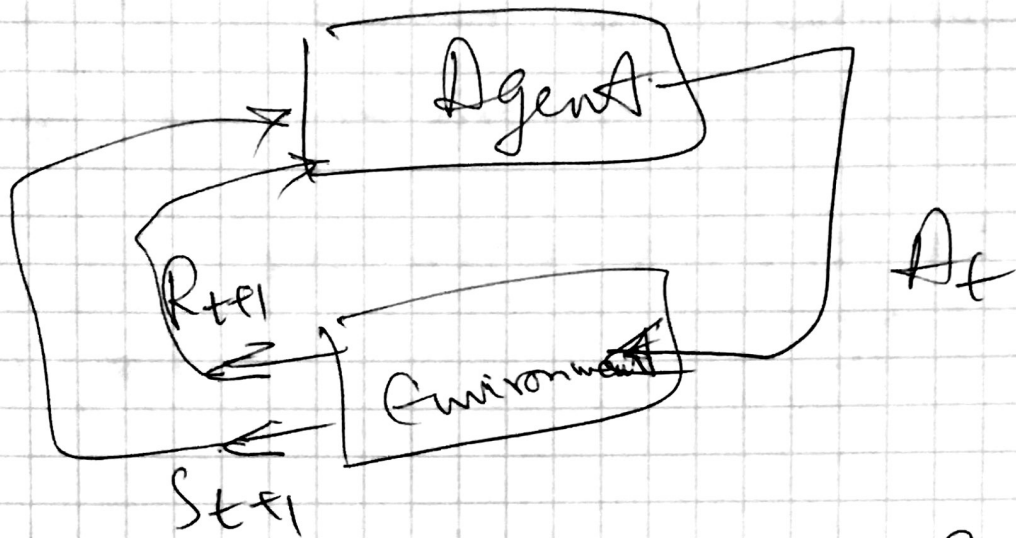


# Ch-3 Notes

(1)



$S_0, A_0, R_1, S_1, A_1, S_2, A_2, R_3, \dots$

Finite MDP: sets of states, actions and rewards all have a finite number of elements

$$p(s', r | s, a) \equiv P\{S_t = s', R_t = r | S_{t-1} = s, A_{t-1} = a\}$$

↳ dynamics, model, plant (control) people.

$$P: S \times R \times S \times A \rightarrow [0, 1]$$

we use  $R_{t+1}$  instead of  $R_t$  to denote reward because it emphasizes that the

next reward and next state  $(R_{t+1}, S_{t+1})$  ⑤  
are jointly determined. Both conventions are  
widely used in literature.

$$\sum_{S' \in S} \sum_{r \in R} p(S', r | S, a) = 1, \text{ for all } S \in S, a \in A(S)$$

\* In an MDP once you know  $p$ ,  
you can completely characterize the whole  
MDP. Where you go next is  
completely determined by where you are  
and actions you take. (Markov property)

\* From the 4-argument dynamics function  $p$   
one can compute state-transition probabilities  
 $p: S \times S \times A \rightarrow [0, 1]$

$$p(S' | S, a) \stackrel{\circ}{=} \Pr(S_t = S' | S_{t-1} = S, A_{t-1} = a) \\ = \sum_{r \in R} p(S', r | S, a)$$

Expected rewards for  $s$ - $a$  pairs. (2)  
 $r: S \times A \rightarrow R.$

$$r(s, a) \doteq E [R_t | S_{t-1} = s, A_{t-1} = a]$$

Want to maximize return. (expected)

$$G_t \doteq R_{t+1} + R_{t+2} + R_{t+3} + \dots + R_T.$$

$T$  = final time step.

Episodic: when agent-environment interaction naturally breaks down into subsequences.

Usually we have some

terminal state.

$T$  can vary from episode to episode.

(4)

\* Continuing tasks:

No natural break down.

Lifelong learning.

This is usually problematic mathematically.  
as  $T = \infty$ , and return could itself  
easily be  $\infty$ .

So we introduce discounting:-

$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots$$
$$= \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$$

$$\textcircled{0} \quad 0 \leq \gamma \leq 1$$

discount rate.

$\{R_k\}$  should be bounded  
as well.



# \* Value functions

⑤

Value of a state  $V(s)$

$\hat{=}$  how good is it for the agent to be in that state.

"how good" = future rewards that can be expected.

But future rewards are consequences of actions taken in the future as well.

$\rightarrow$  Policy  $\rightarrow$  way of acting.

$\rightarrow$  mapping from  $s$  to  $a$ .

$\pi(s) \rightarrow a$ . (Deterministic)

$\pi(a|s) = \text{prob. of taking } a \text{ when in } s$ .

(Stochastic)

$$V_{\pi}(s) \hat{=} E_{\pi} \left[ \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid s_t = s \right]$$

for all  $s \in \mathcal{S}$ .

⑥

$$\underline{q_{\pi}(s, a)} \doteq E_{\pi} \left[ \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid S_t = s, A_t = a \right]$$

↑  
action-value function for  $\pi$ .

$$V_{\pi}(s) \doteq E_{\pi} [G_t \mid S_t = s]$$

$$= E_{\pi} [R_{t+1} + \gamma G_{t+1} \mid S_t = s]$$

$$= \sum_a \pi(a|s) \sum_{s', r} p(s', r|s, a)$$

$$\cdot \left[ r + \gamma E_{\pi} [G_{t+1} \mid S_{t+1} = s'] \right]$$

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

for all  $s \in S$ .

(9)

\* Solving RL means finding a good policy.  
best

Good policy  $\Rightarrow$  good cumulative returns.

$\pi \succ \pi'$  iff  $V_\pi(s) \geq V_{\pi'}(s), \forall s \in S$ .

$\pi_*$   $\hat{=}$  set of all optimal policies  
(there may be more than 1)

$$V_*(s) \hat{=} \max_{\pi} V_\pi(s), \forall s \in S.$$

$$Q_*(s, a) \hat{=} \max_{\pi} Q_\pi(s, a)$$

$$Q_*(s, a) = \mathbb{E} \left[ r(s, a) + \gamma V_*(s') \right]$$

$a \sim \pi^*(a|s)$   
 $s' \sim p(s'|s, a)$

⑧

$$V_{\pi^*}(s) \geq \max_{a \in A(s)} v_{\pi^*}(s, a)$$

$$\geq \max_a E_{\pi^*} [G_t | S_t = s, A_t = a]$$

$$\geq \max_a E_{\pi^*} [R_{t+1} + \gamma G_{t+1} | S_t = s, A_t = a]$$

$$\geq \max_a E_{\pi^*} [R_{t+1} + \gamma V_{\pi^*}(S_{t+1}) |$$

$$S_t = s, A_t = a]$$

Bellman  
Optimality  
Condition

for Value  
Functions.

$$\geq \max_a \sum_{s', r} p(s', r | s, a) [r + \gamma V_{\pi^*}(s)]$$

$$Q_{\pi^*}(s, a) = E [R_{t+1} + \gamma \max_{a'} v_{\pi^*}(S_{t+1}, a') |$$

$$S_t = s, A_t = a]$$

$$\geq \sum_{s', r} p(s', r | s, a) [r + \gamma \max_{a'} v_{\pi^*}(s', a')]$$



\* For finite MDPs, Bellman optimality equation has a unique solution. (9)

System of ' $n$ ' equations in ' $n$ ' unknowns.

(If dynamics  $P$  of the environment are known)

\* If you have  $Q^*$  just take 1 step.  
greedy wrt it. You get optimal policy!

\* If you have  $Q^*$  life is very easy.  
don't even have to do 1-step search

Just pick  $\underset{a}{\operatorname{argmax}} Q^*(s, a)$ .

The  $q$ -function caches the results of all 1-step queries. more memory but allows optimal actions to be selected without having to know dynamics. (possible next states and their values.)

\* Other ways of writing  $V, Q$  (more common) (10)

---

$$* V_{\pi}(s) = r(s, \pi(s)) + \gamma E_{p(s'|s, \pi(s))} [V(s')]$$

$$* Q_{\pi}(s, a) = r(s, a) + \gamma E_{\substack{p(s'|s, a) \\ p(s'|s, a)}} \left[ \max_{a'} Q_{\pi}(s', a') \right]$$

$$= r(s, a) + \gamma E_{p(s'|s, a)} \left[ \max_{a'} Q_{\pi}(s', a') \right]$$

$$* Q(s, a) \leq V(s)$$

$$* V(s) = \max_a Q(s, a)$$

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

$$\pi = \operatorname{argmax}_a Q(s, a)$$