

V, Q, π

π

$S_{10}, A_{10}, R_{11}, S_{11}, A_{11}, \dots, S_{1T}$
 $S_{20}, A_{20}, \dots, S_{2T}$
 $S_{M0}, A_{M0}, \dots, S_{MT}$

Model

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

Experience
 (s, a, r, s')

$$V_{\pi}(s) = E_{\pi} [G_t | S_t = s] \approx \text{Avg}(G_t | S_t = s)$$

Monte Carlo estimation:

- init $\pi, V(s)$

- repeat:

- generate $(s_0, A_0, R_1, s_1, A_1, \dots, s_T, A_T, R_{T+1})$ no π

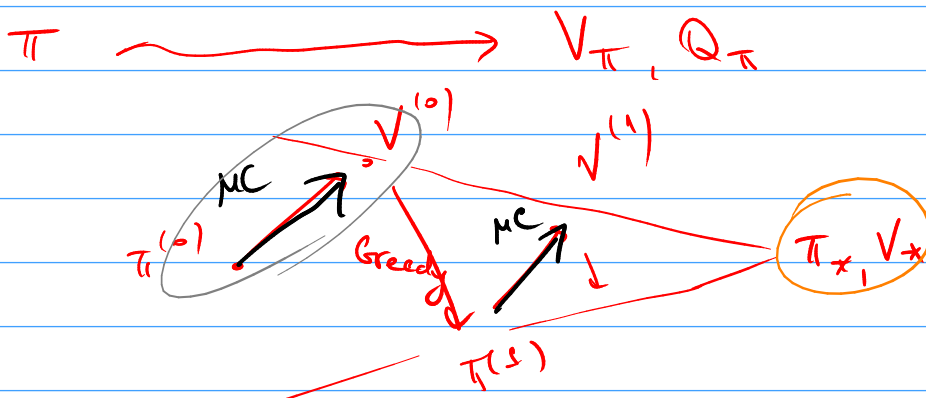
- for $t = T-1, T-2, \dots, 0$:

- $G := \gamma G + R_{t+1}$

- every visit G to $\text{Returns}(S_t)$
 $V(S_t) := \text{Avg}(\text{Returns}(S_t))$

Exploring starts

- G to $\text{Returns}(S_t, A_t)$
 $Q(S_t, A_t) := \text{Avg}(\text{Returns}(S_t, A_t))$



$$\pi'(s) = \arg\max_a Q_{\pi}(s, a)$$

ϵ -soft — $\forall a, p(a|s) \geq \frac{\epsilon}{|A|}$

Policy improv. thm for ϵ -soft strategies

Monte Carlo control

— repeat:

— episode

S_0, A_0

— $\forall t = T-1, \dots, 0$

— get G & returns (S_t, A_t)

— $Q(S_t, A_t) := \text{Avg}(\text{ret} \dots)$

expl. starts

→ od zero Q_π

— $\pi(S_t) := \arg\max_a Q(S_t, a)$

→ wopomy. onyog

— $\pi'(S_t) := \begin{cases} \arg\max_a Q(S_t, a) & \text{c. b. } 1-\epsilon \\ \text{uniform} & \text{c. b. } \epsilon \end{cases}$

Importance sampling

$q(\bar{x})$: even $q(\bar{x})=0$, to u $p(\bar{x})=0$

$$E_{p(\bar{x})}[f] = \int f(\bar{x}) p(\bar{x}) d\bar{x} = \int f(\bar{x}) \frac{p(\bar{x})}{q(\bar{x})} q(\bar{x}) d\bar{x} = E_{q(\bar{x})} \left[f \frac{p}{q} \right]$$

S_t : $A_t, S_{t+1}, A_{t+1}, S_{t+2}, \dots, S_T$

$$p(\text{Traj} | \pi) = p(A_t, \dots, S_T | \pi) = \pi(A_t | S_t) p(S_{t+1} | S_t, A_t) \cdot \pi(A_{t+1} | S_{t+1}) \cdot \dots \cdot p(S_T | S_{T-1}, A_{T-1})$$

$$E_\pi[G_t] = E_{\pi'} \left[G_t \cdot \frac{p(\text{Traj} | \pi)}{p(\text{Traj} | \pi')} \right] = E_{\pi'} \left[G_t \cdot \prod_{k=t}^{T-1} \frac{\pi(A_k | S_k)}{\pi'(A_k | S_k)} \right]$$

$$\frac{p(\text{Traj} | \pi)}{p(\text{Traj} | \pi')} = \frac{\pi(A_t | S_t) p(S_{t+1} | S_t, A_t) \cdot \dots \cdot \pi(A_{T-1} | S_{T-1}) p(S_T | S_{T-1}, A_{T-1})}{\pi'(A_t | S_t) p(S_{t+1} | S_t, A_t) \cdot \dots \cdot \pi'(A_{T-1} | S_{T-1}) p(S_T | S_{T-1}, A_{T-1})}$$

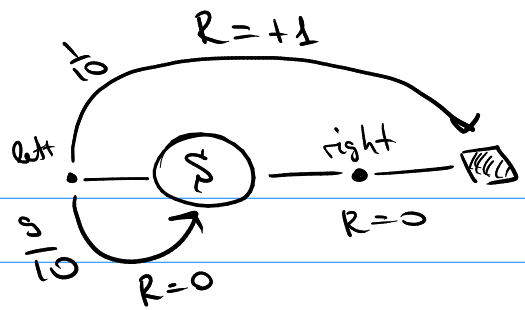
$\sum_{t=T-1}^{\pi, \pi'}$

off-policy estimation
on-policy

$$V_\pi(s) = E_{\pi'} \left[G_t \cdot \sum_{t=T-1}^{\pi, \pi'} \right]$$

$$V_{\pi}(s) = \frac{1}{N} \sum_{t=1}^N p_{t:T-1}^{\pi, \pi'} G_t$$

$$V_{\pi}(s) = \frac{\sum_{t=1}^N p_{t:T-1}^{\pi, \pi'} G_t}{\sum_{t=1}^N p_{t:T-1}^{\pi, \pi'}}$$



$$\pi(\text{left} | s) = 1$$

$$\pi'(\text{left} | s) = 1/2$$

$$V_{\pi}(s) = \mathbb{E}_{\pi'} \left[G \cdot \frac{\pi_k \pi(A_k | S_k)}{\pi'_k \pi'(A_k | S_k)} \right]$$

$$\text{Var}(X) = \mathbb{E}[X^2] - (\mathbb{E}[X])^2$$

$$\mathbb{E}_{\pi'} \left[\left(G \cdot \prod_{k=0}^{T-1} \frac{\pi(A_k | S_k)}{\pi'(A_k | S_k)} \right)^2 \right] = \sum_{\text{Traj}} \pi'(\text{Traj}) \cdot \left(\dots \right)^2 =$$

$$= \frac{1}{2} \cdot \left(\frac{1}{10} \right) \cdot \left(1 \cdot \frac{1}{1/2} \right)^2 +$$

$$+ \frac{1}{2} \cdot \frac{9}{10} \cdot \frac{1}{2} \cdot \left(\frac{1}{10} \right) \cdot \left(1 \cdot \frac{1}{\frac{1}{2} \cdot \frac{1}{2}} \right)^2 +$$

$$+ \frac{1}{2} \cdot \frac{9}{10} \cdot \frac{1}{2} \cdot \frac{9}{10} \cdot \frac{1}{2} \cdot \left(\frac{1}{10} \right) \cdot \left(1 \cdot \frac{1}{\frac{1}{2} \cdot \frac{1}{2} \cdot \frac{1}{2}} \right)^2 + \dots =$$

$$= \frac{1}{10} \sum_{k=0}^{\infty} \left(\frac{9}{10} \right)^k \cdot \left(\frac{1}{2} \right)^{k+1} \cdot 2^{2(k+1)} = \frac{1}{5} \sum_{k=0}^{\infty} \left(\frac{9}{10} \right)^k \cdot 2^k = \frac{1}{5} \sum_{k=0}^{\infty} \left(\frac{5}{2} \right)^k$$

Off-policy MC control

- init

- repeat:

- reset $S_0, A_0, R_1, \dots, L_T, S_T$ no π' , $G=0$, $w=1$

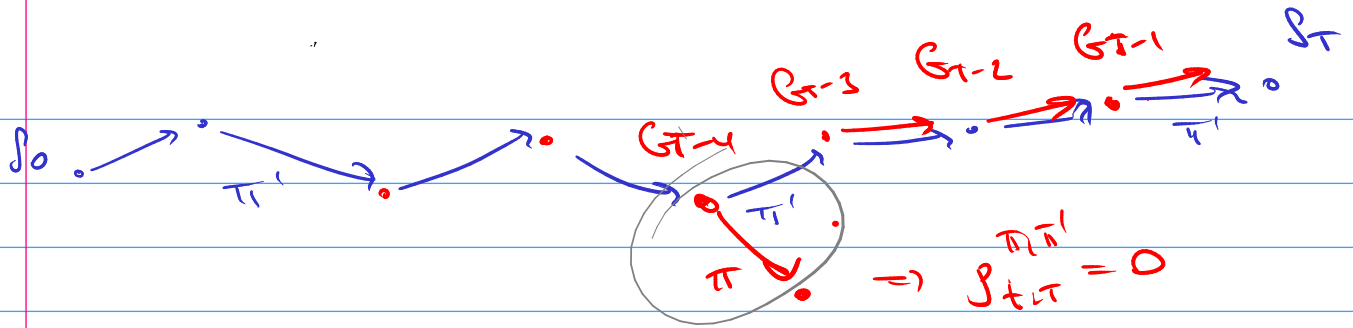
- for $t=T-1, \dots, 0$:

- $G = \gamma G + R_{t+1}$, $W(S_t, A_t) = W(S_t, A_t) + w$

- $Q(S_t, A_t) := Q(S_t, A_t) + \frac{w}{W(S_t, A_t)} [G - Q(S_t, A_t)]$

- $\pi(S_t) := \arg \max_a Q(S_t, a)$, $W := W \cdot \frac{1}{\pi'(A_t | S_t)}$

- if $A_t \neq \pi(S_t)$ then break



TD-learning
 &
 temporal difference
 (bootstrapping)

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

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

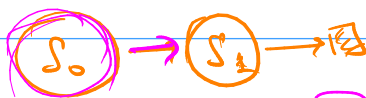
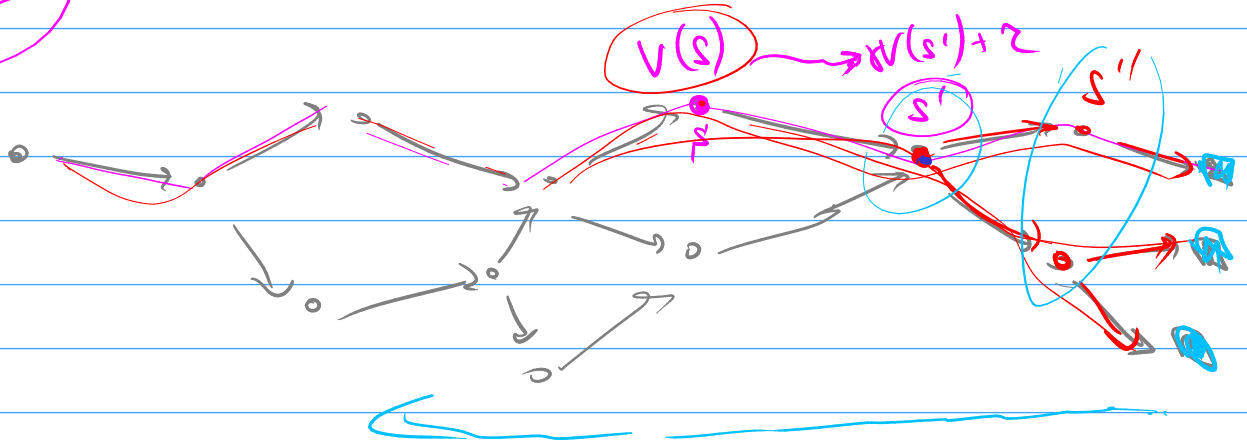
TD(0) for estimation

- init $V_{\pi}(s)$, $s \in \mathcal{S}$
- repeat:

- generate x_0, A_t & ω_{π} - S_t , π - R_{t+1} , S_{t+1}
- observe

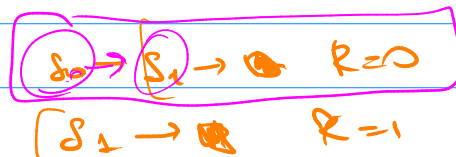
$$V_{\pi}(S_t) := V_{\pi}(S_t) + \alpha (R_{t+1} + \gamma V_{\pi}(S_{t+1}) - V_{\pi}(S_t))$$

$$Q_{\pi}(S_t, A_t) := Q_{\pi}(S_t, A_t) + \alpha (R_{t+1} + \gamma Q_{\pi}(S_{t+1}, A_{t+1}) - Q_{\pi}(S_t, A_t))$$



$\gamma = 1$

- $S_0 \rightarrow \text{box} \quad R=1$
- $S_1 \rightarrow \text{box} \quad R=0$
- $S_2 \rightarrow \text{box} \quad R=1$
- $S_3 \rightarrow \text{box} \quad R=1$



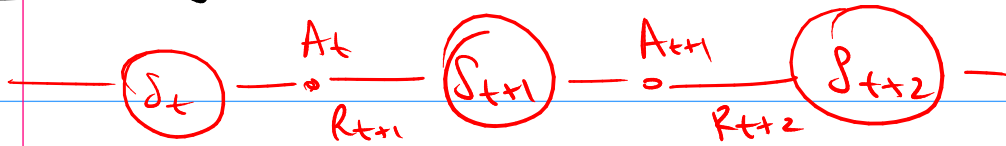
$$V(S_1) = 2/3$$

$$V^{MC}(S_0) = 0 \quad G_1 = 0$$

$$V^{TD}(S_0) = 2/3$$

On-policy TD control

(s, a, r, s')



TD(0):

$$Q(s_t, A_t) := Q(s_t, A_t) + \alpha (R_{t+1} + \gamma Q(s_{t+1}, A_{t+1}) - Q(s_t, A_t))$$

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

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

$$\pi(s) = \begin{cases} \arg\max_a Q(s, a) & 1-\epsilon \\ \text{unit} & \epsilon \end{cases}$$

SARSA

- init Q

- repeat

- generate A no π^Q neurons s', R

- observe A' & s' no π^Q

- observe $Q(s, A)$ no TD(0)

- $s := s', A := A'$

Off-policy TD control

$$\pi^Q = \begin{cases} \arg\max_a Q(s, a), & 1-\epsilon \\ \text{unit} & \epsilon \end{cases}$$

TD(0):

$$Q(s_t, A_t) := Q(s_t, A_t) + \alpha [R_{t+1} + \gamma \max_a Q(s_{t+1}, a) - Q(s_t, A_t)]$$

Q-learning

DQN

