Lecture 20  | Markov Decision Processes |  CS181

Apr
2021

action $a_t$

Agent        World

state $s_{t+1}$        time $t$
reward $r_t$

Data $\quad \mathbb{D} = (s_0, a_0, r_0, s_1, a_1, r_1, \dots)$

Markov Decision Process $\qquad (S, A, r, p)$

states $\quad S = \{1, \dots, |S|\} \quad$ reward function $r(s,a)$

actions $\quad A = \{1, \dots, |A|\}$

transition model $\quad p(s'|s,a)$

next state $\nearrow \quad \uparrow \quad \nwarrow$ action
current state

policy $\quad \pi(a|s)$
prob. of taking action $a$, given in state $s$

| Markov assumption |    | Stationary assumption |

$P_t(s_{t+1}|s_1, s_2, \dots, s_t, a_1, a_2, \dots, a_t) \quad P_t(s_{t+1}|s_t, a_t) = p(s_{t+1}|s_t, a_t)$

$\quad = p_t(s_{t+1}|s_t, a_t)$

# MDP $(S, A, r, p)$

| Planning |
|---|

Input: MDP
Output: Optimal policy

| Reinforcement learning |
|---|

Input: Access to the world

Output: actions
($\rightarrow$ policy)

## 🔲 Objective

**Infinite Horizon**

$$\max_{\pi} \; \mathbb{E}_{\substack{a \sim \pi \\ s \sim p}} \left[ \sum_{t=0}^{\infty} \gamma^t r_t \right]$$

$\gamma \in [0, 1)$

$\gamma^0 3 + \gamma 3 + \gamma^2 3 + \dots$

$3 + \gamma 3 + \gamma^2 3 + \dots$

$= \frac{3}{1 - \gamma}$

$\gamma \to 1$ more patient
$\gamma \to 0$ less patient

**Finite Horizon**

$$\max_{\pi} \; \mathbb{E}_{\substack{a \sim \pi \\ s \sim p}} \left[ \sum_{t=0}^{T} r_t \right]$$

Time horizon $T$

$a_0 \sim \pi(a | s_0)$
$r_0 = r(s_0, a_0)$
$s_1 \sim p(s' | s_0, a_0)$
$a_1 \sim \pi(a | s_1)$
$r_1 = r(s_1, a_1)$
$s_2 \sim \quad \dots$

$\leftarrow$ example reward always 3

# Finite Horizon Planning : Value Iteration

Define $V^*_{(t)}(s)$ = total value from state $s$ under optimal policy with **$t$ steps to go**

## Principle of optimality

An optimal policy consists of :

① An optimal first action

② Followed by an optimal policy from the successor state

## Value Iteration

Base case
$$V^*_{(1)}(s) = \max_a r(s,a)$$

Inductive case
$$V^*_{(t+1)}(s) = \max_a \left[ r(s,a) + \sum_{s' \in S} p(s'|s,a) V^*_{(t)}(s') \right]$$

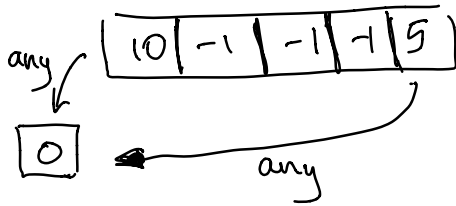$\boxed{\text{Computational complexity}}$

$$O(T \, |S| \, |A| \, L)$$

$L$ : max # states reachable from any state under any action

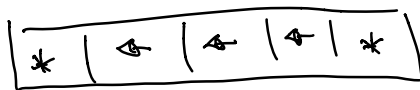$$V^*_{(t+1)}(s) = \max_a \left[ r(s,a) + \sum_{s' \in S} p(s'|s,a) \, V^*_{(t)}(s') \right]$$

$$V^*_{(1)}(s) = \max_a r(s,a)$$

## Example

any ↓

| 10 | -1 | -1 | -1 | 5 |

| 0 |

any

rewards for taking any action in a state

Actions: $\{L, R\}$

Horizon $T = 3$

| * | ← | → | ← | * |

| * |

$V^*_{(1)}$ | 0 | 10 | -1 | -1 | -1 | 5 |

$V^*_{(2)}$ | 0 | 10 | 9 | -2 | 4 | 5 |

action L: $-1 + 10 = 9$
action R: $-1 - 1 = -2$

$V^*_{(3)}$ | 0 | 10 | 9 | 8 | 4 | 5 |

action L: $-1 + 9 = 8$
R: $-1 + 4 = 3$

$V^*_{(4)}$ | 0 | 10 | 9 | 8 | 7 | 5 |

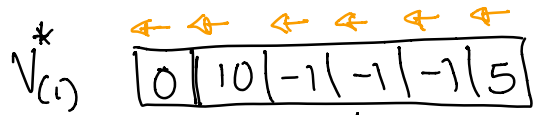action L: $-1 + 8 = 7$
R: $-1 + 5 = 4$

## Policy extraction

Optimal policy $\overline{\Pi}^*_{(t+1)}(s)$ : action to take in state $s$, $t$ steps to go
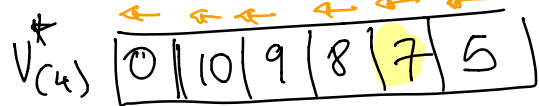
$$\arg \max_a \left[ r(s,a) + \cdots \cdots \right]$$

L = 2 in this grid world

# Infinite Time Horizon    Value Iteration

Assume    deterministic $\pi(s) \in A$

     ( w.l.o.g. )

Define    MDP value function

$$V^{\pi}(s) = \mathbb{E}_{s \sim p} \left[ \sum_{t=0}^{\infty} \gamma^t \, r(s_t, \pi(s_t)) \,\Big|\, s_0 = s \right]$$

       expected discounted value from policy

        $\pi$ in state $s$

Can decompose :

$$V^{\pi}(s) = r(s, \pi(s)) + \gamma \sum_{s' \in S} p(s' \,|\, s, \pi(s)) \, V^{\pi}(s')$$

Define Optimal policy

     $\pi^* \in \arg\max_{\pi} V^{\pi}(s)$    , for all states

Optimal value function

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

       $\longrightarrow$ | Bellman equation + VI .
              Next lecture !