

# Lecture 17 RL : MDPs & Value / Policy Iteration

## 1. Markov Decision Process (MDP)

Parameters:  $(S, A, \{P_{sa}\}, \gamma, R)$

S - Set of states

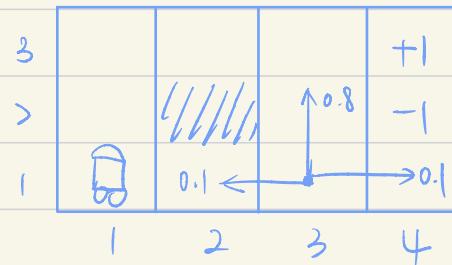
A - Set of actions

$P_{sa}$  - State transition probabilities ( $\sum_{S'} P_{sa}(S') = 1$ )

$\gamma$  - discount factor,  $\gamma \in [0, 1]$

R - reward function

Example: 3



11 states, Action:  $\{N, S, E, W\}$

$$P_{(3,1)N}((3,2)) = 0.8 \quad R((4,3)) = +1 \quad \text{reach destination asap.}$$
$$P_{(3,1)N}((4,1)) = 0.1 \quad R((4,2)) = -1$$

$$P_{(3,1)N}((2,1)) = 0.1 \quad R(S) = -0.02 \text{ for all other states}$$
$$P_{(3,1)N}((3,3)) = 0$$

Process:  $S_0 \rightarrow$  choose action  $a_0 \rightarrow$  get to  $s_1 \sim P_{S_0, a_0} \rightarrow$  choose action  $a_1 \rightarrow$   
 get to  $s_2 \sim P_{s_1, a_1} \dots$

Total payoff:  $R(S_0) + \gamma R(S_1) + \gamma^2 R(S_2) + \dots$

$\gamma = 0.99$  (close to 1)  $\Rightarrow$  encourage to get reward faster

Goal: Choose actions over time to maximize  $E[R(S_0) + \gamma R(S_1) + \gamma^2 R(S_2) + \dots]$

Result: find a policy  $\pi: S \rightarrow A^3$

when in state  $s$ , take action  $\pi(s)$

	→	→	→	+1
↑		↑	-1	
↑	←	←	←	

## 2. Some definition & Preparation

Define:  $V^\pi$ ,  $V^*$ ,  $\pi^*$

- $V^\pi$ : For a policy  $\pi$ ,  $V^\pi: S \rightarrow \mathbb{R}$  is s.t.  $V^\pi(s)$  is the expected total payoff for starting in state  $s$  and executing  $\pi$ .

$\hookrightarrow$  is the value function for policy  $\pi$ .

$$V^\pi(s) = E[R(S_0) + \gamma R(S_1) + \dots | \pi, S_0 = s]$$

Example:

$\rightarrow$	$\rightarrow$	$\rightarrow$	+1
$\downarrow$		$\rightarrow$	-1
$\rightarrow$	$\rightarrow$	$\uparrow$	$\uparrow$

$\pi$

0.52	0.73	0.77	+1
-0.9		-0.82	-1
-0.88	-0.87	-0.85	-1.00

$V^\pi$

- Bellman's equations for  $V^\pi$

$$V^\pi(s) = R(s_0) + \gamma \sum_{s'} P_{\pi(s)}(s') V^\pi(s')$$

<proof>:

$$\begin{aligned} V^\pi(s) &= E[R(s_0) + \gamma R(s_1) + \gamma^2 R(s_2) + \dots | \pi, S_0 = s] \\ &= E[R(s_0) + \gamma (\underbrace{R(s_1) + \gamma R(s_2) + \dots}_{\substack{\text{Expected future rewards} \\ = V^\pi(s_1)}}) | \pi, S_0 = s] \end{aligned}$$

immediate reward

$$\text{let } s = s_0, s' = s_1, V^\pi(s) = E[R(s) + \gamma V^\pi(s')]$$

$\because s' \sim P_{\pi(s)}$  (In state  $s$ , take action  $a = \pi(s)$ )

$$\Rightarrow V^\pi(s) = R(s_0) + \gamma \sum_{s'} P_{\pi(s)}(s') V^\pi(s')$$

Given  $\pi$ , get a linear system of equations in terms of  $V^\pi(s)$

$$\text{e.g. } V^\pi((3,1)) = R((3,1)) + \gamma (0.8V^\pi((3,2)) + 0.1V^\pi((2,1)) + 0.1V^\pi((4,1)))$$

- $V^*$  is the optimal value function

$$V^*(s) = \max_{\pi} V^\pi(s)$$

- Bellman's equations for  $V^*$ : expected future rewards

$$V^*(s) = R(s) + \max_a \gamma \sum_{s'} P_{sa}(s') V^*(s')$$

↗ if take action a

- $\pi^*$  is the optimal policy

$$\pi^*(s) = \arg\max_a \gamma \sum_{s'} P_{sa}(s') V^*(s') \quad (\gamma \text{ can be omitted})$$

strategy: ① Find  $V^* \leftarrow \text{Value Iteration}$

② Use argmax equation to find  $\pi^*$

### 3. Value Iteration

Initialize  $V(s) := 0$  for every  $s$ .

For every  $s$ , update:  $\underset{\text{new estimate of value}}{V(s)} := R(s) + \max_a \gamma \sum_s P_{sa}(s') \underset{\text{old estimate of value}}{V(s')}$  (use synchronous update)

(Bellman backup operator,  $V := B(V)$ )

Value iteration can make  $V$  converge to  $V^*$

#### 4. Policy Iteration

Initialize  $\pi$  randomly.

Repeat:  $\rightarrow V^\pi$  - linear system of equations

Set  $V := V^\pi$ , (i.e. solve Bellman's equations to get  $V^\pi$ )

Set  $\pi(s) := \arg \max_a \sum_s P_{sa}(s') V(s')$

will also converge.

## 5. Pros and Cons of value iteration and policy iteration.

policy iteration : solve  $V^\pi$ , a linear system of equations; when state space is small, then efficient; but if large, then slow; get the optimal value

value iteration : converge to  $V^*$ ; more frequent use.

## 6. What if don't know $P_{sa}$ ?

$$P_{sa}(s') = \frac{\text{\# times took action "a" in state } s \text{ and get to } s'}{\text{\# times took action "a" in state } s}, \text{ or } \frac{1}{|S|} \text{ if above is } 0$$

Putting it together:

Repeat: {

Take actions wrt.  $\pi$  to get experience in MDP.

Update estimates of  $P_{sa}$ . (and possibly  $R$ )

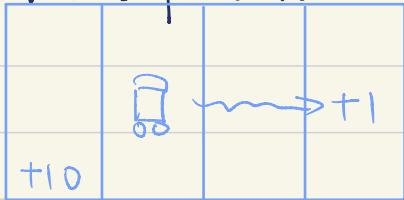
Solve Bellman's equation using value iteration to get  $V$ .

}

$$\text{Update } \pi(s) := \arg\max_a \sum_{s'} P_{sa}(s') V(s')$$

## 7. Exploration vs. Exploitation

example:



if it happens to reach +1 the first time, it may ignore +10, and thinks going to +1 is a good way. "greedy"

Modification:

Repeat: {

0.9 chance wrt  $\pi_t$   
0.1 chance randomly  $\Rightarrow \epsilon$ -greedy

Take actions wrt  $\pi_t$  to get experience in MDP.

Update estimates of  $P_{sa}$ . (and possibly  $R$ )

Solve Bellman's equation using value iteration to get  $V$ .

}

$$\text{Update } \pi_t(s) := \arg \max_a \sum_{s'} P_{sa}(s') V(s')$$