



Temporal Difference Methods for Prediction

Rupam Mahmood

February 28, 2020



(Exercise 6.7 SEB) Design an off-policy version of the TD(0) update that can be used with arbitrary target policy π and covering behavior policy b , using at each step t the importance sampling ratio $\rho_{t:t}$ (5.3).

$$\rho_{t:T-1} \doteq \frac{P(A_t, S_{t+1}, A_{t+1}, \dots, S_T | S_t, A_{t:T-1} \sim \pi)}{P(A_t, S_{t+1}, A_{t+1}, \dots, S_T | S_t, A_{t:T-1} \sim b)} = \frac{\prod_{k=t}^{T-1} \pi(A_k | S_k)}{\prod_{k=t}^{T-1} b(A_k | S_k)}$$

on-policy sample-average MC: $A_t \sim \pi \implies E_\pi [G_t | S_t = s] = v_\pi(s) \rightarrow V(s) \doteq \frac{\sum_{t \in \mathcal{T}(s)} G_t}{|\mathcal{T}(s)|}$

off-policy sample-average MC: $A_t \sim b \implies E_b [\rho_{t:T-1} G_t | S_t = s] = v_\pi(s) \rightarrow V(s) \doteq \frac{\sum_{t \in \mathcal{T}(s)} \rho_{t:T-1} G_t}{|\mathcal{T}(s)|}$

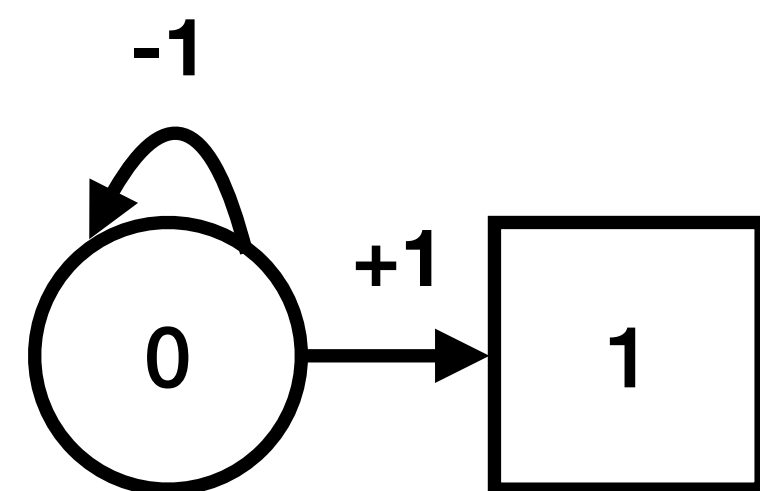
On-policy constant- α MC: $V^{MC}(S_t) \leftarrow V^{MC}(S_t) + \alpha [G_t - V^{MC}(S_t)]$

Off-policy constant- α MC: ?

On-policy TD(0): $V^{TD}(S_t) \leftarrow V^{TD}(S_t) + \alpha [R_{t+1} + \gamma V^{TD}(S_{t+1}) - V^{TD}(S_t)]$

Off-policy TD(0): ?

Live demo of TD updates



Tabular TD(0) for estimating v_π

Input: the policy π to be evaluated

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

Initialize $V(s)$, for all $s \in \mathcal{S}^+$, 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