

Dynamic Programming

February 5, 2020





Rupam Mahmood





A deterministic policy $\pi(s)$ outputs an action $a \in \mathcal{A} = \{a_1, a_2\}$ ally, a policy $\pi(\cdot|s)$ outputs the p204ab2it9es for all actions: $\pi($ How can you write a deterministic policy in this form? Let $\pi(s) = a_i$ and define $\pi(\cdot|s)$. -2.4 -2.9 -3.0 -2.9 k = 3-2.9 | -3.0 | -2.9 | -2.4763.0-2.9 -2.4 -0.0

Example 4.1 Consider the 4×4 gridworld shown below.









optimal policy



	0.0	-1.0	-1.0	-1
Workshaet aue	sti0	1 .0	-1.0	-1.
h - 1	-1.0	-1.0	-1.0	-1
In itorativo poliev ovalu	-1.0	-1.0	-1.0	0.
the Bellman equation r	nany t	times 1	to gen	erat
generate a sequence of a		va¶u₹	fußç	
k=2	+1(\$7	$=20^{2}$	$\pi(2 0)$	\sum_{n}
	-2.0	-2.0	-2.0	<i>s</i> [,] , <i>r</i>
	-2.0	-2.0	-1.7	0.

k = 3

|0.0| - 2.4| - 2.9| - 3.0-2.4 -2.9 -3.0 -2.9 -2.9 -3.0 -2.9 -2.4







Worksheet question

The policy iteration algorithm on page 80 has a subtle bug in that it may never terminate if the policy continually switches between two or more policies that are equally good. This is ok for pedagogy, but not for actual use. Modify the pseudocode so that convergence is guaranteed. Note that there is more than one approach to solve this problem.

Policy Iteration (using iterative policy evaluation) for estimating $\pi \approx \pi_*$

- 1. Initialization
- 2. Policy Evaluation Loop: $\Delta \leftarrow 0$ Loop for each s $v \leftarrow V(s)$ $V(s) \leftarrow \sum_{s'}$ $\Delta \leftarrow \max(\Delta$
- 3. Policy Improvement policy-stable $\leftarrow true$ For each $s \in S$: old-action $\leftarrow \pi($ $\pi(s) \leftarrow \operatorname{arg\,max}_{k}$ If old-action $\neq \eta$ If *policy-stable*, then

 $V(s) \in \mathbb{R}$ and $\pi(s) \in \mathcal{A}(s)$ arbitrarily for all $s \in S$

$$\in S$$
:

$$\sum_{r,r} p(s', r | s, \pi(s)) \left[r + \gamma V(s') \right]$$

$$\Delta, |v - V(s)|)$$

until $\Delta < \theta$ (a small positive number determining the accuracy of estimation)

s)

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

 $\pi(s)$, then *policy-stable* \leftarrow *false*
 n stop and return $V \approx v_*$ and $\pi \approx \pi_*$; else go to 2

Worksheet question

(*Exercise 4.1 S&B*) Consider the 4x4 gridworld below, where actions that would take the agent off the grid leave the state unchanged. The task is episodic with $\gamma = 1$ and the terminal states are the shaded blocks. Using the precomputed values for the equiprobable policy below, what is $q_{\pi}(11, \text{down})$? What is $q_{\pi}(7, \text{down})$?



		1
	1	2
4	5	6
8	9	10
12	13	14





 $R_t = -1$ on all transitions





Worksheet question

(*Exercise 4.1 from S&B*) Suppose in the above gridworld where a new state 15 is added to the gridworld just below state 13, and its actions, left, up, right, and down, take the agent to the states 12, 13, 14, and 15, respectively. Assume that the transitions from the original states are unchanged. What, then is, $v_{\pi}(15)$ for the equiprobable random policy? Now suppose the dynamics of state 13 are also changed, such that action down from state 13 takes the agent to the new state 15. What is $v_{\pi}(15)$ for the equiprobable random policy in this case?



0.0	-14.	-20.	-22.
-14.	-18.	-20.	-20.
-20.	-20.	-18.	-14.
-22.	-20.	-14.	0.0



	1	2	3
	5	6	7
	9	10	11
2	13	14	

 $R_t = -1$ on all transitions



Demo of an MDP in action

