

Temporal Difference Method

Similar to the MC method, temporal difference methods only require experiences, i.e., sample trajectories (sequence of states, actions, rewards).

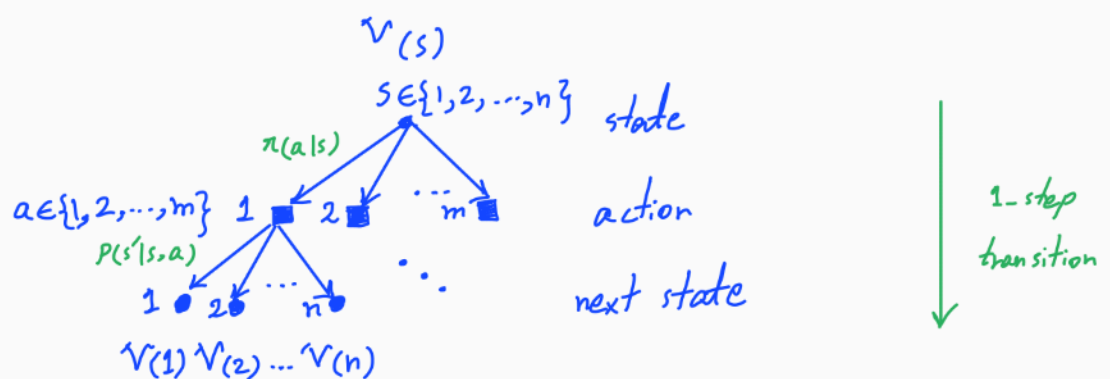
Model-free method

Temporal difference: Using dynamic programming over estimated values to learn online from incomplete episodes.

TD for policy evaluation: Given a stationary policy π , we want to compute

$$V^{\pi}(s) = \mathbb{E}_{\substack{A_t \sim \pi \\ S_{t+1} \sim P(\cdot | S_t, A_t)}} \left[\sum_{t=0}^{\infty} \gamma^t R(S_t, A_t) \mid S_0 = s \right]$$

Idea: Estimate the expected return using the sample immediate reward and the current estimate of the future expected return value function).



Let \hat{V}^π be the current estimate of the value function of π

$$\longrightarrow V^\pi(s) = \mathbb{E}_{\substack{A_t \sim \pi \\ s_{t+1} \sim P(\cdot | s_t, A_t)}} \left[\sum_{t=0}^T \gamma^t R(s_t, A_t) \mid s_0 = s \right]$$

$$= \mathbb{E}_{\substack{A_t \sim \pi \\ s_{t+1} \sim P(\cdot | s_t, A_t)}} \left[R(s_0, A_0) + \gamma \sum_{t=1}^T \gamma^t R(s_t, A_t) \mid s_0 = s \right]$$

Bellman consistency equation

$$V^\pi(s) = \mathbb{E}_{a, s'} [R(s, a) + \gamma V^\pi(s')]$$

$$= \mathbb{E}_{A_0 \sim \pi(s_0)} \left[R(s_0, A_0) + \underbrace{\gamma \mathbb{E}_{\substack{A_t \sim \pi \\ s_{t+1} \sim P(\cdot | s_t, A_t)}} \left[\sum_{t=0}^T \gamma^t R(s_t, A_t) \mid s_0 = s' \right]}_{V^\pi(s')} \right]$$

$$\longrightarrow \hat{V}^\pi(s) = \frac{1}{N} \sum_{i=1}^N (R(s, a^i) + \gamma \hat{V}^\pi(s'^i))$$

samples ← sample i ←

cannot be done online
for $N > 1$
↓
incremental computation
of mean

Recall the MC update:

$$\hat{V}_{N(s)}^\pi(s) = \frac{1}{N(s)} \sum_{i=1}^{N(s)} G_i^s$$

$$\hat{V}_{N(s)+1}^\pi(s) = \hat{V}_{N(s)}^\pi(s) + \frac{1}{N(s)+1} (G_{N(s)+1}^s - \hat{V}_{N(s)}^\pi(s))$$

$$\hat{V}^\pi(s) \leftarrow \hat{V}^\pi(s) + \alpha (G - \hat{V}^\pi(s))$$

for all states visited in the episode learning rate / step size MC target MC error

equivalent for certain choice of $\alpha = \frac{1}{N(s)+1}$

TD(0) update:

$$\hat{V}^\pi(s) \leftarrow \hat{V}^\pi(s) + \alpha (R(s, a) + \gamma \hat{V}^\pi(s') - \hat{V}^\pi(s))$$

for the state just visited learning rate / step size TD target TD error

TD (0) method

Iterative method:

- Initialize \hat{V}^π

- Repeat

- Sample $s_0 \in S$

- Repeat until the episode terminates

- Take action a according to $\pi(a|s)$ at s

- Observe reward $R(s, a)$ and next state s'

- $\hat{V}^\pi(s) \leftarrow \hat{V}^\pi(s) + \alpha [R(s, a) + \gamma \hat{V}^\pi(s') - \hat{V}^\pi(s)]$

- Update the current state $s \leftarrow s'$

$$\begin{bmatrix} \hat{V}^\pi_{(1)} \\ \hat{V}^\pi_{(2)} \\ \vdots \\ \hat{V}^\pi_{(n)} \end{bmatrix}_{|S| \times 1} \xrightarrow{s=2 \xrightarrow{\pi} a=1 \xrightarrow{P} s'=4, r=R(2,1)} \hat{V}^\pi_{(2)} \leftarrow \hat{V}^\pi_{(2)} + \alpha (r + \gamma \hat{V}^\pi_{(4)} - \hat{V}^\pi_{(2)})$$

TD for policy optimization: We want to find an (approximately) optimal policy π^* through iterative policy evaluation and policy improvement.

Idea:

- Evaluate the policy by estimating the Q function using the TD method in an online fashion.
- Improve the policy using the estimated Q function in a greedy manner.

SARSA $\rightarrow s, a, r, s', a'$

Iterative method:

- Initialize \hat{Q} s.t. $\hat{Q}(s, a) = 0$ for terminal states
- Initialize π according to \hat{Q} (e.g., ϵ -greedy)
- Repeat
 - Sample $s_0 \in S$
 - Select action a according to $\pi(a|s)$ at s_0
 - Repeat until the episode terminates

- Take action a

- Observe reward $R(s, a)$ and next state s'

- Select action a' according to $\pi(a|s)$ at s'

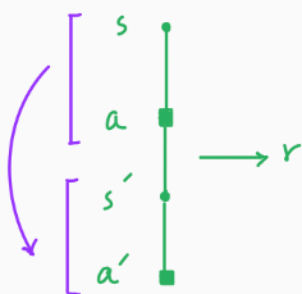
policy evaluation • $\hat{Q}(s, a) \leftarrow \hat{Q}(s, a) + \alpha [\underbrace{R(s, a) + \gamma \hat{Q}(s', a')}_{\text{TD target}} - \hat{Q}(s, a)]$

policy improvement • Build $\pi(a|s)$ according to \hat{Q} (e.g., ϵ -greedy)

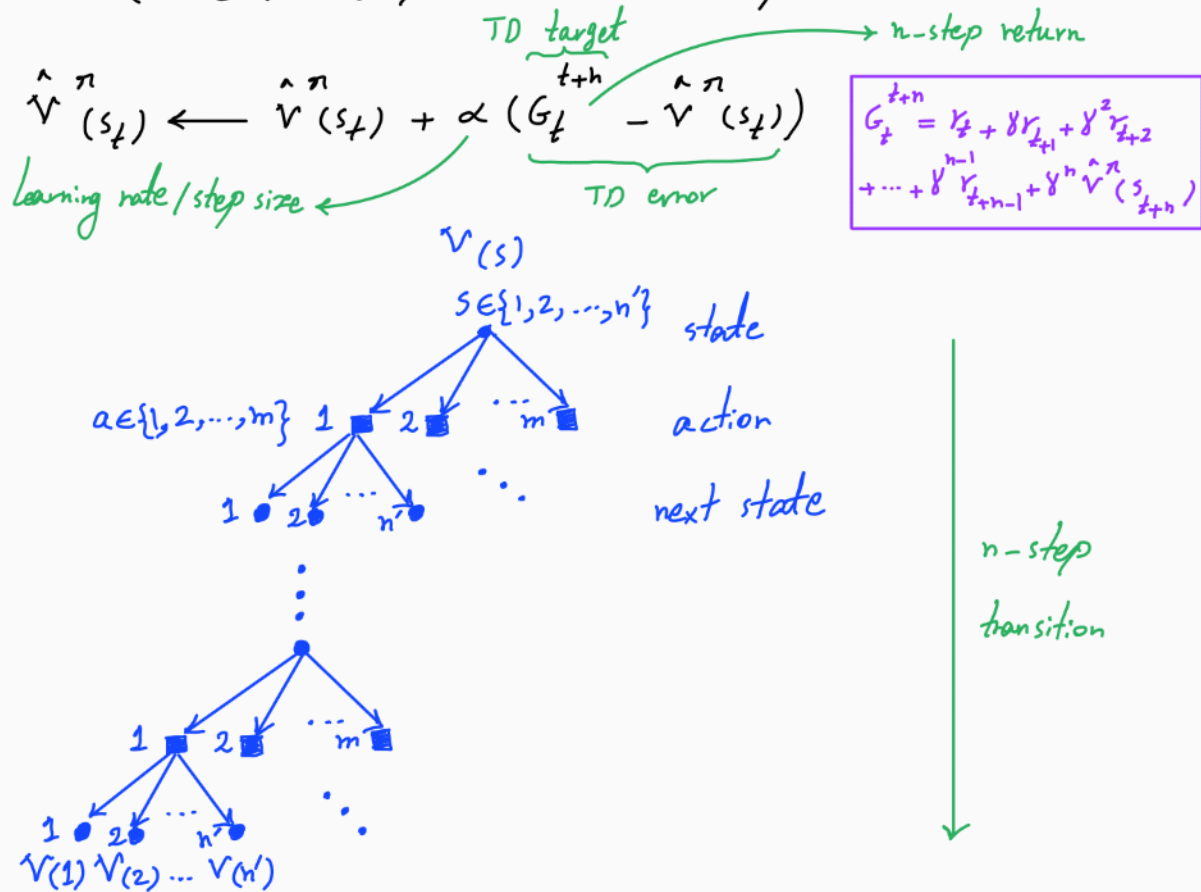
- Update the current state and the current action

$s \leftarrow s'$

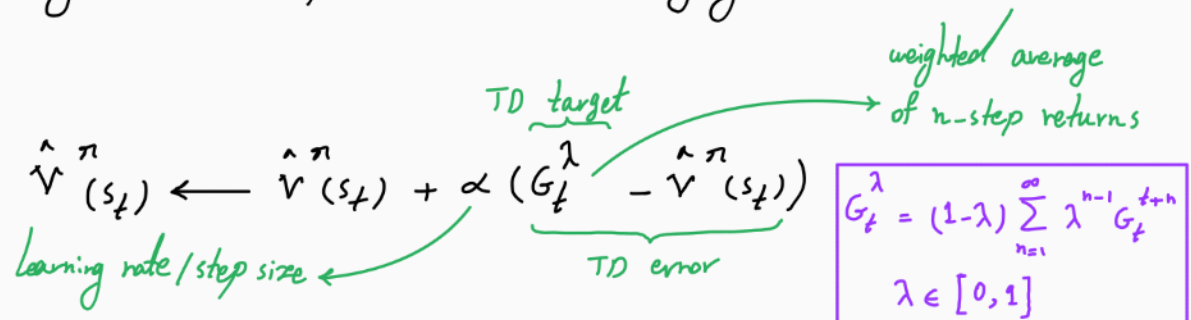
$a \leftarrow a'$



Note [n-step TD method]: Estimate the expected return using n sample immediate rewards and the current estimate of the future expected return (value function) after these n steps.



Note [TD(λ) method]: Updates the expected return using a weighted average of n -step returns over varying n .



* Backward view of TD(λ) using the idea of eligibility traces provides an incremental approach to approximate the forward view.