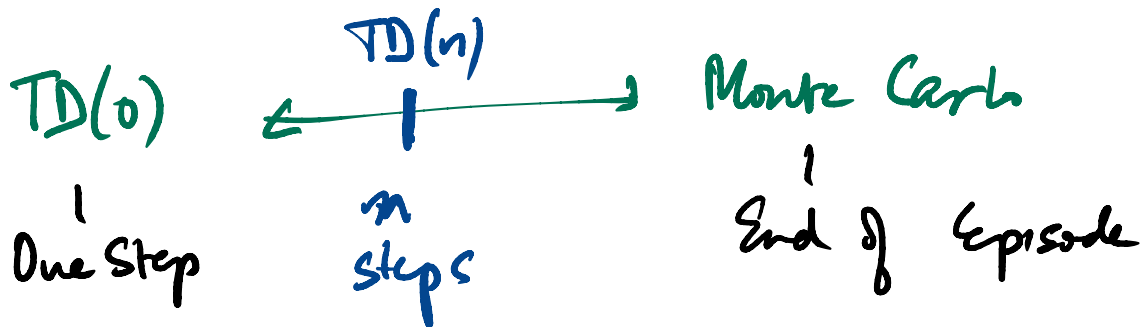## Policy evaluation in MDPs

**Dynamic Programming** — One step ahead, bootstrapped, full model

**Monte Carlo** — Episodes + Mean value

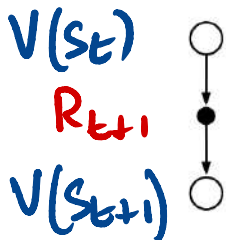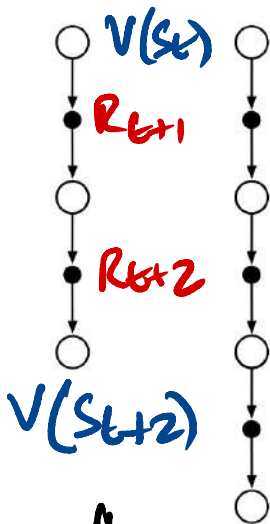**Temporal Difference** — TD(0)

— Episodes, bootstrap one step ahead

TD(0) $\longleftarrow$ TD(n) $\longrightarrow$ Monte Carlo

One Step     n steps     End of Episode

# n-step TD



1-step TD and TD(0)  2-step TD  3-step TD  n-step TD  ∞-step TD and Monte Carlo

$V(S_t)$
$R_{t+1}$
$V(S_{t+1})$

$V(S_t)$
$R_{t+1}$
$R_{t+2}$
$V(S_{t+2})$

$V(S_t) \leftarrow$
$V(S_t) + \gamma(\delta_t)$

$R_{t+1} + \gamma V(S_{t+1})$
$- V(S_t)$

$$V(S_t): R_{t+1} + \gamma R_{t+2} + \gamma^2 V(S_{t+2})$$

$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} \cdots + \gamma^{T-t-1} R_T$$

TD(0)

$$G_t = R_{t+1} + \gamma V_t(S_{t+1})$$

$\hookrightarrow$ V at time t

2 steps

$$G_{t:t+2} = R_{t+1} + \gamma R_{t+2} + \gamma^2 V_{t+1}(S_{t+2})$$

$$G_{t:t+n} = R_{t+1} + \gamma R_{t+2} + \cdots \gamma^{n-1} R_{t+n} + \gamma^n V_{t+n-1}(S_{t+n})$$

$V_t(S_t) \rightarrow$ after n observations $\rightarrow V_{t+n}(S_t)$

$$V_{t+n}(S_t) = V_{t+n-1}(S_t) + \alpha \left[ \underbrace{G_{t:t+n} - V_{t+n-1}(S_t)}_{\delta_{t:t+n}} \right]$$

Algorithm

### $n$-step TD for estimating $V \approx v_\pi$

Input: a policy $\pi$
Algorithm parameters: step size $\alpha \in (0, 1]$, a positive integer $n$
Initialize $V(s)$ arbitrarily, for all $s \in \mathcal{S}$
All store and access operations (for $S_t$ and $R_t$) can take their index mod $n + 1$

Loop for each episode:
    Initialize and store $S_0 \neq$ terminal
    $T \leftarrow \infty$
    Loop for $t = 0, 1, 2, \ldots$ :
        If $t < T$, then:
            Take an action according to $\pi(\cdot|S_t)$
            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 $T \leftarrow t + 1$
        $\tau \leftarrow t - n + 1$   ($\tau$ is the time whose state's estimate is being updated)
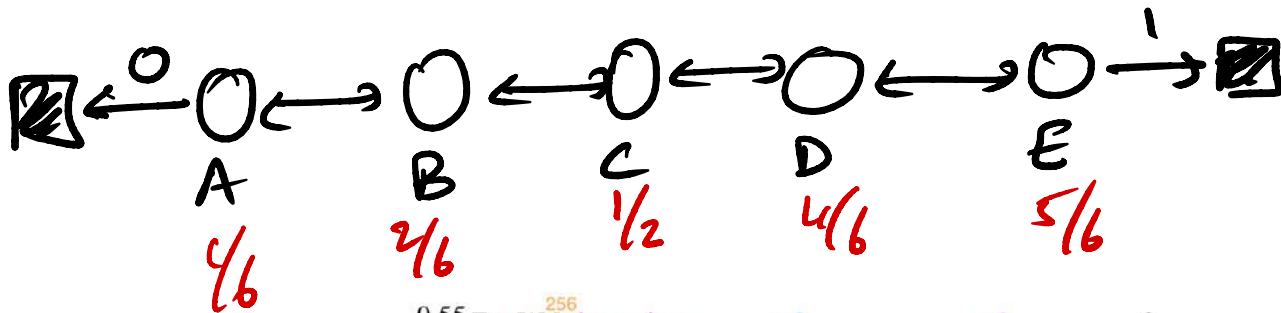        If $\tau \geq 0$:
            $G \leftarrow \sum_{i=\tau+1}^{\min(\tau+n, T)} \gamma^{i-\tau-1} R_i$   ← $G_{b-n:t}$
            If $\tau + n < T$, then: $G \leftarrow G + \gamma^n V(S_{\tau+n})$ ← new $V$   ($G_{\tau:\tau+n}$)
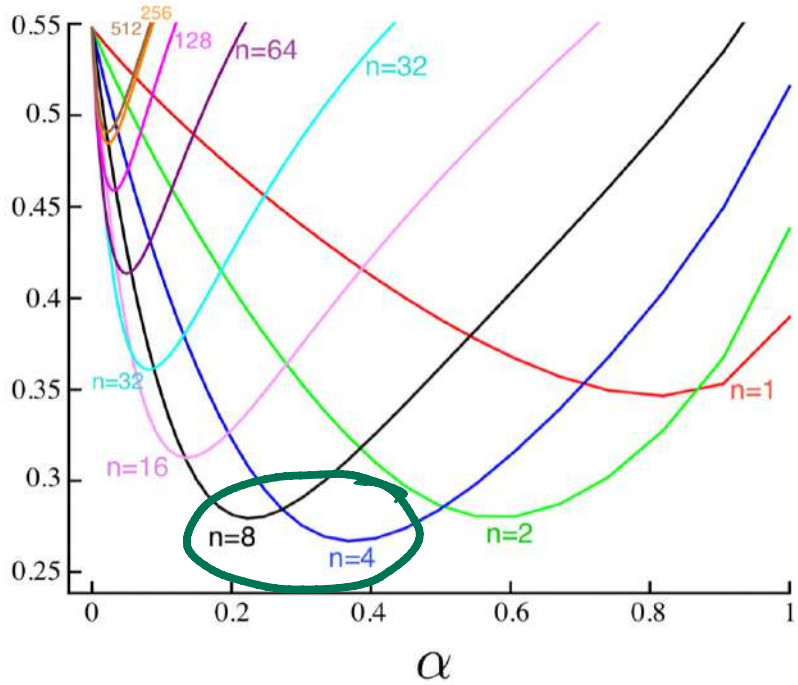            $V(S_\tau) \leftarrow V(S_\tau) + \alpha[G - V(S_\tau)]$ — Update
    Until $\tau = T - 1$

A 4/6    B 2/6    C 1/2    D 4/6    E 5/6

Useful to choose intermediate n

Expand to 19 states

Average RMS error over 19 states and first 10 episodes

Left −1
Right 0

= TD(0)

0.55  512  256  128  n=64  n=32
0.5
0.45  n=32
0.4
0.35  n=16
0.3
0.25  n=8  n=4  n=2  n=1

0   0.2   0.4   0.6   0.8   1
α

## Similarly

### SARSA for policy iteration

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha \left[ R_{t+1} + \gamma Q(S_{t+1}, A_{t+1}) - Q(S_t, A_t) \right]$$

$$(S_t, A_t, R_{t+1}, S_{t+1}, A_{t+1})$$

$$G_{t:t+n} = R_{t+1} + \gamma R_{t+2} + \dots \gamma^{n-1} R_{t+n} + \gamma^n Q_{t+n-1}(S_{t+n}, A_{t+n})$$

$$Q_{t+n}(S_t, A_t) \leftarrow Q_{t+n-1}(S_t, A_t) + \alpha \left[ G_{t:t+n} - Q_{t+n-1}(S_t, A_t) \right]$$

# Expected SARSA

$$G_t = R_{t+1} + \underset{a}{\mathbb{E}} \left[ Q(S_t, A_{t+1}) \right]$$

Extend to n-steps

# Off Policy Learning

Discount each step by $S_k = \dfrac{\pi(A_k | S_k)}{b(A_k | S_k)}$

policy being computed

aux policy

$$V_t(S_t) \leftarrow V_t(S_t) + \alpha \cdot \underline{\rho_t} \left[ G_t - V_t(S_t) \right]$$

Sequence : $\rho_{t:t+n} = \prod_{k=t}^{\min(t+n,T)} \dfrac{\pi(A_k|S_k)}{b(A_k|S_k)}$

$$V_{t+n}(S_t) \leftarrow V_{t+n-1}(S_t) + \alpha \cdot \rho_{t:t+n} \left[ G_{t:t+n} - V_{t+n-1}(S_t) \right]$$
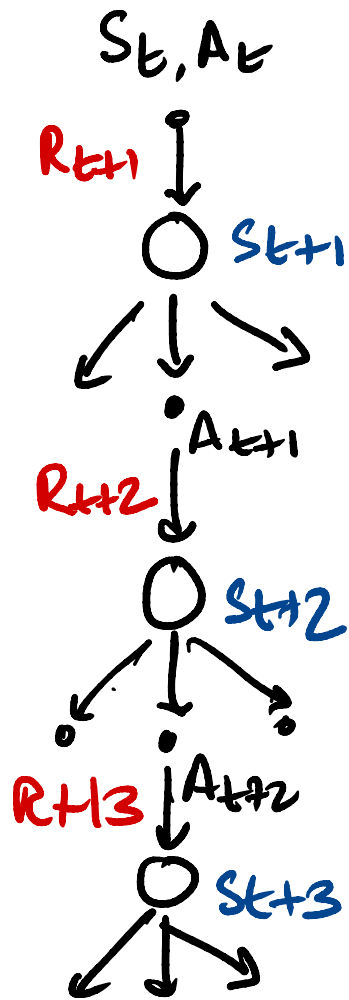
IMPORTANCE SAMPLING

# Off policy without importance sampling

TD(0) — Q learning

SARSA
$Q(S_{t+1}, A_{t+1})$
↓

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha \left[ R_{t+1} + \gamma \max_a Q(S_{t+1}, a) - Q(S_t, A_t) \right]$$

Can we do n-step Q learning?

$$G_{t:t+1} \overset{def}{=} R_{t+1} + \gamma \sum_a \pi(a|S_{t+1}) \cdot Q_t(S_{t+1}, a)$$

EXPECTED SARSA

$$G_{t:t+2} \overset{def}{=} R_{t+1} + \gamma \sum_{a \neq A_{t+1}} \pi(a|S_{t+1}) \, Q_t(S_{t+1}, a)$$

$$+ \gamma \pi(A_{t+1}|S_{t+1}) \Big[ R_{t+2} + \gamma \sum_a \pi(a|S_{t+2}) \, Q_t(S_{t+2}, a) \Big]$$

$$G_{t:t+n} = R_{t+1} +$$

$$\gamma \sum_{a \neq A_{t+1}} \pi(a|S_{t+1}) Q_{t+n-1}(S_{t+1},a)$$

$$+$$

$$\gamma \pi(A_{t+1}|S_{t+1}) \cdot \underline{G_{t+1:t+n}}$$

n-step TD

n-step SARSA, expected SARSA

n-step TREE LEARNING