2);
107 subplot(1,2,1)
108 surf(V_Sarsa,'EdgeColor','none')
109 ylabel('Player Sum')
110 ylim([1 21])
111 xlabel('Dealer showing')
112 xlim([1 10])
113 zlabel('Value')
114 title('Value Function - using SARSA')
115 colormap(gca,coolwarm)
116 hcb1=colorbar;
117
118 h=gcf;
119 set(h,'Position',[50 50 1100 700]);
120 set(h,'PaperOrientation','landscape');
121 print(gcf, '-dpdf', 'Value Function.pdf')
122
123 %%
124 h=gcf;
125 h.PaperPositionMode='auto';
126 set(h,'PaperOrientation','landscape');
127 set(findall(gcf,'-property','FontSize'),'FontSize',12)
128 print(gcf, '-dpdf', 'Value Function.pdf','-fillpage')

```

Appendix E - Matlab codes for Q-Learning of Easy21

```
1 function [policy,action] = ...
   Q_Learning_get_action(Q,policy,state,No,count_state,epsilon_choice)
2
3     [ ~ , action_greedy] = max(Q(state(1),state(2),:));
4     if (Q(state(1),state(2),1)==Q(state(1),state(2),2))
5         action = randi([1 2])-1;
6         policy(state(1),state(2),1)=0.5;
7         policy(state(1),state(2),2)=0.5;
8     else
9         action_greedy=action_greedy-1;
10
11         if epsilon_choice==0
12             epsilon = No/(No+(count_state(state(1),state(2))));
13         else
14             epsilon=epsilon_choice;
15         end
16
17         %Choosing the action
18         p=rand;
19         if p<(1-epsilon/2) %This was deduced from initial form since ...
20             number of possible actions =2
21             action=action_greedy;
22         else
23             action=1-action_greedy;
24         end
25
26         if action_greedy==0
27             policy(state(1),state(2),1)=(1-epsilon/2);
28             policy(state(1),state(2),2)=epsilon/2;
29         elseif action_greedy==1
30             policy(state(1),state(2),2)=(1-epsilon/2);
31             policy(state(1),state(2),1)=epsilon/2;
32         end
33     end
```

```
1 function [Q,Ret_episode,upd_pol,opt_pol,count_s,count_s.a]=...
   Q_Learning_episode(obj,No,init_s,Q,pol,opt_pol,count_s,count_s.a,epsilon)
2
3
4     s=init_s;
5     Terminal=false;
6
7     %Generating the episode
8
9     while ~Terminal
10         %This is basically N(s) - used for step size (?)
11         count_s(s(1),s(2)) = count_s(s(1),s(2)) + 1; %This increments the ...
12             visits to state s(1) and s(2)
13
14         [upd_pol,action] = Q_Learning_get_action(Q,pol,s,No,count_s,epsilon);
15         count_s.a(s(1),s(2),action+1) = count_s.a(s(1),s(2),action+1) + 1;
16
17         [next_state,Terminal,reward,~,~]=step(obj, action);
18
19         %Obtain next action - ONLY IF NOT terminal
20
21         if ~Terminal
```

```

22         [opt_pol,optimal_action] = ...
           Q_Learning_optimal_action(Q,opt_pol,next_state);
23         d=Q(next_state(1),next_state(2),optimal_action+1);
24     else
25         d=0;
26     end
27
28     %Updating Q-value
29     Q(s(1),s(2),action+1) = Q(s(1),s(2),action+1) + ...
        (1/count_s_a(s(1),s(2),action+1))*(reward + d - ...
        Q(s(1),s(2),action+1));
30
31     if ~Terminal
32         s=next_state;
33     else
34         Ret_episode = reward;
35     end
36 end
37 end

```

```

1  clc
2  close all
3  clear all
4
5  % Defining the Game
6
7  %Calling the object into this file
8  import Easy21.v3
9  Game = Easy21.v3;
10
11 %No_vec=[500 300 200 100 10]; %Testing the performance of MC Control with ...
    different No
12 No_vec = 100; %Using the desired No value
13 e = [0,0.01,0.05,0.1,0.2,1]; %Epsilon-greedy
14
15 for i=1:length(e)
16     No=No_vec;
17
18     %Initialisation for SARSA
19
20     Q_QL = zeros(10,21,2); %Q(s,a)
21     c_s_a = zeros(10,21,2); %N(s,a)
22     c_s = zeros(10,21); %N(s)
23
24     max_length = 1000;
25     n_episodes=1000000;
26     b_pol=0.5*ones(10,21,2); %Behavioral policy
27     t_pol = 0.5*ones(10,21,2); %Target policy
28     Rewards=zeros(1,n_episodes);
29
30     index_test=50000;
31     size_test_set=100000;
32     Mean_vect=[];
33     Std_vect=[];
34     xaxis=[];
35
36     for episode=1:n_episodes
37         if mod((episode/n_episodes),0.05)==0
38             episode/n_episodes
39         end
40         init(Game,max_length);

```

```

41         init_s = reset(Game);
42         [Q_QL, Ret_episode, b_pol, t_pol, c_s, c_s_a]=...
43         Q_Learning_episode(Game, No, init_s, Q_QL, b_pol, t_pol, c_s, c_s_a, e(i));
44         Rewards(episode)=Ret_episode;
45
46         %This is when validation testing is done
47         if mod(episode, index_test)==0 || episode==1
48             [mean_test, std_test]=test(b_pol, size_test_set);
49             Mean_vect = [Mean_vect; mean_test];
50             Std_vect = [Std_vect; std_test];
51             xaxis=[xaxis; episode];
52         end
53     end
54
55     figure(1)
56     subplot(1,2,1)
57     hold on
58     grid on
59     plot(xaxis, Mean_vect)
60     ylabel('Mean of Rewards')
61     xlabel('Episode')
62
63     subplot(1,2,2)
64     hold on
65     grid on
66     plot(xaxis, Std_vect)
67     ylabel('Std of Rewards')
68     xlabel('Episode')
69 end
70
71 legend(num2str([0;0.01;0.05;0.1;0.2;1]))
72
73 %%
74 h=gcf;
75 h.PaperPositionMode='auto';
76 set(h, 'PaperOrientation', 'landscape');
77 set(findall(gcf, '-property', 'FontSize'), 'FontSize', 12)
78 print(gcf, '-dpdf', 'Learning Curve .pdf', '-fillpage')
79
80 %%
81 %Plot of policy of hitting
82 figure(2);
83 subplot(1,2,2)
84 h = heatmap(b_pol(:, :, 2));
85 xlabel('Player Sum')
86 xlim([1 21])
87 ylabel('Dealer showing')
88 ylim([1 10])
89 title('Optimal Behaviour policy of action HIT')
90 colormap(gca, 'gray')
91
92 %Histogram of Rewards
93 figure;
94 histogram(Rewards)
95
96 %% Storing Learning Curves mean
97 avg_Rewards_QLearning=Mean_vect;
98 save('avg_Rewards_QLearning.mat', 'avg_Rewards_QLearning')
99 %% Storing Learning Curves std
100 std_Rewards_QLearning=Std_vect;
101 save('std_Rewards_QLearning.mat', 'std_Rewards_QLearning')

```

```

102
103 %% Optimal value plot
104
105 [V_Q-learning, index] = max(Q_QL,[],3);
106
107 figure(2);
108 subplot(1,2,1)
109 surf(V_Q-learning','EdgeColor','none')
110 ylabel('Player Sum')
111 ylim([1 21])
112 xlabel('Dealer showing')
113 xlim([1 10])
114 zlabel('Value')
115 title('Value Function - using Q Learning')
116 colormap(gca,coolwarm)
117 hcb1=colorbar;
118
119 h=gcf;
120 set(h,'Position',[50 50 1100 700]);
121 set(h,'PaperOrientation','landscape');
122 print(gcf, '-dpdf', 'Value Function.pdf')
123
124 %%
125 h=gcf;
126 h.PaperPositionMode='auto';
127 set(h,'PaperOrientation','landscape');
128 set(findall(gcf,'-property','FontSize'),'FontSize',12)
129 print(gcf, '-dpdf', 'Value Function.pdf','-fillpage')

```

Appendix F - Matlab codes/functions used by all methods

```
1 function [mean_test std_test]=test(policy,size_test_set) %validation stage
2     Reward=zeros(1,size_test_set);
3     for i=1:size_test_set
4         Terminal=0;
5         s=[randi(10) randi(10)];
6
7         while ~Terminal
8             p=rand;
9             if p<policy(s(1),s(2),1)
10                 action=0;
11             else
12                 action=1;
13             end
14
15             [s,Terminal,reward]=step(s, action);
16         end
17
18         Reward(i)=reward;
19     end
20
21     mean_test=mean(Reward);
22     std_test=std(Reward);
23
24 end
25
26 function [state,Terminal,reward]=step(s, action)
27     %Action: 1 = hit , 0=stick
28     %colour: 1=black , -1=red
29
30     r=0; %Reward
31
32     player_sum = s(2);
33     dealer_sum=s(1);
34
35     if action==1 %When player hits
36         player_card_val = randi([1 10]);
37
38         x=rand; %Choosing random colour
39         if x<(1/3)
40             colour=-1;
41         else
42             colour=1;
43         end
44         player_card_col=colour;
45         player_sum = player_sum + (player_card_val * ...
46             player_card_col);
47         player_goes_bust = (player_sum>21 || player_sum<1);
48
49         Terminal=0;
50
51         if player_goes_bust == 1
52             r = -1;
53             Terminal = true;
54         end
55
56     end
57
58     if (action==0)%when player sticks
```

```

59         while (dealer_sum>0)&&(dealer_sum<17)
60             dealer_card.val = randi([1 10]);
61             z=rand; %Choosing random colour
62             if z<(1/3)
63                 col=-1;
64             else
65                 col=1;
66             end
67             dealer_card.col=col;
68
69             dealer_sum = dealer_sum + (dealer_card.val * ...
70                 dealer_card.col);
71             dealer_goes_bust = (dealer_sum>21 || dealer_sum<1);
72
73             if dealer_goes_bust == 1
74                 r = 1;
75             end
76         end
77
78         Terminal = true; %game ends when dealer sticks
79         if ~dealer_goes_bust
80             if player_sum>dealer_sum
81                 r=1;
82             elseif player_sum<dealer_sum
83                 r=-1;
84             elseif player_sum==dealer_sum
85                 r=0;
86             end
87         end
88
89         if (Terminal==1) %Game has finished
90             reward=r;
91             state = s;
92         end
93
94         if Terminal==0
95             state(2) = player_sum;
96             state(1) = s(1);
97             reward=r;
98         end
99     end

```

```

1  close all ; clear all ; clc
2
3  %Loading Learning Curves of the three methods
4  load('avg.Rewards_MC.mat');
5  load('std.Rewards_MC.mat');
6
7  load('avg.Rewards_SARSA.mat');
8  load('std.Rewards_SARSA.mat');
9
10 load('avg.Rewards_QLearning.mat');
11 load('std.Rewards_QLearning.mat');
12
13 %construction of the two x-axes to be used
14 n_episodes=1000000;
15 index_test_mean=100000;
16 index_test_std=50000;
17 xaxis_means=[];
18 xaxis_std =[];

```



```

19
20 for episode=1:n_episodes
21     if mod(episode,50000)==0
22         episode
23     end
24
25
26         if mod(episode,index_test_mean)==0 || episode==1
27             xaxis_means=[xaxis_means;episode];
28         end
29
30         if mod(episode,index_test_std)==0 || episode==1
31             xaxis_std=[xaxis_std;episode];
32         end
33     end
34
35 %% Average Reward plots
36 figure(1);
37 subplot(1,2,1)
38 hold on
39 grid on
40 plot(xaxis_means,avg_Rewards_MC);
41 plot(xaxis_means,avg_Rewards_SARSA);
42 plot(xaxis_means,avg_Rewards_QLearning);
43 xlabel('Episode'); ylabel('Average Reward')
44 legend('MC Control','SARSA','Q Learning','Location','best')
45
46 %% Std Plots
47 figure(1);
48 subplot(1,2,2)
49 hold on
50 grid on
51 plot(xaxis_std,std_Rewards_MC);
52 plot(xaxis_std,std_Rewards_SARSA);
53 plot(xaxis_std,std_Rewards_QLearning);
54 xlabel('Episode'); ylabel('Std of Rewards')
55 legend('MC Control','SARSA','Q Learning','Location','best')
56 hold off
57
58 %%
59 h=gcf;
60 h.PaperPositionMode='auto';
61 set(h,'PaperOrientation','landscape');
62 set(findall(gcf,'-property','FontSize'),'FontSize',12)
63 print(gcf, '-dpdf', 'Value Function.pdf','-fillpage')

```