

Dynamic Programming

- Planning by dynamic programming assumes full knowledge of the MDP

for Prediction / Evaluation

- * Input : MDP $\langle S, A, P, R, \gamma \rangle$ and policy π
- * Output : Value function v_π

for Control

- * Input : MDP $\langle S, A, P, R, \gamma \rangle$
- * Output : Optimal value function v^* and optimal policy π^*

Problem

Bellman Equation

Algorithm

Prediction

Bellman Equation Interpretation

Iterative Policy Evaluation

Control

Bellman Expectation +
Greedy Policy Improvement

Policy Iteration

Control

Bellman Optimality Equation

Value Iteration

Example : Grid World

1	2	3	4
5	6 R	7	8
9	10	11	12
13	14	15	16

Action \leftrightarrow
Reward is -1 for all transition

A bot is required to transverse a grid of 4×4 dimensions to reach its goal (1 or 16). Each step is associated with a reward of -1. There are 2 terminal states here : 1 and 16 and 14 non-terminal states given by [2, 3, ..., 15]. Consider a random policy for which at every state, the probability of every action {up, down, left, right} is equal to 0.25.

Initializing V_0 for the Random Policy to all 0s.

0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0

Let's calculate V_1 for all the states of 6:

$$V_1(6) = \sum_{a \in \{u, d, l, r\}} \pi(a|6) \sum_{s' \neq 6} p(s'|6, a) [r + \gamma V_0(s')]$$

$$= \sum_{a \in \{u, d, l, r\}} \underbrace{\pi(a|6)}_{0.25} \sum_{s'} p(s'|6, a) \underbrace{[r + \gamma V_0(s')]}_{=-1} = 0$$

$$= 0.25 * \{-p(2|6, u) - p(10|6, d) - p(5|6, r) \\ - p(7|6, s)\}$$

$$= 0.25 + \{ -1 - 1 - 1 - 1 \}$$

$$= -1$$

$$\Rightarrow V_1(6) = -1$$

Similarly for all non-terminal states

$$V_1(s) = -1$$

For terminal states $p(s'|s,a) = 0$

hence $V_k(1) = V_k(1_b) = 0$ for all k .

So V_1 for the Random policy is given by

0.0	-1.0	-1.0	-1.0
-1.0	-1.0	-1.0	-1.0
-1.0	-1.0	-1.0	-1.0
-1.0	-1.0	-1.0	0.0

For $V_2(s)$, Assume $\gamma = 1$

$$V_2(6) = \sum_{a \in \{u,d,l,r\}} \pi(a|6) \sum_{s'} p(s'|6,a) [0 + \gamma V_1(s')]$$

$= 0.25 \forall a$

$$= 0.25 \times \{ p(2|6,u) [-1 - \gamma] +$$

$$p(10|6,d) [-1 - \gamma] + p(5|6,l) [-1 - \gamma]$$

$$+ p(7|6,r) [-1 - \gamma] \}$$

$$= 0.25 \times \{-2-2-2-2\}$$

$$= -2$$

1 (G)	2	3	4	
5	6	7	8	
9	10	11	12	
13	14	15	16 (G)	

All the states marked in red in above are identical for the purpose of calculating the value function. Hence for all these states:

$$V_2(s) = -2$$

For all the remaining states i.e. 2, 5, 12 and 15 V_2 can be calculate as follows:

$$V_2(2) = \sum_{a \in \{u, d, l, r\}} \pi(a|2) \sum_{s'} p(s'|2, a) [r + \gamma V_1(s')]$$

$$= 0.25$$

$$= 0.25 \times \left[p(2|2,u) [-1-\gamma] + p(6|2,d) [-1-\gamma] \right. \\ \left. + p(1|2,l) [-1-\gamma+0] - p(3|2,r) [-1-\gamma] \right]$$

$$= 0.25 \times \{-2-2-1-2\}$$

$$= -1.75$$

$$\Rightarrow \boxed{V_2(2) = -1.75}$$

V_2 for the random policy

0.0	-1.7	-2.0	-2.0
-1.7	-2.0	-2.0	-2.0
-2.0	-2.0	-2.0	-1.7 ₁₂
-2.0	-2.0	-1.7 ₁₅	0.0

If we repeat this step several times we get V_H

\therefore Note : Check PPT for the solution.
See diagrams.

1. Problem

Frozen Lake Solution

given State Representation

S	F	F	F
F	H	F	H
F	F	F	H
H	F	F	G

1. Using Value Iteration

S = State

F = Frozen

H = Hole

G = Goal

$\gamma = 0.9$

Bellman Equation for Value Iteration

Value function

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

given

G = 1.0 (terminal state)

H = 0 (game over)

other state = 0

Initial Value tables

S	0	0	0
O	H	0	H
O	0	0	H
H	0	0	G

Formula

Step 1: for Iteration Value Update

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

Initial $V(s_1) = 0$

$V(s_1) = 1$

Iteration 1:

$$G = 0$$

State adjacent to the goal update

1. (3,2) But move Right to (3,3)

$$V_{(3,2)} = 0.9 \times 1.0 = 0.9$$

2. (2,3) But move Down to (3,3)

$$V_{(2,3)} = 0.9 \times 1.0 = 0.9$$

3. (3,1) But move Right to (3,2)

$$V_{(3,1)} = 0.9 \times 0.9 = 0.81$$

4. (2,2) But move Right (2,3)

$$V_{(2,2)} = 0.9 \times 0.9 = 0.81$$

5. (2,1) But move Right (2,2)

$$V_{(2,1)} = 0.9 \times 0.81 = 0.729$$

6. (1,2) But move Down to (2,2)

$$V_{(1,2)} = 0.9 \times 0.729 = 0.6561$$

Iteration 2:

Repeating the same steps, value propagate backward

1. (1,1) (avoiding hole):

$$V_{(1,0)} = 0.9 \times 0.97 = 0.96$$

2. (0,2)

$$V_{(0,2)} = 0.9 \times 0.97 = 0.96$$

3. (0,1)

$$V_{(0,1)} = 0.9 \times 0.96 = 0.95$$

4. (0,0) (Start state)

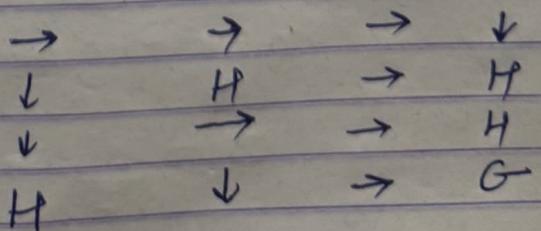
$$V_{(0,0)} = 0.9 \times 0.95 = 0.94$$

S	0.94	0.95	0.96	0.97
0.96	H	0.97	H	0.
0.97	0.98	0.99	H	
H	0.98	0.99	I.	

Extracting the Optimal Policy

(0,0)	→ Right (Toward higher value)
0,1	→ Right
0,2	→ Right
0,3	→ Down
1,3	→ Down
2,3	→ Right
3,2	→ Right
3,3	→ Goal

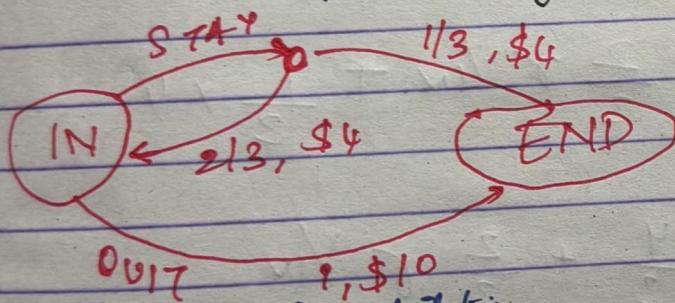
Final Opt



2. Problem

At any time the game is specified by two states: IN and END

When the game is in the state IN, the agent can choose the action from the set $A = \{STAY, QUIT\}$. If QUIT is selected, then the game will go to state END and the agent receive a reward \$10 with a probability of 1. If STAY is selected, the game will stay at IN state and the agent receive a reward \$4 with a probability of $2/3$ or transition to state END, and the agent receive a reward \$4 with a probability of $1/3$.



STATE Transition Probabilistic

Action	Next State	Prob	Reward
QUIT	END	1.0	10
STAY	IN	2/3	4
STAY	END	1/3	4

$$V_{(S)} = \max_a \sum_{s'} P(s'|s, a) [R_{(S, a, s')} + \gamma V(s')]$$

$$\gamma = \frac{b}{b+w} = \frac{1}{1+0.5} = 0.67$$

$$V_{(\text{END})} = 0$$

State = IN given that the agent

for QUIT:

$$V_{\text{QUIT}}(\text{IN}) = 10 + 0.7 \times 0$$

for STAY:

$$V_{\text{STAY}}(\text{IN}) = \left(\frac{2}{3} \times (4 + \gamma V_{(\text{IN})}) + \left(\frac{1}{3} \times (4 + \gamma V_{(\text{END})}) \right) \right)$$

$$= \frac{2}{3} \times 4 + \frac{2}{3} \gamma V_{\text{IN}} + \frac{1}{3} \times 4$$

$$= \frac{8}{3} + \frac{4}{3} + \frac{2}{3} \gamma V_{\text{IN}}$$

$$V_{\text{STAY}}(\text{IN}) = 4 + \frac{2}{3} \gamma V_{\text{IN}}$$

\therefore the agent selects STAY with probability 0.5 and QUIT with 0.5, we take weighted sum

$$V(\text{IN}) = 0.5 \times \left(4 + \frac{2}{3} \gamma V_{\text{IN}} \right) + 0.5 \times 10$$

$$V_{\text{IN}} = 0.5 \times 0.4 + 0.5 \times \frac{2}{3} \gamma V_{\text{IN}} + 5$$

$$= 2 + \frac{1}{3} \gamma V_{\text{IN}} + 5$$

$$= 7 + \frac{1}{3} \gamma V_{\text{IN}} + 5$$

$$V_{IN} - \frac{1}{3}\gamma V_{IN} = 7$$

$$(1 - \frac{1}{3}\gamma) V_{IN} = 7$$

$$V_{IN} = \frac{7}{1 - \frac{1}{3}\gamma}$$

$$\gamma = 0.9$$

$$V_{IN} = \frac{7}{1 - \frac{1}{3} \times 0.9} = [10]$$

$$V_{(IN)} = 0$$

$$\text{Iteration 1} = V_{IN} = 0.5 \times (4 + 0.6 \times 0) + 0.5 \times 10 \\ \Rightarrow V_{IN} = 7$$

$$\text{Iteration 2} = V_{IN} = 0.5 \times (4 + 0.6 \times 7) + 0.5 \times 10 \\ = 9.1$$

$$\text{Iteration 3} = V_{IN} = 0.5 \times (4 + 0.6 \times 9.1) + 0.5 \times 10 \\ = 9.78$$

$$4 \quad V_{IN} = 0.5 \times (4 + 0.6 \times 9.73) + 0.5 \times 10 \\ = 9.919$$

$$5 \quad V_{IN} = 9.9757$$

$$6 \quad V_{IN} = 9.9927$$

$$7 \quad V_{IN} = 9.9978$$

$$8 \quad V_{IN} = 9.999$$