# **Reinforcement Learning - Assignment 2**

Q1

S	a	5'	2	P(8,315,0)
high	search	high	Tsearch	a a
high	Search	0	Bearch	1-8
to low	search	high	-3	1-B
las	secuch	لمن	rseconch	
high	waigh	high	Ywait	1
thigh low	wait	low	rwait	1
low	rechage	hig h	0	1
oses who	٠	b C 3,1	18,0)	
oses who	16 b	p C 3', r	on Eg	r all >0 33 in the boo
ases who This table we have P/	15 b	PC 3, x	(s,c)	33 in the boo
ases who This table we have P/ PI  Average	To be assum  Tigh  Statelow	PC 3, x xxxed red the Styll > h Styllou d for see	on Eg  t # :  Igh   St=  o, Aj= seen  sch = Yse	3.3 in the boo high, A= see th]= 13
Averag	16 b assum high 344=loca 4 rewar	P C 3, x  x sed  Styl = h  Styl=lou  d for sex  d for sex	$  S, c)$ on Eg $  A_1 = Sean$ $  A_2 = Sean$ $  A_3 = Sean$ $  A_4 = Sean$ $  A_4 = Sean$	3.3 in the boo high, A= see th]= 13

### Q2

For this question refer to the file **GridWorld.ipynb**. The file contains solutions to this problem in 2 approaches :

- 1. By solving simultaneous linear equations
- 2. Using policy evaluation

Two approaches have been used as an experiment to compare solutions obtained using iterative means as compared to the exact solution.

Note: The values achieved using policy evaluation have been rounded off to 1 decimal place

## Optimal State Value Function :

3.3	8.8	4.4	5.3	1.5
1.5	3.0	2.3	1.9	0.5
0.1	0.7	0.7	0.4	-0.4
-1.0	-0.4	-0.4	-0.6	-1.2
-1.9	-1.3	-1.2	-1.4	-2.0

Q3

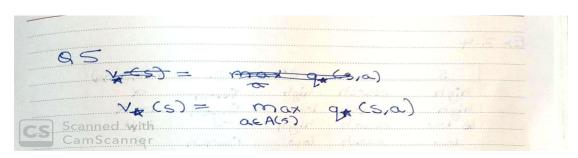
Ex	3:15
	G+= R+11 t/R+2+ = ExxR++1
	If we add a constant 'c' to each reward tray new return Grá 161
******	6; = R+11+c+ Y(R+12+c)+
	= C+ bc+ 82c+ + R+11+ 8R++2
***********	=> Gt = C + Gt
	$V_{\Pi}(s) = E\left[\sum_{t=1}^{\infty} \gamma^{t-1}R_{+} \mid s_{+}=s\right]$
	$V_{\Pi}(s) = E \left[ \sum_{t=1}^{\infty} s^{t-1} (R_{t} + c)   s_{t} = s \right]$
	= E [ = 8 t-1 (R+) 1 s+=s)
	E [ = 1 8+25]
winders	=> Vn(5) = Vn(5) + =
	Thatona = (8-1) 2 constant
Scar	nus the sign of the reward does not matter,
Can	NSCanner WTFSSMTWTFSSMTWTFSSMT

### Q4

For this question refer to the file GridWorld.ipynb.

Optimal State Value Function	Optimal Policy				
[[22. 24.4 22. 19.4 17.5] [19.8 22. 19.8 17.8 16. ] [17.8 19.8 17.8 16. 14.4] [16. 17.8 16. 14.4 13. ] [14.4 16. 14.4 13. 11.7]]	$ \begin{bmatrix} [ \ \ \ \ \ \ \ \ \ \ \ \ \ \ \ \ \ \$				

Q5



#### Q6

For this question refer to the file **DP\_GridWorld.ipynb**.

For **policy iteration** we can see that in the first two iteration the value function has increased for each state with respect to its value in the previous state. In the next iteration it converges to the optimal policy wherein the the value function remains the same as the previous state. Thus policy iteration is demonstrated to improve value function at each iteration.

In the below diagram we also see how the policy is changing on each iteration.

To correct the bug mentioned in Ex 4.4, we could use two possible fixes. This is observing the fact that multiple optimal policies would arise from multiple valid actions in a state have the same action-value function value --> Q(s,a1) = Q(s,a2).

Thus the two possible fixes are:

- 1. In case of multiple optimal actions pick either the first or the last action, by maintaining the same action ordering across states
- 2. We could have a stochastic policy where for a given state the any of the optimal actions is equiprobable to occur while the non optimal actions occur with a zero probability.

```
[[ 0. -14. -20. -22.]
[-14. -18. -20. -20.]
[-20. -20. -18. -14.]
[-22. -20. -14. 0.]]
          ' ' ← ↑
                1
[[ 0. -1. -2. -3.]
[-1. -2. -3. -2.]
[-2. -3. -2. -1.]
[-3. -2. -1. 0.]]
          [[ 0. -1. -2. -3.]
[-1. -2. -3. -2.]
[-2. -3. -2. -1.]
[-3. -2. -1. 0.]]
          ←↑ ' '→↓←↑' '
Dispalying Optimal Policy and Value Function
[[ 0. -1. -2. -3.]
[-1. -2. -3. -2.]
[-2. -3. -2. -1.]
[-3. -2. -1. 0.]]
          1
```

Below are the value functions after every iteration of **Value Iteration**. At the end you can see the optimal value function and policy. Notice that it matches the one generated by policy iteration.

```
[[ 0. -1. -1. -1.]
[-1. -1. -1. -1.]
[-1. -1. -1. -1.]
[-1. -1. -1. 0.]]
[[ 0. -1. -2. -2.]
 [-1. -2. -2. -2.]
 [-2. -2. -2. -1.]
[-2. -2. -1. 0.]]
[[ 0. -1. -2. -3.]
[-1. -2. -3. -2.]
 [-2. -3. -2. -1.]
[-3. -2. -1. 0.]]
[[ 0. -1. -2. -3.]
 [-1. -2. -3. -2.]
[-2. -3. -2. -1.]
[-3. -2. -1. 0.]]
Optimal Policy and Value Function
 [[ 0. -1. -2. -3.]
 [-1. -2. -3. -2.]
[-2. -3. -2. -1.]
 [-3. -2. -1. 0.]]
             1
       1
             \ ^{1}\ ^{1}\rightarrow\downarrow\leftarrow\uparrow\ ^{1}\ ^{1}\ \rightarrow\downarrow
      1
                                                       ']]
```

#### **Q7**

Refer to the file Jacks\_Car\_Rental.ipynb for the solution.

Below are the heatmaps for the policies after every iteration until the converge. This heatmap is for the modified Jack's Car Rental decribed in Ex4.7. Prior to running for the Ex 4.7 I have also verified the correctness of the code by running on the original Jack's Car Rental Example. Necessary comments are placed in the code to adapt current code to run the original Jack's Car Rental problem

