

When does TD have lower variance than MC?

Sam Greydanus
sgrey@google.com

September 11, 2018

Abstract

The goal of this writeup is to clarify the connection between Monte Carlo (MC) and Temporal Difference (TD) value estimation. I will also show that TD value estimates always have equal or lower variance than MC estimates.

1 Setup

Monte Carlo (MC). Monte Carlo methods estimate the value of a state s_t by computing the average return $G(s_t)$ at that state. A simple every-visit MC update would look like Equation 1 where $\alpha = \frac{1}{k}$ and k is the visit count of s_t .

$$V(s_t) \leftarrow V(s_t) + \alpha[G(s_t) - V(s_t)] \quad (1)$$

Temporal Differences (TD). Another way to estimate the value of a state is with *temporal difference* (TD) learning. The idea of TD learning is to express the value of s_t in terms of the values of its successor state s_{t+1} along a trajectory, plus the change in reward due to the $s_t \rightarrow s_{t+1}$ transition.

$$V(s_t) \leftarrow V(s_t) + \alpha[r_{t+1} + \gamma V(s_{t+1}) - V(s_t)] \quad (2)$$

This approach is often called *bootstrapping* because each state uses the next state's value to update its own.

Implementation notes. Looking at Equations 1 and 2, we can see that the core difference between MC and TD updates is the substitution $G(s_t) = r_{t+1} + \gamma V(s_{t+1})$. Unfortunately, the classic algorithms for MC and TD value estimation are very different. In the textbook, *Reinforcement Learning: An Introduction* [1], the MC and TD algorithms are given on pages 92 and 120 respectively. For MC learning, the authors compute $G_0 \dots G_N$ across each episode and average across episodes. For TD learning, the authors update $V(s_t)$ after each state transition in an episode.

In spite of these textbook differences, most modern implementations of MC and TD are strikingly similar. The value updates happen at the end of each episode and proceed in reverse order from the last episode in the trajectory to the first. We choose to follow this format (and thus depart slightly from textbook definitions) for two reasons. First, doing so enables us to make tighter analogies between MC and TD. Second, most real-world applications of RL algorithms use this format.

Pseudocode for our versions of MC and TD value estimation can be found in Figure 1.

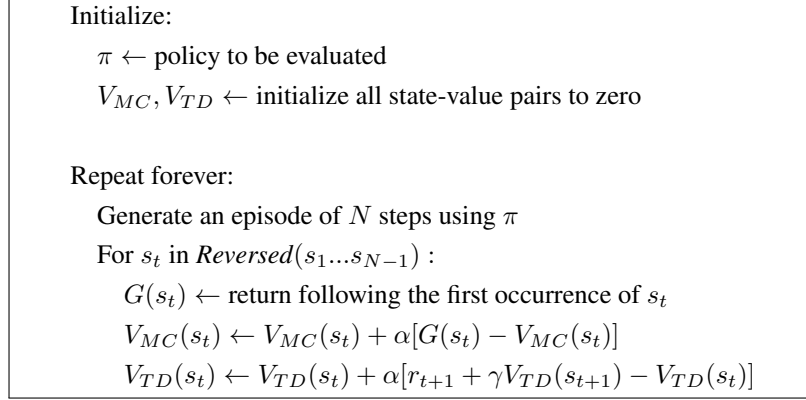


Figure 1: Estimating value using Monte Carlo and Temporal Difference methods.

2 Proof

In this section, we will show why TD value estimates always have equal or lower variance than MC estimates. Note that this proof only works for tabular agents/environments; when state approximators are introduced, all bets are off.

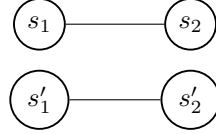


Figure 2: Adjacent states along two different trajectories.

Consider two adjacent states s_1 and s_2 along trajectory T , and two adjacent states s'_1 and s'_2 along trajectory T' (see Figure 2). We intend to perform a value update using first T and then T' .

Now, consider all cases for which Equation 3 holds:

$$V_{TD}(s_2) = G(s_2) \quad \text{and} \quad V_{TD}(s'_2) = G(s'_2) \quad \text{assumption of Lemma 1} \quad (3)$$

Lemma 1. *Wherever Equation 3 is true, TD and MC updates are identical.* The way to see this is by expressing the TD updates along trajectory T' using Equation 2:

$$V_{TD}(s'_1) = V_{TD}(s_1) + \alpha(r_{12'} + \gamma V_{TD}(s'_2) - V_{TD}(s_1)) \quad \text{TD update rule} \quad (4)$$

The next step is to recognize that the MC and TD value estimates after the first value update - the one along T - are identical:

$$V_{TD}(s_1) = 0 + (r_{12} + \gamma V_{TD}(s_2) - 0) \quad \text{initial value } V(s_1) = 0 \quad (5)$$

$$= r_{12} + \gamma G(s_2) = G(s_1) \quad \text{definition of return} \quad (6)$$

$$= V_{MC}(s_1) \quad \text{definition of first MC update} \quad (7)$$

Substituting Equation 7 into Equation 4, we have

$$V_{TD}(s'_1) = V_{MC}(s_1) + \alpha(r_{12'} + \gamma V_{TD}(s'_2) - V_{MC}(s_1)) \quad (8)$$

$$= V_{MC}(s_1) + \alpha(r_{12'} + \gamma G(s'_2) - V_{MC}(s_1)) \quad (9)$$

$$= V_{MC}(s_1) + \alpha(G(s'_1) - V_{MC}(s_1)) \quad (10)$$

$$= V_{MC}(s'_1) \quad (11)$$

Equations 8-11 tell us that MC and TD updates along T and T' are identical whenever Equation 3 is true.

Lemma 2. When $s_1 \neq s'_1$ and $s_2 = s'_2$, the TD update has lower variance than the MC update. To see this, we must first recall that value estimates are simply averages. The variance of these averages is proportional to $\frac{\sigma^2}{n}$ where n is the number of items being averaged over. Next, we need to introduce a count function, $N(s)$ which measures the number of times the agent has visited state s . In the case we are considering,

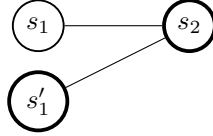


Figure 3: In Lemma 2, we consider the case where s_2 and s'_2 are the same state but s_1 and s'_1 are not.

$$N(s_1) = N(s'_1) = 1 \quad \text{and} \quad N(s_2) = N(s'_2) = 2 \quad (12)$$

$$\text{Var}[G(s_1)] = \text{Var}[G(s'_1)] = \sigma^2 \quad (13)$$

$$\text{Var}[G(s_2)] = \text{Var}[G(s'_2)] = \frac{\sigma^2}{2} \quad (14)$$

Now recall that the difference between a MC update and a TD update is given by the substitution $G(s_t) = r_{t+1} + \gamma V(s_{t+1})$. Let's compare the variances of these quantities at state s_1

$$\text{Var}[G(s_1)] \stackrel{?}{>} \text{Var}[r_{12} + \gamma V_{TD}(s_2)] \quad \text{variance of both terms} \quad (15)$$

$$\stackrel{?}{>} \text{Var}[r_{12}] + \text{Var}[\gamma V_{TD}(s_2)] \quad \text{assume independence} \quad (16)$$

$$> 0 + \gamma \text{Var}[G(s_2)] \quad \text{assume deterministic rewards and } \gamma < 1 \quad (17)$$

Proof. Whenever s_2 or s'_2 is a terminal node, the assumption in Equation 3 is true. By induction, this assumption holds for each preceding state, except for the case where $s_1 \neq s'_1$ and $s_2 = s'_2$.

In this case, shown in Figure 2, $V_{TD}(s'_2)$ gets updated during the first value update - the one along T - but this value update never gets reflected in the return $G(s'_1)$. Thus $V_{TD}(s'_1)$ and $V_{MC}(s'_1)$ go out of sync. But we've shown in Lemma 2 that TD updates always have lower variance than MC updates in this case. This means that, regardless of how trajectories T and T' overlap, TD updates will always have equal or lower variance than MC updates. This argument extends to an arbitrary number of trajectories and updates.

References

- [1] Richard Sutton and Andrew Barto. *Reinforcement Learning: An Introduction*. MIT Press, Cambridge, 2 edition, 2018.