

Input  $\pi$  to eval

Initialize:

$V(s) \leftarrow R$  arbitrarily for all  $s \in S$   
 $\text{Returns}(s) \leftarrow \text{empty}$  list  $\forall s \in S$

loop forever (for each episode)

gen an ep following  $\pi: S_0, A_0, R, S_1, A_1, \dots, S_{T-1}, A_{T-1}, R_T$

$G \leftarrow 0$

loop for each state in the ep,  $t = T-1, T-2, \dots, 0$

$G \leftarrow \gamma G + R_{t+1}$

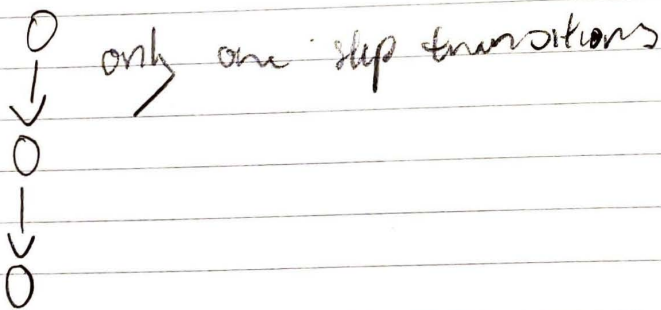
Unless  $S_t$  in  $S_0, S_1, \dots, S_{t-1}$

Append  $G$  to  $\text{Returns}(S_t)$

$V(S_t) \leftarrow \text{average}(\text{Returns}(S_t))$

Simulate many games using a policy with an MC approach then average the returns for each state

No need to pre-comp probs



policy improvement - use greedy

$$q_{\pi_{k+1}}(s, \pi_{k+1}(s)) = q_{\pi_k}(s, \arg \max_a q_{\pi_k}(s, a))$$

$$= \max_a q_{\pi_k}(s, a)$$

$$\geq q_{\pi_k}(s, \pi_k(s))$$

$$\geq V_{\pi_k}(s)$$

$$\pi_0 \rightarrow q_{\pi_0} \rightarrow \pi_1 \rightarrow q_{\pi_1} \rightarrow \dots \rightarrow \pi_* \rightarrow q_*$$

Int  
 $a(s) \in A(s)$  arbitrarily  $\forall s \in S$   
 $Q(s, a) \in \mathbb{R}$   
 $\text{Returns}(s, a) \subseteq \text{empty list}$   $\forall s, a \in A(s)$

Loop Forever (for each ep)

Choose  $S_0 \in S, A_0 \in A(S_0)$  randomly ( $p(S_0, A_0) > 0 \forall A_0, S_0$ )  
 generate an ep  $\pi: S_0, A_0, R_1, \dots, S_T, A_T, R_T$   
 $C \leftarrow 0$

Loop from each step  $t = T-1, \dots, 0$

$C \leftarrow \gamma C + R_{t+1}$

unless the pair  $S_t, A_t$  in our ep

Append  $C$  to  $\text{Returns}(S_t, A_t)$

$Q(S_t, A_t) \leftarrow \text{average}(\text{Returns}(S_t, A_t))$

$\pi(S_t) \leftarrow \text{argmax}_a Q(S_t, a)$

On policy limit visit

$A^* \leftarrow \text{argmax}_a Q(S_t, a)$   
 $\forall a \in A(S_t)$

$\pi(a | S_t) \leftarrow \begin{cases} 1 - \epsilon + \epsilon / |A(S_t)| \\ \epsilon / |A(S_t)| \end{cases}$

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

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

$$\geq \frac{\epsilon}{|A(s)|} \sum_a q_{\pi}(s, a) + (1-\epsilon) \sum_a \frac{\pi(a|s) - \frac{\epsilon}{|A(s)|}}{1-\epsilon} q_{\pi}(s, a)$$

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

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



# Off policy MC

Need to behave non-optimally to learn  
but want to act optimally

Two policies

- Target
- behavior

High var low convergence more power

$b$  (behavior),  $\pi$  (target)

$$\pi(a|s) > 0 \Rightarrow b(a|s) > 0$$

$$\Pr\{A_t, S_{t+1}, A_{t+1}, \dots, S_T | S_t, A_{t:T-1} \sim \pi\}$$

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

$$= \prod_k \pi(A_k | S_k) p(S_{k+1} | S_k, A_k)$$

importance sampling ratio - estimating expected values  
of one distro from another

$$P_{t:T-1} = \frac{\prod_{k=t}^{T-1} \pi(A_k | S_k) p(S_{k+1} | S_k, A_k)}{\prod_{k=t}^{T-1} b(A_k | S_k) p(S_{k+1} | S_k, A_k)} = \prod_{k=t}^{T-1} \frac{\pi(A_k | S_k)}{b(A_k | S_k)}$$

$$E[G_+ | S_{t=0}] = V_b(s) \quad - \text{wont give us } \pi$$

$$E[P_{t:T-1} G_+ | S_{t=0}] = V_\pi(s)$$

$$V(t) = \frac{\sum_{t \in T(t)} P_t \cdot T(t) - 1 \cdot C_t}{|T(t)|}$$

$\pi$  given later

weighted version

$$V(t) = \frac{\sum_{t \in T(t)} P_t \cdot T(t) - 1 \cdot C_t}{\sum_{t \in T(t)} P_t \cdot T(t) - 1}$$

Want

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

$$V_{n+1} = V_n + \frac{W_n}{C_n} [C_n - V_n] \quad (C_n \text{ is weight of weights})$$

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

# MCTS

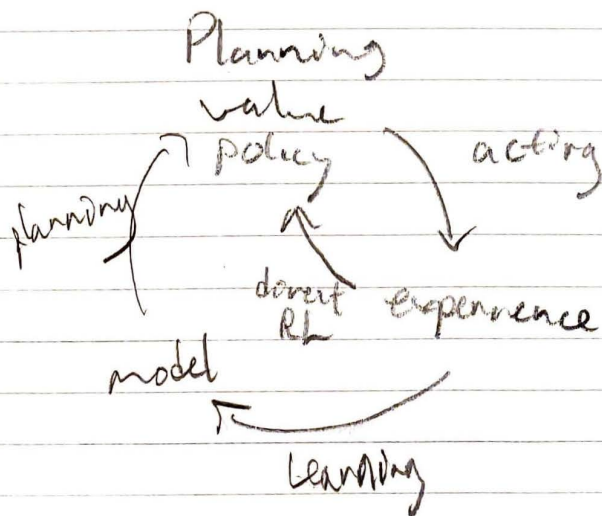
- 1.) Selection: Follow the tree policy to a leaf
- 2.) Expansion: expand the tree by adding a new leaf node to the tree
- 3.) Simulation: From the selected node or new child run a MC episode with the rollout policy
- 4.) Backup: Return the generated episode and its return

## Rollout Algos

Estimate action values by averaging the rewards of many trajectories that start with a given action.

Unlike MC we don't estimate  $q^*$  or  $q_\pi$

Instead they only get values for only the actions associated with the current state.





# Tab Dyna Q

Init  $Q(S, a)$ ,  $Model(S, a) \forall S \in S, a \in A(S)$

loop:

$S \leftarrow$  current state

$A \leftarrow \epsilon$ -greedy( $S, Q$ )

Take action  $A$ : observe resultant reward  $R, S'$

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

$Model(S, A) \leftarrow R, S'$  (assumes deterministic)

loop  $n$  times:

$S \leftarrow$  random previous state

$A \leftarrow$  random action previously taken in  $S$

$R, S' \leftarrow Model(S, A)$

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