

1. In Chapter 6 we noted that the Monte Carlo error can be written as the sum of TD errors (6.6) if the value estimates don't change from step to step. Show that the  $n$ -step error used in (7.2) can also be written as a sum of TD errors (again if the value estimates don't change) generalizing the earlier result.

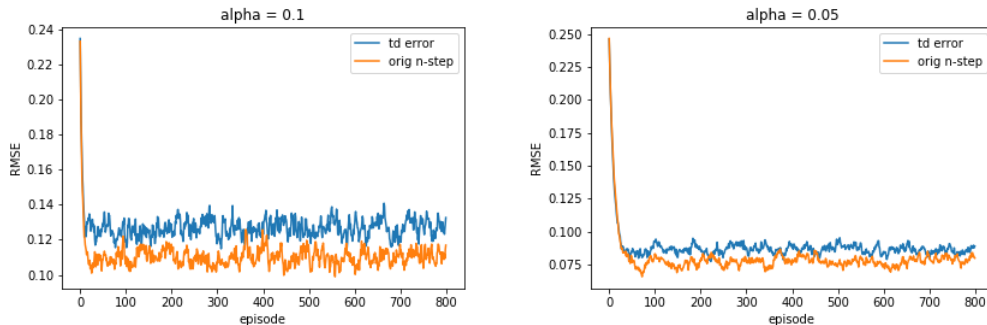
$$\begin{aligned}
 G_{t:t+n} - V(S_t) &= R_{t+1} + \gamma G_{t+1:t+n} - V(S_t) + \gamma V(S_{t+1}) - \gamma V(S_{t+1}) \\
 &= \delta_t + \gamma(G_{t+1:t+n} - V(S_{t+1})) \\
 &= \delta_t + \gamma\delta_{t+1} + \gamma^2\delta_{t+2} + \cdots + \gamma^{n-1}\delta_{t+n-1} + \gamma^n V(S_{t+n}) - \gamma^n V(S_{t+n}) \\
 &= \sum_{k=t}^{t+n-1} \gamma^{k-t} \delta_k
 \end{aligned}$$

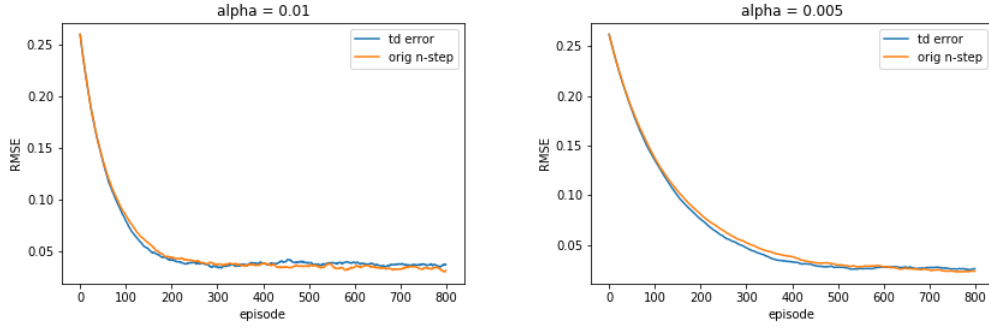
2. *Programming.* With an  $n$ -step method, the value estimates do change from step to step, so an algorithm that used the sum of TD errors (see previous exercise) in place of the error in (7.2) would actually be a slightly different algorithm. Would it be a better algorithm or a worse one? Devise and program a small experiment to answer this question empirically.

There is only a difference in the two methods if we update state estimates for some states from  $\{S_{t+1}, \dots, S_{t+n}\}$  between time steps  $t$  and  $t+n$ . Because of this reason, I would like to try the algorithms on a problem where we may return to the same state shortly. My expectation is that the  $n$ -step update works better, because it uses a more recent estimation of  $V(S_t)$ . It does not use estimations of  $V(S_{t+i})$  (for  $i = 1 \dots < n$ ) either, but I am not sure if it is for the better or the worse. Let's try it on a random walk example, like in Example 6.2.

The code can be found at [https://github.com/hannagabor/SBRL/blob/master/7.2/random\\_walks\\_td.ipynb](https://github.com/hannagabor/SBRL/blob/master/7.2/random_walks_td.ipynb).

The results show that if alpha is bigger, then the TD-error update is worse than the original  $n$ -step method. If alpha is small enough, then there is no difference. (At least on this problem.)





3. Why do you think a larger random walk task (19 states instead of 5) was used in the examples of this chapter? Would a smaller walk have shifted the advantage to a different value of  $n$ ? How about the change in left-side outcome from 0 to -1 made in the larger walk? Do you think that made any difference in the best value of  $n$ ?

If  $n$  is big, but the random walk has only a few states, then there is a higher chance that we update states that are far from the end state. That sounds bad.

As for the left-side outcome: I can't see any change.

4. Prove that the  $n$ -step return of Sarsa (7.4) can be written exactly in terms of a novel TD error, as

$$G_{t:t+n} = Q_{t-1}(S_t, A_t) + \sum_{k=t}^{\min(t+n, T)-1} \gamma^{k-t} (R_{k+1} + \gamma Q_k(S_{k+1}, A_{k+1}) - Q_{k-1}(S_k, A_k))$$

$$\begin{aligned} G_{t:t+n} &= R_{t+1} + \gamma G_{t+1:t+n} + Q_{t-1}(S_t, A_t) - Q_{t-1}(S_t, A_t) \\ &\quad + \gamma Q_t(S_{t+1}, A_{t+1}) - \gamma Q_t(S_{t+1}, A_{t+1}) \\ &= Q_{t-1}(S_t, A_t) + (R_{t+1} + \gamma Q_t(S_{t+1}, A_{t+1}) - Q_{t-1}(S_t, A_t)) \\ &\quad + \gamma (G_{t+1:t+n} - Q_t(S_{t+1}, A_{t+1})) \\ &\quad + \gamma^2 Q_{t+1}(S_{t+2}, A_{t+2}) - \gamma^2 Q_{t+1}(S_{t+2}, A_{t+2}) \\ &= Q_{t-1}(S_t, A_t) + (R_{t+1} + \gamma Q_t(S_{t+1}, A_{t+1}) - Q_{t-1}(S_t, A_t)) \\ &\quad + \gamma (R_{t+2} + \gamma^2 Q_{t+1}(S_{t+2}, A_{t+2}) - Q_t(S_{t+1}, A_{t+1})) \\ &\quad + \gamma^2 (G_{t+2:t+n} - Q_{t+1}(S_{t+2}, A_{t+2})) \\ &\quad + \gamma^3 Q_{t+2}(S_{t+3}, A_{t+3}) - \gamma^2 Q_{t+2}(S_{t+3}, A_{t+3}) \\ &\quad \dots \\ &= Q_{t-1}(S_t, A_t) + \sum_{k=t}^{\min(t+n, T)-1} \gamma^{k-t} (R_{k+1} + \gamma Q_k(S_{k+1}, A_{k+1}) - Q_{k-1}(S_k, A_k)) \end{aligned}$$