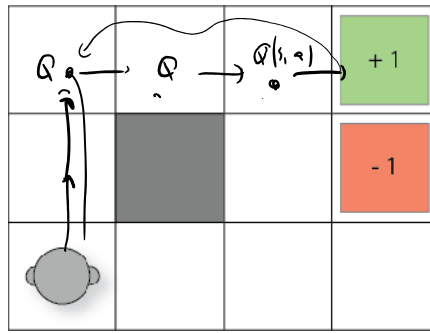
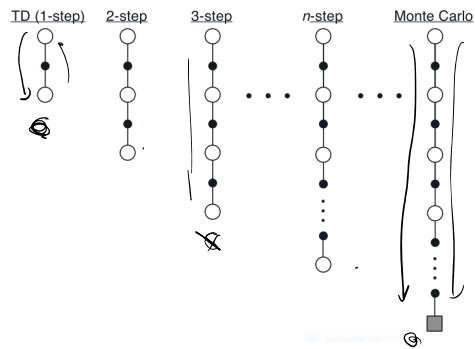


n-step temporal difference learning

Discounted future rewards

$$\begin{aligned} G_t &= r_t + \gamma V(s') \\ &= r_t + \gamma r_{t+1} + \gamma^2 V(s'') \\ &= r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \gamma^3 V(s''') \end{aligned}$$

Truncated future rewards



$\beta = 0.4$

$$Q(s_t) \leftarrow Q(s_t) + \alpha [r + \gamma V(s') - Q(s_t)]$$

$$\text{SARSA: } Q(s,a) := Q(s,a) + \alpha [r + \gamma Q(s',a') - Q(s,a)]$$

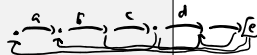
Change the update:

$$\begin{aligned} G_t &\leftarrow r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \dots + \gamma^n Q(s',a') \\ Q(s,a) &\leftarrow Q(s,a) + \alpha [G_t + \gamma Q(s',a') - Q(s,a)] \end{aligned}$$

n-step Sarsa for estimating $Q \approx q_\pi$, or $Q \approx q_\pi$ for a given π

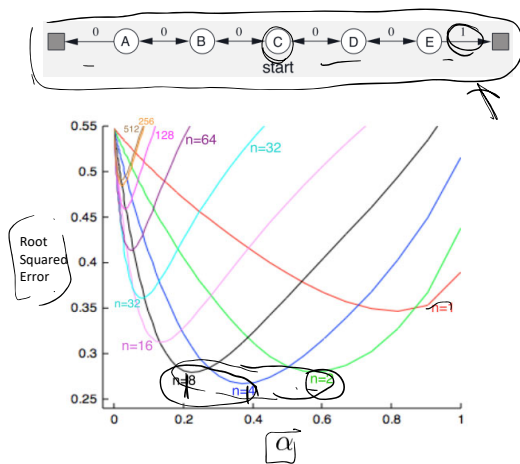
Initialize $Q(s,a)$ arbitrarily, for all $s \in S, a \in A$
Initialize π to be ϵ -greedy with respect to Q , or to a fixed given policy
Parameters: step size $\alpha \in (0,1)$, small $\epsilon > 0$, a positive integer n
All store and access operations (for S_t, A_t , and R_t) can take their index mod n

Repeat (for each episode):
Initialize and store $S_0 \neq \text{terminal}$
Select and store an action $A_0 \sim \pi(\cdot|S_0)$
 $T \leftarrow \infty$
For $t = 0, 1, 2, \dots$:
If $t < T$, then:
Take action A_t
Observe and store the next reward as R_{t+1} and the next state as S_{t+1}
If S_{t+1} is terminal, then:
 $T \leftarrow t + 1$
else:
Select and store an action $A_{t+1} \sim \pi(\cdot|S_{t+1})$
 $\tau \leftarrow t - n + 1$ (τ is the time whose estimate is being updated)
If $\tau \geq 0$:
 $G \leftarrow \sum_{i=\tau}^{t+n-1} \gamma^{t-i} R_{i+1}$
If $\tau + n < T$, then $G \leftarrow G + \gamma^n Q(S_{\tau+n}, A_{\tau+n})$ ($G_{\tau+\tau+n}$)
 $Q(S_\tau, A_\tau) \leftarrow Q(S_\tau, A_\tau) + \alpha [G - Q(S_\tau, A_\tau)]$
If π is being learned, then ensure that $\pi(\cdot|S_\tau)$ is ϵ -greedy wrt Q
Until $\tau = T - 1$



With 1-step learning					With 5-step learning				
State	Action				State	Action			
	North	South	East	West		North	South	East	West
(0,0)	0	0	0	0	(0,0)	0	0	0	0
(0,1)	0	0	0	0	(0,1)	0	0	0	0
(0,2)	0	0	0	0	(0,2)	0	0	<u>0.2953</u>	0
...					...				
(1,2)	0	0	0	0	(1,2)	0	0	<u>0.3281</u>	0
(2,1)	0	0	0	0	(2,1)	<u>0.405</u>	0	<u>0</u>	0
(2,2)	0	0	<u>0.45</u>	0	(2,2)	<u>0</u>	<u>0.3645</u>	<u>0.45</u>	0
(2,3)	0	0	0	0	(2,3)	0	<u>0</u>	<u>0</u>	0
...					...				

Example: Random walk



MCTS + Reinforcement learning

AlphaZero:

