

4) MONTE CARLO
 $S_0, A_0, R_0, S_1, A_1, R_1, \dots$ (episodes) \rightarrow π^*

model based $\left\{ \begin{array}{l} p(s'|s, a) - \text{Probability of next state} \\ p(r|s, a) - \end{array} \right.$

model free $\{ V(s), Q(s, a) \text{ (value)} -$
 $\pi(s, a) \text{ (policy)} -$
 $\pi^* \text{ (control)} \}$

Monte Carlo

(π^* \rightarrow π_{MC}) \rightarrow $\pi_{MC} \approx \pi^*$

$(S_0, A_0, R_0, \dots, S_T, A_T, R_T, S_T)$: episode \rightarrow $\pi_{MC} \approx \pi^*$
 \rightarrow $\pi_{MC} \approx \pi^*$

$$G_t = \sum_{k=t}^T \gamma^{k-t} \cdot r_k$$

$U(s_t) = U(s_t) \cup \{G_t\}$ \rightarrow $U(s_t)$ $\approx \pi^*$

(first visit OR every visit)

$V^\pi(s) = \text{average}(U(s))$ $\approx \pi^*$

\rightarrow π^* \rightarrow V^π \rightarrow π^* \rightarrow V^π \rightarrow π^*

$\forall s, a \quad \pi_k \rightarrow Q^{\pi_k}(s, a)$: Evaluation ①

$\forall s \quad \pi_{k+1}(s) = \arg \max_a Q^{\pi_k}(s, a)$: Improvement ②

? π_{k+1} ? π_k ? π^*

E-Greedy Policy

π_{k+1} \rightarrow π_k עבוק $Q^{\pi_k}(s, a)$ מילוי
 $\arg \max_a Q^{\pi_k}(s, a)$ מילוי + E-הימורון נסוב
 סימולציה של אction a מילוי E-הימורון

$$\begin{aligned}
 Q^{\pi_k}(s, \pi^{k+1}(s)) &= \sum_a \pi_{k+1}(a|s) Q^{\pi_k}(s, a) \\
 &= \frac{\epsilon}{|A|} \sum_a Q^{\pi_k}(s, a) + (1-\epsilon) \max_a Q^{\pi_k}(s, a) \\
 &\stackrel{?}{=} \frac{\epsilon}{|A|} \sum_a Q^{\pi_k}(s, a) + (1-\epsilon) \cdot \sum_a \frac{\pi_k(a|s) - \epsilon/|A|}{1-\epsilon} \cdot Q^{\pi_k}(s, a) \\
 &= \sum_a \pi_k(a|s) \cdot Q^{\pi_k}(s, a) = V^{\pi_k}(s)
 \end{aligned}$$

$$Q^{\pi_k}(s, \pi^{k+1}(s)) \stackrel{?}{=} V^{\pi_k}(s) \quad <=$$

• PI - δ $\sim \mathcal{U}(0, \infty)$ $\sim N(\mu_N)$

3

Temporal Difference

first step to also to MC style how π^* finds

plus has GT values $V(S_t) + V(S_t) + \alpha [G_t - V(S_t)]$
 instead of π^* values "learning rate" α \leftarrow $\frac{\text{learning rate}}{\text{error}}$

: TD(0) approximate $\gamma \hat{V}$. plus signs? punishment etc

$$V(S_t) \leftarrow V(S_t) + \alpha [R_t + \gamma \hat{V}(S_{t+1}) - V(S_t)]$$

$\underbrace{R_t}_{\text{actual reward}}$ $\underbrace{\gamma \hat{V}(S_{t+1})}_{\text{predicted value}}$
 $\underbrace{\text{TD-error}}_{\delta_t} \quad \delta_t$

$$G_T - V(S_t) = \sum_{k=t}^T \gamma^{k-t} \delta_k \quad : \text{Bellman Eq}$$

: episodes 8 . B-1 A and ye !/and/3

B, I	A, O, B, O
B, I	B, I
B, I	B, I
B, O	B, I

$V(B)$? $V(A)$? now

Minimize MSE on train Monte Carlo : $V(A)=0$

Now we want \hat{P} prob Maximum Likelihood $\text{TD}(0)$: $V(A)=V(B)=\frac{3}{4}$

TD prediction

- מילוי אמצעי TD
- משלב predict & fit
- מילוי אמצעי TD

MC prediction

- MC over one episode
- episodes of 300 -
- מילוי אמצעי MC

On policy SARSA

(on policy) policy- π influenced control methods

$$Q_{t+1}(s, a) = Q_t(s, a) + \alpha_t(r(s, a) + \gamma Q_t(s', a') - Q_t(s, a))$$

$a' = \pi(s')$

Off policy Q-Learning

$$Q_{t+1}(s, a) = Q_t(s, a) + \alpha_t[\pi(s, a) + \max_a Q_t(s', a') - Q_t(s, a)]$$

behavior policy \neq target policy

$$Q_{k+1}(s, a) = r(s, a) + \gamma \sum_{s'} p(s'|a, s) \max_{a'} Q_k(s', a'): Q-\delta \text{ reward function}$$

new policy update

⑤ Q^* - δ (�້າງວ່າ Q^* ດີວ່າ) ອົດນ. Q-Learning : ເລັກ

: ($t=1$ ນີ້) $s' = g(s, a)$ ເລັກທີ່ມີກຳນົດໃຫຍ່

$$\Delta_t = \|Q_t - Q^*\|_\infty = \max_{s, a} |Q_t(s, a) - Q^*(s, a)|$$

$$|Q_{t+1}(s, a_t) - Q^*(s, a_t)| = |r(s_t, a_t) + \gamma \max_{a'} Q_t(s_{t+1}, a') - r(s_t, a_t) - \gamma \max_{a'} Q^*(s_{t+1}, a')|$$

(ສິ່ງທີ່ມີກຳນົດໃຫຍ່)

$$= \gamma \left| \max_{a'} Q_t(s_{t+1}, a') - \max_{a'} Q^*(s_{t+1}, a') \right|$$

$$\leq \gamma \max_a |Q_t(s_{t+1}, a) - Q^*(s_{t+1}, a)|$$

$$\leq \gamma \max_a \max_s |Q_t(s, a) - Q^*(s, a)| \leq \gamma \Delta_t$$

ກີ່າວ່າ $Q^*(s, a) = \lim_{n \rightarrow \infty} Q^n(s, a)$, ແລ້ວ $Q^*(s, a) \rightarrow \text{ກຳນົດ}$