

single-  
agent  
RL

Markov Decision Process (MDP) ✖

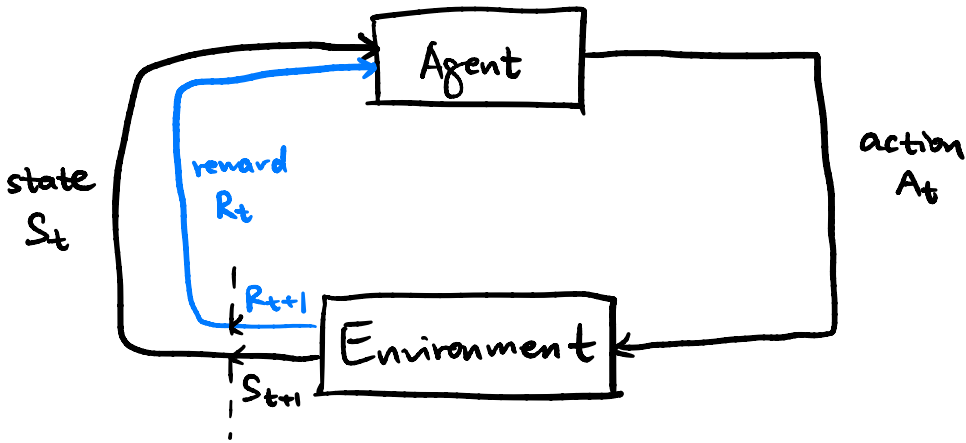
Dynamic Programming (DP)

Monte Carlo (MC)

Temporal Difference and Q-learning (TD) ✖

Policy Gradient

# Theoretical framework: MDP



logic: environment has state, agent observes state and make an action, he receives reward based on his state-action pair. However, his action also changes the state so at the next time step his new reward is based on his new state etc.

The agent's objective: maximize aggregate reward at all time steps.

★  
↓ Main Difficulty in RL:

**Greedy strategy fails!** If too shortsighted, only choose the action that brings highest current reward, might fall into a very bad state and get very low future rewards!

e.g.: Consider bandit problem, RL problem with no state transition, agent just continuously make actions and receive reward (same as bandit in the casino).

At time  $t$ , make action  $A_t$ , assume reward

$$R_t | A_t = a \sim \mathcal{N}(\underbrace{q_*(a)}_{\substack{\text{mean reward} \\ \text{for action } a \\ \text{unknown to agent}}}, 1)$$

assume 4 bandits,  $A_t$  can take value  $\{1, 2, 3, 4\}$  (which to select)



$$q_*(1) = 0$$



$$q_*(2) = 3$$



$$q_*(3) = 2$$



$$q_*(4) = 10$$

obviously, optimal strategy is to always select bandit 4 but if we use completely greedy strategy,

$$A_0 = 1 \Rightarrow R_0 = -0.5 \quad (\text{realization of } \mathcal{N}(0, 1))$$

$$A_1 = 2 \Rightarrow R_1 = 3.1 \quad (\text{from } \mathcal{N}(3, 1))$$

we will think that choosing bandit 2 is a lot better, so we stick to bandit 2 and miss bandit 4.

even with no state transition, greedy is not a good strategy (ε-greedy instead)

trade-off { exploitation (maximize the reward)  
exploration (know about what happens for other actions)

{ If always exploit, might miss a better action  
If always explore, might have low total reward



motivation of a lot of  
algorithms and concepts

# Basic Setting of MDP:

State  $s \in S$ , action  $a \in A(s)$ , reward  $r \in R \subseteq \mathbb{R}$   
 $\uparrow$   
the set of available actions might depend on the state

Infinite time-horizon

Finite MDP for simplicity:  $|S| < \infty$ ,  $\forall s \in S, |A(s)| < \infty$ ,

denote  $A = \bigcup_{s \in S} A(s)$  is still finite  $|R| < \infty$

(the set of all possible actions regardless of the state)

dynamics

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

(time-homogeneous)

given the action  $\{A_t\}$ ,  $\{S_t\}$  is Markov, transition from  $S_{t-1}$  to  $S_t$  only depends on the value of  $A_{t-1}$ .  
(given  $S_t, A_t, \text{future } R_{t+1}, S_{t+1}, A_{t+1}, \dots$  past  $S_{t-1}, A_{t-1}, R_{t-1}, \dots$ )

seeing current state  $s$ , making action  $a$ , the prob that next state is  $s'$  and get immediate reward  $r$

Calculations:

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

$$= \sum_r p(s', r | s, a)$$

state-transition prob

$$r(s, a) \triangleq \mathbb{E}(R_t | S_{t-1} = s, A_{t-1} = a)$$

expected reward for  
state-action pair

$$= \sum_r r \cdot p(r | s, a)$$

$$= \sum_r r \cdot \sum_{s'} p(s', r | s, a)$$

$$r(s, a, s') \triangleq \mathbb{E}(R_t | S_{t-1} = s, A_{t-1} = a, S_t = s')$$

expected reward for  
state-action-next state  
triple

$$= \sum_r r \cdot p(r | s, a, s')$$

$$= \sum_r r \cdot \frac{p(s', r | s, a)}{p(s' | s, a)} = \frac{\sum_r r \cdot p(s', r | s, a)}{\sum_r p(s', r | s, a)}$$

so dynamics provides expression for everything  
we are interested in.

Modelling the environment

# Modelling the agent

Agent obj: maximize the sum of rewards, in infinite horizon setting there's convergence problems so use discounting

$$G_t \triangleq \sum_{k=0}^{\infty} \gamma^k \cdot R_{t+k+1}$$

↓  
discount rate

return at time  $t$   
(aggregating all rewards  $\geq$  time  $t$ )

notice that  $R_{t+1}$  is the reward received based on  $S_t, A_t$ .

Always write  $G_t = R_{t+1} + \gamma \cdot G_{t+1}$ .

Agent obj: maximize  $\mathbb{E}G_0$

The agent makes decision based on a policy

$$\pi(a|s) = \mathbb{P}(A_t = a | S_t = s) \quad (\text{time-homogeneous})$$

a conditional distribution on  $A$ , meaning on seeing current state  $s$ , react with action  $a$  with some probability.

Why not always deterministically?  
Because of exploration!

At this point, it's possible to compute

$$\begin{aligned} E(R_{t+1} | S_t) &= \sum_r r \cdot P(R_{t+1} = r | S_t) \quad \text{policy connects state and action} \\ &= \sum_r r \cdot \sum_a P(A_t = a | S_t) \cdot P(R_{t+1} = r | S_t, A_t = a) \\ &= \sum_r r \cdot \sum_a \pi(a | S_t) \cdot \sum_{s'} p(s', r | S_t, a) \end{aligned}$$

Value Function:

$$V_\pi(s) \triangleq E_\pi(G_t | S_t = s)$$

state value func  
for policy  $\pi$

$$q_\pi(s, a) \triangleq E_\pi(G_t | S_t = s, A_t = a)$$

state-action value func  
for policy  $\pi$

connection of  $V_\pi, q_\pi$ ?

$$\begin{aligned} \text{Obviously, } V_\pi(s) &= \sum_a P(A_t = a | S_t = s) \cdot E_\pi(G_t | S_t = s, A_t = a) \\ &= \sum_a \pi(a | s) \cdot q_\pi(s, a) \end{aligned}$$

On the other hand,

$$\begin{aligned} q_\pi(s, a) &= E_\pi(R_{t+1} + \gamma \cdot G_{t+1} | S_t = s, A_t = a) \\ &= \underbrace{r(s, a)}_{\text{calculated!}} + \gamma \cdot E_\pi(G_{t+1} | S_t = s, A_t = a) \end{aligned}$$



$$= \sum_r r \cdot \sum_{s'} p(s', r | s, a) + \gamma \cdot \sum_{s'} \underbrace{\mathbb{E}_{\pi} \left( \overset{R_{t+2}, R_{t+3}, \dots}{G_{t+1} | S_{t+1}=s', S_t=s, A_t=a} \right)}_{\substack{\text{Markov} \\ \mathbb{E}_{\pi}(G_{t+1} | S_{t+1}=s') \\ \parallel \\ V_{\pi}(s')}} \cdot \mathbb{P}(S_{t+1}=s' | S_t=s, A_t=a)$$

$$= \sum_r r \cdot \sum_{s'} p(s', r | s, a) + \gamma \cdot \sum_{s'} V_{\pi}(s') \cdot \sum_r p(s', r | s, a)$$

$$= \sum_{r, s'} (r + \gamma \cdot V_{\pi}(s')) \cdot p(s', r | s, a)$$

MDP time-homogeneous Markov + policy time-homogeneous  
+ infinite time horizon

↓

$Q_{\pi}, V_{\pi}$  time-homogeneous (not depend on  $t$ )

that's why consider inf time horizon  
(recall stochastic control)

# Bellman Consistency Equation:

Describe what condition  $V_\pi, q_\pi$  has to satisfy,

$\star$ : consider  $S_t \rightarrow S_{t+1}$  then  
can write back as  
value func!

$$\begin{aligned} V_\pi(s) &= \mathbb{E}(R_{t+1} | S_t = s) + \gamma \cdot \mathbb{E}_\pi(G_{t+1} | S_t = s) \\ \text{future reward} \uparrow & \\ \text{following} & \\ \text{policy } \pi \text{ with} & \\ \text{current state } s & \\ &= \sum_r r \cdot \mathbb{P}(R_{t+1} = r | S_t = s) + \gamma \cdot \sum_{s', a} p(s', a | s) \cdot \mathbb{E}_\pi(G_{t+1} | S_t = s, A_t = a, S_{t+1} = s') \\ & \quad \text{consider action} \quad \text{Markov!} \\ &= \sum_r r \cdot \sum_a \mathbb{P}(A_t = a | S_t = s) \cdot p(r | s, a) + \\ & \quad \gamma \cdot \sum_{s', a} p(s', a | s) \cdot \mathbb{E}_\pi(G_{t+1} | S_{t+1} = s') \\ & \quad \quad \quad = \pi(a | s) \cdot p(s' | a, s) \\ &= \sum_r r \cdot \sum_a \pi(a | s) \sum_{s'} p(s', r | s, a) + \\ & \quad \gamma \cdot \sum_{s', a} \pi(a | s) \cdot \sum_r p(s', r | s, a) \cdot V_\pi(s') \\ &= \underbrace{\sum_a \pi(a | s)}_{\text{average w.r.t. policy based on current state}} \underbrace{\sum_{s', r} p(s', r | s, a)}_{\text{average w.r.t. dynamics}} \cdot \underbrace{[r + \gamma V_\pi(s')]}_{\substack{\text{immediate} \\ \text{reward}} \quad \text{discounted future reward}} \end{aligned}$$

Remark: One naturally thinks about if it's possible to solve out  $V_\pi$  from Bellman consistency equation, actually one can prove that solution  $\exists$  and is unique. (Exercise)

However, one would need a given policy  $\pi$  and the dynamics  $p$ , so it depends on knowledge of the model.

Bellman consistency equation for  $q_\pi$ :

$$\begin{aligned}
 q_\pi(s, a) &= \mathbb{E}_\pi(G_t | S_t = s, A_t = a) \\
 &= \mathbb{E}(R_{t+1} | S_t = s, A_t = a) + \gamma \cdot \mathbb{E}_\pi(G_{t+1} | S_t = s, A_t = a) \\
 &= \sum_r r \cdot p(r | s, a) + \gamma \cdot \sum_{s', a'} p(S_{t+1} = s', A_{t+1} = a' | S_t = s, A_t = a) \\
 &\quad \cdot \underbrace{\mathbb{E}_\pi(G_{t+1} | S_t = s, A_t = a, S_{t+1} = s', A_{t+1} = a')}_{q_\pi(s', a')} \\
 &= \sum_r r \cdot \sum_{s'} p(s', r | s, a) + \gamma \sum_{s', a'} p(s' | s, a) \cdot \pi(a' | s') \cdot q_\pi(s', a') \\
 &= \underbrace{\sum_{r, s'} p(s', r | s, a)}_{\text{average w.r.t. dynamics}} \left[ \underbrace{r}_{\text{immediate reward}} + \gamma \sum_{a'} \underbrace{\pi(a' | s')}_{\text{discounted future reward}} q_\pi(s', a') \right]
 \end{aligned}$$

# Optimal Value Function

Since the obj is to maximize  $IEG_0$ , if a policy has higher state value  $V_\pi(s)$  for  $\forall s \in S$ , it's a better policy. Does there exist the optimal policy?

$$\begin{cases} \forall s \in S, V_*(s) \triangleq \sup_{\pi} V_\pi(s) \\ \forall s \in S, a \in A, q_*(s) \triangleq \sup_{\pi} q_\pi(s, a) \end{cases}$$

optimal value function

as pointwise sup

Thm:  $\exists$  deterministic policy  $\pi_*$  (not necessarily unique) s.t.

$$\forall s \in S, \forall a \in A, V_{\pi_*}(s) = V_*(s), q_{\pi_*}(s, a) = q_*(s, a)$$

value of optimal policy  $\rightarrow$   $V_{\pi_*}(s)$   
optimal value  $\rightarrow$   $V_*(s)$   
optimal value  $\rightarrow$   $q_{\pi_*}(s, a)$   
optimal value  $\rightarrow$   $q_*(s, a)$

Surprisingly, pointwise sup of value function is just local, for each  $s \in S$ , sup may be approx by different  $\pi$  the value function of optimal policy!  
global optimal

Pf:

$$\pi_*(s) \triangleq \arg \sup_a \mathbb{E} \left[ R_{t+1} + \gamma \cdot V_*(S_{t+1}) \mid S_t = s, A_t = a \right]$$

means take this action w.p. 1 (deterministic policy)

whenever at state  $s$ , put all prob mass on the action that maximizes expectn of sum of immediate and discounted future reward.

$\forall s \in S, V_{\pi^*}(s) \leq V_*(s)$  by def of  $V_*$ .

On the other hand,  $\forall s \in S,$

$$V_*(s) = \sup_{\pi} \left\{ \mathbb{E}_{\pi} [R_{t+1} + \gamma G_{t+1} | S_t = s] \right\}$$

lower property

$$= \sup_{\pi} \left\{ \mathbb{E}_{\pi} \left[ R_{t+1} + \gamma \cdot \mathbb{E}_{\pi} (G_{t+1} | S_t = s, A_t, S_{t+1}) \right] \right\}$$

change to  $\pi'$

$$\leq \sup_{\pi} \left\{ \mathbb{E}_{\pi} \left[ R_{t+1} + \gamma \cdot \underbrace{\sup_{\pi'} \mathbb{E}_{\pi'} (G_{t+1} | S_t = s, A_t, S_{t+1})}_{V_*(S_{t+1})} \right] \right\}$$

Markov  $\Pi$

$$= \sup_{\pi} \left\{ \mathbb{E}_{\pi} \left[ R_{t+1} + \gamma \cdot V_*(S_{t+1}) | S_t = s \right] \right\}$$

$$= \sup_{\pi} \left\{ \mathbb{E} \left[ R_{t+1} + \gamma \cdot V_*(S_{t+1}) | S_t = s, A_t \sim \pi(\cdot | s) \right] \right\}$$

def of  $\tilde{\pi}$

$$\leq \mathbb{E}_{\pi^*} \left[ R_{t+1} + \gamma \cdot V_*(S_{t+1}) | S_t = s \right]$$

↑  
policy appears here

$$\Downarrow$$

$$V_*(S_t) \leq \mathbb{E}_{\pi^*} \left[ R_{t+1} + \gamma \cdot V_*(S_{t+1}) | S_t \right]$$

iteratively

$$V_*(S_t) \leq \mathbb{E}_{\pi^*} \left( R_{t+1} + \gamma \cdot \underbrace{\mathbb{E}_{\pi^*} [R_{t+2} + \gamma \cdot V_*(S_{t+2}) | S_{t+1}]}_{\text{upper bound of } V_*(S_{t+1})} \right)$$

$$= \mathbb{E}_{\pi^*} \left[ R_{t+1} + \gamma R_{t+2} + \gamma^2 V_*(S_{t+2}) | S_t \right]$$

$$S_0 \quad V_*(S_t) \leq \dots \leq \mathbb{E}_{\pi_*} [R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots | S_t] \\ \leq V_{\pi_*}(S_t)$$

proves  $\forall s \in S, V_*(s) = V_{\pi_*}(s)$ .

and  $\pi_*$  is deterministic optimal policy.

For  $q$ : recall  $q_{\pi}(s, a) = \sum_{r, s'} [r + \gamma V_{\pi}(s')] \cdot p(s', r | s, a)$   
 (connection of  $q_{\pi}, V_{\pi}$ )

take sup w.r.t.  $\pi$  on both sides PW

$$q_*(s, a) \leq \sum_{r, s'} [r + \gamma V_*(s')] \cdot p(s', r | s, a) \\ \text{(sup goes inside sum)} \quad \parallel \\ \sum_{r, s'} [r + \gamma \cdot V_{\pi_*}(s')] \cdot p(s', r | s, a) \\ \parallel \\ q_{\pi_*}(s, a)$$

proves that it's also optimal w.r.t.  $q$ .

## Bellman Optimality Equation:

Plug in  $\pi = \pi_*$  in Bellman consistency equation

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

$$V_*(s) = \sum_a \pi_*(a|s) \sum_{s', r} p(s', r|s, a) \cdot [r + \gamma \cdot V_*(s')]$$

$$q_{\pi}(s, a) = \sum_{s', r} p(s', r|s, a) [r + \gamma \cdot \sum_{a'} \pi(a'|s') \cdot q_{\pi}(s', a')] \\ \Downarrow$$

$$q_*(s, a) = \sum_{s', r} p(s', r|s, a) [r + \gamma \cdot \sum_{a'} \pi_*(a'|s') \cdot q_*(s', a')]$$

☆:

but it's not a good form since we don't know  $\pi_*$  in prior!

To get rid of  $\pi_*$ , discover relationship between  $V_*$ ,  $q_*$

$$V_*(s) = \sup_{\pi} V_{\pi}(s) = \sup_{\pi} \sum_a \pi(a|s) \cdot q_{\pi}(s, a)$$

$$\leq \sup_{\pi} \underbrace{\sum_a \pi(a|s)}_{\text{weighted average}} \cdot q_*(s, a)$$

$$= \max_a q_*(s, a)$$

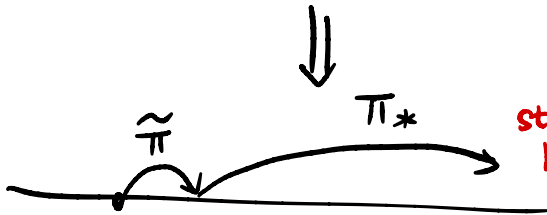
but  $V_*(s) = V_{\pi_*(s)} < \max_a q_*(s, a)$  by contradiction

consider  $\tilde{\pi}(s) \triangleq \operatorname{argmax}_a q_*(s, a)$ ,

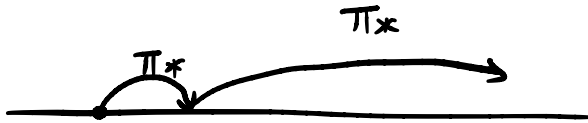
$$q_{\tilde{\pi}_*}(s, \tilde{\pi}(s)) = \max_a q_*(s, a) > V_{\pi_*(s)}$$

||

$$q_{\pi_*(s)}(s, \pi_*(s))$$



strictly better, contradiction with  $\pi_*$  being optimal policy!



$$\tilde{\pi} = \pi_*$$

⇓

$$V_*(s) = \max_a q_*(s, a)$$

Policy Improvement Theorem  
(proof not mentioned)

☆: compact form of Bellman optimality equation



Useful Bellman Optimality Equation:

$$V_*(s) = \max_a q_*(s, a)$$

(HJB type)

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

recall

$$q_{\pi}(s, a) = \sum_{r, s'} (r + \gamma \cdot V_{\pi}(s)) \cdot p(s', r | s, a)$$

$$q_*(s, a) = \sum_{r, s'} [r + \gamma \cdot V_*(s')] \cdot p(s', r | s, a)$$

$$= \sum_{r, s'} [r + \gamma \cdot \max_{a'} q_*(s', a')] \cdot p(s', r | s, a)$$

$$\pi_*(s) = \operatorname{argmax}_a q_*(s, a) \text{ is optimal policy}$$
$$= \operatorname{argmax}_a \sum_{r, s'} [r + \gamma \cdot V_*(s')] \cdot p(s', r | s, a)$$

Correspondence:

①: Stochastic control (cts time, finite horizon, non-randomized policy)

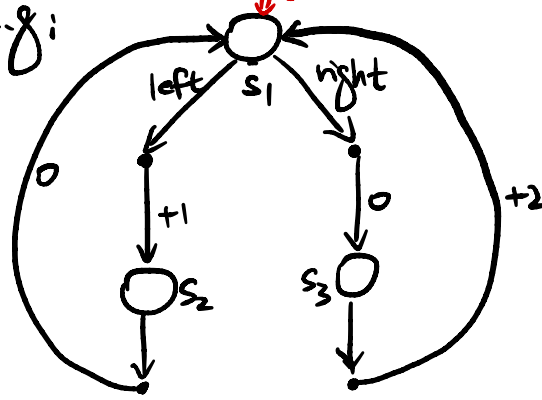
②: Statistical decision theory  
(no time evolution, recall the way deriving Bayes estimator, similar to optimal policy)

information theory  
methods for optimality  
in RL

Solve control  
with RL

only decision to make here

e.g:



rewards shown are deterministic  
two deter ~ policy  $\pi_{left}, \pi_{right}$ ,  
which is optimal?

df:  $S = \{s_1, s_2, s_3\}$   $A = \{left, right\}$ ,

Bellman optimality equation:

$$\begin{cases} V^*(s_2) = \gamma V^*(s_1) \\ V^*(s_3) = 2 + \gamma V^*(s_1) \\ V^*(s_1) = \max\{1 + \gamma V^*(s_2), \gamma V^*(s_3)\} \end{cases}$$

$$\Downarrow$$

$$V^*(s_1) = \max\{1 + \gamma^2 V^*(s_1), 2\gamma + \gamma^2 V^*(s_3)\}$$

$$\Downarrow$$

$$\begin{cases} \gamma \in [0, \frac{1}{2}), & \pi_l \text{ optimal} \\ \gamma = \frac{1}{2}, & \pi_l, \pi_r \text{ optimal} \\ \gamma \in (\frac{1}{2}, 1), & \pi_r \text{ optimal} \end{cases}$$

(intuitively match!)

# DP: (model-based)

**Policy iteration:** current policy  $\pi_k$ , ① do policy evaluation

(get  $V_{\pi_k}$  or  $Q_{\pi_k}$  through Bellman consistency equation as fixed point iteration),

② policy improvement get better  $\pi_{k+1}$

repeat until convergence

$$\pi_{k+1}(s) = \underset{a}{\operatorname{argmax}} Q_{\pi_k}(s, a)$$

**Value iteration:** directly solve  $V^*$  or  $Q^*$  through fixed point iteration of Bellman consistency equation and construct  $\pi^*$

Check:  $\mathcal{J}v(s) = \max_a \mathbb{E}[R_{t+1} + \gamma \cdot v(S_{t+1}) | S_t = s, A_t = a]$   
is contraction mapping

$$\exists 0 \leq k < 1, \forall v, v', \|\mathcal{J}v - \mathcal{J}v'\|_{\infty} \leq k \|v - v'\|_{\infty}$$

↑ actually the discount rate  $\gamma$

$\mathcal{J}$  is Bellman optimality operator

pros: easy to implement, efficient, fit with small problems

cons: model-based, can't deal with cts state/action space,

# MC

Natural since  $v, q$  are conditional expectation, policy evaluation done by MC.

## First-visit MC ES:

- ①: Exploring start, any  $(s, a)$  pos probability of selected as initial state (maintain exploration!)
- ②: For fixed policy  $\pi$ , experience:  $S_0, A_0, R_1, \dots$   
Calculate returns at each time, find out the time of first visit to  $(s, a)$  and add the return at this time of first visit into  $\text{list}(s, a)$   
ensure i.i.d. MC samples (can be every-visit also)  
(maintain a list for each state-action pair)
- ③: Update estimate for  $q_\pi(s, a)$  as sample average of all numbers in  $\text{list}(s, a)$
- ④: Construct greedy deter $\sim$  policy (after enough experience gained)  
 $\pi'(s) = \underset{a}{\operatorname{argmax}} q_\pi(s, a)$  as policy improvement
- ⑤: Iterate until  $\pi \rightarrow \pi_*$ ,  $q \rightarrow q_*$

pros: model-free (pure experience)

cons: Inefficient, deter $\sim$  policy, Xots state action space has to wait  $\uparrow$  until end of episode to calculate return

# TD:

$$\left\{ \begin{array}{l} \text{MC: } V_{\pi}(s) = \mathbb{E}_{\pi}(G_t | S_t = s), \text{ sample } G_t \text{ with condition} \\ \text{TD: } V_{\pi}(s) = \mathbb{E}_{\pi}(R_{t+1} + \gamma V_{\pi}(S_{t+1}) | S_t = s), \text{ expect} \\ \text{to see } \mathbb{E}_{\pi}(R_{t+1} + \gamma V_{\pi}(S_{t+1}) - V_{\pi}(S_t) | S_t = s) = 0. \end{array} \right.$$

temporal difference error  $\delta_t$

$$\left\{ \begin{array}{l} 0 \Rightarrow V_{\pi}(s) \text{ is good} \\ > 0 \Rightarrow V_{\pi}(s) \text{ too small} \\ < 0 \Rightarrow V_{\pi}(s) \text{ too large} \end{array} \right.$$

**TD(0) (one-step TD):** only policy evaluation

- ①: For fixed policy  $\pi$ , initialize state  $S$ , figure out action  $A$  given by  $\pi$  and  $S$
- ②: Take action  $A$ , observe reward  $R$ , next state  $S'$

learn the guess from the guess, bootstrap!

$$\textcircled{3}: V(S) \leftarrow V(S) + \alpha \cdot [R + \gamma V(S') - V(S)]$$

↑ learning rate, parameter

$S \leftarrow S'$  go to next state

- ④: loop until episode ends

better than MC

online

pros: model-free, much more efficient (experience while update)

cons:  $\times$  cts state/action space

TD(0) policy evaluation is proved to converge to  $V_{\pi}$  w.p. 1 if stochastic approx scheme holds, i.e.

$$\begin{cases} \sum_n \alpha_n = \infty & \text{large enough, overcome fluctuation} \\ \sum_n \alpha_n^2 < \infty & \text{small enough, guarantee convergence} \end{cases} \text{ for } \alpha_n \text{ as learning rate at time } n$$

↓ Idea of TD(0) applied on estimating  $q_*$ ,  $Q \approx q_*$

## SARSA (state-action-reward-state-action)

- ①: Init state  $S$ . Generate action  $A$  based on  $S$  and  $\epsilon$ -greedy policy derived from  $Q \approx q_*$
- ②: Take action  $A$ , observe reward  $R$ , next state  $S'$
- ③: Generate  $A'$  based on  $S'$  and  $\epsilon$ -greedy policy derived from  $Q$ . (Bellman consistency equation)
- ④:  $Q(S, A) \leftarrow Q(S, A) + \alpha \cdot [R + \gamma Q(S', A') - Q(S, A)]$   
 $q_{\pi}(s, a) = \mathbb{E}_{\pi} [R_{t+1} + \gamma \cdot q_{\pi}(S_{t+1}, A_{t+1}) \mid S_t = s, A_t = a]$   
 $S \leftarrow S', A \leftarrow A'$  (next time step) TD error for  $q_{\pi}$
- ⑤: Loop until episode ends actually following this action

on-policy: learning of  $Q$  based on  $A'$ , generated by the policy which is constructed based on  $Q$  and we are actually following it!

**Q-learning:** ( $Q \approx q_*$ )

①: Init state  $S$ . Generate action  $A$  based on  $S$  and  $\epsilon$ -greedy policy derived from  $Q$ . *actually following the action here*

②: Take action  $A$ , reward  $R$ , next state  $S'$

③:  $Q(S, A) \leftarrow Q(S, A) + \alpha \cdot [R + \delta \cdot \max_a Q(S', a) - Q(S, A)]$   
*(Bellman optimality equation)*

$S \leftarrow S'$  (next time step)

④: Loop until episode ends

TD error

$$\begin{aligned} q_*(s, a) &= \sum_{r, s'} [r + \delta \cdot \max_{a'} q_*(s', a')] \cdot p(s', r | s, a) \\ &= \mathbb{E}[R_{t+1} | S_t = s, A_t = a] + \delta \cdot \mathbb{E} \left[ \max_{a'} q_*(S_{t+1}, a) \mid S_t = s, A_t = a \right] \\ &= \mathbb{E} \left[ R_{t+1} + \delta \cdot \max_a q_*(S_{t+1}, a) \mid S_t = s, A_t = a \right] \end{aligned}$$

*write back as expectation*

**Off-policy:** uses the next action in the way of  $\max_a Q(s', a)$ , actually it's the greedy action w.r.t.  $Q$ . We are just estimating the return assuming a greedy policy were followed but not actually following this policy.

On/off-policy depends on if the action you use in TD error is exactly what you follow!



## Policy Gradient Method:

Directly update policy  $\pi$ , can deal with problems where estimation of value function is infeasible.

Parametrize  $\pi(a|s, \theta)$  with parameter  $\theta$ , obj of agent is to get the optimal policy  $\pi$  to maximize

$$J(\theta) = \mathbb{E}_{\pi}(G_0 | S_0 = s_0)$$

### Thm (Policy Gradient):

$$\nabla J(\theta) \propto \sum_{s \in S} \mu(s) \sum_{a \in A} q_{\pi}(s, a) \cdot \nabla \pi(a|s, \theta)$$

where  $\mu$  is on-policy dist as prob meas. on  $S$ .

$$g(s) \triangleq \sum_{k=0}^{\infty} \gamma^k \cdot \mathbb{P}_{\pi}(S_k = s | S_0 = s_0), \quad \mu(s) \triangleq \frac{g(s)}{\sum_{s' \in S} g(s')}$$

(in continuing case,  $\mu$  is stationary dist)

↓  
proportion of time MDP has spent in state  $s$

Naturally, shall update  $\theta$  in the direction of  $\nabla J(\theta)$

$$\theta \leftarrow \theta + \alpha \cdot \sum_s \mu(s) \sum_a q_{\pi}(s, a) \cdot \nabla \pi(a|s, \theta)$$

how to calculate  $\mu(s)$  and  $q_{\pi}(s, a)$ ?

Turns out we don't need to calculate those but can write them in terms of expectation.

Assume  $S_0 \sim \mu$  (stationary distribution)

assume  $\gamma=1$ , then

$S_t \xrightarrow{\alpha} \mu$  ( $t \rightarrow \infty$ )  
possible to argue if  
( $S_t$ ) nice enough

$$\sum_S \mu(s) \sum_a q_{\pi}(s, a) \cdot \nabla \pi(a|s, \theta)$$

$$= \mathbb{E}_{\pi} \left[ \sum_a q_{\pi}(S_t, a) \cdot \nabla \pi(a|S_t, \theta) \right]$$

$(q_{\pi}(S_t, a) = \mathbb{E}_{\pi}(G_t | S_t, A_t = a))$

$$= \mathbb{E}_{\pi} \sum_a \pi(a|S_t, \theta) \cdot \frac{\nabla \pi(a|S_t, \theta)}{\pi(a|S_t, \theta)} \cdot q_{\pi}(S_t, a)$$

$$= \mathbb{E}_{\pi} \mathbb{E}_{A_t \sim \pi(\cdot | S_t, \theta)} \left[ \frac{\nabla \pi(A_t | S_t, \theta)}{\pi(A_t | S_t, \theta)} \cdot q_{\pi}(S_t, A_t) \mid S_t \right]$$

$$= \mathbb{E}_{\pi} \left[ \frac{\nabla \pi(A_t | S_t, \theta)}{\pi(A_t | S_t, \theta)} \cdot q_{\pi}(S_t, A_t) \right]$$

$$= \mathbb{E}_{\pi} \left[ \nabla \log \pi(A_t | S_t, \theta) \cdot \mathbb{E}_{\pi}(G_t | S_t, A_t) \right]$$

$$= \mathbb{E}_{\pi} \left[ G_t \cdot \nabla \log \pi(A_t | S_t, \theta) \right]$$

gradient has form of expectation!

$\Downarrow$   
SGD

$$\star: \theta \leftarrow \theta + \alpha \cdot G_t \cdot \nabla \log \pi(A_t | S_t, \theta)$$

## REINFORCE:

- ①: Experience  $S_0, A_0, R_1, \dots$  following current policy under parameter  $\theta$
- ②: At each time  $t$ , build up  $G_t$  and perform 
$$\theta \leftarrow \theta + \alpha \cdot \gamma^t \cdot G_t \cdot \nabla \log \pi(A_t | S_t, \theta)$$

combine with DL:  $\pi$  can be approx by NN with softmax output layer.

## Actor-Critic:

- actor: approx policy, generating actions
- critic: approx state value func to assess the (TD error) action taken

①: Current state  $S$ , generate action  $A \sim \pi(\cdot | S, \theta)$   
take  $A$ , get reward  $R$ , next state  $S'$

②: TD error

$$\delta \leftarrow R + \gamma \hat{v}(S', w) - \hat{v}(S, w)$$

approx state value func, parametrized by parameter  $w$

$$\textcircled{3}: \begin{cases} w \leftarrow w + \alpha^w \cdot \delta \cdot \nabla_w \hat{V}(S_t, w) & \text{update critic} \\ \theta \leftarrow \theta + \alpha^\theta \cdot \gamma^t \cdot \delta \cdot \nabla_\theta \log \pi(A_t | S_t, \theta) & \text{update actor} \end{cases}$$

learning rate

$$\textcircled{4}: S \leftarrow S' \text{ (next time step)}$$

DL: approx policy & value func with NN, so organize 2 NNs and  $\nabla_w, \nabla_\theta$  can be derived easily numerically.

- Pros: model-free, cts state action space, enough randomized policy, online
- Cons: parametric form, time-consuming training

weighted mean-square error of value approx:

$$\overline{VE}(w) = \sum_s \mu(s) [V_\pi(s) - \hat{V}(s, w)]^2$$

replace  $G_t$  with  $R + \gamma \hat{V}(S', w) \Leftarrow G_t \cdot \nabla_w \hat{V}(S_t, w)$

(TD idea)

again use SGD idea avoid calculation of expectation