

Lecture - 05

Temporal Difference Learning

- Dynamic programming . (bootstrapping)
- Monte Carlo . (does not requ. dynamics fn)

- TD Learning → Comb. of MC & DP .

TD prediction:-

Simple every-visit Monte Carlo.

Constant - α

MC

$$\underline{\text{MC}} : V(S_{t+1}) = V(S_t) + \alpha [G_t - V(S_t)]$$

↑
Sample return.

$$\underline{\text{TD}} : V(S_{t+1}) = V(S_t) + \alpha [R_{t+1} + \gamma V(S_{t+1}) - V(S_t)]$$

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

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

Tabular TD for estimating v_π

Input: the policy π to be evaluated

Algorithm parameter: step size $\alpha \in (0, 1]$

Initialize $V(s)$, for all $s \in S^+$, arbitrarily except that $V(\text{terminal}) = 0$

Loop for each episode:

 Initialize S

 Loop for each step of episode:

$A \leftarrow$ action given by π for S

 Take action A , observe R, S'

$$V(S) \leftarrow V(S) + \alpha [R + \gamma V(S') - V(S)]$$

$$S \leftarrow S'$$

 until S is terminal

$$v_\pi(s) = E_\pi [G_t | s_t = s] \quad MC$$

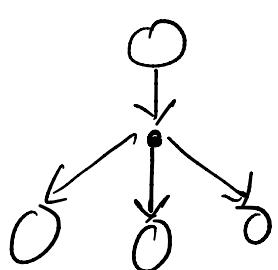
$$= E_\pi [R_{t+1} + \gamma G_{t+1} | s_{t+1} = s]$$

$$= E_\pi [R_{t+1} + \gamma v_\pi(s_{t+1}) | s_{t+1} = s]$$

TD.

Backup diagram of TD :-

DP



MC



TD



$$\text{TD error: } S_t = \frac{R_{t+1} + \gamma V(S_{t+1})}{\underline{\hspace{1cm}}} - \frac{V(S_t)}{\underline{\hspace{1cm}}}$$

$$G_t - V(S_t) = R_{t+1} + \gamma G_{t+1} - V(S_t)$$

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

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

$$+ \gamma G_{t+1} - \gamma^2 V(S_{t+2})$$

$$= \delta_t + \gamma (G_{t+1} - V(S_{t+1}))$$

$$= \delta_t + \gamma \delta_{t+1} + \gamma^2 (G_{t+2} - V(S_{t+2}))$$

$$= \delta_t + \gamma \delta_{t+1} + \gamma^2 \delta_{t+2}$$

$$+ \dots + \gamma^{T-t-1} \delta_{T-1} + \gamma^{T-t} (\cancel{G_T} - \cancel{V(S_T)})$$

$$= \sum_{k=t}^{T-1} \gamma^{k-t} \delta_k$$

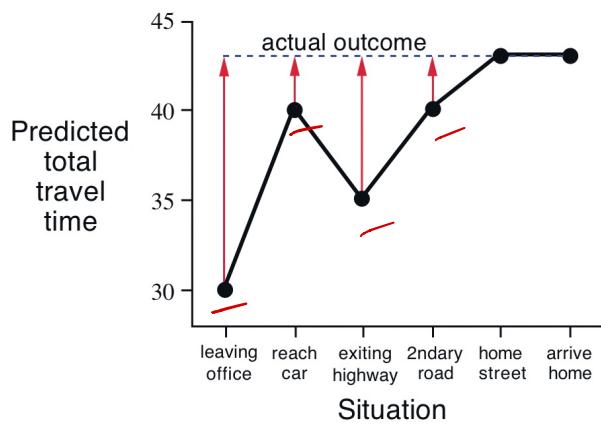
Driving home:

State	Elapsed Time (minutes)	Predicted Time to Go	Predicted Total Time
leaving office, friday at 6	0	<u>30</u>	<u>30</u>
reach car, raining	5	<u>35</u>	<u>40</u>
exiting highway	20	<u>15</u>	<u>35</u>
2ndary road, behind truck	30	<u>10</u>	<u>40</u>
entering home street	40	<u>3</u>	<u>43</u>
arrive home	43	<u>0</u>	<u>43</u>

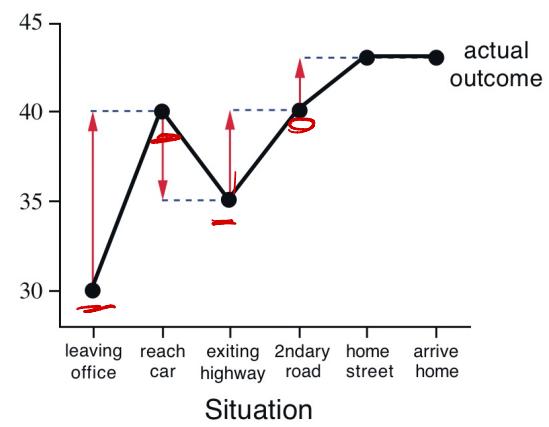
Rewards — elapsed time on each leg
of the journey.

$\gamma = 1$ No discount.

return → actual time to go for the
next step



MC



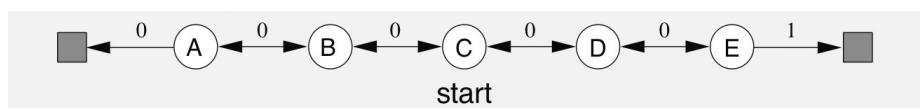
TD

Advantages of TD prediction methods!:-

Random walk :-

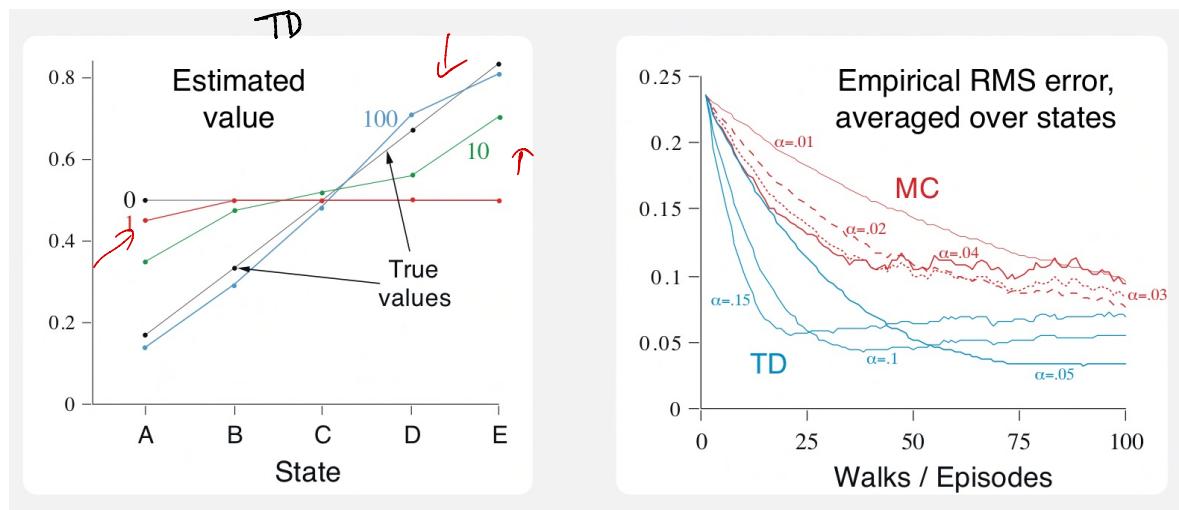
undiscounted

Markov Reward
process
(MRP).



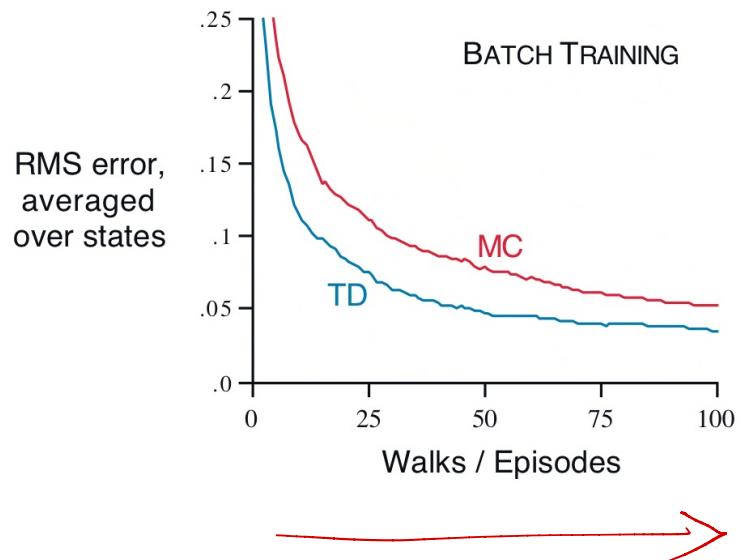
$\frac{1}{6}$ $\frac{2}{6}$ $\frac{3}{6}$ $\frac{4}{16}$ $\frac{5}{6}$
 \uparrow
 0.5

Initial
State
Value of
0.5
for all s .



Batch updating :-

- 1 ep.
- 2 epis.
- 3 epis.



A, 0, C, -5

1)	A, 0, B, 0	B, 1 B, 1 B, 1 B, 0
2)	B, 1	
3)	B, 1	
4)	B, 1	

$v(A), v(B)$

$$v(B) = \frac{6}{8} = \frac{3}{4}.$$

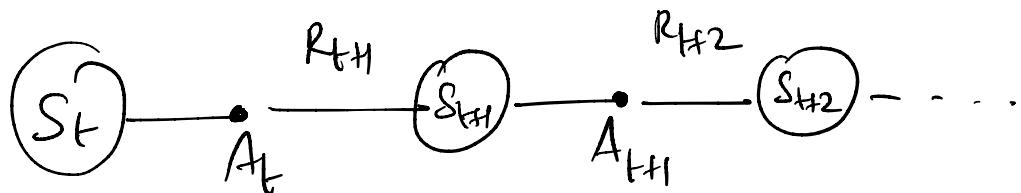
$$v(A) = ? \quad \frac{3}{4} ?$$

Batch MC \rightarrow estimate the min.
MSE on the training set.

Batch TD \rightarrow Certainty-equivalence estimate

On-Policy TD Control: Sarsa.

for term. S.t.
 $Q(S_{t+1}, A_{t+1}) = 0$



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

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

SARSA.



Sarsa (on-policy TD control) for estimating $Q \approx q_*$

Algorithm parameters: step size $\alpha \in (0, 1]$, small $\varepsilon > 0$

Initialize $Q(s, a)$, for all $s \in \mathcal{S}^+, a \in \mathcal{A}(s)$, arbitrarily except that $Q(\text{terminal}, \cdot) = 0$

Loop for each episode:

 Initialize S

 Choose A from S using policy derived from Q (e.g., ε -greedy)

 Loop for each step of episode:

 Take action A , observe R, S'

 Choose A' from S' using policy derived from Q (e.g., ε -greedy)

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

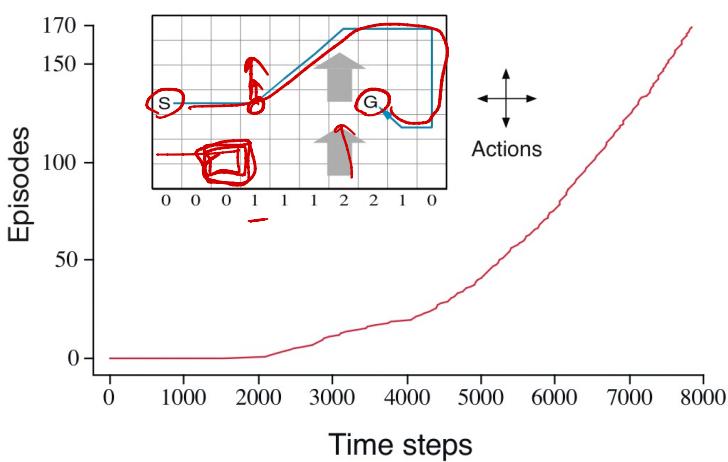
$S \leftarrow S'; A \leftarrow A';$

 until S is terminal

Windy gridworld:

undiscounted.

episode:



$$\epsilon = 0.1, \alpha = 0.05 \quad \text{init} \quad Q = 0.$$

Off-Policy TD Control:

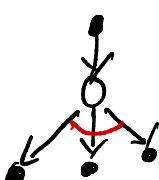
Q-learning [Watkins, 1989]:

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha$$

$$\left[R_{t+1} + \gamma \max_a Q(s_{t+1}, a) - Q(s_t, a_t) \right]$$

α_{V_s} .

Sarsa



Q-learning (off-policy TD control) for estimating $\pi \approx \pi_*$

Algorithm parameters: step size $\alpha \in (0, 1]$, small $\varepsilon > 0$

Initialize $Q(s, a)$, for all $s \in \mathcal{S}^+, a \in \mathcal{A}(s)$, arbitrarily except that $Q(\text{terminal}, \cdot) = 0$

Loop for each episode:

 Initialize S

 Loop for each step of episode:

 Choose A from S using policy derived from Q (e.g., ε -greedy)

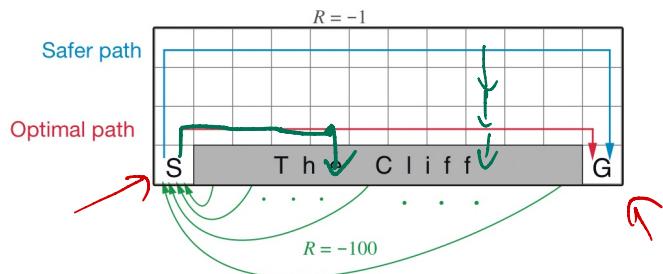
 Take action A , observe R, S'

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

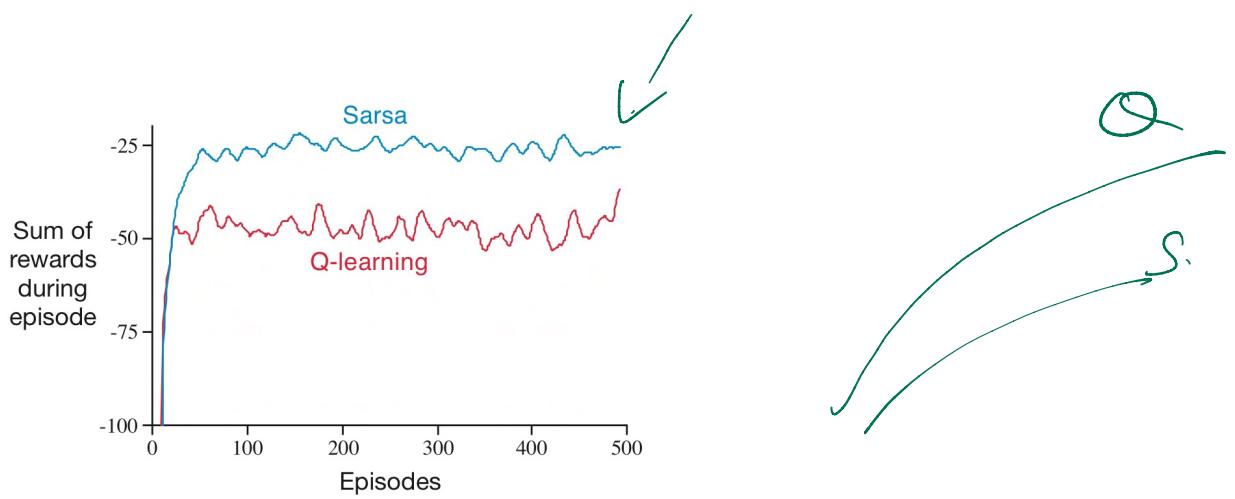
$S \leftarrow S'$

 until S is terminal

Cliff walking:



- undiscounted,
episode,



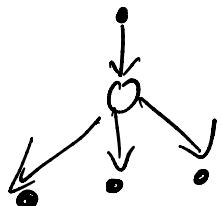
Expected Sarsa :-

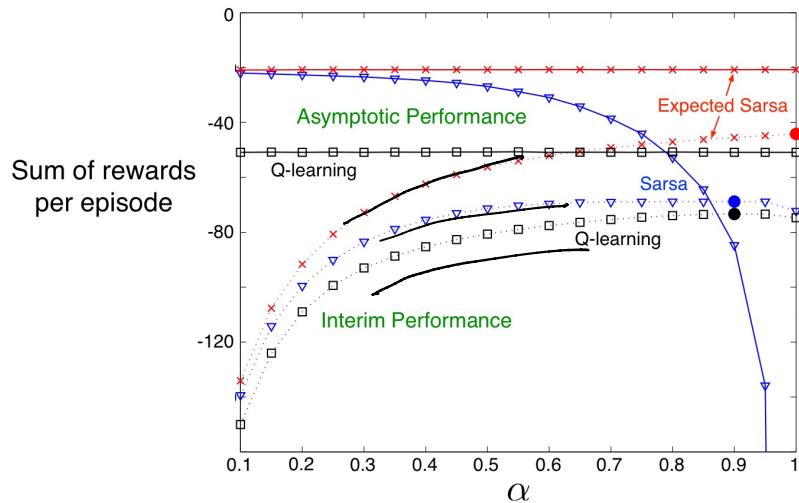
$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha [R_{t+1} +$$

$$\gamma E_{\pi} [Q(s_{t+1}, a_{t+1}) (s_{t+1}) - Q(s_t, a_t)]$$

$$\leftarrow Q(s_t, a_t) + \alpha [R_{t+1} + \gamma \sum_a \pi(a|s_{t+1}) Q(s_{t+1}, a) - Q(s_t, a_t)]$$

Expected Sarsa .





π -greedy
beh. pol - ϵ -greedy } Expected Sarsa

Maximization bias :-

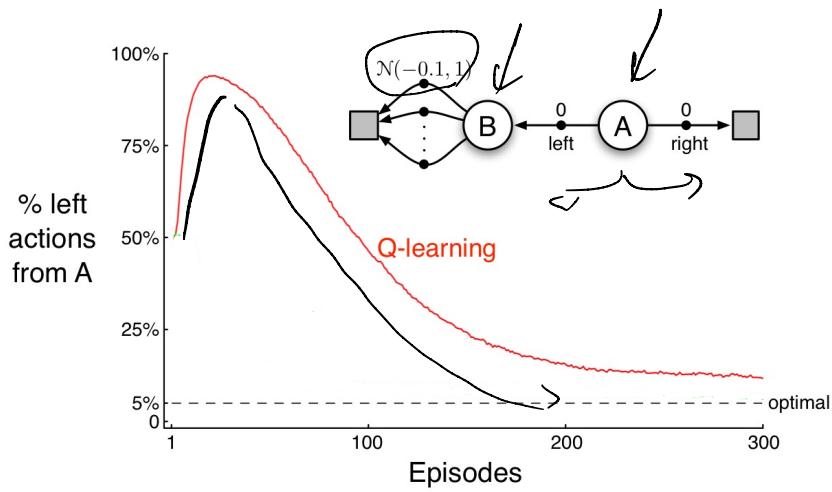
Q-learning : greedy :

Sarsa ! < greedy.

Positive bias - / Maximization bias

s. $q(s,a) = 0.$

$Q(s,a)$



How to avoid max. bias?

Bandits $Q_1(a)$ $Q_2(a)$

Use Q_1 for max: $A^* = \underset{a}{\operatorname{argmax}} Q_1(a)$

Use Q_2 for estm: $Q_2(A^*) = Q_2(\underset{a}{\operatorname{argmax}} Q_1(a))$

$$\mathbb{E}[Q_2(A^*)] = \bar{q}(A^*)$$

$$Q_1(\underset{a}{\operatorname{argmax}} Q_2(a))$$

Double learning:

~~Double Q-learning:~~

$$Q_1(S_t, A_t) = Q_1(S_t, A_t) + \alpha [R_{t+1} + \gamma Q_2(S_{t+1}, \underset{a}{\operatorname{argmax}} Q_1(A_{t+1}, a))]$$
$$- Q_1(S_t, A_t)$$

Double Q-learning, for estimating $Q_1 \approx Q_2 \approx q_*$

Algorithm parameters: step size $\alpha \in (0, 1]$, small $\varepsilon > 0$

Initialize $Q_1(s, a)$ and $Q_2(s, a)$, for all $s \in \mathcal{S}^+$, $a \in \mathcal{A}(s)$, such that $Q(\text{terminal}, \cdot) = 0$

Loop for each episode:

 Initialize S

 Loop for each step of episode:

 Choose A from S using the policy ε -greedy in $Q_1 + Q_2$

 Take action A , observe R, S'

 With 0.5 probability:

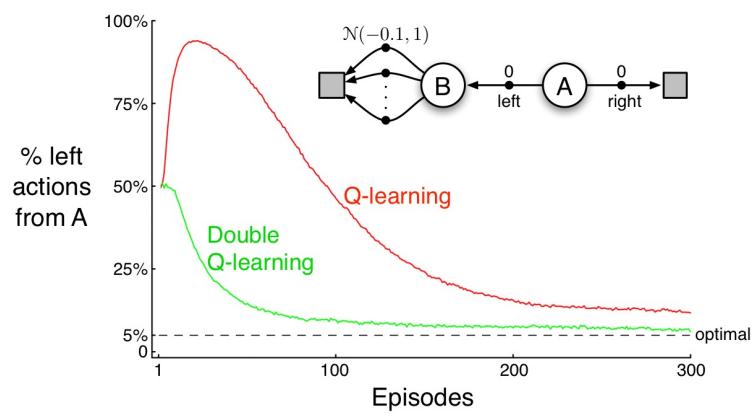
$$Q_1(S, A) \leftarrow Q_1(S, A) + \alpha \left(R + \gamma Q_2(S', \arg \max_a Q_1(S', a)) - Q_1(S, A) \right)$$

 else:

$$Q_2(S, A) \leftarrow Q_2(S, A) + \alpha \left(R + \gamma Q_1(S', \arg \max_a Q_2(S', a)) - Q_2(S, A) \right)$$

$S \leftarrow S'$

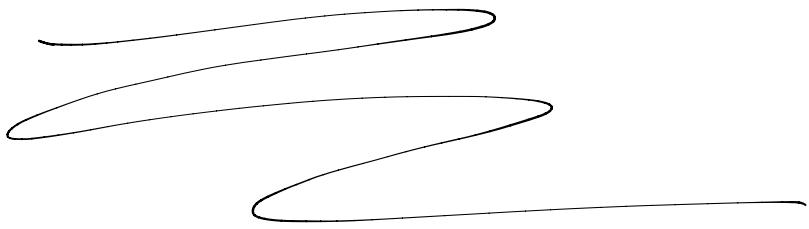
 until S is terminal



MC



n-step
bootstrapping



1-step
TD

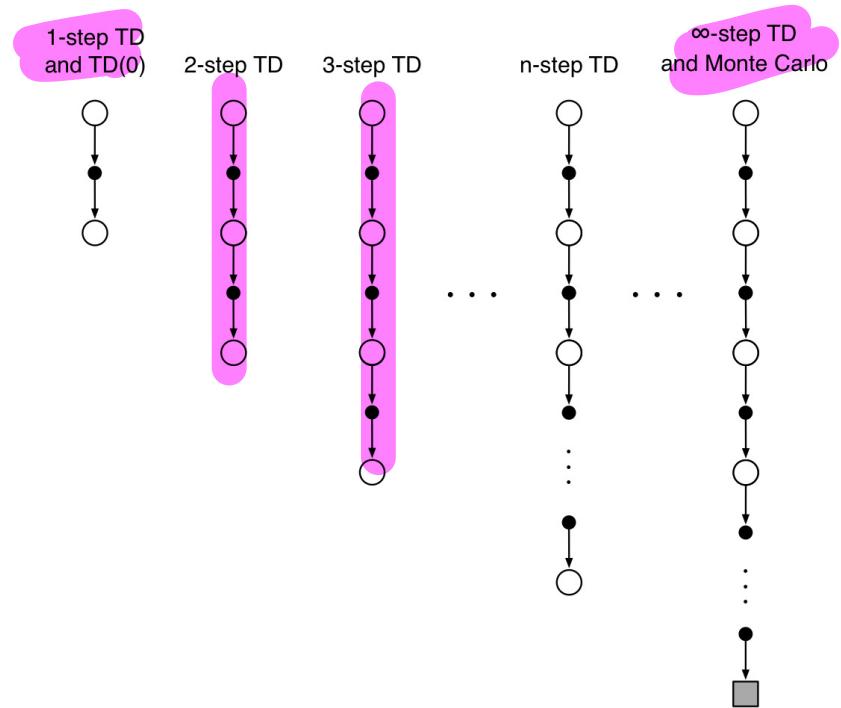


n-step TD

MC

1-step TD

n-step TD Prediction:



$S_t, R_{t+1}, S_{t+1}, R_{t+2}, \dots, R_T, S_T$

$$\underline{\text{MC}}: G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots + \gamma^{T-t-1} R_T$$

$$\underline{\text{1-step TD}}: G_{t:t+1} = R_{t+1} + \gamma V_t(S_{t+1})$$

$$\underline{\text{2-step TD}}: G_{t:t+1} = R_{t+1} + \gamma R_{t+2} + \gamma^2 V_{t+1}(S_{t+2})$$

$$\underline{\text{n-step TD}}: G_{t:t+n} = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{n-1} R_{t+n} + \gamma^n V_{t+n-1}(S_{t+n})$$

$t+n > T$, missy — zero.

$$V_{t+n}(s_t) = V_{t+n-1}(s_t) + \alpha [G_{t:t+n} - V_{t+n-1}(s_t)]$$

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

Input: a policy π

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 $t = 0, 1, 2, \dots :$

 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 $T \leftarrow t + 1$

$\tau \leftarrow t - n + 1$ (τ is the time whose state's estimate is being updated)

 If $\tau \geq 0$:

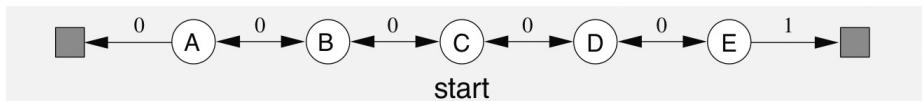
$G \leftarrow \sum_{i=\tau+1}^{\min(\tau+n, T)} \gamma^{i-\tau-1} R_i$

 If $\tau + n < T$, then: $G \leftarrow G + \gamma^n V(S_{\tau+n})$

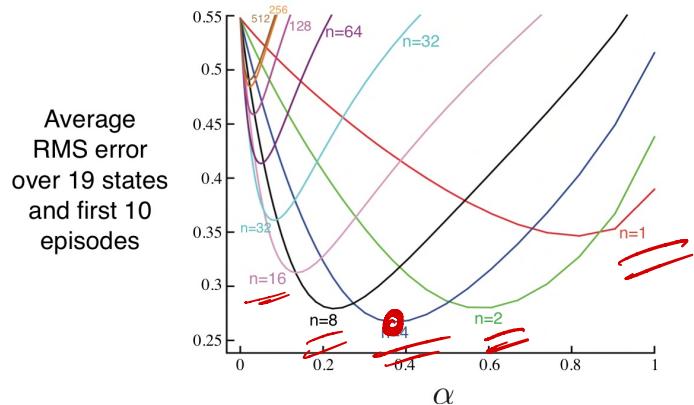
$V(S_\tau) \leftarrow V(S_\tau) + \alpha [G - V(S_\tau)]$

$(G_{\tau:\tau+n})$

 Until $\tau = T - 1$



19 states instead of 5, with -1 outcome on left, all values initialized to 0.



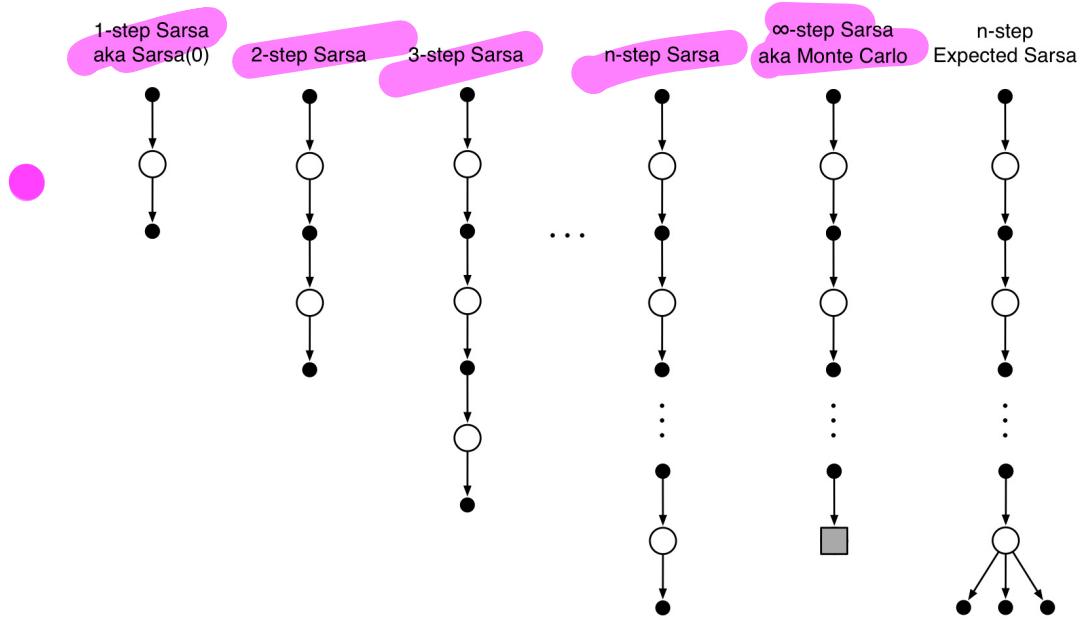
n -step
hyperparam

n -step Sarsa :-

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

$$\text{with } G_{t:t+n} = G_t \text{ if } t+n > T.$$

$$Q_{t+n}(s_t, a_t) = Q_{t+n-1}(s_t, a_t) + \alpha \left[G_{t:t+n} - Q_{t+n-1}(s_t, a_t) \right]$$



***n*-step Sarsa for estimating $Q \approx q_*$ or q_π**

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

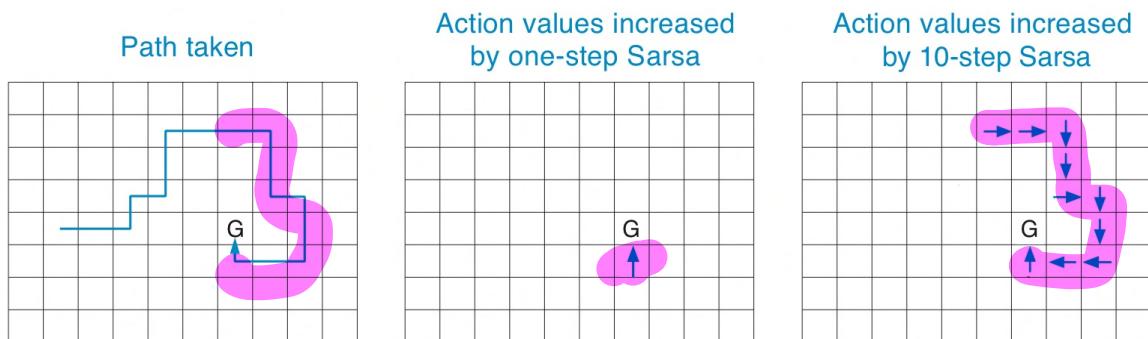
Loop for each episode:

```

    Initialize and store  $S_0 \neq$  terminal
    Select and store an action  $A_0 \sim \pi(\cdot | S_0)$ 
     $T \leftarrow \infty$ 
    Loop 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  $A_{t+1} \sim \pi(\cdot | S_{t+1})$ 
             $\tau \leftarrow t - n + 1$  ( $\tau$  is the time whose estimate is being updated)
            If  $\tau \geq 0$ :
                 $G \leftarrow \sum_{i=\tau+1}^{\min(\tau+n, T)} \gamma^{i-\tau-1} R_i$ 
                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  $\pi$  is being learned, then ensure that  $\pi(\cdot | S_\tau)$  is  $\varepsilon$ -greedy wrt  $Q$ 
    Until  $\tau = T - 1$ 

```

Illustration! 1-step Sarsa vs. 10-step Sarsa.



n-step expected Sarsa

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

$$\widehat{V}_t(s) = \sum_a \pi(a|s) Q_t(s, a).$$

n-step Off-Policy learning:

Importance Sampling:

$$V_{t+n}(s_t) = V_{t+n-1}(s_t) + \alpha P_{t:t+n-1} \left[G_{t:t+n} - V_{t+n-1}(s_t) \right]$$

where : $P_{t:h} = \prod_{k=t}^{\min(h, T-1)} \frac{\pi(A_k | s_k)}{b(A_k | s_k)}$

Off-Poly n-step Sarsa:

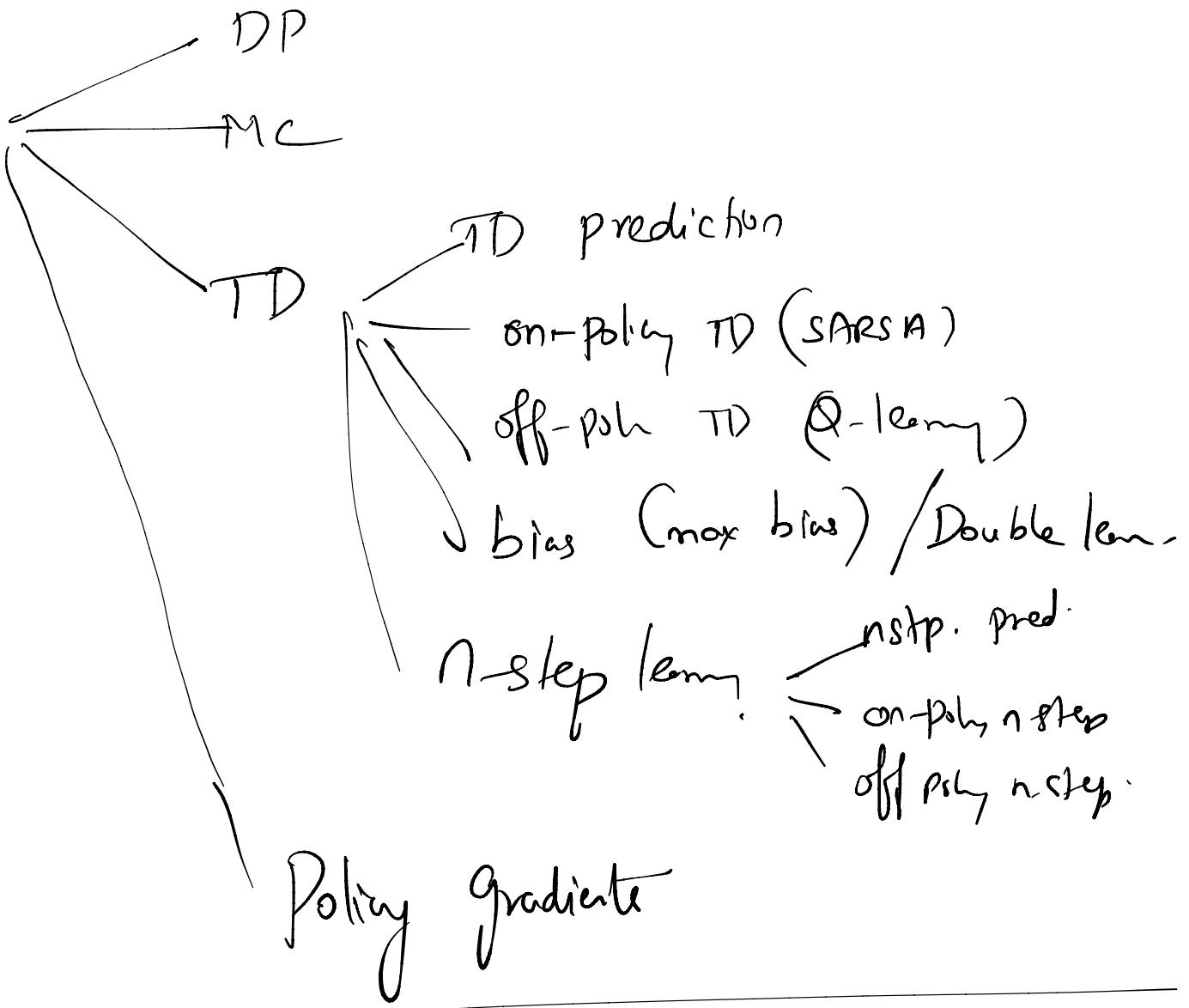
$$Q_{t+n}(s_t, a_t) = Q_{t+n-1}(s_t, a_t) + \alpha P_{t:t+n} \left[G_{t:t+n} - Q_{t+n+1}(s_{t+n}, a_{t+n}) \right]$$

Off-policy n -step Sarsa for estimating $Q \approx q_*$ or q_π

Input: an arbitrary behavior policy b such that $b(a|s) > 0$, for all $s \in \mathcal{S}, a \in \mathcal{A}$
 Initialize $Q(s, a)$ arbitrarily, for all $s \in \mathcal{S}, a \in \mathcal{A}$
 Initialize π to be greedy with respect to Q , or as a fixed given policy
 Algorithm parameters: step size $\alpha \in (0, 1]$, a positive integer n
 All store and access operations (for S_t , A_t , and R_t) can take their index mod $n + 1$

Loop for each episode:

- Initialize and store $S_0 \neq$ terminal
- Select and store an action $A_0 \sim b(\cdot|S_0)$
- $T \leftarrow \infty$
- Loop 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 $A_{t+1} \sim b(\cdot|S_{t+1})$
 - $\tau \leftarrow t - n + 1$ (τ is the time whose estimate is being updated)
 - If $\tau \geq 0$:
 - $\rho \leftarrow \prod_{i=\tau+1}^{\min(\tau+n, T-1)} \frac{\pi(A_i|S_i)}{b(A_i|S_i)}$ $(\rho_{\tau+1:\tau+n})$
 - $G \leftarrow \sum_{i=\tau+1}^{\min(\tau+n, T)} \gamma^{i-\tau-1} R_i$
 - If $\tau + n < T$, then: $G \leftarrow G + \gamma^n Q(S_{\tau+n}, A_{\tau+n})$ $(G_{\tau:\tau+n})$
 - $Q(S_\tau, A_\tau) \leftarrow Q(S_\tau, A_\tau) + \alpha \rho [G - Q(S_\tau, A_\tau)]$
 - If π is being learned, then ensure that $\pi(\cdot|S_\tau)$ is greedy wrt Q
- Until $\tau = T - 1$



— Function approximation