

Problem 1: Value Iteration

a), b)

Gamma = 1.0 gives the following value function and policy.

0.812	0.868	0.918	1.000	→	→	→	
0.762		0.660	-1.000	T		T	
0.705	0.655	0.611	0.388	T	←	←	←

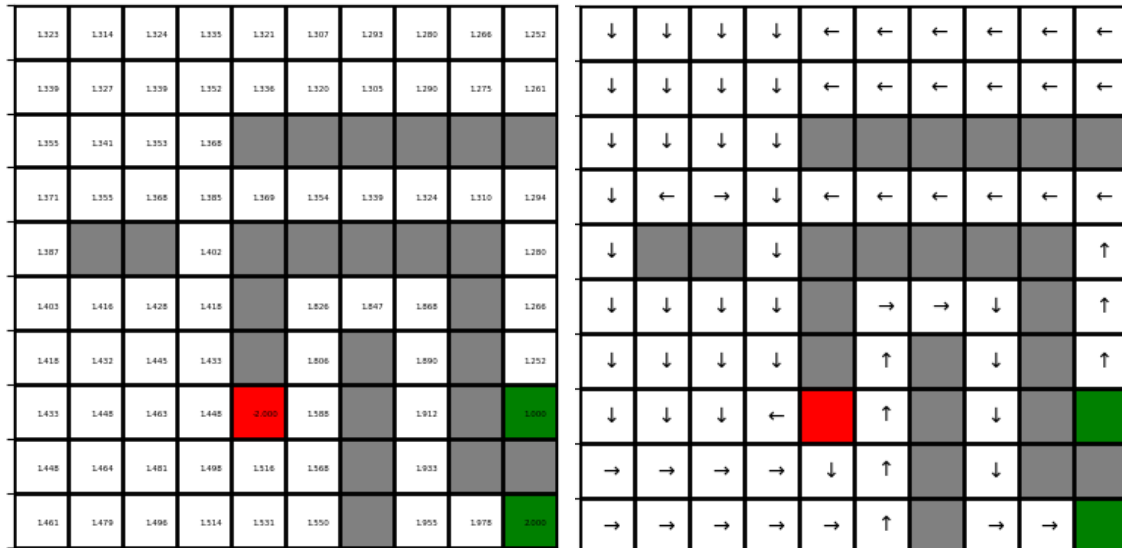
c)

See figures. With greater discount factor (bigger gamma) it is easier to see the +2 terminal state even when close to the +1 state. Thus, gamma = 0.99 finds a better policy.

Gamma = 0.9

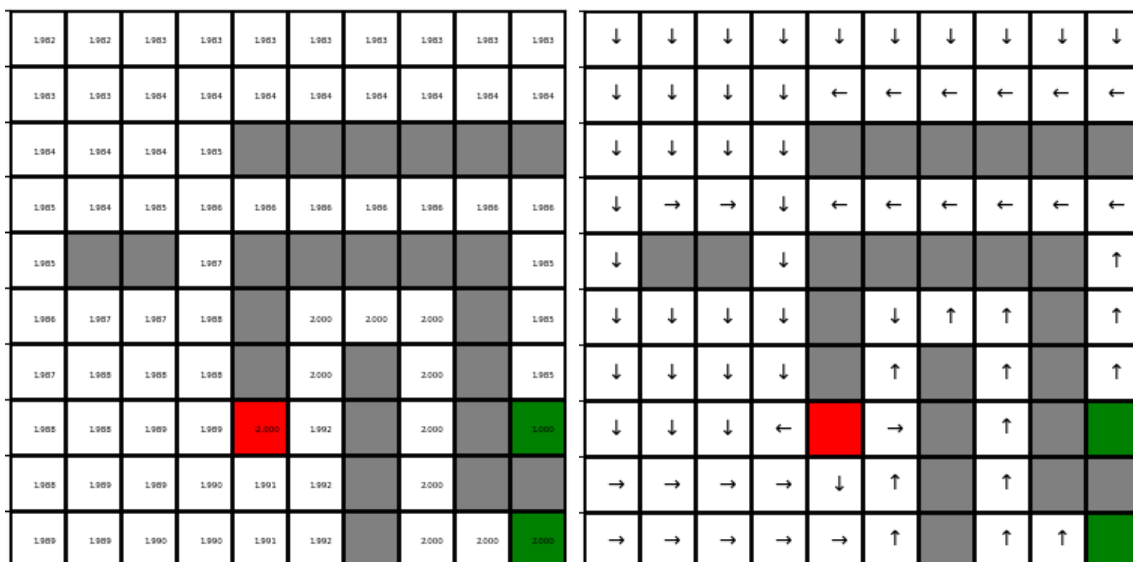
0.154	0.172	0.192	0.211	0.189	0.169	0.151	0.135	0.121	0.108	↓	→	↓	↓	←	←	←	←	←	←
0.172	0.192	0.216	0.239	0.212	0.188	0.168	0.148	0.131	0.116	↓	↓	↓	↓	←	←	←	←	←	←
0.192	0.216	0.242	0.272							→	→	↓	↓						
0.212	0.242	0.272	0.307	0.349	0.392	0.440	0.495	0.556	0.624	→	→	→	→	→	→	→	→	→	↓
0.189			0.273						0.705	↑			↑						↓
0.169	0.190	0.214	0.242		0.775	0.876	0.984		0.792	→	→	→	↑		→	→	↓		↓
0.155	0.173	0.193	0.214		0.690		1.111		0.890	→	→	→	↑		↑		↓		↓
0.170	0.190	0.211	0.193	2.000	0.491		1.249		2.000	↓	↓	↓	←		↑		↓		
0.190	0.213	0.240	0.269	0.305	0.431		1.403			→	→	→	→	↓	↑		↓		
0.208	0.235	0.265	0.299	0.338	0.381		1.576	1.780	2.000	→	→	→	→	→	↑		→	→	

Gamma = 0.99



d)

With gamma = 1.0 we should from any state converge towards the +2 reward. Thus all states should have a reward close to 2 and the policy becomes kind of random in the cases where neighboring states are equally good. This is confirmed by the following results. As shown, the policy is unable to reach any terminal state.



Problem 2: Policy Iteration

a)

I choose Iterative Policy Evaluation because it was the first algorithm in the assignment. Well, it also looks easier :P

b)

gamma = 1.0. Identical with 1a).

0.812	0.868	0.918	1.000	→	→	→	
0.762		0.660	-1.000	T		T	
0.705	0.655	0.611	0.388	T	←	←	←

c)

Policy iteration can change policy to an equally good one for each iteration, thus never satisfy the termination requirement of two consecutive identical policies.

d)

Value iteration uses 13 iterations and its biggest error (except from changing terminal states from 0 to +1/-1) was 0.792 on the tiny grid. For policy iteration the corresponding numbers are 7 iterations and 8.156 error.

Problem 3

Adding a negative reward of -0.01 for every non-terminal state finds a useful and optimal policy. As can be seen below, we are guaranteed to end in the +2 reward even when next to the +1 reward.

1501	1493	1501	1510	1499	1489	1475	1467	1457	1446
1512	1504	1512	1522	1510	1499	1487	1475	1464	1452
1524	1514	1523	1535						
1536	1524	1535	1546	1535	1524	1513	1502	1491	1479
1545			1559						1465
1559	1569	1577	1570		1909	1921	1932		1457
1570	1580	1589	1581		1895		1944		1446
1581	1591	1601	1592	2505	1684		1955		1580
1591	1602	1614	1625	1635	1671		1966		
1601	1612	1624	1636	1645	1659		1977	1989	2000

↓	↓	↓	↓	←	←	←	←	←	←
↓	↓	↓	↓	←	←	←	←	←	←
↓	↓	↓	↓						
↓	←	→	↓	←	←	←	←	←	←
↓			↓						↑
↓	↓	↓	↓		→	→	↓		↑
↓	↓	↓	↓		↑		↓		↑
↓	↓	↓	←		→		↓		
→	→	→	→	↓	↑		↓		
→	→	→	→	→	↑		→	→	