

② weight average method.

1). upper confidence bound

(UCB) action selection

(see 6).

$$A_t = \underset{a}{\operatorname{argmax}} \left(Q_t(a) + c \sqrt{\frac{\log t}{N_t(a)}} \right),$$

$c > 0$ controls the degree of exploration.

12. gradient bandit alg

12.1 $H_t(a)$: preference of actions

$$P_t(A_t=a) = \frac{e^{H_t(a)}}{\sum_{i=1}^k e^{H_t(i)}} = \pi_t(a)$$

12.3 $\pi_t(a)$ is the probability of taking action a at time t .

$$Q_{t+1} = Q_t + \alpha(R_t - Q_t) \quad \alpha \in (0, 1] \text{ is constant}$$

$$H_{t+1}(a) = H_t(a) + \alpha(R_t - \bar{R}_t)(1 - \pi_t(a)) \quad \text{if } A_t=a$$

$$H_{t+1}(a) = H_t(a) - \alpha(R_t - \bar{R}_t)\pi_t(a) \quad = (1-\alpha)^n Q_t + \sum_{i=1}^n \alpha(1-\alpha)^{n-i} R_i \quad \text{fixed, then if } A_t \neq a$$

$\alpha > 0$ is step size, $R_t \in \mathbb{R}$ is the t . estimate values of actions $R_t(a) \sim \mathcal{N}(\hat{q}_t(a), 1)$

8. A simple bandit alg

① Initialize

$$Q(a) \leftarrow 0, \quad N(a) \leftarrow 0,$$

for $a=1$ to k

② Repeat forever

$$A \leftarrow \begin{cases} \operatorname{argmax}_a Q(a) & \text{with } p = 1-\epsilon \\ \text{random action, with } p = \epsilon \end{cases}$$

③ Action A

$$R \leftarrow \text{bandit}(A)$$

$$N \leftarrow N+1$$

$$Q(A) \leftarrow Q(A) + \frac{1}{N+1} (R - Q(A))$$

$$A_t = \operatorname{argmax}_a Q_t(a)$$

④ Greedy action selection: $Q_t(a) = E(R_t | A_t=a)$

⑤ Greedy action selection: $Q_t(a) = \operatorname{argmax}_a Q_t(a)$

⑥ Action selection rule: $Q_t(a) = \operatorname{argmax}_a Q_t(a)$

⑦ Greedy action selection: $Q_t(a) = \operatorname{argmax}_a Q_t(a)$

⑧ Greedy action selection: $Q_t(a) = \operatorname{argmax}_a Q_t(a)$

⑨ Greedy action selection: $Q_t(a) = \operatorname{argmax}_a Q_t(a)$

⑩ Greedy action selection: $Q_t(a) = \operatorname{argmax}_a Q_t(a)$

⑪ Greedy action selection: $Q_t(a) = \operatorname{argmax}_a Q_t(a)$

⑫ Greedy action selection: $Q_t(a) = \operatorname{argmax}_a Q_t(a)$

k -armed bandit problem

1). k -armed bandit alg

1). k -armed bandit alg

2). k -armed bandit alg

3). k -armed bandit alg

4). k -armed bandit alg

5). k -armed bandit alg

6). k -armed bandit alg

7). k -armed bandit alg

8). k -armed bandit alg

9). k -armed bandit alg

10). k -armed bandit alg

11). k -armed bandit alg

12). k -armed bandit alg

13). k -armed bandit alg

14). k -armed bandit alg

15). k -armed bandit alg

16). k -armed bandit alg

17). k -armed bandit alg

18). k -armed bandit alg

k -armed bandit problem

1). k -armed bandit problem

1). k -armed bandit problem

2). k -armed bandit problem

3). k -armed bandit problem

4). k -armed bandit problem

5). k -armed bandit problem

6). k -armed bandit problem

7). k -armed bandit problem

8). k -armed bandit problem

9). k -armed bandit problem

10). k -armed bandit problem

11). k -armed bandit problem

12). k -armed bandit problem

13). k -armed bandit problem

14). k -armed bandit problem

15). k -armed bandit problem

16). k -armed bandit problem

17). k -armed bandit problem

18). k -armed bandit problem

$$\textcircled{4} \therefore \frac{\partial E(R_t)}{\partial H_t(a)}$$

$$= \sum_b \pi_t(b) (q_*(b) - X_t) \frac{\partial \pi_t(b)}{\partial H_t(a)} \cdot \frac{1}{\pi_t(b)}$$

B. gradient bandit alg as
stochastic gradient ascent

$$\textcircled{3.1} \quad H_{t+1}(a) \doteq H_t(a) + \alpha \frac{\partial E(R_t)}{\partial H_t(a)}$$

(exact gradient ascent)

$$E(R_t) \doteq \sum_b \pi_t(b) q_*(b)$$

$$\textcircled{5} \therefore E_{R_t \sim q_*(A_t)} (R_t | A_t) = q_*(A_t)$$

B.2: the updates of gradient bandit alg (B.1)
are equal to B.1 in expected value.

$$\therefore \frac{\partial E(R_t)}{\partial H_t(a)} = E_{A_t \sim \pi_t} \left[(R_t - X_t) \frac{\partial \pi_t(A_t)}{\partial H_t(a)} \cdot \frac{1}{\pi_t(A_t)} \right]$$

(2)

, gradient bandit alg is an instance
of stochastic gradient ascent.

$$\textcircled{6} \therefore \frac{\partial \pi_t(b)}{\partial H_t(a)} = \pi_t(b) (1_{a=b} - \pi_t(a))$$

(v) 13.3 proof:

$$\textcircled{1} \quad \frac{\partial E(R_t)}{\partial H_t(a)} = \frac{\partial}{\partial H_t(a)} \sum_b \pi_t(b) q_*(b)$$

$$\text{choose } X_t = \bar{R}_t$$

$$\therefore \frac{\partial E(R_t)}{\partial H_t(a)} = E_{A_t \sim \pi_t} \left[(R_t - \bar{R}_t) (1_{a=A_t} - \pi_t(a)) \right] = \sum_b q_*(b) \frac{\partial \pi_t(b)}{\partial H_t(b)}$$

$$\textcircled{2} \quad \because \sum_b \pi_t(b) = 1 \quad \therefore \sum_b \frac{\partial \pi_t(b)}{\partial H_t(b)} = 0$$

$$\textcircled{3} \quad \therefore \frac{\partial E(R_t)}{\partial H_t(a)} = \sum_b (q_*(b) - X_t) \frac{\partial \pi_t(b)}{\partial H_t(a)}$$

X_t does not depend on b .

4.2 finite MDP

6. 6.1 one step dynamics of the environment
 \Leftrightarrow MDP with finite state and action sets.

$$P(s'|s, a)$$

$$\hat{=} P(S_{t+1} = s' | S_t = s, A_t = a)$$

6.2 the 6.1 complete

specify the dynamics of a finite MDP.

6.3 Given 6.1, one can compute $P(r, s', A_t = a | S_t = s, A_t = a)$

anything else one might want \Rightarrow the state signal has the anything about the environment. Markov property.

\Leftrightarrow the environment and task as a whole have the Markov property.

7.

7.1 expected rewards for state-action pairs.

$$r(s, a) \hat{=} E(R_{t+1} | S_t = s, A_t = a)$$

$$= \sum_{s' \in S} r \sum_{a'} p(s', r | s, a)$$

7.2 state-transition probability

$$p(s' | s, a) \hat{=} Pr(S_{t+1} = s' | S_t = s, A_t = a) = \sum_{r \in R} p(s', r | s, a)$$

3. Markov property:

All that matters is in the current state signal.

$$Pr(\text{future} | \text{now}, \text{history})$$

$$= Pr(\text{future} | \text{now}).$$

3.2 like linear, it's a simple

property. A fully understanding of the theory of Markov case

is an essential foundation for extending it to the more complex and realistic non-Markov case.

2. ① return: sum of rewards

$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots$$

② unified notation for episodic and continuing tasks

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

4. A reinforcement learning task that satisfies the Markov property is called a Markov decision process, or MDP.

Finite MDP are all we need to understand 90% of modern RL.

but often nearly so.

Finite MDP

$$= E_{A \sim P(A|S, R)} (S, R) | P(S_t, R_t | S_t, A_t)$$

11. (agent, environment, action)
 \Leftrightarrow (controller, controlled system, state-action-next-state triples, control signal).

$$E_{\pi} (A_t | S_t = s, A_t, S_{t+1}, R_{t+1})$$

$$= \frac{1}{\pi} \sum_{s, r} \cancel{P(s, r | S_t, A_t)}$$

$$= E_{A \sim P(A|S, R)} (S, R) | P(S_t, R_t | S_t, A_t)$$

$$\cancel{E_{\pi} (A_t | S_t, A_t)}$$

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

$$13.6 \text{ the value function } v_\pi \text{ is } \quad \textcircled{2} \quad i. E_{\substack{X \sim P(X) \\ X \sim P(X)}} f(X)$$

the unique solution to its Bellman equation.

$$14. \text{ discrete-time optimization problem:}$$

$$\text{Bellman equation: } \quad = E_{X \sim P(X)} E_{X \sim P(X)} f(X)$$

(2) continuous-time optimization problem:
 Hamilton-Jacobi-Bellman equation: $= E_{\substack{X \sim P(X) \\ A \sim P(A | S=s)}} E_{\substack{X \sim P(X) \\ A \sim P(A | S=s)}} f(X, Y, Z)$
 (a partial differential equation).

7.3 expected rewards for state-action-next-state triples,
 $\hat{E}(R_{t+1} | S_t = s, A_t = a, A_{t+1} = a')$

$$r(s, a, s')$$

12. value functions v_π and q_π
 can be estimated from experience. $= \frac{\sum r_p(s, t, s, a)}{\sum p(s | s, a)}$

13.1 Bellman equation for v_π . $\hat{v}_\pi(s) = \sum_{a, s'} \pi(a|s) p(s, s' | s, a) (r + \gamma v_\pi(s'))$

$$= \int_X \left(\sum_a \pi(a|x) \right) p(x, x') f(x) dx$$

$$= \sum_a \pi(a|x) \int_X f(x) dx$$

$$= \sum_a \pi(a|x) f(x)$$

$$13.2 \hat{v}_\pi(s) = \int_X \int_Y f(x, y) p(x, y) dy dx$$

$$= \int_X p(x) f(x) dx$$

9. state-value function
 for policy π .

$$v_\pi(s) = E_{\pi} \left(A_t | S_t = s \right)$$

$$= \frac{1}{\pi} \sum_{a, s} P(a | s)$$

10. action-value function for policy π .
 $q_\pi(s, a) = E_{\pi} (A_t | S_t = s, A_t = a)$

$$E_{\substack{X \sim P(X) \\ Y \sim P(Y)}} f(X, Y, Z)$$

$$(X, Y, Z \sim P(X, Y, Z))$$

$$= E_{\substack{X \sim P(X) \\ Y \sim P(Y)}} E_{\substack{Z \sim P(Z) \\ X \sim P(X)}} f(X, Y, Z)$$

15. Along a reinforcement learning task means, roughly, finding a policy that achieves a lot of reward over the long run.
16. $\pi \geq \pi' \Leftrightarrow V_{\pi}(s) \geq V_{\pi'}(s) \quad \forall s \in S$
17. Bellman optimality equation for V_{π} .
- $$V_{\pi}(s) = \max_{a \in A(s)} Q_{\pi}(s, a)$$
18. Bellman optimality equation for $Q_{\pi}(s, a)$
- $$= E(R_{t+1} + \gamma V_{\pi}(S_{t+1}) | S_t = s, A_t = a)$$
19. If you have the optimal value function V_{π^*} , then the actions that appear best after a one-step search will be optimal actions
20. With Q_{π^*} , the optimal action is the one that maximizes $Q_{\pi^*}(s, a)$.
21. The nature of RL makes it possible to approximate optimal policies in ways that put more effort into learning to make good decisions for frequently encountered states. This is one key property that distinguishes RL from other approaches to approximately solving MDPs.
22. Bellman optimality equations \Rightarrow optimal value functions \Rightarrow an optimal policy.
23. The RL problem is deeply indebted to the idea of MDPs from the field of optimal control. RL adds to MDPs a focus on approximation and incomplete info for realistically large environments dynamics.
- 15.5 (L-see|18) has a unique solution independent of the policy.
- 16.2 $\max_{\pi^*} E(Q_{\pi^*}(s, a) | S_t = s, A_t = a)$ optimal policy π^* exists. (2)
- $= \max_{\pi^*} E(R_{t+1} + \gamma V_{\pi^*}(S_{t+1}) | S_t = s, A_t = a)$ It may not be unique.
- 16.3 optimal state-value function
- $$= \max_{\pi^*} E(V_{\pi^*}(s') | S_{t+1} = s')$$
- 16.4 optimal action-value function
- $$= \max_{\pi^*} \sum_{s' \in S} P(s', r | s, a) (r + \gamma V_{\pi^*}(s'))$$
- $$V_{\pi^*}(s) = \max_{\pi^*} \sum_{s' \in S} \sum_{a \in A(s')} P(s', r | s, a) (r + \gamma V_{\pi^*}(s'))$$
- $$Q_{\pi^*}(s, a) = \max_{\pi^*} \sum_{s' \in S} \sum_{a \in A(s')} P(s', r | s, a) Q_{\pi^*}(s', a)$$

Moreover, if there is strict inequality/strict inequality of (1) at any state, then there must be done in a sweep through π in S a system of $|S|$ linear equations in $|S|$ unknowns.

9. update the values in place
 - \Leftrightarrow with each new backup value immediately overwriting the old one.

$$\text{proof: } v_n(s) \leq q_n(s, \pi^*(s))$$

$$= E_{\pi^*}(R_{t+1} + \gamma v_n(s_{t+1}) | S_t = s)$$

$$\leq E_{\pi^*}(R_{t+1} + \gamma q_n(s_{t+1}, \pi(s_{t+1})) | S_t = s)$$

$$= E_{\pi^*}(R_{t+1} + \gamma E_{\pi^*}(R_{t+2} + \gamma v_n(s_{t+2})) | S_t = s)$$

$$= E_{\pi^*}(R_{t+1} + \gamma R_{t+2} + \gamma^2 v_n(s_{t+2}) | S_t = s)$$

$$\leq \dots \leq E_{\pi^*}(R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{n-1} R_{t+n} + \gamma^n v_n(s) | S_t = s)$$

$$= v_n(s)$$

is sufficiently small.
13. Policy improvement: the process of making a new policy that improves on an original policy, by making it greedy with respect to the value function.

the process of improving the process of making a new policy that improves on an original policy, by making it greedy with respect to the value function.

12. policy improvement theorem: $\pi^*(s, \pi'(s)) \geq v_n(s), \forall s \in S$

11. policy improvement theorem: $\pi^*(s) = \sum_a \pi(a|s) \sum_{s'} \pi(s'|s, a) v_n(s')$

1. dynamic programming (DP) refers to a collection of algs that can be used to compute optimal policies given a perfect model of the environment as a MDP.

2. we usually assume the environment is a finite MDP.

$$\begin{aligned} v_H &= E_{\pi^*}(R_{t+1} + \gamma v_H(s_{t+1}) | S_t = s) \\ &\leq E_{\pi^*}(R_{t+1} + \gamma v_n(s_{t+1}) | S_t = s) \end{aligned}$$

3. $v_n(s) = E_{\pi^*}(R_{t+1} + \gamma v_{n+1}(s_{t+1}) | S_t = s)$

4. condition A satisfied

$\Rightarrow v_n \rightarrow v_H$, as $n \rightarrow \infty$

$$\begin{aligned} &v_n(s) = E_{\pi^*}(R_{t+1} + \gamma v_{n+1}(s_{t+1}) | S_t = s) \\ &\text{use the in-place version.} \end{aligned}$$

5. for DP algos, we usually use the in-place version.

6. if condition A is satisfied

$\Rightarrow v_n \rightarrow v_H$, as $n \rightarrow \infty$

$$\begin{aligned} &v_n(s) = E_{\pi^*}(R_{t+1} + \gamma v_{n+1}(s_{t+1}) | S_t = s) \\ &\text{the initial } v_0 \text{ is chosen arbitrarily.} \end{aligned}$$

7. iterative policy evaluation:

8. Each iteration of iterative policy evaluation backs up the value of every state once to produce

$\Rightarrow v_n \rightarrow v_H$, so we call it a full backup.

9. or eventual termination is guaranteed from all states under the policy π^* .

10. A satisfied condition

$\Rightarrow v_n \rightarrow v_H$, unique.

19.8 condition A satisfied \Rightarrow

value iteration converges.

$\underline{\pi}$ denotes policy improvement.
 $\overline{\pi}$ denotes policy improvement.

20. asynchronous DP algs are in place iterative DP algs that back up states in an arbitrary order.

21. GPI: generalized policy iteration. We use GPI to refer to the general idea of letting policy evaluation and policy improvement processes interact, independent of the granularity and other details of the two processes.

21.2 Almost all RL methods are well described as GPI.

22. DP is exponentially faster than any direct search in policy space could be.

23 bootstrapping: update estimates on the basis of other estimates.

14. greedy policy $\pi'(s) = \arg\max_a \pi_\pi(s, a)$
 $\forall s \in S$

18. 'a finite MDP has a finite/15. Policy improvement must give us a number of policies, i.e. policy iteration strictly better policy except when the original policy is already optimal. must converge to an optimal policy proof: $\pi'(s) = \arg\max_a \pi_\pi(s, a)$
= $\arg\max_a \sum_{s', r} p(s', r|s, a)(r + \gamma \pi_\pi(s'))$

19. value iteration:

$$\underline{v}_\pi(s) = \max_a \pi_\pi(s, a)$$
$$\underline{v}_\pi(s) = \max_a \sum_{s', r} p(s', r|s, a)(r + \gamma \underline{v}_\pi(s'))$$
$$\text{if } \underline{v}_\pi = \bar{v}_\pi$$

19.2 policy evaluation:

$$\bar{v}_\pi = E_{\pi}(\underline{R}_{t+1} + \gamma \bar{v}_\pi(S_{t+1}) | S_t = s)$$

then $\bar{v}_\pi(s) = \max_a \sum_{s', r} p(s', r|s, a)(r + \gamma \bar{v}_\pi(s'))$

'this is the same as the Bellman

19.3 value iteration is identical optimality equation

19.4 value iteration is obtained for the stochastic policies. simply by turning the Bellman 17 policy iteration: $\pi_0 \rightarrow \pi_1 \rightarrow \pi_2 \rightarrow \dots \rightarrow \pi_k \rightarrow \pi_{k+1}$

Monte Carlo Methods

7. advantages Monte Carlo methods can have over DP methods :

- ① ability to learn from actual experience.
- ② ability to learn from simulated experience.
- ③ estimating the value of a single state is independent of the number of states.
- ④ attractive when one requires the value of only one or a subset of states.
- ⑤ just as good as π^* .
- ⑥ policy iteration converges.

10.3 We made two unlikely assumptions in order to obtain this guarantee of convergence for the Monte Carlo methods:

- ① the episodes have exploring starts.
- ② policy evaluation can be done with an infinite number of episodes. (To avoid the infinite number of episodes, we can use the idea of API).

11. Monte Carlo ES (Monte Carlo with exploring starts) use API on an episode-by-episode basis. Monte Carlo ES converges, but has not yet been formally proved.

1. Monte Carlo methods are ways of solving the problem based on averaging sample returns
2. Monte Carlo methods are our first learning method for estimating value functions and discovering optimal policies.
3. rule → data

If we know rule, data will be useless.

4. We compute value functions from knowledge of MDP, or we learn them from sample returns with the MDP.
5. policy evaluation:

first-visit MC, every-visit MC

Both converge quadratically to $v_{\pi}(s)$ as the number of visits (or first visits) to s goes to infinity. (not proved).

6. the estimates for each state are independent. ▷ Monte Carlo methods do not bootstrap.

12. On-policy methods attempt to evaluate or improve the policy that is used to make decisions, whereas off-policy methods evaluate or improve a policy different from that used to generate the data.

16. coverage: $\pi(a|s) > 0 \Rightarrow \mu(a|s) > 0$
where π is target policy, μ is behavior policy.

14.3 define new environment
with π^{old} , with $\mu = 1 - \varepsilon$.

$$17.2 \quad \mathbb{E}_f p^T = \mathbb{E}_f p_f^T = \int p^T f = \frac{1}{2} \int f$$

$$17.3 \quad \mathbb{E}_f Q_\pi = \mathbb{E}_f p_f^T Q_\pi, \quad p_f^T = \frac{\pi^T \pi(\pi(s)) p(\pi(s), f)}{\pi^T \pi(\pi(s)) p(\pi(s), f)}$$

replics action at random with equal pr,
and then behave like old environment,

with $\mu = \pi$.

17.4. Monte Carlo ES is an on-policy method with exploring starts.

17.4. policy iteration works for ε -soft policies.

$$\text{If } q_\pi(s, \pi'(s)) = v_\pi(s), \forall s, \quad \text{i.e. } \pi' = \pi,$$

(π is no longer improved)

$$\Rightarrow v_\pi(s) = q_\pi(s, \pi'(s))$$

$$= (1 - \varepsilon) \max_a q_\pi(s, a) + \frac{\varepsilon}{|A(s)|} \sum_a q_\pi(s, a)$$

$$= (1 - \varepsilon) \max_a \sum_{s', r, \pi(s, a)} p(s', r | s, a) (r + \gamma v_\pi(s'))$$

$$+ \frac{\varepsilon}{|A(s)|} \sum_a p(s', r | s, a) (r + \gamma v_\pi(s'))$$

$$= (1 - \varepsilon) \max_a \tilde{q}_\pi(s, a) + \frac{\varepsilon}{|A(s)|} \sum_a \tilde{q}_\pi(s, a)$$

$$\text{if } \pi(s) = \tilde{q}_\pi(s) \quad \text{then}$$

$$\text{proof: } q_\pi(s, \pi'(s)) = \frac{1}{\alpha} \sum_a \pi(a|s) q_\pi(s, a)$$

$$= \frac{\varepsilon}{|A(s)|} \sum_a q_\pi(s, a) + (1 - \varepsilon) \max_a q_\pi(s, a)$$

$$\text{if } \pi(s) \neq \tilde{q}_\pi(s) \quad \text{then}$$

$$= \frac{\varepsilon}{|A(s)|} \sum_a \pi(a|s) q_\pi(s, a) + (1 - \varepsilon) \max_a q_\pi(s, a)$$

$$\leq \frac{\varepsilon}{|A(s)|} (\pi(a|s) - \frac{\varepsilon}{|A(s)|}) q_\pi(s, a)$$

$$\leq \frac{\varepsilon}{|A(s)|} (\pi(a|s) - \frac{\varepsilon}{|A(s)|}) \max_a q_\pi(s, a)$$

$$= \max_a (1 - \varepsilon) q_\pi(s, a)$$

$$\therefore v_\pi(s) \leq q_\pi(s, \pi'(s)) \wedge \text{by policy improvement}$$

theorem, $\pi' \geq \pi$.

14.3 define new environment

with π^{old} , with $\mu = 1 - \varepsilon$.

$$17.4. \quad \text{weighted importance sampling}$$

$$V(s) = \frac{\sum_a p(s, a) R_a}{\sum_a p(s, a)}$$

$$17.5. \quad \text{return-specific importance sampling}$$

$$V(s) = \frac{\sum_a p(s, a) R_a}{\sum_a p(s, a)}$$

$$18.1. \quad \text{flat partial returns: } \tilde{G}_t^* = R_t + \gamma R_{t+1} + \gamma^2 R_{t+2} + \dots + \gamma^{T-t} R_T$$

$$0 \leq t < T$$

$$= (1 - \gamma) \sum_{h=t+1}^T \gamma^{h-t} \tilde{G}_h^* + \gamma^{T-t} \tilde{G}_T^*$$

$$18.2. \quad \text{discounting-aware importance sampling estimator:}$$

$$V(s) = \frac{\sum_a ((1-\gamma)^T \sum_{h=t+1}^T \gamma^{h-t} p(s, a) + \gamma^{T-t} R_T)}{|A(s)|}$$

$$18.3. \quad G_t = R_t + \gamma R_{t+1} + \gamma^2 R_{t+2} + \dots + \gamma^{T-t} R_T$$

$$= (1 - \gamma) \sum_{h=t+1}^T \gamma^{h-t} G_h + \gamma^{T-t} G_T$$

$$18.4. \quad \text{per-reward importance sampling}$$

$$V(s) = \frac{\sum_a ((1-\gamma)^T \sum_{h=t+1}^T \gamma^{h-t} p(s, a) + \gamma^{T-t} R_T)}{|A(s)|}$$

$$18.5. \quad \text{per-reward importance sampling}$$

$$V(s) = \frac{\sum_a G_a}{|A(s)|}, \quad \tilde{G}_a = p_a^T R_a + \gamma p_a^T R_{a+1} + \gamma^2 p_a^T R_{a+2} + \dots + \gamma^{T-a} p_a^T R_T$$

19. incremental implementation

19.2. weighted importance sampling

$$C_{n+1} = C_n + W_{n+1}$$

$$W_n = p^t \pi(t)$$

$$V_{n+1} = \frac{\sum_{k=1}^n W_k C_k}{\sum_{k=1}^n W_k}$$

$$\therefore V_{n+1} C_n = \sum_{k=1}^n W_k C_k$$

$$V_n C_{n-1} = \sum_{k=1}^{n-1} W_k C_k$$

$$\therefore V_{n+1} C_n = V_n C_{n-1} + W_n C_n$$

$$\therefore V_{n+1} C_n = V_n C_n + W_n (C_n - V_n)$$

$$\therefore V_{n+1} = V_n + \frac{W_n}{C_n} (C_n - V_n), n \geq 1$$

Temporal-Difference Learning

1. TD learning is a combination of Monte Carlo ideas and dynamic programming ideas.
 2. constant- α MC : $V(S_t) \leftarrow V(S_t) + \alpha(A_t - V(S_t))$, α is constant.
 3. TD(0) : $V(S_t) \leftarrow V(S_t) + \alpha(R_{t+1} + \gamma V(S_{t+1}) - V(S_t))$
 4. $v_\pi(s) \doteq E_\pi(A_t | S_t = s)$ (1)
 $= E_\pi\left(\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} | S_t = s\right) = E_\pi(R_{t+1} + \gamma \sum_{k=0}^{\infty} \gamma^k R_{t+k+2} | S_t = s)$
 $= E_\pi(R_{t+1} + \gamma v_\pi(S_{t+1}) | S_t = s)$ (2)
- Roughly speaking, MC methods use an estimate of (1) as a target, whereas PP methods use an estimate of (2) as a target.
- MC: a sample return is used in place of the real expected return.
- DP: $v_\pi(S_{t+1})$ is not known and $V(S_{t+1})$ is used instead.
- TD: combine the sampling of MC with the bootstrapping of DP.
5. TD, MC : sample backups.
DP : full backups.
 6. TD error : $\delta_t \doteq R_{t+1} + \gamma V(S_{t+1}) - V(S_t)$
 - 6.2 MC error can be written as a sum of TD errors.
- $$\begin{aligned}
 A_t - V(S_t) &= R_{t+1} + \gamma A_{t+1} - V(S_t) \\
 &= R_{t+1} + \cancel{\gamma A_{t+1}} + \gamma V(S_{t+1}) - V(S_t) + \gamma(A_{t+1} - V(S_{t+1})) \\
 &= \delta_t + \gamma(A_{t+1} - V(S_{t+1})) \\
 &= \dots \\
 &= \delta_t + \gamma \delta_{t+1} + \dots + \gamma^{T-t-1} \delta_{T-1} + \gamma^{T-t}(A_T - V(S_T)) \\
 &= \sum_{k=0}^{T-t-1} \gamma^k \delta_{t+k}
 \end{aligned}$$

7. advantages of TD prediction methods.

- ① TD have an advantage over DP in that they do not require a model of the environment.
- ② TD over MC: they are naturally implemented in an online, fully incremental ~~faster~~ fashion.
- ③ for any fixed policy π , TD has been proved to converge to v_π under mild conditions.

8. under batch updating, TD(0) converges deterministically to a single answer independent of the step-size parameter, α , as long as α is chosen to be sufficiently small. The constant- α MC method also converges ~~deterministically~~ deterministically under the same condition, but to a different answer.

8.2 Under normal updating, the methods do not move all the way to their respective batch answers, but in some sense they take steps in these directions.

8.3 How is it that batch TD MC is optimal only in a limited way, and that TD is optimal in a way that is more relevant to predicting returns.

8.4 Batch MC always find the estimates that minimize mse error on the training set, whereas batch TD(0) always finds the estimates that would be exactly correct for the maximum-likelihood model of the Markov process.

8.5 batch TD(0) converges to the certainty-equivalence estimate.

8.6 certainty-equivalence estimate is the estimate of the value function that would be exactly correct if the model were exactly correct. It ~~is assumed~~ that the estimate of the underlying process was known with certainty rather than being approximated. This is equivalent to assuming

9. Sarsa is an on-policy TD control algorithm.

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha(R_{t+1} + \gamma Q(S_{t+1}, A_{t+1}) - Q(S_t, A_t))$$

The rule uses $(S_t, A_t, R_{t+1}, S_{t+1}, A_{t+1})$ that make up a transition from one state-action pair to the next. This quintuple gives rise to the name Sarsa for the algorithm.

10. One of the early breakthroughs in RL was the development of an off-policy TD control algorithm known as Q-learning (1989),

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha(R_{t+1} + \gamma \max_a Q(S_{t+1}, a) - Q(S_t, A_t))$$

10.2. Q directly approximates q_* , independent of the policy being followed.

11. Any method guaranteed to find optimal behavior in the general case must have a minimal requirement that all ~~possible~~ ~~action~~ state-action pairs continue to be updated.

12. Expected Sarsa :

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha(R_{t+1} + \gamma E(Q(S_{t+1}, A_{t+1}) | S_{t+1}) - Q(S_t, A_t))$$

Given S_{t+1} , this algorithm moves deterministically in the same direction as Sarsa moves in expectation, and accordingly it is called expected Sarsa.

12.2 Expected Sarsa can be an on-policy or off-policy algorithm.

Except for the same additional computational cost, Expected Sarsa may completely dominate both Sarsa and Q-learning.

B. A maximum over estimated values is used as an estimate of the maximum value, which can lead to a significant positive bias. We call this maximization bias. (L-)

14. doubled learning : $A^* = \underset{a}{\operatorname{argmax}} Q_1(a)$, $Q_2(A^*)$ is unbiased in the sense that $E(Q_2(A^*)) = q(A^*)$. Similarly, we can yield a second unbiased estimate $Q_1(\underset{a}{\operatorname{argmax}} Q_2(a))$

here Q_1, Q_2 are independent estimate of the true value $q(a)$, $\forall a \in A$.
14.2 doubled learning seems to eliminate the harm caused by maximization bias.

15. TD methods are alternatives to MC methods for solving the prediction problem. In both cases, the extension to the control problem is via the idea of generalized policy iteration (GPI) that we abstracted from dynamic programming. This is the idea that approximate policy and value functions should interact in such a way that they both move toward their optimal values.

16. TD methods are general methods for learning to make long-term predictions about dynamical systems.

Chapter 7
Multi-step Bootstrapping

1. Multi-step TD methods generalize both MC methods and one-step TD methods so that one can switch from one to the other smoothly.
2. n-step TD methods are still TD methods because they still change an earlier estimate based on how it differs from a later estimate.

3.1 The target of Monte Carlo backups :

$$G_t \doteq R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots + \gamma^{T-t-1} R_T.$$

3.2 the target of one-step backups:

$$G_t^{(1)} \doteq R_{t+1} + \gamma V_t(S_{t+1})$$

3.3 the target of n-step backup:

$$G_t^{(n)} \doteq R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{n-1} R_{t+n} + \gamma^n V_{t+n}(S_{t+n}), \quad n \geq 1, \quad 0 \leq t < T_n$$

$$G_t^{(n)} \doteq G_t, \quad \text{if } t+n \geq T.$$

4. n-step TD algorithm :

$$V_{t+n}(S_t) \doteq V_{t+n-1}(S_t) + \alpha (G_t^{(n)} - V_{t+n-1}(S_t)), \quad 0 \leq t < T$$

$$V_{t+n}(s) = V_{t+n-1}(s), \quad \forall s \neq S_t.$$

5. error reduction property of n-step returns ~~(n-step Sarsa)~~

$$\max_s |E_\pi(G_t^{(n)}| S_t=s) - V_n(s)| \leq \gamma^n \max_s |V_{t+n}(s) - V_n(s)|$$

6.1 n-step returns in terms of estimated action values:

~~$$G_t^{(n)} \doteq R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{n-1} R_{t+n} + \gamma^n Q_{t+n-1}(S_{t+n}, A_{t+n}), \quad n \geq 1, \quad 0 \leq t < T_n$$~~

$$Q_t^{(n)} \doteq Q_t \text{ if } t+n \geq T.$$

~~6.2 n-step Sarsa algorithm :~~

$$Q_{t+n}(S_t, A_t) \doteq Q_{t+n-1}(S_t, A_t) + \alpha (G_t^{(n)} - Q_{t+n-1}(S_t, A_t)), \quad 0 \leq t < T$$

$$Q_{t+n}(s, a) = Q_{t+n-1}(s, a), \quad \forall s, a \text{ such that } s \neq S_t \text{ or } a \neq A_t.$$

~~$$G_t^{(n)} = Q_{t-1}(S_t, A_t) + \sum_{k=t}^{\min(t+n-1, T-1)} \gamma^{k-t} (R_{k+1} + \gamma Q_k(S_{k+1}, A_{k+1}) - Q_{k-1}(S_k, A_k))$$~~

n-step

7. the target of Expected Sarsa.

$$G_t^{(n)} = R_{t+1} + \dots + \gamma^{n-1} R_{t+n} + \gamma^n \sum_a \pi(a|S_{t+n}) Q_{t+n-1}(S_{t+n}, a), n \geq 1,$$

8. ~~off-policy~~ off-policy version of n-step TD: $0 \leq t \leq T-1$.

$$V_{t+n}(S_t) = V_{t+n-1}(S_t) + \alpha p_t^{t+n} (G_t^{(n)} - V_{t+n-1}(S_t)), 0 \leq t < T$$

where p_t^{t+n} is called importance sampling ratio.

$$p_t^{t+n} = \prod_{k=t}^{\min(t+n-1, T-1)} \frac{\pi(A_k|S_k)}{\mu(A_k|S_k)}$$

9. off-policy form of n-step Sarsa:

$$Q_{t+n}(S_t, A_t) = Q_{t+n-1}(S_t, A_t) + \alpha p_{t+1}^{t+n+1} (G_t^{(n)} - Q_{t+n-1}(S_t, A_t)), 0 \leq t < T$$

10. off-policy form of n-step Expected Sarsa. ~~($\pi \neq \pi^{t+n}$)~~

$$Q_{t+n}(S_t, A_t) = Q_{t+n-1}(S_t, A_t) + \alpha p_{t+1}^{t+n+1} (G_t^{(n)} - Q_{t+n-1}(S_t, A_t)), 0 \leq t < T$$

11. It is probably inevitable that off-policy training is slower than on-policy training—after all, the data is less relevant to what you are trying to learn.

12. n-step tree backup algorithm.

12.1 off-policy learning without importance sampling.

12.2 expected action value:

$$V_t = \sum_a \pi(a|S_t) Q_{t-1}(S_t, a)$$

12.3 TD. error: $\delta_t = R_{t+1} + \gamma V_{t+1} - Q_{t-1}(S_t, A_t)$

12.4 $G_t^{(1)} = R_{t+1} + \gamma V_{t+1} = Q_{t-1}(S_t, A_t) + \delta_t$

$$G_t^{(2)} = R_{t+1} + \gamma V_{t+1} = \gamma \pi(A_{t+1}|S_{t+1}) Q_{t-1}(S_{t+1}, A_{t+1}) + \gamma \pi(A_{t+1}|S_{t+1})(R_{t+2} + \gamma V_{t+2})$$

$$= Q_{t-1}(S_t, A_t) + \delta_t + \gamma \pi(A_{t+1}|S_{t+1}) \delta_{t+1}$$

$$\therefore G_t^{(n)} = Q_{t-1}(S_t, A_t) + \sum_{k=t}^{\min(t+n-1, T-1)} \gamma^k \pi(A_k|S_k) \delta_{k+1}$$

$$12.5 \quad Q_{t+n}(S_t, A_t) = Q_{t+n-1}(S_t, A_t) + \alpha (Q_t^{(n)} - Q_{t+n-1}(S_t, A_t)), \quad 0 \leq t < T$$

$$Q_{t+n}(s, a) = Q_{t+n-1}(s, a), \quad \forall s, a \text{ such that } s \neq S_t \text{ or } a \neq A_t.$$

13. A unifying algorithm: n-step $Q(\sigma)$

13.2 n-step returns of $Q(\sigma)$

$$\textcircled{1} \quad Q_t^{(1)} = R_{t+1} + \gamma (\sigma_{t+1} Q_t(S_{t+1}, A_{t+1}) + (1-\sigma_{t+1}) V_{t+1})$$

$$= S_t + Q_{t-1}(S_t, A_t)$$

$$\text{where } S_t = R_{t+1} + \gamma (\sigma_{t+1} Q_t(S_{t+1}, A_{t+1}) + (1-\sigma_{t+1}) V_{t+1}) - Q_{t-1}(S_t, A_t)$$

$$\textcircled{2} \quad Q_t^{(2)} = R_{t+1} + \gamma (\sigma_{t+1} Q_t(S_{t+1}, A_{t+1}) + (1-\sigma_{t+1}) V_{t+1})$$

$$- \gamma (1-\sigma_{t+1}) \pi(A_{t+1}|S_{t+1}) Q_t(S_{t+1}, A_{t+1})$$

$$+ \gamma (1-\sigma_{t+1}) \pi(A_{t+1}|S_{t+1}) (R_{t+2} + \gamma (\sigma_{t+2} Q_t(S_{t+2}, A_{t+2}) + (1-\sigma_{t+2}) V_{t+2}))$$

$$- \gamma \sigma_{t+1} Q_t(S_{t+1}, A_{t+1})$$

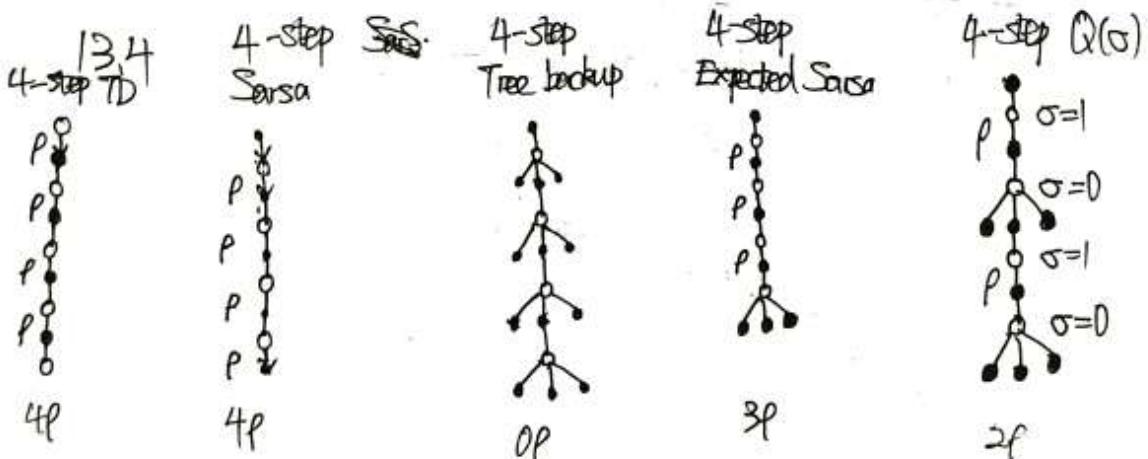
$$+ \gamma \sigma_{t+1} (R_{t+2} + \gamma (\sigma_{t+2} Q_t(S_{t+2}, A_{t+2}) + (1-\sigma_{t+2}) V_{t+2}))$$

$$= Q_{t-1}(S_t, A_t) + S_t + \gamma ((1-\sigma_{t+1}) \pi(A_{t+1}|S_{t+1}) + \sigma_{t+1}) S_{t+1}$$

$$\textcircled{3} \quad \therefore Q_t^{(n)} = Q_{t-1}(S_t, A_t) + \sum_{k=t}^{\min(t+n-1, T-1)} \sum_{i=1}^K \pi(A_i|S_k) \gamma ((1-\sigma_i) \pi(A_i|S_i) + \sigma_i)$$

13.3 for off-policy:

$$\rho_t^{t+n} = \prod_{k=t}^{\min(t+n-1, T-1)} \left(\sigma_k \frac{\pi(A_k|S_k)}{\mu(A_k|S_k)} + 1 - \sigma_k \right)$$



Chapter 8

Planning and learning with tabular methods

difficult to apply efficiently to

1. MC and TD can be unified by n-step methods. the stochastic optimal control
~~simply~~ We can develop a unified view of problems that are the focus in rl.
planning and learning methods.
4. all state-space planning methods

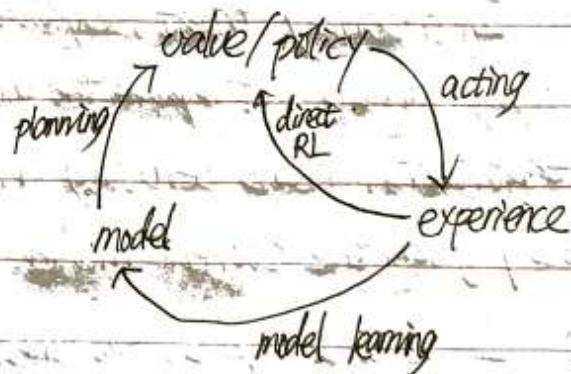
2. By a model of the environment we share a common structure:
mean anything that an agent can use model → simulated backups → values → policy
to predict how the environment will respond to its actions
5. The heart of both learning and planning methods is the estimation of value functions by backup operations.
- 2.2. distribution model
- 2.3. sample model.
6. random-sample one-step tabular
- 2.4. model can be used to simulate Q-planning is a planning method the environment and produce simulated based on one-step tabular experience.
- Q-learning and on random samples

3. We refer use planning to refer to any from a sample model. computational process that takes a model as input and produces or improves a policy for interacting with the modeled environment.
7. relationship between experience, model, values, and policy

3.2 state-space planning

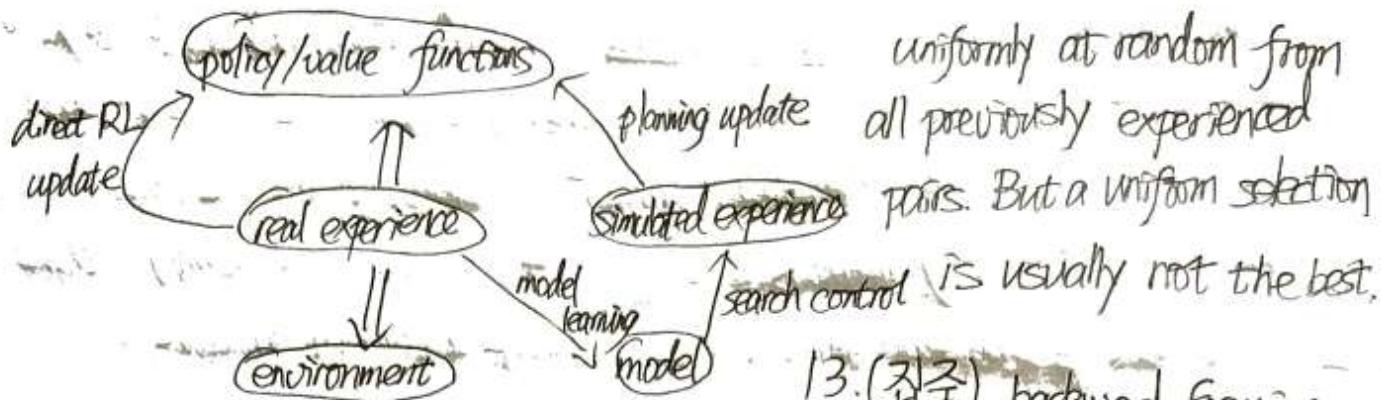
3.3 plan-space planning : evolutionary methods and "partial-order planning".

3.4 plan-space planning methods are



8. Dyna architecture

创造环境, 然后让环境告诉你 RL.



uniformly at random from all previously experienced pairs. But a uniform solution is usually not the best.

13. (집중) backward focusing

9. Dyna-Q is a simple architecture of Dyna. Of planning computations:

9.2 Dyna-Q includes planning, acting, model-learning, work backward from and direct RL — all ~~occurring~~ occurring continually. arbitrary states that have changed in value, either tabular Q-planning method. The direct RL method performing useful backups is one-step tabular Q-learning. The model-learning or terminating the method is also table-based and assumes the propagation. world is deterministic.

14. (선택) prioritized sweeping

10. Dyna-Q+ is Dyna-Q with an exploration bonus : prioritize the backups that encourages exploration.

according to a measure of their urgency, and perform them in order of priority.

11. In a planning context, exploration means trying actions that improve the model, whereas exploitation means behaving in the optimal way given the current model.

15. prioritized sweeping for a deterministic environment (선택과 집중) (c)

12. In Dyna-Q agents, simulated transitions are started in state-action pairs selected

16. prioritized sweeping uses sample backups about planning as part of in deterministic environments but uses action selection. (L-)
full backups in stochastic environments. 20. The classical state-space
- 16.2 One of prioritized sweeping's limitations planning methods in AI is that it uses full backups, which in are planning-as-part-of-action-select stochastic environments may waste lots of action methods collectively computation on low-probability transitions known as heuristic search.
17. all kinds of state-space planning can be 20.2 heuristic search can viewed as sequences of backups, varying be viewed as an extension only in the type of backup, full or sample, of the idea of a greedy large or small, and in the order in which policy beyond a single step (L-) the backups are done.
- 20.3 The great focusing
18. forward focusing: focus on states according of memory and computational to how easily they can be reached from resources on the current the states that ~~are~~ are visited frequently decision is presumably the under the current policy. reason why heuristic search
19. there tends to be two ways of thinking can be so effective. about planning. One ^{way} conceives of planning 21. Monte Carlo Tree Search as the gradual improvement of a policy or (MCTS). value function that is good in all states 21.2 MCTS is one of the generally rather than focused on any simplest examples of planning particular state. The other way thinks as part of the policy. It is

most often used when the model of the world is completely known and cheap to compute.

21.3 MCTS typically involves no approximate value functions or policies that are retained from one time step to the next; these are computed on each step and then discarded.

21.4 MCTS incrementally builds a partial game tree to select each of a game playing program's moves.

21.5 each iteration of MCTS proceeds in four stages:

- ① selection
- ② expansion
- ③ simulation
- ④ backpropagation

22 It is straightforward to integrate learning and planning processes simply by allowing both to update the same estimated value function.

23 The smaller the backups, the more incremental the planning methods can be. Among the smallest backups are one-step sample backups, as in Dyna.

24 backward focusing : prioritized sweeping
forward focusing : when planning is done exclusively as part of action selection, as in classical heuristic search and MCTS.

Chapter 9

On-policy Prediction with Approximation TD(0) backup: $S_t \rightarrow R_{t+1} + \gamma \hat{v}(S_{t+1}, \theta)$

1. The kind of generalization we require is often called function approximation because it takes examples from a desired function and attempts to generalize from them to construct an approximation of the entire function.
n-step TD backup: $S_t \rightarrow G_t^{(n)}$
DP policy-evaluation backup:
 $s \rightarrow E_{\pi}(R_{t+1} + \gamma \hat{v}(S_{t+1}, \theta_t) | t=s)$
here in DP, an arbitrary state s is backed up, whereas in other cases the state encountered

- 1.2 function approximation is an instance of supervised learning.
in actual experience, S_t , is backed up.

2. we write $\hat{v}(s, \theta) \approx v_t(s)$ for the approximated value of state s given weight vector θ .
We use function approximation methods for value

3. All of the prediction methods covered in this book have been described as backups, that is, as updates to an estimated value function that shifts its value at particular states toward a "backed-up value" function they produce for that state.
simply by passing to then the $s \rightarrow g$ of each backup as a training example. We then interpret the approximate as an estimated value function.

- 3.2 we refer to an individual backup by the notation $s \rightarrow g$, where s is the state backed up and g is the backed-up value, or target, that s 's estimated value is shifted toward.
5. In the tabular case a continuous measure of prediction quality was not necessary because the

- 3.3 MC backup: $S_t \rightarrow G_t$
learned value function

could come to equal the true value function exactly.

6. MSVE : mean squared value error

$$\text{MSVE}(\theta) = \sum_{s \in S} d(s) (v_\pi(s) - \hat{v}(s, \theta))^2$$

6.2 $d(s)$ is the fraction of time spent in s under the target policy π .
This is called the on-policy distribution.

6.3 In continuing tasks, the on-policy distribution is the stationary distribution under π .

6.4 In an episodic task, the on-policy distribution is a little different in that it depends on how the initial states of episodes are chosen.

6.5 It is not completely clear that MSVE is the right performance objective for rl.

7. SGD methods are among the most widely used of all function approximation methods and are particularly well suited to online reinforcement learning.

7.2 the approximate value function $\hat{v}(s, \theta)$ is a differentiable function of θ for all $s \in S$.

8. SGD: $\theta_{t+1} = \theta_t - \frac{1}{2} \alpha \nabla (v_\pi(s_t) - \hat{v}(s_t, \theta_t))^2$
 $= \theta_t + \alpha (v_\pi(s_t) - \hat{v}(s_t, \theta_t)) \nabla \hat{v}(s_t, \theta_t)$

8.2 the convergence results for SGD methods assume that α decreases over time.

9. We do not seek or expect to find a value function that has zero error for all states, but only an approximation that balances the errors in different states.

10. SGD method for state-value prediction: i. bootstrapping methods do
 $\theta_{t+1} = \theta_t + \alpha(U_t - \hat{v}(S_t, \theta_t)) \nabla \hat{v}(S_t, \theta_t)$ not converge as robustly as
here U_t is target value, not the true gradient methods.
value $v_t(S_t)$. $S_t \rightarrow U_t$ (one does not obtain the same

10.2 If U_t is an unbiased estimate, guarantee).

that is, if $E(U_t) = v_t(S_t)$, for each t , then θ_t is

11. condition: $\sum_{n=1}^{\infty} \alpha_n = \infty$, $\sum_{n=1}^{\infty} \alpha_n^2 < \infty$. of changing the weight vector θ_t

12. If U_t is an unbiased estimate, that is, if $E(U_t) = v_t(S_t)$, for each t , then θ_t is guaranteed to converge to a local optimum under the usual stochastic approximation conditions II for decreasing α . They include only a part of its effect on the target. c. we call them semi-gradient

13. i. $E(U_t) = E(\theta_t) = v_t(S_t)$ methods.

i. the gradient-descent version of MC state-value prediction is guaranteed to find a locally optimal solution. here θ_t is MC target. In some cases such as the linear case semi-gradient methods do converge reliably.

14.4 advantages

14. i. bootstrapping target isn't independent of θ_t .
ii. bootstrapping methods will not produce a true gradient-descent method.

- ① significantly faster to learn.
- ② They enable learning to be continual and online.

15. At time t , $\hat{v}_t(S_t, \theta_t)$ At time $t+1$, $\hat{v}_{t+1}(S_{t+1}, \theta_{t+1})$

we can choose state A or state B,
but we may never have this choice
to choose B after time t.

(6) semi-gradient TD(0)

$$U_t = R_{t+1} + \gamma \hat{v}(S_{t+1}, \theta).$$

7. $\nabla(y - f(x, \theta))^2 = -2(y - f(x, \theta))f'(x, \theta)$ MSVE under linear function
- ① for approximate methods, we update θ . approximation if α is reduced
- ② for tabular methods, we update f directly, over time according to the
 $\therefore f'(x, \theta) = 1$.

8. If $v(s, \theta)$ is not continuous about s .
in some sense, then $\hat{v}(s, \theta)$ cannot be converges to a point near learnt.
Same for all AI problem. (ε). the local optimum under linear

9. $\hat{v}(s, \theta) = \theta^T \phi(s) = \sum \theta_i \phi_i(s)$ function approximation.

~~if~~ $\phi_i: S \rightarrow \mathbb{R}$ are called basic functions. 19.8.2

19.2 $\hat{v}(\cdot, \theta)$ is a linear function of θ . $\theta_{t+1} = \theta_t + \alpha(R_{t+1} + \gamma \theta_t^T \phi_{t+1} -$

19.3 s may not be continual, but $\phi(s)$ $\theta_t^T \phi_t) \phi_t$

can still be meaningful if $\phi(s)$ is here $\phi_t = \phi(S_t)$.

continuous in some sense.

$$\therefore \theta_{t+1} = \theta_t + \alpha(R_{t+1} \phi_t - \phi_t)$$

$$\phi_t \rightarrow (\phi_{t+1})^T \phi_t).$$

19.4 $\nabla \hat{v}(s, \theta) = \phi(s)$.

$$\therefore E(\theta_{t+1} | \theta_t) = \theta_t + \alpha(b - A\theta_t)$$

19.5 the linear SGD case is one of the most favorable for mathematical analysis.

$$b = E(R_{t+1} \phi_t)$$

$$A = E(\phi_t (\phi_t - \gamma \phi_{t+1})^T)$$

i. if the system converges,
 $\Rightarrow b - A\theta_{TD} = 0$
 $\Rightarrow \theta_{TD} = A^\top b.$

θ_{TD} is called TD fixpoint.

19.8.3 proof of convergence
of linear TD(0)

19.8.4 In the continuing case:

$$MSVE(\theta_{TD}) \leq \frac{1}{1-\gamma} \min_{\theta} MSVE(\theta).$$

19.9 n-step semi-gradient TD

20. One needs to introduce features
for combinations of feature values
when using linear function approximation
methods.

21. for multi-dimensional continuous state
spaces, function approximation for RL
has much in common with the tasks
of interpolation and regression.

22. represent highly-complex interactions
among the state variables.

\Leftrightarrow allow for more accurate approximations
of more complicated functions.

$$\text{23. } s \in \mathbb{R}^d, \phi_i(s) = \prod_{j=1}^d s_j^{c_{ij}}$$

ϕ_i is polynomial basis, $c_{ij} \in \{0, 1, \dots, N\}, N \geq 0$. (here the function has to

\Rightarrow define $(N+1)^d$ distinct functions
for dimension d. They make up
the order-N polynomial basis.

24. Fourier series express periodic
functions as a weighted sum of
sine and cosine basis functions
of different frequencies.

24.2 f is periodic $\Leftrightarrow f(x) = f(x+T)$
for all x and some period T .

24.3 With enough basis functions
essentially any function can
be approximated as accurately
as desired.

24.4 even function

\Leftrightarrow function that is symmetric
about origin.

odd function \Leftrightarrow function
that ~~is~~ anti-symmetric
about the origin.

24.5 we can represent

any even function with just
the cosine basis functions.

be mathematically well-behaved.)

24.6 The one-dimensional order- N Fourier cosine basis consists of the $N+1$ functions

$$\phi_i(s) = \cos(i\pi s), \quad s \in [0, 1]$$
$$i = 0, 1, \dots, N.$$

24.7 Unlike polynomial basis functions, Fourier basis functions are always bounded and do not require exponentiation.

24.8 For a state space that is the d -dimensional unit hypercube with the origin in one corner, $s \in \mathbb{R}^d$, $s_i \in [0, 1]$.

$$\phi_i(s) = \cos(\pi c^i \cdot s).$$

ϕ_i is order- N Fourier cosine basis, $c^i = (c_1^i, \dots, c_d^i)^T$, $c_j^i \in \{0, \dots, N\}$, $\forall j = 1, \dots, d$, $i = 0, \dots, (N+1)^d$

~ This defines a function for each of the $(N+1)^d$ possible integer vectors c^i .

24.9 Not surprisingly, Fourier basis functions have trouble with discontinuities. which binary features are present

24.10 As is true for polynomial approximation, the number of basis functions in the order- N Fourier cosine basis grows exponentially with the state space dimension.

24.11 advantages of Fourier basis functions are that it is easy to select functions by setting the c^i vectors to account for suspected interactions among the state variables, and by limiting the values in the c^i vectors so that the approximation can filter out high frequency components.

considered to be noise.

24.12 We do not recommend using the polynomial basis

for online learning.

25. ~~Given a state,~~

indicate within which page.

circles the state lies, and thus coarsely code for its location.

Representing a state with features that overlap in this way is known as coarse coding.

25.2 the shape of the features will determine the nature of the generalization.

25.3 Features with large receptive fields

give broad generalization, and are able to make discriminations much finer than the width of the receptive fields, because ultimately the finest discrimination is controlled more by the total number of features, though initial generalization

is indeed controlled by the size and shape of the receptive fields.

26. tile coding works with partitions, ∵ the overall number of features that are active at one time is the same for any state.

26.4 tile coding also gains computational advantages from its use of binary feature vectors.

26.5 generalization occurs to states other than the one trained if those states fall within any of the same tiles, proportional to the number of tiles in common.

26.6 Hashing produces tiles regions randomly spread throughout the state space, but that for multi-dimensional continuous spaces still form an exhaustive partition.

that is flexible and computationally efficient.

26.7 hashing frees us from the curse of dimensionality

26.2 It may be the most practical feature representation for modern sequential digital computers. In the sense that

memory requirements need not be exponential in the number of dimensions, but need merely match the real demands of the task.

realm of nonlinear function approximators. Nonlinear methods may be able to fit target functions much more precisely.

27. RBF

radial basis functions (RBFs) are the natural generalization of coarse coding to continuous-valued features.

27.2 a typical RBF feature :

$$\phi_i(s) = \exp\left(-\frac{\|s - c_i\|^2}{2\sigma_i^2}\right)$$

27.3 The primary advantage of RBFs over binary features is that they produce approximate functions that vary smoothly and are differentiable. Although this is appealing, in most cases it has no practical significance.

27.4 An RBF network is a linear function approximator using RBFs for its features.

27.5 some learning methods for RBF networks exclusively on change the centers and widths of the features as well, bringing them into the

28. ANNs

28.2 ANNs are widely used for nonlinear function approximation. 28.3 both feedforward and recurrent ANNs have been used in rl.

28.4 an ANN with a single

hidden layer having a large enough finite number of sigmoid units can approximate any continuous function on a compact

region of the network's input space to any degree of accuracy.

Here nonlinearity is essential.

28.5 creating the hierarchical representations without relying

hand-crafted features has been an enduring

challenge for AI. This is why learning algorithms for ANNs with hidden layers have received so much attention over the years.

28.6 backpropagation algorithm consists of alternating forward and backward passes through the network. Each forward pass computes the activation of each unit given the current activations of the network's input units. After each forward pass, a backward pass efficiently computes the partial derivative for each weight.

28.7 We can train ANNs using RL principles instead of backpropagation. These methods are less efficient than bp algorithm, but they may be closer to how real neural networks learn.

28.8 bp algorithm does not work well for deeper ANNs. It requires computation proportional to the square of the number of weights,

28.9 A particularly effective method for reducing overfitting by deep ANNs is the dropout method. whereas all other methods are of complexity linear in the number of weights.

28.10 Because of its special architecture, a deep convolutional network can be trained by bp algorithm.

28.11 Units in the same feature map share the same weights. In CNN. This means that a feature map detects the same feature no matter where it is located in the input array.

28.12 The subsampling layers reduce the network's sensitivity to the spatial locations of the features detected, that is, they help make the CNN's responses spatially invariant.

29. LSTD.

29.2 Least-Squares TD algorithm of linear TD(0), but it is also much more expensive computationally.

Chapter 10

有的人先聪明，有的人后聪明。

On-policy control with approximation. $\hat{Q}_t^{(n)} = Q_t$, if $t+n \geq T$.

1. $S_t, A_t \rightarrow U_t$

2.3 n-step update equation:

The target U_t can be any approximation. $\theta_{t+1} = \theta_{t+1} + \alpha(\hat{Q}_t^{(n)} - \hat{q}(S_t, A_t, \theta_{t+1}))$ of $q_\pi(S_t, A_t)$.
 $\nabla \hat{q}(S_t, A_t, \theta_{t+1})$, $0 \leq t < T$.

1.2 the general gradient-descent update 3. three classical settings for
for action-value prediction is formulating the goal in MDP.

$\theta_{t+1} = \theta_t + \alpha(U_t - \hat{q}(S_t, A_t, \theta_t)) \nabla \hat{q}(S_t, A_t, \theta_t)$ ① episodic setting

② discounted setting

$\theta_{t+1} = \theta_t + \alpha(R_{t+1} + \gamma \hat{q}(S_{t+1}, A_{t+1}, \theta_t) - \hat{q}(S_t, A_t, \theta_t)) \nabla \hat{q}(S_t, A_t, \theta_t)$ ③ average reward

we call this method episodic semi-gradient one-step Sarsa. setting.

1.4 To form control methods, we need to 4. In the average-reward setting,
couple such action-value prediction methods the quality of a policy π
with techniques for policy improvement and is defined as the average
action selection. Suitable techniques applicable rate of reward while following
to continuous actions, or to actions from that policy.

(large discrete sets, are a topic of ongoing research with as yet no clear resolution.)

4.2 $\eta(\pi) = \lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T E(R_t | A_{0:t-1} \sim \pi)$
 $= \lim_{T \rightarrow \infty} E(R_t | A_{0:t-1} \sim \pi)$
 $= \sum_s d_\pi(s) \sum_a \pi(a|s) \sum_{s',r} p(s',r|s,a)r$

2. n-step semi-gradient Sarsa

2.2 n-step return.

4.3 $d_\pi(s)$ is the steady-state distribution

$G_t^{(n)} = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{n-1} R_{t+n} + \gamma^n \hat{q}(S_{t+n}, A_{t+n}, \theta_{t+n-1})$ which is

$n \geq 1$, $0 \leq t < T-n$

$d_\pi(s) = \lim_{T \rightarrow \infty} \Pr_{\pi}(S_T = s | A_{0:T-1} \sim \pi)$,

今日，初次见面，人比画美。

which is assumed to exist and to be independent of s_0 . This property is known as ergodicity. It means that where the MDP starts or any early decision made by the agent can have only a temporary effect; in the long run your expectation of being in a state depends only on the policy and the MDP transition probabilities.

4.4 Ergodicity is sufficient to guarantee the existence of the limits in the equations in 4.2.

4.5 We consider all policies that attain the maximal value of $\eta(\pi)$ to be optimal.

4.6 steady state distribution is the special distribution under which, if you select actions according to π , you remain in the same distribution.

$$\sum_s d_\pi(s) \sum_a \pi(a|s, \theta) p(s'|s, a) = d_\pi(s')$$

4.7 differential return:

$$A_t = R_{t+1} - \bar{R}_t + \hat{q}(S_{t+1}, A_{t+1}, \theta)$$

4.8 differential value functions:

$$V_\pi(s) \doteq E_\pi(A_t | S_t = s)$$

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

4.12 differential semi-gradient Sarsa for control:

Learn how to play and ~~speak~~ talk. Improve myself.

→ 3 days, still keep in touch.

$$\delta \leftarrow R - \bar{R} + \hat{q}(S', A', \theta) - \hat{q}(S, A, \theta) \quad 6.5. \quad \bar{R} \leftarrow \bar{R} + \beta S,$$

$$\bar{R} \leftarrow \bar{R} + \beta S,$$

why use δ , maybe $\delta \approx R - \bar{R}$?

$$\theta \leftarrow \theta + \alpha \delta \nabla \hat{q}(S, A, \theta)$$

8.3.2. In the approximate case most

here $\bar{R} \leftarrow \bar{R} + \beta S$ may be wrong, (2)
it should be $\bar{R} \leftarrow \bar{R} + \beta(R - \bar{R})$. by a value function.

5. deprecating the discounted setting

8.

8.2. tabular cases { episodic case
continuing case

$$5.2. J(\pi) = \sum_S d_\pi(s) v_\pi^T(s)$$

$$= \sum_S d_\pi(s) \sum_a \pi(a|s) \sum_{S'} \sum_r p(s', r|s, a) (r + \gamma v_\pi^T(s')) \quad \text{approximate case} \quad \text{episodic case}$$

$$= \eta(\pi) + \gamma \sum_S v_\pi^T(s') \sum_S d_\pi(s) \sum_a \pi(a|s) p(s'|s, a) \quad \text{continuing case}$$

$$= \eta(\pi) + \gamma \sum_S v_\pi^T(s') d_\pi(s')$$

$$= \eta(\pi) + \gamma J(\pi)$$

$$\therefore J(\pi) = \frac{1}{1-\gamma} \eta(\pi) \quad \begin{matrix} \downarrow \text{return} \\ \rightarrow \text{time} \end{matrix} \quad \dots$$

$\therefore J(\pi)$ orders policies identically to the undiscounted objective.

8.3.1 for the continuing case we have

to introduce the discounted formulation cannot be carried over to control in the presence of approximations.

6. n-step differential semi-gradient Sarsa

8.3.3 The arbitrary policies

$$6.2. C_t^{(n)} \doteq R_{t+1} - \bar{R} + R_{t+2} - \bar{R} + \dots + R_{t+n} - \bar{R} \quad \text{that remain need to be ranked,}$$

$$+ \hat{q}(S_{t+n}, A_{t+n}, \theta), \quad t+n < T, n \geq 1. \quad \text{and the } \eta(\pi) \text{ provides an}$$

$$C_t^{(n)} \doteq C_t, \quad t+n \geq T.$$

effective way to do this.

6.3 n-step TD error

$$\delta_t \doteq C_t^{(n)} - \hat{q}(S_t, A_t, \theta).$$

7. The ordering of all policies in the average discounted return setting would be exactly the same as in the average-reward setting.
page 3.

$$6.4. \quad \theta_{t+1} \doteq \theta_t + \alpha \delta_t \nabla \hat{q}(S_t, A_t, \theta_t)$$

Chapter 11

今日仍和计算机系

off-policy methods with approximation where $\hat{S}_t = R_{t+1} + \gamma \sum_a \pi(a|S_{t+1}) \hat{q}(S_{t+1}, a, \theta_t)$

1. explore the convergence problem. $-\hat{q}(S_t, A_t, \theta_t)$

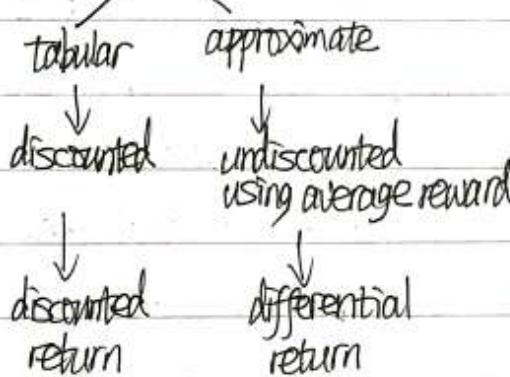
2. tabular case corresponds to a special (episodic)

case of function approximation.

or

$S_t = R_{t+1} - \bar{R}_t + \sum_a \pi(a|S_{t+1}) \hat{q}(S_{t+1}, a, \theta_t)$
 $-\hat{q}(S_t, A_t, \theta_t)$ (continuing)

3. episodic \ continuing



4. semi-gradient off-policy TD(0)

$$\theta_{t+1} = \theta_t + \alpha \delta_t \nabla \hat{v}(S_t, \theta_t)$$

$$\text{where } \delta_t = \frac{\pi(A_t|S_t)}{M(A_t|S_t)}$$

$$S_t = R_{t+1} + \gamma \hat{v}(S_{t+1}, \theta_t) - \hat{v}(S_t, \theta_t)$$

when episodic and potentially discounted

5.2 This algorithm does not use importance sampling. In the tabular case it is clear that this is appropriate, L, but with function approximation it is a judgement call. (ε)

6. n-step version of semi-gradient Expected Sarsa

$$\theta_{t+n} = \theta_{t+n-1} + \alpha \delta_{t+n} \nabla \hat{v}(S_t, \theta_{t+n-1})$$

$$(C_t^{(n)} - \hat{q}(S_t, A_t, \theta_{t+n-1})) \nabla \hat{q}(S_t, A_t, \theta_{t+n-1})$$

$$S_t = R_{t+1} - \bar{R}_t + \hat{v}(S_{t+1}, \theta_t) - \hat{v}(S_t, \theta_t) \text{ where } C_t^{(n)} = R_{t+1+n} + \gamma^n R_{t+n}$$

when continuing and undiscounted using average reward.

$$+ \gamma^n \sum_a \pi(a|S_{t+n}) \hat{q}(S_{t+n}, a, \theta_{t+n-1})$$

(episodic)

5. semi-gradient Expected Sarsa

$$C_t^{(n)} = R_{t+1} - \bar{R}_{t+n} + R_{t+n} - \bar{R}_{t+n-1}$$

$$+ \sum_a \pi(a|S_{t+n}) \hat{q}(S_{t+n}, a, \theta_{t+n-1})$$

(continuing)

6.2 for episodes, $p_t = 1$, if $t \geq T$,
 $\hat{Q}_t^{(n)} = Q_t$, if $t+n > T$.

7. n-step semi-gradient tree-backup algorithm

$$\theta_{t+n} = \theta_{t+n-1} + \alpha (\hat{Q}_t^{(n)} - \hat{q}(s_t, a_t, \theta_{t+n-1})) \cdot \nabla \hat{q}(s_t, a_t, \theta_{t+n-1}),$$

where $\hat{Q}_t^{(n)} = \hat{q}(s_t, a_t, \theta_n) + \sum_{k=t}^{t+n-1} \gamma^k \sum_{i=t}^k r_i(a_i | s_i)$ in this case, but there has been no theoretical analysis.

with \hat{q} as defined for Expected Sarsa (5), 10.4 stability is guaranteed for $\hat{Q}_t^{(n)} = Q_t$ if $t+n > T$, and $\gamma = 1$ in the function approximation methods continuing case.

8. n-step semi-gradient $Q(\pi)$

9. episode terminates on all transitions with 1% probability, much like a $\gamma = 0.99$ discount rate

10. Baird's counterexample

10.2 Even the simplest combination

of bootstrapping and function approximation can be unstable if the backups are not done according to the on-policy distribution.

10.3 It may be possible to guarantee convergence of Q-learning as long as

the behavior policy is sufficiently

11.3 Any two of these three is fine; the danger arises only in the presence of all three.

close to the estimation policy, for example, when it is the ϵ -greedy policy. To the best of our knowledge, Q-learning has never been found to diverge.

that do not extrapolate from the observed targets. These methods, called averagers, include nearest neighbor methods and local weighted regression, but not popular methods such as tile coding and backpropagation.

11. The deadly triad.

11.2 ① training on a distribution of transitions other than that naturally generated by the process whose expectation is being estimated. (e.g., off-policy learning)

② scalable function approximation. (e.g., linear semi-gradient).

③ bootstrapping (e.g., DP, TD learning)

Chapter 12

IB 和 TT 教案

eligibility traces

of the episode. In addition learning

1. TD(λ), λ refers to the use of an eligibility trace.

can occur and affect behavior immediately after a state is

2. eligibility traces unify and generalize TD and MC methods, producing a family delayed n steps.

of methods spanning a spectrum with MC at one end ($\lambda=1$) and one-step forward views: looking forward from the updated state.

TD at the other ($\lambda=0$).

backward views: looking backward

2.2 eligibility traces also provide a way to recently visited states ~~using~~.

of implementing MC online and on continuing problems without episodes. somewhat complex to implement.

3. ① short-term memory vector, 6. tabular and state aggregation
the eligibility trace $e_t \in \mathbb{R}^n$. ~~cases~~ are special cases of

② long-term weight vector $\theta_t \in \mathbb{R}^n$ linear function approximation.

③ trace-decay parameter $\lambda \in [0, 1]$.

7.

4. The primary computational advantage of 7.2 n -step return
eligibility traces over n -step methods is $G_t^{(n)} = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{n-1} R_{t+n}$

that only a single trace vector is required $+ \gamma^n \hat{v}(S_{t+n}, \theta_{t+n-1})$,

rather than a store of the last n feature $0 \leq t \leq T-n$

vectors. Learning also occurs continually $n \geq 1$.

and uniformly in time rather than being 7.3 average of n -step returns
delayed and then catching up at the end possesses an error reduction

property similar to that of individual n-step returns and thus can be used to construct backups with guaranteed convergence properties.

7.4 averaging produces a substantial new range of algorithms. For example,

one could average one-step and infinite-step returns to obtain another way of

interrelating TD and MC. In principle, $\Rightarrow f \uparrow \downarrow$, if $\lambda \uparrow$, n fixed.

one could even average experience-based backups with DP backups to get

a simple combination of experience-based and model-based methods.

7.5 ensemble methods \longleftrightarrow averaging ε .

7.6 A backup that averages simpler component backups is called a compound backup.

8. TD(λ) algorithm can be understood as one particular way of averaging n-step backups.

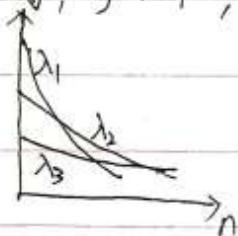
8.2 λ -return

$$G_t^\lambda = (1-\lambda) \sum_{n=1}^{\infty} \lambda^{n-1} G_t^{(n)}$$

8.3 If $\lambda=0$, then the overall backup reduces to its first component, the one-step TD backup, whereas if $\lambda=1$, then the overall backup reduces to its last component, the MC backup.

$$8.4 f(\lambda, n) = (1-\lambda)\lambda^n$$

$$\lambda \in [0, 1], n \geq 1, n \in \mathbb{Z}$$



$$\lambda_1 < \lambda_2 < \lambda_3$$

8.5 If we want, we can separate these post-termination terms from the main sum, yielding

$$G_t^\lambda = (1-\lambda) \sum_{n=1}^{T-1} \lambda^{n-1} G_t^{(n)} + \lambda^{T-1} G_T$$

9. a rate factor such as λ is sometimes an awkward way of characterizing the speed of the decay.

Lily 唱的生日歌 너무 좋다.

For some purposes it is better to specify a time constant, or half-life, T_λ , the time by which the weighting sequence will have fallen to half of its initial value.

10. off-line λ -return algorithm.

10.2 As an off-line algorithm, it makes no changes to the weight vector during the episode. Then, at the end of the episode, a whole sequence of off-line updates are made according to our usual semi-gradient rule, using the λ -return as the target:

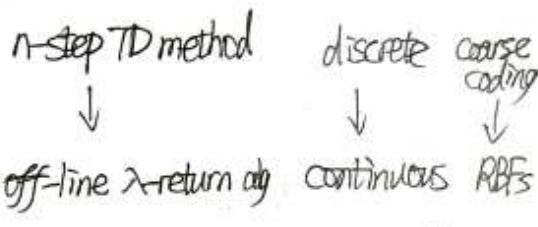
$$\theta_{t+1} = \theta_t + \alpha (G_t^\lambda - \hat{v}(s_t, \theta_t)) \nabla \hat{v}(s_t, \theta_t), \quad t=0, \dots, T-1.$$

10.3 The λ -return gives us an alternative way of moving smoothly between MC and one-step TD methods that can be compared with the n-step TD way.

Chapter 10.4

the overall performance of the offline λ -return algorithms is comparable to that of the n-step algorithms.

10.5



10.6 The approach is forward view of a learning algorithm. We decide how to update each state by looking forward to future rewards and states.

11. TD(λ) is one of the oldest and most widely used algorithms in rl. It was the first algorithm for which a formal relationship was shown between a more theoretical forward view and a

more computational congenial backward view using eligibility traces.

11.2 TD(λ) improves over the off-line λ -return algorithm in

three ways. First it updates the weight vector on every

step of an episode rather than only at the end, ~~this~~. Second, its computations are equally distributed in time rather than all at the end of the episode. And third, it can be applied to continuing problems rather than just episodic problems.

$$11.3 \quad e_0 = 0, \quad e_t = \nabla \hat{v}(s_t, a_t) + \gamma \lambda e_{t-1}$$

where γ is the discount rate.

11.4 the trace is said to indicate the eligibility of each component of the weight vector for undergoing learning changes should a reinforcing event occur. The reinforcing events are TD errors, one step

$$\delta_t = R_{th} + \gamma \hat{v}(s_{th}, a_t) - \hat{v}(s_t, a_t)$$

$$11.5 \quad \theta_{t+1} = \theta_t + \alpha \delta_t e_t$$

11.6 eligibility traces \leftrightarrow momentum.
It has the same idea ^{in form} as momentum, but it has further explanation.

11.7 TD(0) is oriented backward in time. At each moment we look at the current TD error and assign it backward to each prior state according to how much that state contributed to the current eligibility trace at that time.

$$11.8 \quad \text{if } \lambda=0, \Rightarrow e_t = \nabla \hat{v}(s_t, a_t)$$

\therefore 11.5 reduces to the one-step semi-gradient TD update.

\therefore we call that TD(0).

11.9 if $\lambda=1$, then the credit given to earlier states falls only by γ per step. This turns out to be just the right thing to do to achieve MC behavior.

11.10 If $\lambda=1$ and $\gamma=1$, then the eligibility traces do not decay at all with time. In this case, the method behaves like MC for undiscounted, episodic task.

11.11 TD(1) is a way of implementing MC that is more general and that significantly increases their range of applicability. Whereas the earlier MC were limited to episodic tasks, TD(1) can be applied to discounted continuing tasks as well. Moreover, TD(1) can be performed incrementally and on-line.

11.12 If α is chosen larger, λ -return alg is only a little worse whereas TD(λ) is much worse and may even be unstable. This is a weakness of TD(λ).

12. online λ -return algorithm

12.2 h-truncated λ -return

$$C_t^{(\lambda)h} \doteq (1-\lambda) \sum_{n=1}^{h-t+1} \lambda^n C_t^{(n)} + \lambda^{h-t+1} C_t^{(h-t)}$$

$$0 \leq t < h \leq T$$

$$12.3 \quad \theta_{t+1}^h \doteq \theta_t^h + \alpha (C_t^{(\lambda)h} - \hat{v}(s_t, \theta_t^h)) \nabla \hat{v}(s_t, \theta_t^h), \text{ B. true online TD}(\lambda)$$

$$\forall 0 \leq t < h \leq T$$

$$12.4 \quad \theta_0^h \doteq \theta_0, \quad \theta_h \doteq \theta_h^h, \quad \theta_T \doteq \theta_T^h, \quad \theta_t \doteq \theta_t^h$$

12.5 h is horizon.

12.6 the online λ -return alg is fully online, determining a new weight vector θ_t at each step t during an episode, using only info available at time t . Its main drawback is that it is computationally complex, passing over the entire episode so far on every step. L

Note that it is strictly more complex than the off-line λ -return alg, which passes through all the steps at the time of termination but does not make any updates during the episode.

The online alg can be expected to perform better than the off-line one.

13.2 It is "truer" to the

idea of the online TD(λ) alg,

triver even than the TD(λ) alg itself. 13.5 This alg has been
 13.3 The sequence of weight vectors proven to produce exactly
 produced by the on-line λ -return alg the same sequence of weight
 can be arranged in a triangle:
 θ_0^0
 $\theta_0^1 \theta_1^1$
 $\theta_0^2 \theta_1^2 \theta_2^2$
 \dots
 $\theta_0^T \theta_1^T \theta_2^T \dots \theta_T^T$

Really only the weight vectors
 on the diagonal, the θ_t^t , need to
 be produced by the alg. The strategy
 is to find a compact, efficient way
 of computing each θ_t^t from the one
 before.

13.4 for linear case in which $\hat{v}(s, \theta) = \theta^T \phi(s)$, the episode end that there is
 the true ~~online~~ TD(λ) alg is:

$$\theta_{t+1} = \theta_t + \alpha \delta_t e_t + \alpha (\theta_t^T \phi_t - \phi_{t-1}^T \phi_t) (e_t - \phi_t)$$

where $\phi_t = \phi(S_t)$,

$$e_t = r \lambda e_{t-1} + (1 - \alpha r \lambda) e_{t-1}^T \phi_t \phi_t$$

15.5 a_t and e_t are updated on
 each time step $t < T$ and then,
 at time T when G is observed,
 they are used to compute θ_T . Page 6.

14. eligibility trace in the true
 online TD(λ) is called a
 dutch trace.

The trace used in TD(λ) is
 called an accumulating trace.

15. Dutch traces in MC learning
 15.2 $\theta_{t+1} = \theta_t + \alpha (G - \theta_t^T \phi_t) \phi_t$
 $0 \leq t < T$

We assume return G is a single
 reward received at the end of

$$15.3 \quad \theta_T = \theta_{T-1} + \alpha (G - \theta_{T-1}^T \phi_{T-1}) \phi_{T-1}$$

$$= F_{T-1} \phi_{T-1} + \alpha G \phi_{T-1}$$

where $F_t = I - \alpha \phi_t \phi_t^T$ is a forgetting,

$$15.4 \quad \theta_T = a_{T-1} + \alpha G e_{T-1}$$

$$\text{where } a_{T-1} = F_{T-1} F_{T-2} \cdots F_0 \theta_0$$

$$G_{T-1} = \sum_{k=0}^{T-1} F_{T-1} F_{T-2} \cdots F_{k+1} \phi_k = e_{T-1} + (1 - \alpha e_{T-1}^T \phi_{T-1}) \phi_{T-1}$$

Chapter 13

Policy Gradient Methods

a state-value function.

1.① methods that learn the values of actions 2.4 episodic case: $\eta(\theta) = v_{\pi_\theta}(s_0)$ and then select actions based on their continuing case: $\eta(\theta) = r(\theta)$ estimated action values. where $v_{\pi_\theta}(s_0)$ is the value of the start state under π_θ .

② methods that learn a parameterized policy that can select actions without consulting a value function. $r(\theta)$ is average reward rate.

1.3 $\theta \in \mathbb{R}^n$ policy weight vector 3.1 policy approximation.
ω : value functions weight, as in $\hat{v}(s, a)$ generally require that

2. policy gradient methods. $\pi(a|s, \theta) \in (0, 1), \forall a, s, \theta$.

$\theta_{t+1} = \theta_t + \alpha \widehat{\nabla \eta(\theta_t)}$ 3.2 policy-based methods offer useful ways of dealing with where $\widehat{\nabla \eta(\theta_t)}$ is a stochastic estimate whose expectation approximates the gradient of the performance measure with respect to its argument θ_t . continuous action spaces.

2.2 $\eta(\theta)$ is performance measure with respect to the policy weights. 3.3 If the action space is discrete and not too large, then a natural kind of parameterization is to form

2.3 methods that learn approximations to both policy and value functions are often called actor-critic methods, where 'actor' is a reference to the learned policy, and 'critic' refers to the learned value function, usually

$$3.5 \quad \pi(a|s, \theta) = \frac{\exp(h(s, a, \theta))}{\sum_b \exp(h(s, b, \theta))}$$

- 3.6 ① An immediate advantage of selecting actions according to the softmax in action preferences (3.5) is that the approximate policy can approach determinism, whereas with ϵ -greedy action selection over action values there is always an ϵ probability of selecting a random action.
- 3.6 ② Perhaps the simplest advantage that policy parameterization may have over action-value parameterization is that the policy may be a simpler function to approximate.
- 3.6 ③ Action-value methods have no natural way of finding stochastic optimal policies, whereas policy approximating methods can. This is a third significant advantage of policy-based methods.
- 3.6 ④ the choice of policy parameterization is sometimes a good way of injecting prior knowledge about the desired form of the policy into the RL system.
4. Policy Gradient Theorem
- 4.1 assume that every episode starts in some particular (non-random) state s_0 . Then, in the episodic case we define performance as $v(\theta) \doteq v_{\pi_\theta}(s_0)$
- 4.2 policy gradient theorem:
- $$\nabla v_\pi(s) = \sum_a d_\pi(s) q_\pi(s, a) \nabla_\theta \pi(a|s, \theta)$$
- 4.4 proof:
- $$\begin{aligned} \nabla v_\pi(s) &= \nabla \left(\sum_a \pi(a|s) q_\pi(s, a) \right) \\ &= \sum_a (\nabla \pi(a|s) q_\pi(s, a) + \pi(a|s) \nabla q_\pi(s, a)) \\ &= \sum_a (\nabla \pi(a|s) q_\pi(s, a) + \pi(a|s) \nabla \sum_{s', r} p(s', r|s, a) (r + \gamma v_\pi(s'))) \\ &= \sum_a \nabla \pi(a|s) q_\pi(s, a) + \sum_a \pi(a|s) \sum_{s', r} p(s'|s, a) \nabla v_\pi(s') \\ &= \sum_a \nabla \pi(a|s) q_\pi(s, a) + \gamma \sum_{s'} p(s \rightarrow s') \sum_a \nabla \pi(a'|s') q_\pi(s', a') \\ &\quad + \gamma^2 \sum_{s''} p(s \rightarrow s'') \sum_{a''} \nabla \pi(a''|s'') q_\pi(s'', a'') \\ &\quad + \dots \end{aligned}$$
- $$\begin{aligned} &= \sum_{x \in S} \sum_{k=0}^{\infty} \gamma^k p(s \rightarrow x | k) \sum_a \nabla \pi(a|x) q_\pi(x, a) \\ &= \sum_{x \in S} d_\pi(x) \sum_a \nabla \pi(a|x) q_\pi(x, a) \end{aligned}$$
- OK

顶天立地，立地之选择，立地→顶天

5. REINFORCE

$$5.2 \nabla \eta(\theta) = \sum_s d_\pi(s) \sum_a q_\pi(s, a) \nabla_\theta \pi(a|s, \theta)$$

as the performance gradient.
(Not proved).

$$\therefore \nabla \eta(\theta) = E_\pi \left(r^t \sum_a q_\pi(s_t, a) \nabla_\theta \pi(a|s_t, \theta) \right) \quad (2)$$

$$\begin{aligned} \therefore \nabla \eta(\theta) &= E_\pi \left(r^t \sum_a \pi(a|s_t, \theta) q_\pi(s_t, a) \frac{\nabla_\theta \pi(a|s_t, \theta)}{\pi(a|s_t, \theta)} \right) \text{ thus slow to learn.} \\ &= E_\pi \left(r^t q_\pi(s_t, A_t) \frac{\nabla_\theta \pi(A_t|s_t, \theta)}{\pi(A_t|s_t, \theta)} \right) \\ &= E_\pi \left(r^t A_t \frac{\nabla_\theta \pi(A_t|s_t, \theta)}{\pi(A_t|s_t, \theta)} \right) \\ (\because E_\pi(A_t|s_t, A_t) &= q_\pi(s_t, A_t)) \end{aligned}$$

5.3 REINFORCE:

$$\theta_{t+1} = \theta_t + \alpha r^t A_t \frac{\nabla_\theta \pi(A_t|s_t, \theta)}{\pi(A_t|s_t, \theta)}$$

5.4 REINFORCE uses the complete return from time t . In this sense REINFORCE is a MC algorithm and is well defined only for the episodic case.

5.5 $\frac{\nabla_\theta \pi(A_t|s_t, \theta)}{\pi(A_t|s_t, \theta)}$ has been given several names and notations in the literature. We will refer to it simply as eligibility vector.

5.6. the expected update over an episode is in the same direction

5.7 As a MC method it may be of high variance and

6. REINFORCE with Baseline.

$$6.2 \therefore \sum_a b(s) \nabla_\theta \pi(a|s, \theta) = b(s) \nabla_\theta b = 0, \quad \forall s \in S.$$

$$\therefore \nabla \eta(\theta) = \sum_s d_\pi(s) \sum_a (q_\pi(s, a) - b(s)) \nabla_\theta \pi(a|s, \theta)$$

$$6.3 \therefore \theta_{t+1} = \theta_t + \alpha (A_t - b(s_t)) \frac{\nabla_\theta \pi(A_t|s_t, \theta)}{\pi(A_t|s_t, \theta)}$$

6.4 In general, the baseline leaves the expected value of the update unchanged, but it can have a large effect on its variance.

6.5 One natural choice for the baseline is an estimate of the state value,

$\hat{v}(s, w)$, where $w \in \mathbb{R}^m$ is a second learned weight vector.

7. Although the REINFORCE-with-baseline method learns both a policy and a state-value function, we do not consider it to be an actor-critic method because its state-value function is

used only as a baseline, not as a critic. $q_\pi(s, a) \doteq \sum_{t=1}^{\infty} E[R_{t+h} - r(\pi)] | S_t = s, A_t = a$
 That is, it is not used for bootstrapping, 9.5 proof of policy gradient theorem
 but only as a baseline for the state
 being updated. This is a useful
 distinction, for only through bootstrapping
 do we introduce bias and an asymptotic
 dependence on the quality of the
 function approximation.

$$\begin{aligned} \text{if } \nabla v_\pi(s) &= \nabla \left(\sum_a \pi(a|s) q_\pi(s, a) \right), \forall s \in S \\ &= \sum_a (\nabla \pi(a|s) q_\pi(s, a) + \pi(a|s) \nabla \sum_{s', r} p(s'|s, a) (r - r(\theta) + v_\pi(s'))) \end{aligned}$$

$$\begin{aligned} \text{then } \nabla r(\theta) &= \sum_a \nabla \pi(a|s) q_\pi(s, a) + \\ &\quad \sum_a \pi(a|s) \sum_{s'} p(s'|s, a) \nabla v_\pi(s') \\ &\quad - \nabla v_\pi(s), \quad \forall s \in S. \end{aligned}$$

8. Actor-Critic Methods

8.2 One-step actor-critic (episodic)

$$\theta_{t+1} \doteq \theta_t + \alpha (a_t^{(t)} - \hat{v}(s_t, w)) \frac{\nabla \pi(A_t | S_t, \theta)}{\pi(A_t | S_t, \theta)}$$

$$= \theta_t + \alpha (R_{t+1} + \gamma \hat{v}(S_{t+1}, w) - \hat{v}(S_t, w)) \frac{\nabla \pi(A_t | S_t, \theta)}{\pi(A_t | S_t, \theta)} \wedge \nabla \eta(\theta) = \nabla r(\theta)$$

The natural state-value function learning (no γ^t)
 method to pair with this is semi-gradient TD(0).

8.3 actor-critic with eligibility traces (episodic)

9. for continuing problems without episode boundaries,

$$\eta(\theta) \doteq r(\theta) = \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{t=0}^n E(R_t | \theta_0 = \theta_1 = \dots = \theta_{t-1} = \theta)$$

$$= \lim_{t \rightarrow \infty} E(R_t | \theta_0 = \theta_1 = \dots = \theta_{t-1} = \theta)$$

$$= \sum_s d_{\pi_\theta}(s) \sum_a \pi_\theta(a|s, \theta) \sum_{s', r} p(s', r|s, a) r$$

$$= \sum_s d_{\pi_\theta}(s) \sum_a \nabla \pi(a|s) q_{\pi_\theta}(s, a)$$

$$+ \sum_{s'} d_{\pi_\theta}(s') \nabla v_{\pi_\theta}(s')$$

$$- \sum_s d_{\pi_\theta}(s) \nabla v_{\pi_\theta}(s)$$

$$= \sum_s d_{\pi_\theta}(s) \sum_a \nabla \pi(a|s) q_{\pi_\theta}(s, a)$$

9.3 d_π is steady-state distribution under π . 10. policy-based methods offer

$$d_\pi(s) \doteq \lim_{t \rightarrow \infty} P(S_t = s | A_0, t \sim \pi)$$

practical ways of dealing with large

$$9.4 v_\pi(s) \doteq \sum_{t=1}^{\infty} E[R_{t+h} - r(\pi)] | S_t = s$$

actions spaces, even continuous spaces
 with an infinite number of actions.