

Module 4 – Self-Check

As noted in the video lectures, it is possible to solve value iteration for a 1-d Markov Decision Process in a spreadsheet by explicitly modeling the π , V , and Q arrays. You have been provided with a starter spreadsheet (module-4-self-check.xlsx)

For the self-check, you have been provided with a starter spreadsheet (module-4-self-check.xlsx). Fill in the formulas for Value Iteration (stochastic version) so that it looks like the following:

Q	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V	W	X	Y	Z	AA	AB	AC	AD	AE	AF	AG	AH
1	epsilon	0.01																		1	2	3	4	5	6	7		1	2	3	4	5	6	7
2	discount rate	0.90																			5.00	0.00	0.00	0.00	0.00	0.00	0.00		0.00	0.00	0.00	0.00	0.00	5.00
3	% planned	0.90																																
4	% surprise	0.10																																
5																																		
6		t=	0	1	2	3	4	5	6	7																								
7																																		
8			1	x	?	?	?	?	?	x																								
9			2	x	<	?	?	?	?	>	x																							
10			3	x	<	<	?	?	?	>	x																							
11			4	x	<	<	?	?	?	>	x																							
12			5	x	<	<	?	?	?	>	x																							
13			6	x	<	<	?	?	?	>	x																							
14			7	x	<	<	?	?	?	>	x																							
15			8	x	<	<	?	?	?	>	x																							
16			9	x	<	<	?	?	?	>	x																							
17			10	x	<	<	?	?	?	>	x																							
18																																		
19																																		

1. the cell B1 is named “discount_rate”, you can refer to it by that name in cell formulas, for example, “=discount_rate*10”
2. the cell B2 is named “planned”, you can refer to it by that name in formulas.
3. The cell B3 is named “surprise”, you can refer to it by that name in formulas.
4. if(condition,then,else) can be used to test values in cells. It can also be nested (the first one will need an =).
5. =max(cellref1, cellref2) will return the maximum value.
6. You do not need to fill in formulas for the goal states 1 or 7 except in the case of $V(s)$.
7. Use “<” for “left”, “?” for “pick random” and “>” for “right” in your policy formula.
8. There isn’t a good way to control the number of iterations automatically, so epsilon is there to let you know when you should stop copying.

You should get the same results as above and converge in 10 iterations.