

- MDP
 $s \xrightarrow{\pi} s' \xrightarrow{\pi} s'' \xrightarrow{\pi} \dots$ discount factor

- Bellman equation $V_\pi(s) = E[\gamma r_{t+1} + V_\pi(s') | S=s]$

- DP
 $= \sum_a \pi(a|s) \sum_{s'} p(s', r(s, a)) [r + V_\pi(s')] \rightarrow \text{Bellman equation}$

 $\pi^* = \arg \max_\pi V_\pi(s)$
 $V_{\pi^*}(s) = \max_a q_{\pi^*}(s, a) = \max_a \sum_{s'} p(s', r(s, a)) [r + V_{\pi^*}(s')] \quad \left. \begin{array}{l} \text{Bellman optimality equation} \\ q_{\pi^*}(s, a) = \sum_{s'} p(s', r(s, a)) [r + \max_a f_{\pi^*}(s', a)] \end{array} \right.$

- Model-based method $\xrightarrow{P^I}$ env response: $p(s'|s, a), r$

model-free method \xrightarrow{MC} \xrightarrow{TD}

- Policy iteration $\xrightarrow{\text{PE (estimation)}}$ $\xrightarrow{\text{PI (improvement)}}$

• initial π_0, V_{π_0} for $s \in S$

Loop for $s \in S$:

PE $\left\{ \begin{array}{l} V_{\pi_0}(s) = \sum_a \pi_0(a|s) \sum_{s'} p(s', r(s, a)) [r + V_{\pi_0}(s')] \\ |V_{\pi_0}(s) - V_{\pi_0}(s)| < \epsilon, \rightarrow V_{\pi_0} \text{ is convergent} \end{array} \right.$

• for $s \in S$:

Pi $\left\{ \begin{array}{l} a^* = \arg \max_a q_{\pi_0}(s, a) \\ \text{if } \pi_0(a|s) \neq 1, \rightarrow \pi_0 \text{ not optimal} \rightarrow \pi_0(a|s) = 1 \text{ else } \pi_0(a|s) = 0 \\ \text{if } \pi_0 \text{ is not optimal, go to PE again} \end{array} \right.$

PE \rightarrow Pi \rightarrow PE \rightarrow Pi \dots



- MC

① X env respond $p(s'|s, a)$

[generate experience + avg]

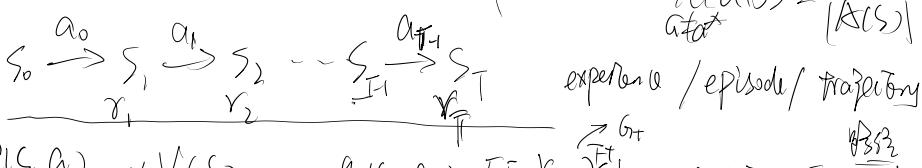
② $S \uparrow, V_{\pi^*} \text{ cost } \uparrow,$

$p(s_t)$, randomly assign start port $(s_0, a_0) \rightarrow$ MC with start port

$\xrightarrow{\text{exploitation, greedy action}} \epsilon = 0.1$

$\xrightarrow{\text{exploration}}$

greedy $a^* = \arg \max_a q(s, a)$
 $\pi(a|s) = 1 - \epsilon + \frac{\epsilon}{|\mathcal{A}(s)|}$
 $\pi(a|s) = \frac{\epsilon}{|\mathcal{A}(s)|}$



$\underline{g_T}$

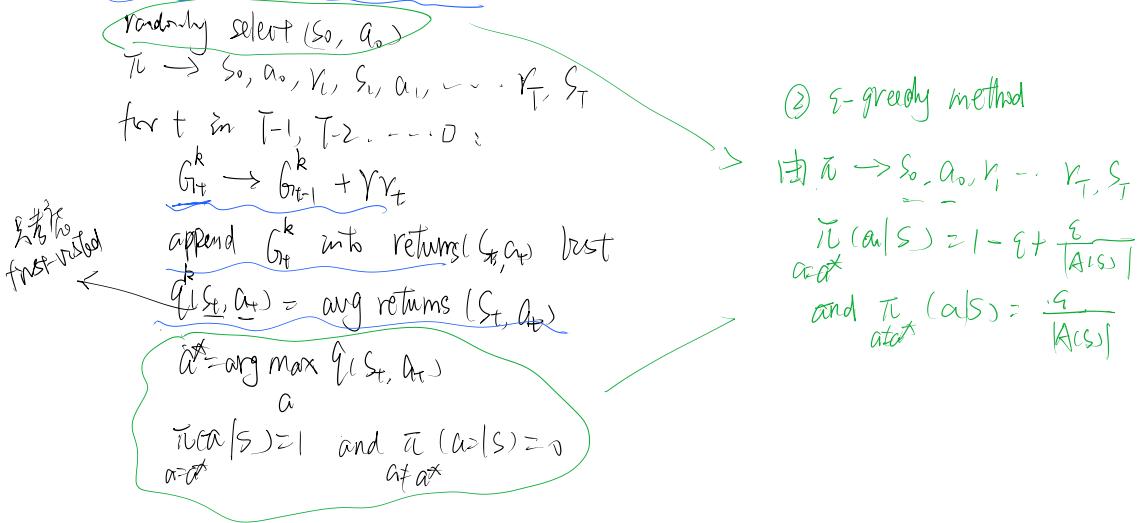
$\underline{g_T}$

$\underline{g_T}$

$$V(S) \rightarrow V(S), \quad Q(S_t, A_t) = E[\underbrace{R_{t+1} + \gamma V(S_{t+1})}_{G_t} | S_t, A_t] = \bar{E}[G_t | S_t, A_t]$$

Version 1 w/ start port

- initialize $\pi, q(s,a)$, returns(s, a) empty list for $s \in S$
- for episode k in $1, 2, \dots, n$:



- incremental implementation

$$\begin{aligned} \text{Same } (S_t, A_t) & \quad G_t^1, G_t^2, \dots, G_t^n \\ q^n(S_t, A_t) &= \frac{G_t^1, G_t^2, \dots, G_t^n}{\sum G_t^k} = \frac{\sum W_k G_t^k}{\sum W_k} \quad \xrightarrow{\text{long}} \\ q^{n+1}(S, a) &= \frac{\sum W_k G_t^k}{\sum W_k} \rightarrow q^{n+1}(S, a) = \frac{q^n(S, a) \sum W_k + G_t^n}{\sum W_k} = \frac{q^n \sum W_k + G_t^n + G_t^n - q^n}{\sum W_k} \\ &= q^n + \frac{W_k}{\sum W_k} (G_t^n - q^n) \end{aligned}$$

Version 2 (w/ incremental implementation)

- initialize $\pi, q(s,a)$, returns(s, a) empty list for $s \in S$
- for episode k in $1, 2, \dots, n$:

$$\pi \rightarrow S_0, A_0, R_1, \dots, S_T, A_T$$

for t in $T-1, T-2, \dots, 1$:

$$G_t^k \rightarrow G_{t+1}^k + \gamma r_t$$

$$q^{k+1}(S_t, A_t) = q^k + \frac{W_k}{\sum W_k} (G_t^k - q^k(S_t, A_t))$$

- TD learning

$$\begin{aligned} q(S_t, A_t) &\leftarrow q(S_t, A_t) + \alpha [G_t - q(S_t, A_t)] \\ V(S_t) &\leftarrow V(S_t) + \alpha [G_t - V(S_t)] \end{aligned} \quad \left. \begin{array}{l} \\ \end{array} \right\}$$

$$V(S_t) \leftarrow V(S_t) + \alpha [R_{t+1} + \gamma V(S_{t+1})] \quad \text{1-step learning} \rightarrow \text{TD}(0)$$

Proof: $V(S_t) = E[G_t | S=S_t] = E[R_{t+1} + \gamma G_{t+1} | S=S_t] \rightarrow \text{TD error}$
 $= E[R_{t+1} + \gamma V(S_{t+1}) | S=S_t] \rightarrow \text{Comb of MC and DP}$

generalization of TD error:

$$G_t - V(S_t) = R_{t+1} + \gamma G_{t+1} - V(S_t) + V(S_{t+1}) - V(S_{t+1}) \\ = \delta_t + \gamma(G_{t+1} - V(S_{t+1})) = \delta_t + \gamma S_{t+1} + \gamma(G_{t+2} - V(S_{t+2})) = \dots = \sum_{k=t}^T \gamma \delta_k$$

where $\delta_t = R_{t+1} + \gamma G_{t+1} - V(S_t)$

property

TD learning variance \downarrow , bias \uparrow compared to MC
 $S_t, S_{t+1}, G_{t+1}, V_{t+1} \rightarrow$ fewer var

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

- initialize $\pi, q(s,a)$ for $s \in S, a \in A, q(s_t, a_t) = 0$
- for every episode:

initialize s_0

$\pi(\cdot | s_0) \rightarrow a_0 \rightarrow (\epsilon\text{-greedy considered})$

for t in $0, 1, 2, \dots, T-1$:

take action $a_t \rightarrow r_{t+1}, s_{t+1}$

$\pi(\cdot | s_{t+1}) \rightarrow a_{t+1}$

$q(s_t, a_t) \leftarrow q(s_t, a_t) + \alpha [r_{t+1} + \gamma q(s_{t+1}, a_{t+1}) - q(s_t, a_t)] \rightarrow \approx \text{Value update}$

$a^* = \arg \max_a q(s_t, a)$

update $\pi(a_t | s_t)$ based on ϵ -greedy

wait until n -step to update $q(s_t, a_t)$

$\rightarrow \approx \text{Value update}$

$\rightarrow \approx \text{Policy update}$

Q-learning

- initialize $\pi, q(s,a)$ for $s \in S, a \in A, q(s_t, a_t) = 0$
- for every episode:

initialize s_0

$\pi(\cdot | s_0) \rightarrow a_0$

for t in $0, 1, 2, \dots, T-1$:

take action $a_t \rightarrow r_{t+1}, s_{t+1}$

$\pi(\cdot | s_{t+1}) \rightarrow a_{t+1}$

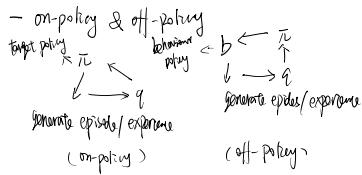
$q(s_t, a_t) \leftarrow q(s_t, a_t) + \alpha [r_{t+1} + \max_a q(s_{t+1}, a_{t+1}) - q(s_t, a_t)]$

update $\pi(a_t | s_t)$ based on ϵ -greedy

$\approx n$ step in (uninitial).

11-step Q-learning

$$q(s_t, a_t) = r_{t+1} + \gamma r_{t+2} + \dots + \gamma^{n-1} r_{t+n} + \gamma^n q(s_{t+n}, a_{t+n}) - q(s_t, a_t)$$



- importance sampling

$$\frac{\prod_{t=0}^T \pi(a_t | s_t) p(s_{t+1} | s_t, a_t)}{\prod_{t=0}^T b(a_t | s_t) p(s_{t+1} | s_t, a_t)} = \frac{\prod_{t=0}^T \pi(a_t | s_t)}{\prod_{t=0}^T b(a_t | s_t)} = P \text{ (importance sampling ratio)}$$