

Reference 1:

Value Iteration

Initialize array v arbitrarily (e.g., $v(s) = 0$ for all $s \in \mathcal{S}^+$)

Repeat

$\Delta \leftarrow 0$

For each $s \in \mathcal{S}$:

$temp \leftarrow v(s)$

$v(s) \leftarrow \max_a \sum_{s'} p(s'|s, a) [r(s, a, s') + \gamma v(s')]$

$\Delta \leftarrow \max(\Delta, |temp - v(s)|)$

until $\Delta < \theta$ (a small positive number)

Output a deterministic policy, π , such that

$$\pi(s) = \arg \max_a \sum_{s'} p(s'|s, a) [r(s, a, s') + \gamma v(s')]$$

Source: <https://moodle.bath.ac.uk/mod/page/view.php?id=974187>

Reference 2:

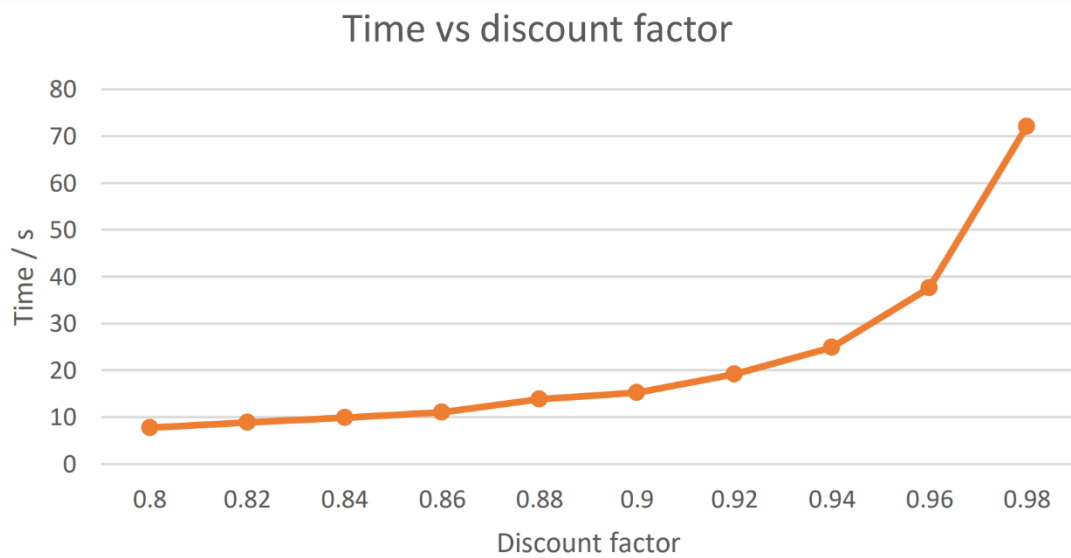
$$V_{k+1}(s) = \max_a \sum_{s', r} p(s', r|s, a) (r + \gamma V_k(s'))$$

Source: <https://moodle.bath.ac.uk/mod/page/view.php?id=974187>

Reference 3:

With constant theta = 0.001, n = 10,000 tests

Discount factor	Time taken / s	Avg. score
0.8	7.75	12.1072
0.82	8.9375	12.2835
0.84	9.9375	12.6719
0.86	11	12.742
0.88	13.8125	12.9469
0.9	15.2656	13.0513
0.92	19.1562	13.2906
0.94	24.9062	13.3043
0.96	37.5469	13.1885
0.98	72.0781	12.6026
1 Infinite	Never found	

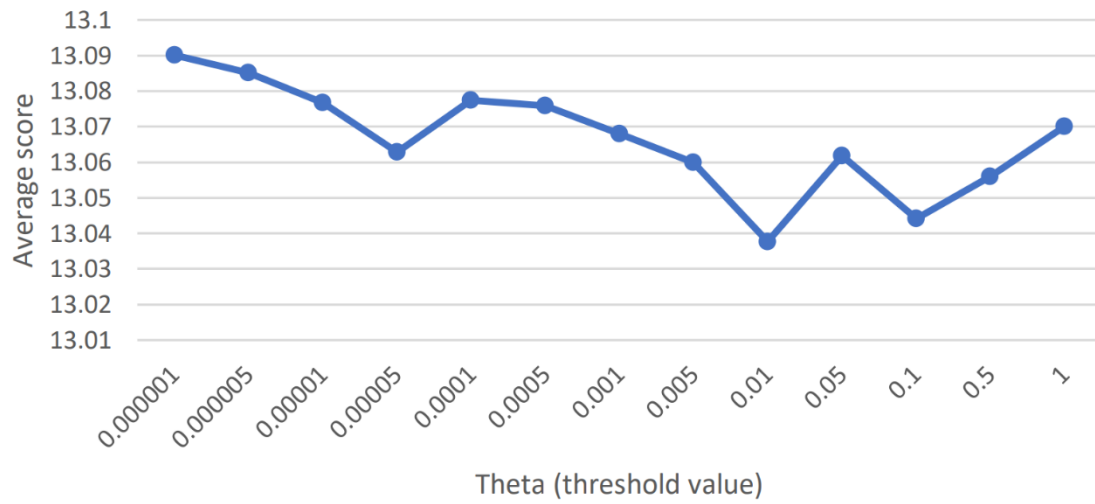


Reference 4:

Constant discount factor = 0.9, n = 10,000 tests

Theta	Time taken / s	Avg. score
0.000001	24.5781	13.0901
0.000005	22.4219	13.0851
0.00001	21.6719	13.0768
0.00005	19.4062	13.0628
0.0001	18.1562	13.0774
0.0005	16.0469	13.0759
0.001	15.0625	13.068
0.005	13.2656	13.06
0.01	12.2656	13.0377
0.05	9.9844	13.0618
0.1	9.1562	13.0442
0.5	7.1094	13.056
1	6.203	13.07

Avg. score vs theta



Time vs theta

