

Department of Computing

MLNC – Machine Learning & Neural Computation – Dr Aldo Faisal

Lab Assignment 1

**Lab assignment: Understanding of MDPs**

Consider a stair climbing MDP made of 7 states  $\mathcal{S} = \{s_1, \dots, s_7\}$ . The top of the states and the bottom of the states are absorbing states  $s_1, s_7$ . Reaching the bottom yields you +100 reward, reaching the top yields you -100 reward. Each step up gives you -10 reward (effort), each step down gives you 10 reward. The two possible actions  $\mathcal{A} = \{Down, Up\}$  are deterministic. We do not specify a specific  $\gamma$ .

1. Draw the MDP as a graph (by hand)
2. Using below Matlab code as template for specifying an MDP, specify above Stair Climbing MDP (You will have to fix the code).
3. Implement code that takes a deterministic policy (i.e.  $\pi(s) = a$ ) and computes the value function for this policy for  $\gamma = \frac{1}{2}$
4. Compute for the "Always Down" policy and the "Always Up" policy the two value functions (code it so that it works for an arbitrary  $\gamma \in [0, 1]$ ).
5. Plot how varying  $\gamma$  will make an "always up" policy be more rewarding than an "always down" if you start in  $s_4$ .
6. Think about a way to find a better deterministic policy (e.g. by choosing actions so as to be greedy on the value function).

```

1  %% Basic code to specify an MDP
2  %% Learning in Autonomous Systems coursework
3  %% Aldo Faisal (2015), Imperial College London
4  function [S, A, T, R, StateNames, ActionNames, Absorbing] = StairClimbingMDP()
5  % States are: { s1 <-- s2 <=> s3 <=> s4 <=> s5 <=> s6 --> s7 };
6  S = 7;
7  StateNames = ['s1'; 's2'; 's3'; 's4'; 's5'; 's6'; 's7'];
8
9  % Actions are: {L,R} --> {1, 2 }
10 A = 2;
11 ActionNames = ['L'; 'R'];
12
13 % Matrix indicating absorbing states
14 Absorbing = [
15 %P  1   2   3   4   5   6   7   G   <-- STATES
16      1   0   0   0   0   0   0   1
17 ];
18
19 % load transition
20 T = transition_matrix()
21

```

```

22 % load reward matrix
23 R = reward_matrix(S,A)
24
25 %-----
26
27 % the transition subfunction
28 function prob = transition_function(priorState, action, postState)
29 % reward function (defined locally)
30 T = transition_matrix()
31 prob = T(postState,priorState,action)
32
33 % get the transition matrix
34 function T = transition_matrix()
35 TL = [
36 % MODIFY HERE
37 % 1 ...7      <-- FROM STATE
38 1   1   0   0   0   0   0 ; % 1 TO STATE
39 0   0   1   0   0   0   0 ; % .
40 0   0   0   1   0   0   0 ; % .
41 0   0   0   0   1   0   0 ; % .
42 0   0   0   0   0   1   0 ; % .
43 0   0   0   0   0   0   0 ; % .
44 0   0   0   0   0   0   1 ; % 7
45 ];
46 TR = [
47 % MODIFY HERE
48 % 1 ...7      <-- FROM STATE
49 1   0   0   0   0   0   0 ; % 1 TO STATE
50 0   0   0   0   0   0   0 ; % .
51 0   1   0   0   0   0   0 ; % .
52 0   0   1   0   0   0   0 ; % .
53 0   0   0   1   0   0   0 ; % .
54 0   0   0   0   1   0   0 ; % .
55 0   0   0   0   0   1   1 ; % 7
56 ];
57 T = cat(3, TL, TR); %transition probabilities for each action
58
59
60 %-----
61
62 % the locally defined reward function
63 function rew = reward_function(priorState, action, postState)
64 % reward function (defined locally)
65 % MODIFY HERE
66 if ((priorState == 2) && (action == 1) && (postState == 1))
67     rew = -1;
68 elseif ((priorState == 6) && (action == 2) && (postState == 7))
69     rew = 1;
70 elseif (action == 1)
71     rew = 0;
72 else
73     rew = 0;
74 end
75
76 % get the reward matrix
77 function R = reward_matrix(S, A)
78 % i.e. 11x11 matrix of rewards for being in state s,
79 %performing action a and ending in state s'
80 R = zeros(S, S, A);
81 for i = 1:S
82     for j = 1:A
83         for k = 1:S
84             R(k, i, j) = reward_function(i, j, k);
85         end

```

```
86     end
87 end
```