

Challenge Question Solution

Recall in 6.1, we compared TD target and MC target, where the MC target G_t is the actual return following time t , and the TD target being $R_{t+1} + \gamma V(S_{t+1})$, in which R_{t+1} is the observed reward at time t and $V(S_{t+1})$ is our value estimate of S_{t+1} .

Given $S_t = s$, and S_{t+1} , we can rewrite G_t as $R_{t+1} + \gamma G_{t+1}$.

Therefore,

$$\begin{aligned} E[G_t | S_t = s, S_{t+1}] \\ = E[R_{t+1} + \gamma G_{t+1} | S_t = s, S_{t+1}] \end{aligned}$$

Since reward R_{t+1} is observed and γ is a constant discount rate,

$$= R_{t+1} + \gamma E[G_{t+1} | S_{t+1}] \quad \text{under the markov assumption}$$

Recall the definition of the value function in 3.5, $v_{\pi}(s) = E_{\pi}[G_t | S_t = s]$

$$= R_{t+1} + \gamma v_{\pi}(S_{t+1})$$

Because we assume in the question that we have found a V that exactly equals v_{π} , we have

$$E[G_t | S_t = s, S_{t+1}] = R_{t+1} + \gamma V(S_{t+1}) \quad (\text{TD target})$$

Provided the variance decomposition equation of C_t , we have

$$\begin{aligned} \text{Var}(C_t | S_{t=s}) &= E[\text{Var}(C_t | S_{t=s}, S_{t+1})] + \text{Var}(E[C_t | S_{t=s}, S_{t+1}] | S_{t=s}) \\ \text{MC Target} &= \underbrace{E[\text{Var}(C_t | S_{t=s}, S_{t+1})]}_{\geq 0} + \text{Var}(\underbrace{R_{t+1} + rV(S_{t+1})}_{\text{TD target}} | S_{t=s}) \end{aligned}$$

Since the variance of C_t is always greater or equal to 0, so does its expectation. Hence, we've shown that the Monte Carlo target has at least as high of variance as the TD target, whose variance is just the second term in our final result.