

4.1 Solution

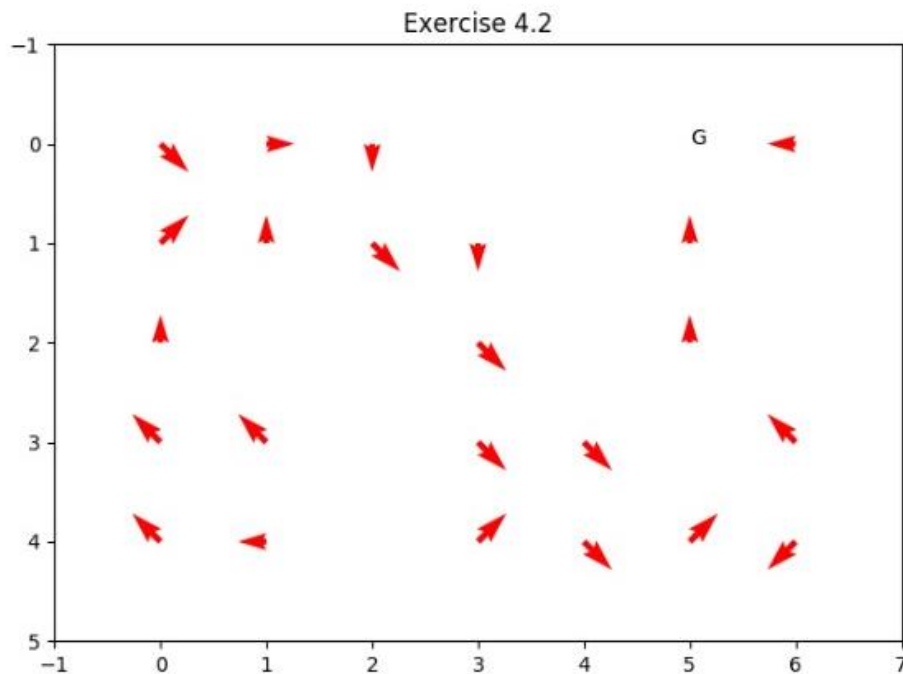
TD(0)								
0	-0.97	-2.38	-13.05	W	W	G	0.0	
1	-9.72	-18.64	-14.01	-9.52	W	0.0	W	
2	-55.14	W	W	-4.84	W	0.0	W	
3	S	-90.74	W	-1.0	-0.05	W	0.0	
4	-63.3	-87.86	W	-0.01	-0.01	-0.01	-1.0	
5								
	0	1	2	3	4	5	6	7

4.2 Solution

Value grid:

0	-13.32	-15.75	-16.21	W	W	G	0.0	
1	-15.07	-19.05	-14.1	-22.42	W	0.0	W	
2	-19.11	W	W	0.0	W	0.0	W	
3	S	-16.19	W	-0.05	-0.05	W	-5.05	
4	-12.94	-13.1	W	0.0	0.0	-0.05	0.0	
5								
	0	1	2	3	4	5	6	7

The optimal policy:



We got this policy which looks meaningful and precise after a number of tries.

Note: As the policy selection is non-deterministic in both the cases, the value grid or the policy changes with every try and the solution (in case of 4.2) is not always optimal or meaningful. Also, the convergence of the algorithm depends on the parameters alpha and gamma. We have used the values for which we got the fastest convergence in the presented solution.