

Lecture 5. Model free Control

Prediction = $V_{\pi}(s)$

Control = improve policy $\pi^i(s) \geq V_{\pi}(s)$
 $\pi' \geq \pi$

S.1 Model free

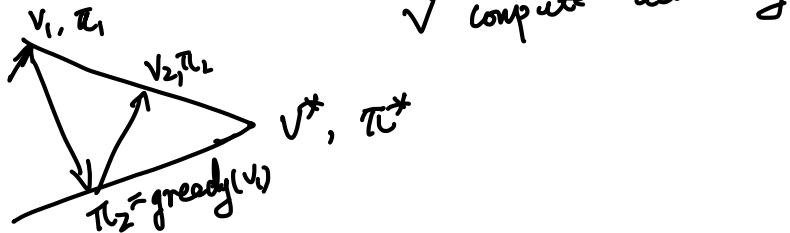
① boat, wave, wind, complicate

② go: know every detail, space too large

S.2 on-policy: learn on job
 learn about π from samples from π

off-policy: learn policy $\pi \neq$ sample policy μ

S.3 use V to impolicy
 V compute accurately



policy iteration

① MC. to evaluate V

② TD to evaluate V

MDP

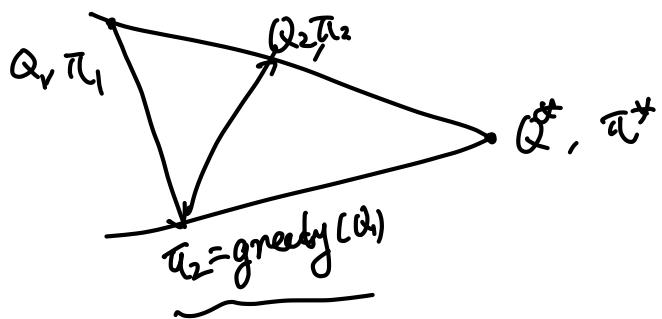
$$\pi'(s) = \arg \max_a R_s^a + \gamma \sum_{s'} P_{ss'}^a V(s')$$

Model free

$$\pi'(s) = \arg \max_a Q(s, a)$$

state-action value function

5.4 Policy iteration with action-value function



① MC

② TD

① 估計
② 面包

example



巨大問題

Greedy:

left Reward 0
 $V(\text{left}) = 0$

right $R: +1$
 $V(\text{right}) = 1$

right $R: +3$
 $V(\text{right}) = \frac{1+3}{2} = 2$

⋮
rights

5.5 ε-greedy exploration

$$\pi(a|s) = \begin{cases} 1-\varepsilon + \frac{\varepsilon}{m} & \text{if } a^* = \underset{a}{\operatorname{argmax}} Q(s, a) \\ \frac{\varepsilon}{m} & \text{others} \end{cases}$$

$1-\varepsilon$ greedy

ε uniform

Theorem ε -greedy policy improvement

ε -greedy π , π' , ε -greedy w.r.t $q_{\bar{\pi}}$

$$\Rightarrow V_{\bar{\pi}'}(s) \geq V_{\pi}(s)$$

$$\text{Pf: } V_{\bar{\pi}'}(s) \geq q_{\bar{\pi}}(s, \bar{\pi}'(s)) = \underbrace{\sum_a \bar{\pi}'(a|s) q_{\bar{\pi}}(s, a)}$$

Lect 3 DP

$$\boxed{V_{\pi}(s) \leq q_{\pi}(s, \pi'(s))} = \mathbb{E}_{\pi'} [R_{t+1} + \gamma v_{\pi}(S_{t+1}) \mid S_t = s] \\ \leq \mathbb{E}_{\pi'} [R_{t+1} + \gamma q_{\pi}(S_{t+1}, \pi'(S_{t+1})) \mid S_t = s] \\ \leq \mathbb{E}_{\pi'} [R_{t+1} + \gamma R_{t+2} + \gamma^2 q_{\pi}(S_{t+2}, \pi'(S_{t+2})) \mid S_t = s] \\ \leq \mathbb{E}_{\pi'} [R_{t+1} + \gamma R_{t+2} + \dots \mid S_t = s] = \underline{v_{\pi'}(s)} \\ V_{\bar{\pi}'}(s) \leq \underline{V_{\bar{\pi}'}(s)}$$

$$V_{\bar{\pi}'}(s) \geq q_{\bar{\pi}}(s, \bar{\pi}'(s)) = \sum_a \bar{\pi}'(a|s) q_{\bar{\pi}}(s, a)$$

$$= \frac{\varepsilon}{m} \sum_a q_{\bar{\pi}}(s, a) + (1-\varepsilon) \max_a q(s, a)$$

$$\sum_a \bar{\pi}(a|s) q_{\bar{\pi}}(s, a) \geq \frac{\varepsilon}{m} \sum_a q_{\bar{\pi}}(s, a) + (1-\varepsilon) \sum_a \frac{\bar{\pi}(a|s) - \frac{\varepsilon}{m}}{1-\varepsilon} q_{\bar{\pi}}(s, a)$$

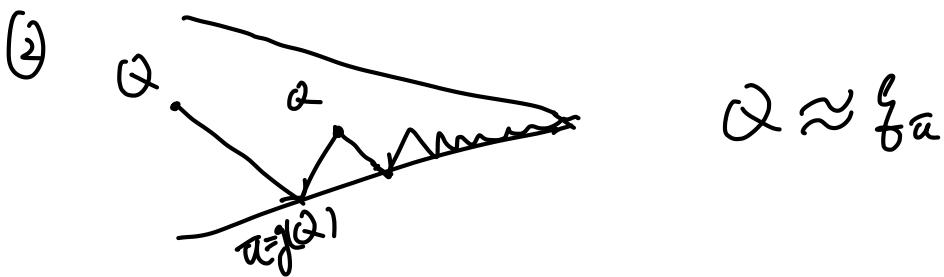
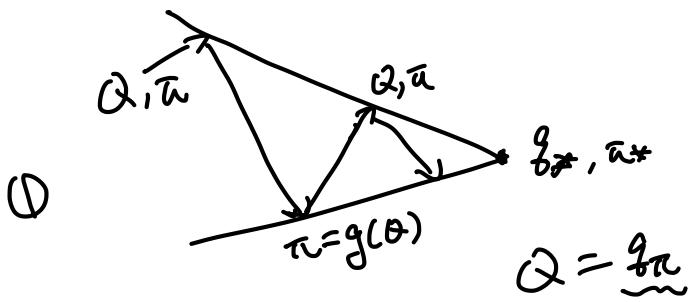
$$\sum_a \frac{\bar{\pi}(a|s) - \frac{\varepsilon}{m}}{1-\varepsilon}$$

$$= \sum_a \bar{\pi}(a|s) q_{\bar{\pi}}(s, a) = V_{\bar{\pi}}(s)$$

$$\stackrel{?}{=} \frac{1-\varepsilon}{1-\varepsilon}$$

$$= 1$$

5.6. MC. policy iteration



5.6.1 Greedy in the limit with infinite exploration
 Definition (GLIE)

① $N_K(s, a) = +\infty$ K : episode
 $K \rightarrow \infty \uparrow \pi$

② $\lim_{K \rightarrow \infty} \pi_K(a|s) = 1$ ($a = \arg \max_{a'} Q_K(s, a')$)

5.6.2 GLIE MC control

① Sample k th episode using $\pi = \{S_1, A_1, R_1, S_2, A_2, R_2\}$

② S_t, A_t .

$$N(S_t, A_t) \leftarrow N(S_t, A_t) + 1$$

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{G_t - Q(S_t, A_t)}{N(S_t, A_t)}$$

③

$$\varepsilon \leftarrow \frac{1}{K}$$

$$\pi \leftarrow \varepsilon\text{-greedy } (\mathbb{Q})$$

Theorem: GLIE MC control converges to the optimal action-value function

$$\mathbb{Q}(s, a) \rightarrow \mathbb{q}_x(s, a)$$

5.7. MC vs TD control

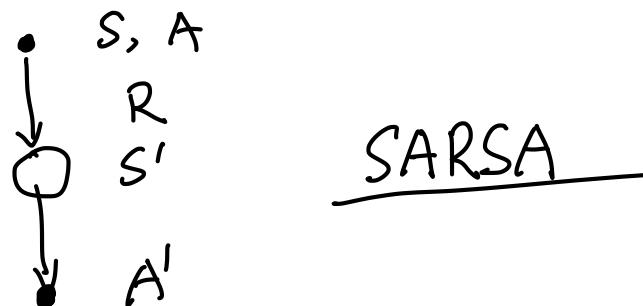
TD: lower variance

= online

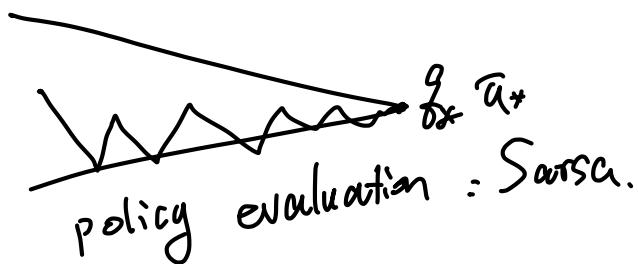
= incomplete sequences

use TD to replace MC

5.7.1



$$Q(s, a) \leftarrow Q(s, a) + \alpha (R + \gamma Q(s', a') - Q(s, a))$$



Sarsa Algorithm for On-Policy Control

```

Initialize  $Q(s, a), \forall s \in \mathcal{S}, a \in \mathcal{A}(s)$ , arbitrarily, and  $Q(\text{terminal-state}, \cdot) = 0$ 
Repeat (for each episode):
    Initialize  $S$ 
    Choose  $A$  from  $S$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)
    Repeat (for each step of episode):
        Take action  $A$ , observe  $R, S'$ 
        Choose  $A'$  from  $S'$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)
         $Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma Q(S', A') - Q(S, A)]$ 
         $S \leftarrow S'; A \leftarrow A'$ ;
    until  $S$  is terminal

```

Theorem convergence of Sarsa

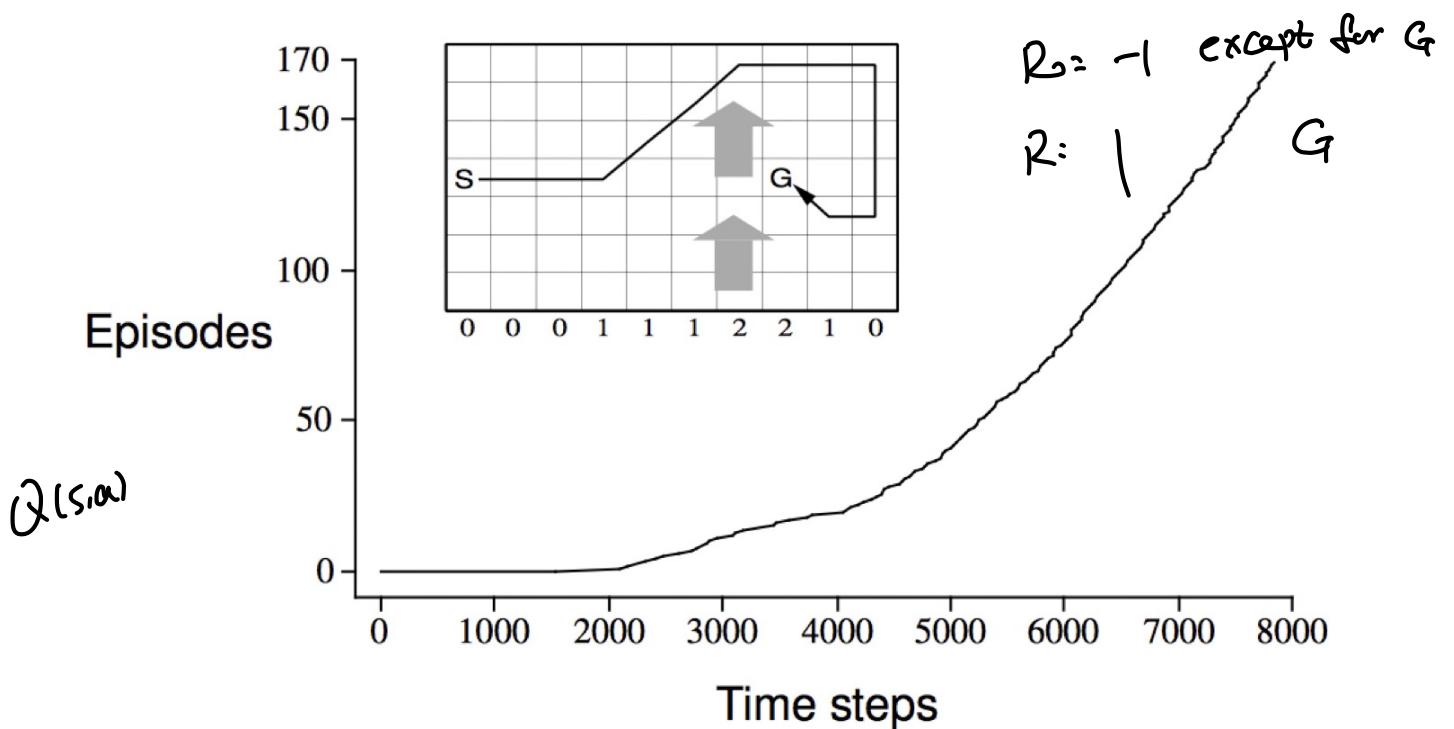
Condition:

① GLIE

$$\textcircled{2} \quad \sum_{t=1}^{\infty} \alpha_t < \infty$$

$$\sum_{t=1}^{\infty} \alpha_t^2 < \infty$$

Sarsa on the Windy Gridworld



5.8. n-step Sarsa.

$$n=1 \text{ (Sarsa)} \quad q_t^{(1)} = R_{t+1} + \gamma Q(S_{t+1})$$

$$n=2 \quad q_t^{(2)} = R_{t+1} + \gamma R_{t+2} + \gamma^2 Q(S_{t+2})$$

⋮

$$n=\infty \text{ MC} \quad q_t^{(\infty)} = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{T-1} R_T$$

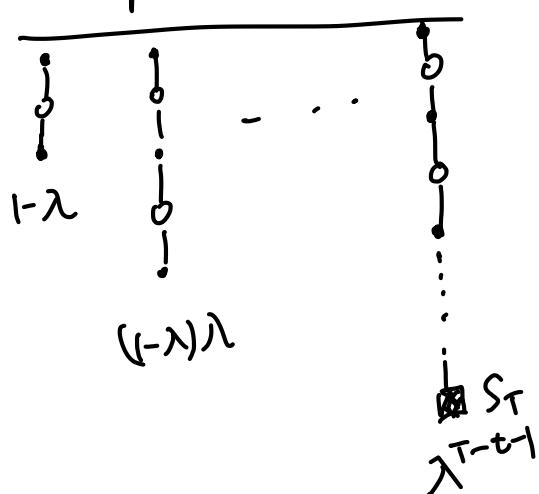
Define n-step Q-return

$$q_t^{(n)} = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{n-1} R_{t+n} + \gamma^n Q(S_{t+n})$$

n-step Sarsa

$$Q(S_t, A_t) = Q(S_t, A_t) + \alpha (q_t^{(n)} - Q(S_t, A_t))$$

5.9. Forward view Sar(a|λ)



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

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha (q_t^{(\lambda)} - Q(S_t, A_t))$$

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

5.10 Backward View Sarsa(λ)

Sarsa(λ) Algorithm

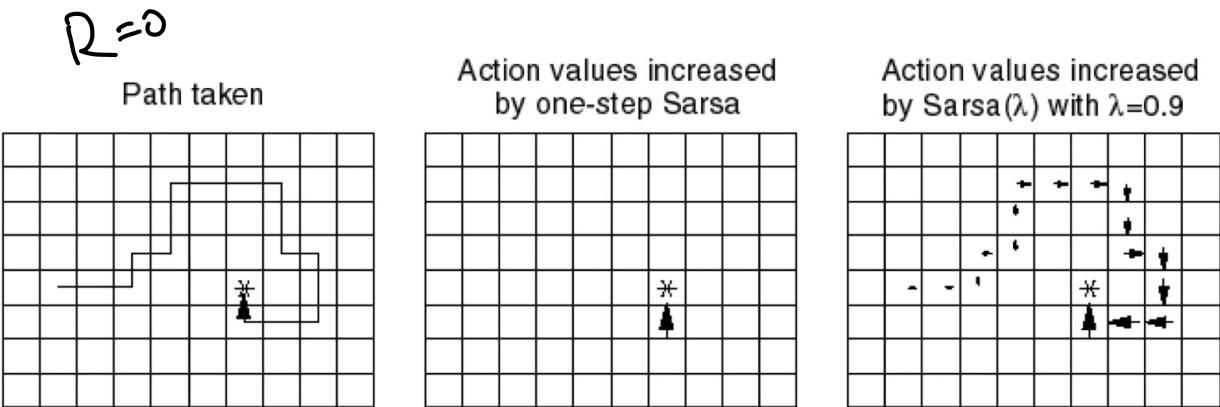
Eligibility Trace

```

Initialize  $Q(s, a)$  arbitrarily, for all  $s \in \mathcal{S}, a \in \mathcal{A}(s)$ 
Repeat (for each episode):
     $E(s, a) = 0$ , for all  $s \in \mathcal{S}, a \in \mathcal{A}(s)$ 
    Initialize  $S, A$ 
    Repeat (for each step of episode):
        Take action  $A$ , observe  $R, S'$ 
        Choose  $A'$  from  $S'$  using policy derived from  $Q$  (e.g.,  $\varepsilon$ -greedy)
         $\delta \leftarrow R + \gamma Q(S', A') - Q(S, A)$ 
         $E(S, A) \leftarrow E(S, A) + 1$ 
        For all  $s \in \mathcal{S}, a \in \mathcal{A}(s)$ :
             $Q(s, a) \leftarrow Q(s, a) + \alpha \delta E(s, a)$ 
             $E(s, a) \leftarrow \gamma \lambda E(s, a)$ 
         $S \leftarrow S'; A \leftarrow A'$ 
    until  $S$  is terminal

```

Sarsa(λ) Gridworld Example



S.11 off-policy learning

target policy $\pi(a|s) \rightarrow V_\pi(s)$ or $Q_\pi(s, a)$

$$\{S_1, A_1, R_2, S_2, A_2, \dots, S_T\} \sim M$$

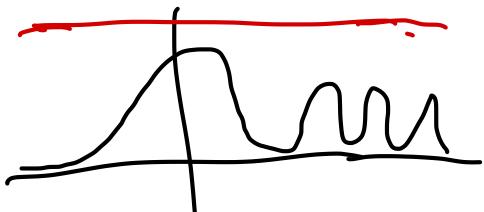
- ① Learn from observing other agents
- ② Re-use experience
- ③ Learn about optimal policy while following exploratory policy
- ④ Train multiple policies

S.12. Importance sampling

$$E_{X \sim p}[f(x)] = \sum p(x)f(x)$$

$$= \sum f(x) \frac{p(x)}{q(x)} q(x)$$

$$= E_{X \sim q(x)} \left[f(x) \frac{p(x)}{q(x)} \right]$$



S.13. Imp. sampling for off-policy

$$G_t \quad \mu(A_t|S_t)$$

$$\mu(A_{t+1}|S_{t+1}) \quad \dots$$

$$G_t^{\text{imp}} = \frac{\pi(A_t|S_t)}{\mu(A_t|S_t)} \frac{\pi(A_{t+1}|S_{t+1})}{\mu(A_{t+1}|S_{t+1})} \dots \frac{\pi(A_T|S_T)}{\mu(A_T|S_T)} G_t$$

$$M_C \quad V(S_t) \leftarrow V(S_t) + \alpha (G_t^{\text{imp}} - V(S_t))$$

$$TD \quad V(S_t) \leftarrow V(S_t) + \alpha \left(\frac{\pi(A_t | S_t)}{d(A_t | S_t)} (R_{t+1} + \gamma V(S_{t+1}) - V(S_t)) \right)$$

5.14 $\pi(S_{t+1}) = \arg \max_{a'} Q(S_{t+1}, a')$

$$R_{t+1} + \gamma Q(S_{t+1}, A')$$

$$= R_{t+1} + \gamma Q(S_{t+1}, \arg \max_{a'} Q(S_{t+1}, a'))$$

$$= R_{t+1} + \gamma \max_a Q(S_{t+1}, a')$$

Q-learning Control: (no policy explicitly)

$$(Q(S, A) \leftarrow Q(S, A) + \alpha(R + \gamma \max_a Q(S', a) - Q(S, A))$$

Q-Learning Algorithm for Off-Policy Control

Initialize $Q(s, a), \forall s \in \mathcal{S}, a \in \mathcal{A}(s)$, arbitrarily, and $Q(\text{terminal-state}, \cdot) = 0$
 Repeat (for each episode):

 Initialize S

 Repeat (for each step of episode):

 Choose A from S using policy derived from Q (e.g., ϵ -greedy)

 Take action A , observe R, S'

$Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) - Q(S, A)]$

$S \leftarrow S'$;

 until S is terminal

Relationship Between DP and TD (2)

Full Backup (DP)	Sample Backup (TD)
Iterative Policy Evaluation $V(s) \leftarrow \mathbb{E}[R + \gamma V(S') s]$	TD Learning $V(S) \xleftarrow{\alpha} R + \gamma V(S')$
Q-Policy Iteration $Q(s, a) \leftarrow \mathbb{E}[R + \gamma Q(S', A') s, a]$	Sarsa $Q(S, A) \xleftarrow{\alpha} R + \gamma Q(S', A')$
Q-Value Iteration $Q(s, a) \leftarrow \mathbb{E} \left[R + \gamma \max_{a' \in \mathcal{A}} Q(S', a') s, a \right]$	Q-Learning $Q(S, A) \xleftarrow{\alpha} R + \gamma \max_{a' \in \mathcal{A}} Q(S', a')$ 值极 .

where $x \xleftarrow{\alpha} y \equiv x \leftarrow x + \alpha(y - x)$