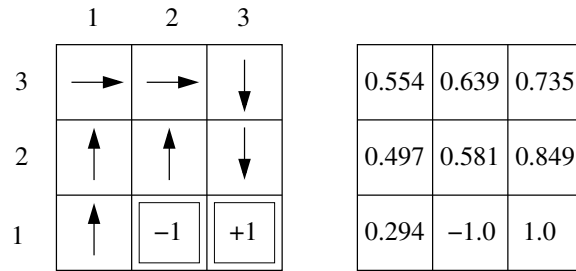


In the mini grid world shown below, there are two terminal states: state (2,1) with a negative reward of -1, and (3,1) with a positive reward of +1. The transition model is the same as the grid world in the course slides: an action succeeds with probability 0.8, and goes to the left or right with probability 0.1, respectively. However, moves into the wall are not allowed. The optimal policy π is in the left figure, and the correct utility function the optimal policy is in the right figure.



Below is a series of three trials (1, 2, and 3 left to right) in this environment. Starting in state (1,3), actions were taken according to the fixed policy π above, and ended once a terminal state is reached. The trials are as follows:

<i>S</i>	<i>A</i>	<i>R</i>	<i>S</i>	<i>A</i>	<i>R</i>	<i>S</i>	<i>A</i>	<i>R</i>
(1,3)	E	0	(1,3)	E	0	(1,3)	E	0
(2,3)	E	0	(2,3)	E	0	(1,2)	N	0
(3,3)	S	0	(2,2)	N	0	(1,3)	E	0
(3,2)	S	1	(2,3)	E	0	(2,3)	E	0
(3,1)			(3,3)	S	0	(3,3)	S	0
			(3,2)	S	1	(3,2)	S	1
			(3,1)			(3,1)		

1. Estimate the transition function $T(s, a, s')$ as much as possible given these limited trials.

	(1,1)	(1,2)	(1,3)	(2,1)	(2,2)	(2,3)	(3,1)	(3,2)	(3,3)
(1,2),N									
(1,3),E									
(2,2),N									
(2,3),E									
(3,2),S									
(3,3),S									

2. Perform policy evaluation with your estimated MDP to find the state values. Assume $\gamma = 0.9$. Set up the linear equations to solve exactly without iteration.

$$V^\pi((1, 2)) =$$

$$V^\pi((1, 3)) =$$

$$V^\pi((2, 2)) =$$

$$V^\pi((2, 3)) =$$

$$V^\pi((3, 1)) =$$

$$V^\pi((3, 2)) =$$

$$V^\pi((3, 3)) =$$

3. Perform TD learning for the three trials left to right. Use $\alpha(n) = 1/n$, where n is the trial number. Only write the non-trivial updates.

Trial 1:

$$V^\pi(x, y) = \dots \text{ (put in appropriate } x \text{ and } y, \text{ and add as necessary)}$$
$$V^\pi(x, y) = \dots$$

Trial 2:

$$V^\pi(x, y) = \dots \text{ (put in appropriate } x \text{ and } y, \text{ and add as necessary)}$$
$$V^\pi(x, y) = \dots$$

Trial 3:

$$V^\pi(x, y) = \dots \text{ (put in appropriate } x \text{ and } y, \text{ and add as necessary)}$$
$$V^\pi(x, y) = \dots$$