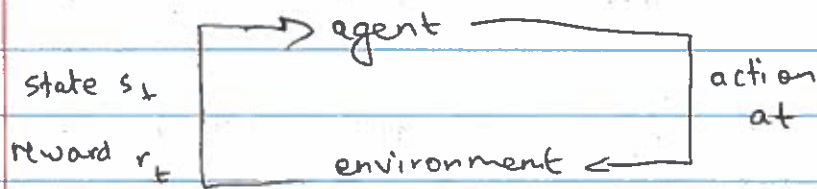


①

Reinforcement learning

07/21



## \* Challenges

- uncertainty
- explore or exploit
- delayed vs immediate feedback
- evaluative vs instructive feedback
- complex worlds, computational guarantees

Markov decision processes (MDPs)

## \* Definition

- state space  $\mathcal{S}$  with states  $s \in \mathcal{S}$
- action space  $\mathcal{A}$  with actions  $a \in \mathcal{A}$
- transition probabilities



for all state action pairs  $(s, a)$ ,

$$P(s' | s, a) = P(s_{t+1} = s' | s_t = s, a_t = a)$$

prob moving from state  $s$  to state  $s'$   
after taking action  $a$ .

#### \* Assumptions:

- time-independent

$$P(s_{t+1} = s' | s_t = s, a_t = a)$$

$$= P(s_{t+\Delta+1} = s' | s_{t+\Delta} = s, a_{t+\Delta} = a)$$

- Markov condition

$$P(s_{t+1} | s_t, a_t) = P(s_{t+1} | \overbrace{s_t, a_t, s_{t-1}, a_{t-1}, s_{t-2}, \dots}^{\text{additional history}})$$

conditional independence.

#### \* Defn (contd)

- reward function.

$R(s, s', a)$  = real-valued reward signal after taking action  $a$  in state  $s$  and moving to state  $s'$ .

- Simplifications for LSE 150

- discrete, finite state space  $\mathcal{S}$
  - discrete, finite action space  $\mathcal{A}$
  - reward function  $R(s, s', a) = R(s) = R_s$   
only depends on current state
  - bounded, deterministic rewards,  $\max_s |R_s| < \infty$
- (Note: The first two items are grouped by a bracket in the original image as 'vs continuous, infinite'.)*

Ex: board game with dice (eg. backgammon, monopoly, ...)

$\mathcal{S}$  = board position and roll of dice (right before agent decides to move)

(2)

$P(s'|s, a)$  = how state changes due to agent's move,  
opponent rolls dice, opponent's move, agent rolls dice

rewards :  $R(s) = \begin{cases} +1 & \text{win} \\ -1 & \text{lose} \\ 0 & \text{otherwise (at any other move)} \end{cases}$

$$\text{MDP} = \{S, A, P(s'|s, a), R(s)\}$$

\* Decision-making

policy : deterministic mapping from state to actions

$$\pi : S \rightarrow A$$

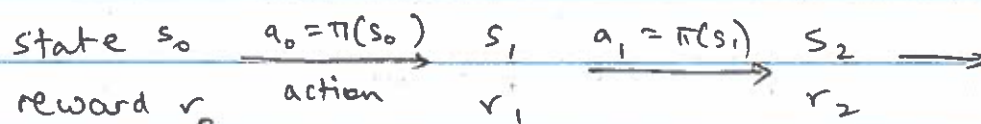
# policies =  $|A|^{|S|}$  exponential in # states

- dynamics under policy  $\pi$

$$P(s'|s, \pi(s))$$

↪ action taken by agent in state  $s$

- experience under policy  $\pi$



+ How to measure accumulated rewards over time?

long-term discounted return

discount factor  $0 \leq \gamma < 1$

$$\text{return} = \sum_{t=0}^{\infty} \gamma^t r_t$$

Possibilities

$\gamma = 0 \rightarrow$  only immediate reward at  $t=0$  matters

$\gamma < 1 \rightarrow$  near-sighted agents

$\gamma \approx 1 \rightarrow$  far-sighted agents

Intuitively :  $\gamma < 1$

- near future weighted more heavily than distant future
- confidence in MDP as model of real world may diminish with far off predictions

Also, mathematically convenient : leads to recursive algorithms whose convergence is calculated by  $\gamma$ .

State value function

$V^{\pi}(s)$  = expected discounted return following policy  $\pi$  from initial state  $s$ .

$$V^{\pi}(s) = E^{\pi} \left[ \sum_{t=0}^{\infty} \gamma^t R(s_t) \mid s_0 = s \right]$$

↑  
expectation operator  
(Computing mean)



(3)

Relating value function in different states :

$$V^{\pi}(s) = E^{\pi} [R(s_0) + \gamma R(s_1) + \gamma^2 R(s_2) + \dots | s_0 = s]$$

$$= R(s) + \gamma E^{\pi} [R(s_1) + \gamma R(s_2) + \gamma^2 R(s_3) + \dots | s_0 = s]$$

$$= R(s) + \gamma \sum_{s'=1}^n P(s'|s, \pi(s)) E^{\pi} [R(s_1) + \gamma R(s_2) + \dots | s = s']$$

$$V^{\pi}(s) = R(s) + \gamma \sum_{s'=1}^n P(s'|s, \pi(s)) V^{\pi}(s')$$

Bellman  
equation

### Action-value function

$Q^{\pi}(s, a)$  = expected return from initial state  $s$ , taking action  $a$ , then following policy  $\pi$ .

$$= E^{\pi} \left[ \sum_{t=0}^{\infty} \gamma^t R(s_t) | s_0 = s, a_0 = a \right]$$

$\rightarrow a$  may equal  $\pi(s)$ ,  
or not.

$$Q^{\pi}(s, a) = R(s) + \gamma \sum_{s'} P(s'|s, a) V^{\pi}(s')$$

### Optimality in MDPs

Thm : there is always (at least) one optimal policy  $\pi^*$  for which  $V^{\pi^*}(s) \geq V^{\pi}(s)$  for all policies  $\pi$  and states  $s \in \mathcal{S}$ .

Goal: how to compute optimal policy  $\pi^*$ ?

\* Optimal state-value function

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

\* Optimal action-value function

$$Q^*(s, a) = Q^{\pi^*}(s, a)$$

$$\text{Note: } V^*(s) = \max_a Q^*(s, a)$$

There may be multiple optimal policies, but optimal value functions are unique.

\* Relations - given MDP, how to recover  $\pi^*$  from  $V^*$ ,  $Q^*$ ?

$$\pi^*(s) = \operatorname{argmax}_a [Q^*(s, a)]$$

$$= \operatorname{argmax}_a \left[ R(s) + \gamma \sum_{s'} P(s'|s, a) V^*(s') \right]$$

$$= \operatorname{argmax}_a \left[ \sum_{s'} R(s'|s, a) V^*(s') \right]$$

(4)

Planning

Assume complete world of environment as

MDP =  $\{S, A, P(s'|s, a), R(s)\}$ , also  $\gamma < 1$

how to compute  $\pi^*(s)$ , or equivalently  $V^*(s)$  or  $Q^*(s, a)$

i) Policy evaluation

How to compute  $V^\pi(s)$ ?

From Bellman eqn:

$$V^\pi(s) = R(s) + \gamma \sum_{s'=1}^N P(s'|s, \pi(s)) V^\pi(s')$$

for  $s = 1, 2, \dots, N$  where  $N = \# \text{ states of MDP}$

This is a system of linear equations:  $N$  equations,  
 $N$  unknowns  $V^\pi(s)$ ,  $s = 1, \dots, N$

~~Put~~ Put all unknowns on left side:

$$V^\pi(s) - \gamma \sum_{s'=1}^N \{P(s'|s, \pi(s)) V^\pi(s')\} = R(s)$$

$$\sum_{s'=1}^N \{[I(s, s') - \gamma P(s'|s, \pi(s))] V^\pi(s')\} = R(s)$$

↑  
indicator  
function

Can rewrite this:

$$(I - \gamma P) V = R$$

Diagram annotations:

- $I$ :  $n \times n$  identity matrix
- $P$ :  $n \times n$  matrix
- $R$ :  $n \times 1$  vector of knowns
- $V$ :  $n \times 1$  vector of unknowns

Ex: states  $s \in \{0, 1\}$   
transitions  $P(s' | s, \pi(s))$   
rewards  $R(s) = \begin{pmatrix} r_0 \\ r_1 \end{pmatrix}$

state value function  $V^\pi(s) = \begin{pmatrix} v_0 \\ v_1 \end{pmatrix}$

Solve:

$$\left[ \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} - \gamma \begin{pmatrix} p_{00} & p_{01} \\ p_{10} & p_{11} \end{pmatrix} \right] \begin{pmatrix} v_0 \\ v_1 \end{pmatrix} = \begin{pmatrix} r_0 \\ r_1 \end{pmatrix}$$

2 eqns  
2 unknowns

## 2) Policy improvement

\* How to compute  $\pi'$  such that  $V^{\pi'}(s) \geq V^\pi(s)$  for all states  $s$ ?

\* Recall  $Q^\pi(s, a)$  = expected return from state  $s$ , follow action  $a$ , then revert to policy  $\pi$ .

How to compute  $Q^\pi(s, a)$ ?

Evaluate policy (step #1) to compute  $V^\pi(s)$ .

$$\text{Then: } Q^\pi(s, a) = R(s) + \gamma \sum_{s'} P(s' | s, a) V^\pi(s')$$



\* Define "greedy" policy.

$$\begin{aligned}\pi'(s) &= \text{greedy}[\pi] \\ &= \underset{a}{\operatorname{argmax}} [Q^\pi(s, a)] \\ &= \underset{a}{\operatorname{argmax}} \left[ \sum_{s'} P(s'|s, a) V^\pi(s') \right]\end{aligned}$$

Thm greedy policy  $\pi'$  everywhere performs better or equal to original policy  $\pi$

$$V^{\pi'}(s) \geq V^\pi(s) \text{ for all } s$$

Intuition: if better to choose action  $a$  over  $\pi(s)$  in state  $s$  before following policy  $\pi$ , it's always better to choose action  $a$  in state  $s$  (for all times)

Proof:

$$\begin{aligned}V^\pi(s) &= Q^\pi(s, \pi(s)) \\ &\leq \max_a Q^\pi(s, a) \\ &= Q^\pi(s, \pi'(s)) \\ &= R(s) + \gamma \sum_{s'} P(s'|s, \pi'(s)) V^\pi(s')\end{aligned}$$

So far: better to (take one step under  $\pi'$ , then resort to  $\pi$ ), than to (always ~~res~~ follow  $\pi$ )

"One-step" inequality.

$$V^\pi(s) \leq R(s) + \gamma \sum_{s'} P(s'|s, \pi'(s)) V^\pi(s')$$

Apply one-step inequality to itself on the right side:

$$V^{\pi}(s) \leq R(s) + \gamma \sum_{s'} P(s' | s, \pi'(s)) \left[ R(s') + \gamma \sum_{s''} P(s'' | s', \pi'(s')) V^{\pi}(s'') \right]$$

Better to take two steps under  $\pi'$ , then revert to  $\pi$ , than to always follow policy  $\pi$ .

Apply this one-step inequality to itself  $t$ -times to get  $t$ -step inequality: better to take  $t$  steps under  $\pi'$  (then revert to  $\pi$ ), than to always follow  $\pi$ .

Let  $t \rightarrow \infty$ : always better to follow  $\pi'$  than  $\pi$ .

$V^{\pi}(s) \leq V^{\pi'}(s)$  for all states  $s$  since right side converges for  $\gamma < 1$ .