

---

# True Online TD( $\lambda$ )

---

Harm van Seijen  
Richard S. Sutton

HARM.VANSEIJEN@UALBERTA.CA  
SUTTON@CS.UALBERTA.CA

Department of Computing Science, University of Alberta, Edmonton, Alberta, T6G 2E8, Canada

## Abstract

TD( $\lambda$ ) is a core algorithm of modern reinforcement learning. Its appeal comes from its equivalence to a clear and conceptually simple forward view, and the fact that it can be implemented online in an inexpensive manner. However, the equivalence between TD( $\lambda$ ) and the forward view is exact only for the off-line version of the algorithm (in which updates are made only at the end of each episode). In the online version of TD( $\lambda$ ) (in which updates are made at each step, which generally performs better and is always used in applications) the match to the forward view is only approximate. In a sense this is unavoidable for the conventional forward view, as it itself presumes that the estimates are unchanging during an episode. In this paper we introduce a new forward view that takes into account the possibility of changing estimates and a new variant of TD( $\lambda$ ) that *exactly* achieves it. Our algorithm uses a new form of eligibility trace similar to but different from conventional accumulating and replacing traces. The overall computational complexity is the same as TD( $\lambda$ ), even when using function approximation. In our empirical comparisons, our algorithm outperformed TD( $\lambda$ ) in all of its variations. It seems, by adhering more truly to the original goal of TD( $\lambda$ )—matching an intuitively clear forward view even in the online case—that we have found a new algorithm that simply improves on classical TD( $\lambda$ ).

## 1. Why True Online TD( $\lambda$ ) Matters

Temporal-difference (TD) learning is a core learning technique in modern reinforcement learning (Sutton, 1988; Kaelbling, Littman & Moore; Sutton & Barto, 1998; Szepesvári, 2014). One of the main challenges in rein-

forcement learning is to make predictions, in an initially unknown environment, about the (discounted) sum of future rewards, the *return*, based on currently observed feature and a certain behavior policy. With TD learning it is possible to learn good estimates of the expected return quickly by bootstrapping from other expected-return estimates. TD( $\lambda$ ) is a popular TD algorithm that combines basic TD learning with *eligibility traces* to further speed learning. The popularity of TD( $\lambda$ ) can be explained by its simple implementation, its low computational complexity, and its conceptually straightforward interpretation, given by its *forward view*.

The forward view of TD( $\lambda$ ) (Sutton & Barto, 1998) is that the estimate at each time step is moved toward an update target known as the  $\lambda$ -return; the corresponding algorithm is known as the  $\lambda$ -return algorithm. The  $\lambda$ -return is an estimate of the expected return based on both subsequent rewards and the expected-return estimates at subsequent states, with  $\lambda$  determining the precise way these are combined. The forward view is useful primarily for understanding the algorithm theoretically and intuitively. It is not clear how the  $\lambda$ -return could form the basis for an online algorithm, in which estimates are updated on every step, because the  $\lambda$ -return is generally not known until the end of the episode. Much clearer is the off-line case, in which values are updated only at the end of an episode, summing the value corrections corresponding to all the time steps of the episode. For this case, the overall update of the  $\lambda$ -return algorithm has been proven equal to the off-line version of TD( $\lambda$ ), in which the updates are computed on each step, but not used to actually change the value estimates until the end of the episode, when they are summed to produce an overall update (Sutton & Barto, 1998). The overall per-episode updates of the  $\lambda$ -return algorithm and of off-line TD( $\lambda$ ) are the same, even though the updates on each time step are different. By the end of the episode they must sum to the same overall update.

One of TD( $\lambda$ )’s most appealing features is that at each step it only requires temporally local information—the immediate next reward and next state; this enables the algorithm to be applied online. The online updates will be slightly dif-

ferent from those of the off-line version of TD( $\lambda$ ) because the estimates change during the episode and these participate in determining the  $\lambda$ -returns. Thus the overall effect of online TD( $\lambda$ ) by the end of the episode will not equal that of the off-line  $\lambda$ -return algorithm. They may be close if the step-size parameter is small and the episode is short, but in general the overall updates will differ. On balance this is a good thing (for online TD( $\lambda$ )) because online TD updates are generally better than off-line ones (see Sutton & Barto, Sections 7.1–3). What is needed is not an online way to produce the same results as the (off-line)  $\lambda$ -return algorithm, but rather a new forward view for the online case, with a new online  $\lambda$ -return-like algorithm, and then a new version of TD( $\lambda$ ) that achieves the same result when applied online.

This is exactly what we provide in this paper. First we present a new online forward view, based on a modified  $\lambda$ -return, which we call the *truncated  $\lambda$ -return*. Then we present a new version of TD( $\lambda$ ) that obtains exact step-by-step equivalence with our online forward view. We call this variation *true online TD( $\lambda$ )* because it is truer to the idea of TD( $\lambda$ ) as matching a  $\lambda$ -based forward view. Empirically, we demonstrate that true online TD( $\lambda$ ) also performs better on three benchmark problems.

## 2. Problem Setting and Notation

In this section, we formally present our learning framework, the various versions of TD( $\lambda$ ), and the conventional forward view. As a convention, we indicate random variables by capital letters (e.g.,  $S_t, R_t$ ), vectors by bold letters (e.g.,  $\theta, \phi$ ), functions by lowercase letters (e.g.,  $v$ ), and sets by calligraphic font (e.g.,  $\mathcal{S}, \mathcal{A}$ ).

### 2.1. Markov Reward Processes

We focus in this paper on discrete-time *Markov reward processes* (MRPs), which can be described as 4-tuples of the form  $\langle \mathcal{S}, p, r, \gamma \rangle$ , consisting of  $\mathcal{S}$ , the set of all states;  $p(s'|s)$ , the transition probability function, giving for each state  $s \in \mathcal{S}$  the probability of a transition to state  $s' \in \mathcal{S}$  at the next step;  $r(s, s')$ , the reward function, giving the expected reward after a transition from  $s$  to  $s'$ .  $\gamma$  is the discount factor, specifying how future rewards are weighted with respect to the immediate reward.

The *return* at time step  $t$  is the discounted sum of rewards observed after time step  $t$ :

$$G_t = \sum_{i=1}^{\infty} \gamma^{i-1} R_{t+i}.$$

An MRP can contain *terminal states*, dividing the sequence of state transitions into *episodes*. When a terminal state is

reached the current episode ends and the state is reset to the initial state. The return for an episodic MRP is defined as:

$$G_t = \sum_{i=1}^{T-t} \gamma^{i-1} R_{t+i},$$

where  $T$  is the time step that the terminal state is reached.

We are interested in learning the value-function  $v$  of an MRP, which maps each state  $s \in \mathcal{S}$  to the expected value of the return:

$$v(s) = \mathbb{E}\{G_t | S_t = s\}.$$

### 2.2. Function Approximation

In the general case, the state-space can be huge or even continuous. In such cases it is not possible to represent the value function  $v$  exactly. In this case, or when data-generalization is desired, function approximation is used to represent  $v$ . The approximate value-function  $\hat{v}(s, \theta)$  gives the approximate value of state  $s$ , given a weight vector  $\theta$ . Our learning algorithms produce a sequence of weight vectors  $\theta_t$ , and as a shorthand we use  $\hat{v}_t(s) = \hat{v}(s, \theta_t)$ .

A common approach is to use linear function approximation, in which case the value of a state is the inner product between the weight vector  $\theta$  and a state-dependent feature vector  $\phi(s)$ . In this case, the value of state  $s$  at time step  $t$  is approximated by:

$$\hat{v}_t(s) = \theta_t^\top \phi(s) = \sum_i \theta_{i,t} \phi_i(s),$$

where  $\phi_i(s)$  is the value of feature  $i$  in state  $s$ , and  $\theta_{i,t}$  is the weight of that feature at time step  $t$ . If  $s$  is a terminal state, then by definition  $\phi(s) = 0$ , so that  $\hat{v}_t(s) = 0$ . We also use the shorthand  $\phi_t = \phi(S_t)$ .

Although TD( $\lambda$ ) is not strictly speaking a gradient-descent method (Barnard, 1993) it is still useful to view it in these terms (as in Sutton & Barto, 1998). Stochastic gradient descent involves incremental updates of the weight vector  $\theta$  at each time step in the direction of the gradient of the time step's error. The updates can generally be written as

$$\theta_{t+1} = \theta_t + \alpha [U_t - \hat{v}_t(S_t)] \nabla_{\theta_t} \hat{v}_t(S_t), \quad (1)$$

where  $U_t$  is the *update target* and  $\nabla_{\theta_t} \hat{v}_t(S_t)$  is the gradient of  $\hat{v}$  with respect to  $\theta_t$ . In the case of linear function approximation, this gradient has a particularly simple form:

$$\nabla_{\theta_t} \hat{v}_t(S_t) = \phi_t.$$

There are many ways to construct an update target  $U_t$  from observed states and rewards. For example, *Monte Carlo* methods use the full return as the update target:

$$U_t = G_t.$$

TD methods use an update target that is based on value estimates of other states. For example, the TD(0) update target is:

$$U_t = R_{t+1} + \gamma \hat{v}_t(S_{t+1}). \quad (2)$$

The update (1) is an *online* update, meaning that the weight vector is updated at every time step  $t$ . Alternatively, the weight vector can be updated *off-line*. With off-line updating, the weight vector stays constant during an episode, and instead the weight corrections are collected on the side. Let the weight vector for (all of) episode  $k$  be  $\theta_k$ . The weight correction for time step  $t$  of this episode is:

$$\Delta_t = \alpha [U_t - \hat{v}_k(S_t)] \nabla_{\theta_k} \hat{v}_k(S_t).$$

After the episode has terminated, the weight vector is updated by adding all weight corrections collected during the episode:

$$\theta_{k+1} = \theta_k + \sum_{t=0}^{T-1} \Delta_t.$$

Online updating not only has the advantage that it can be applied to non-episodic tasks, but it will in general also produce better value-function estimates under temporal-difference learning (see Sutton & Barto, 1998, Sections 7.1–3).

### 2.3. Linear TD( $\lambda$ )

Under linear function approximation, the TD(0) method, based on update target (2), performs at time step  $t + 1$  the following update:

$$\theta_{t+1} = \theta_t + \alpha \delta_t \phi_t,$$

where

$$\delta_t = R_{t+1} + \gamma \theta_t^\top \phi_{t+1} - \theta_t^\top \phi_t \quad (3)$$

is called the *TD error*. This update only affects the weights  $\theta_i$  for which  $\phi_i$  is non-zero at time step  $t$ . The idea behind TD( $\lambda$ ) is to update weights for which  $\phi_i$  was non-zero in the (near) past as well. This is implemented by means of an eligibility vector  $e$ , which reflects how much each feature is ‘eligible’ for the current TD error. TD( $\lambda$ ) performs an update of each  $\theta_i$ , with the current TD error, proportional to its trace value  $e_i$ :<sup>1</sup>

$$\theta_{t+1} = \theta_t + \delta_t e_t. \quad (4)$$

The original linear TD( $\lambda$ ) (Sutton, 1988) used *accumulating traces*, in which  $e_t$  is initialized to the zero vector,  $\mathbf{0}$ , and then updated according to:

$$e_t = \gamma \lambda e_{t-1} + \alpha \phi_t, \quad (5)$$

where  $\lambda \in [0, 1]$ . Note that for  $\lambda = 0$ , the TD( $\lambda$ ) update reduces to the TD(0) update. Algorithm 1 shows pseudocode for linear TD( $\lambda$ ) with accumulating traces.

<sup>1</sup>We fold the step-size  $\alpha$  into the eligibility vector.

---

#### Algorithm 1 linear TD( $\lambda$ ) with accumulating traces

---

```

initialize  $\theta$  arbitrarily
loop {over episodes}
  initialize  $e = \mathbf{0}$ 
  initialize  $S$ 
  repeat {for each step in the episode}
    generate reward  $R$  and next state  $S'$  for  $S$ 
     $\delta \leftarrow R + \gamma \theta^\top \phi(S') - \theta^\top \phi(S)$ 
     $e \leftarrow \gamma \lambda e + \alpha \phi(S)$ 
     $\theta \leftarrow \theta + \delta e$ 
     $S \leftarrow S'$ 
  until  $S$  is terminal
end loop

```

---

Singh and Sutton (1996) introduced an interesting variation on the traces of TD( $\lambda$ ) called *replacing traces*. These traces do not increase with repeated visits to a state (or state feature) but are rather reset to a constant value (generally 1) each time the state (or state feature) is visited. The trace due to the re-visit does not add to the existing trace, it *replaces* it. Replacing traces result in faster and more reliable learning on many problems (Sutton, 1996). Replacing traces are usually defined only for discrete states or linear function approximation with binary features (that are either 1 or 0, present or not present). We generalize the update slightly to allow the features to take on any values, and include  $\alpha$  in the trace, as follows:

$$e_{i,t} = \begin{cases} \gamma \lambda e_{i,t-1} & \text{if } \phi_{i,t} = 0 \\ \alpha \phi_{i,t} & \text{if } \phi_{i,t} \neq 0 \end{cases} \quad \text{for all } i.$$

This update, with (4), is exactly the same as TD( $\lambda$ ) with replacing traces as it has been explored in the literature, and also allows us to use the same idea with features that take on non-binary values, as in the experiments later in this paper.

### 2.4. The Conventional Forward View

The conventional forward view relates TD( $\lambda$ ) (with accumulating traces) to the  $\lambda$ -return algorithm. This algorithm performs at each time step a standard update (1) with the  $\lambda$ -return as update target. The  $\lambda$ -return  $G_t^\lambda$  is an estimate of the expected return based on a combinations of rewards and other estimates:

$$G_t^\lambda = (1 - \lambda) \sum_{n=1}^{\infty} \lambda^{n-1} G_{t,\theta}^{(n)},$$

with  $\theta$  the weight vector corresponding to the current episode, and  $G_{t,\theta}^{(n)}$  the  $n$ -step return, defined by:

$$G_{t,\theta}^{(n)} = \sum_{i=1}^n \gamma^{i-1} R_{t+i} + \gamma^n \theta^\top \phi_{t+n}. \quad (6)$$

For an episodic task, the  $\lambda$ -return is defined as:

$$G_t^\lambda = (1 - \lambda) \sum_{n=1}^{T-t-1} \lambda^{n-1} G_{t,\theta}^{(n)} + \lambda^{T-t-1} G_{t,\theta}^{(T-t)}, \quad (7)$$

where  $T$  is the time step that the terminal state is reached. Because the value of a terminal state is always 0, the last  $n$ -step return in this equation is equal to the full return:

$$G_{t,\theta}^{(T-t)} = \sum_{i=1}^{T-t} \gamma^{i-1} R_{t+i} = G_t.$$

Note that for  $\lambda = 0$ ,  $G_t^\lambda$  is equal to the TD(0) update target (2). On the other hand, for  $\lambda = 1$ ,  $G_t^\lambda$  is equal to the full return  $G_t$ .

It has been shown that off-line TD( $\lambda$ ) is equal to the (off-line)  $\lambda$ -return algorithm (Sutton & Barto, 1998; Sutton, 1988). However, online TD( $\lambda$ ) is only approximately equal to the online  $\lambda$ -return algorithm. In fact, no version of online TD( $\lambda$ ) can be constructed that matches the online  $\lambda$ -return algorithm exactly, because the  $\lambda$ -return uses rewards and states beyond the current time step. In Section 3, we present a new forward view that can be implemented online.

## 2.5. Other Variations on TD( $\lambda$ )

Several variations on TD( $\lambda$ ) other than those treated in this paper have been suggested in the literature. Shapire and Warmuth (1996) introduced a variation of TD( $\lambda$ ) for which upper and lower bounds on the accuracy of the value estimate of the current state can be derived and proven. Maei, Sutton, and others (Maei, 2011; Maei & Sutton, 2010; Sutton et al, in preparation; Sutton et al., 2009) have explored generalizations of TD( $\lambda$ )-like algorithms to *off-policy learning*, in which the behavior policy (generating the data) and the target policy (whose value function is being learned) are allowed to be different. Although in this paper we consider only the on-policy case, where these two policies must be the same, extending the ideas of true online TD( $\lambda$ ) to the off-policy case would be an interesting and natural avenue of future work.

## 3. An Online Forward View

The conventional forward view assumes that the value estimates are constant during an episode. Consequently, the online TD( $\lambda$ ) algorithm matches the forward view only approximately. In order to construct a version of online TD( $\lambda$ ) that matches the forward view exactly, we have to think more carefully about what we mean by the term ‘online’. In Section 2.2 we specified that online learning means that the value estimates are updated at every time step. This definition does not rule out the possibility that

the update at time step  $t$ , uses information beyond time step  $t$ . One could for example record a sequence of state transitions and rewards, and at the end of an episode use the full sequence to construct update targets for all visited states. When the  $\lambda$ -return is used as the basis for an online forward view, this broad interpretation of the term online is required, because the  $\lambda$ -return is an update target that is based on the full experience sequence.

On the other hand, if we look at online TD( $\lambda$ ) the term online seems to mean something more restrictive. Specifically, online TD( $\lambda$ ) uses for the update at time step  $t$  only information up to time step  $t$ . This restriction is what makes online TD( $\lambda$ ) very generally applicable. For example, it can be applied to non-episodic tasks, and can be easily generalized to control tasks. We will call the more restrictive definition of online, as applies to online TD( $\lambda$ ), *strict-online*. The goal of this section is to define a strict-online forward view.

To construct a strict-online forward view, we start by defining a version of the  $\lambda$ -return that is truncated at a specific time step. Let  $t'$  be the time step the  $\lambda$ -return is truncated. The *truncated  $\lambda$ -return*  $G_t^{\lambda|t'}$  is defined as:

$$G_t^{\lambda|t'} = (1 - \lambda) \sum_{n=1}^{t'-t-1} \lambda^{n-1} G_{t,\theta_{t+n-1}}^{(n)} + \lambda^{t'-t-1} G_{t,\theta_{t'-1}}^{(t'-t)}. \quad (8)$$

Note that this definition is closely related to the definition of regular  $\lambda$ -return for episodic tasks (7). Basically,  $T$  is replaced by  $t'$ , and the weight vectors used in the  $n$ -step returns have now a specific time index.

The idea behind the new strict-online forward view is in essence a simple one: at each time step, the truncated  $\lambda$ -returns from all previous time steps are updated, such that they are now truncated at the current time step; the weight vector of the current time step is determined by sequentially performing TD backups, using the updated  $\lambda$ -returns, starting from the initial weight vector,  $\theta_{init}$ . The weight vectors for the first three time steps are shown below. We use a second index for  $\theta$  to indicate at which time step the

$\lambda$ -returns are truncated.<sup>2</sup>

$$\boldsymbol{\theta}_{0,0} : \boldsymbol{\theta}_{0,0} = \boldsymbol{\theta}_{init},$$

$$\begin{aligned} \boldsymbol{\theta}_{1,1} : \boldsymbol{\theta}_{1,0} &= \boldsymbol{\theta}_{init} \\ \boldsymbol{\theta}_{1,1} &= \boldsymbol{\theta}_{1,0} + \alpha_0 \left[ G_0^{\lambda|1} - \boldsymbol{\theta}_{1,0}^\top \boldsymbol{\phi}_0 \right] \boldsymbol{\phi}_0, \end{aligned}$$

$$\begin{aligned} \boldsymbol{\theta}_{2,2} : \boldsymbol{\theta}_{2,0} &= \boldsymbol{\theta}_{init} \\ \boldsymbol{\theta}_{2,1} &= \boldsymbol{\theta}_{2,0} + \alpha_0 \left[ G_0^{\lambda|2} - \boldsymbol{\theta}_{2,0}^\top \boldsymbol{\phi}_0 \right] \boldsymbol{\phi}_0, \\ \boldsymbol{\theta}_{2,2} &= \boldsymbol{\theta}_{2,1} + \alpha_1 \left[ G_1^{\lambda|2} - \boldsymbol{\theta}_{2,1}^\top \boldsymbol{\phi}_1 \right] \boldsymbol{\phi}_1, \end{aligned}$$

$$\begin{aligned} \boldsymbol{\theta}_{3,3} : \boldsymbol{\theta}_{3,0} &= \boldsymbol{\theta}_{init} \\ \boldsymbol{\theta}_{3,1} &= \boldsymbol{\theta}_{3,0} + \alpha_0 \left[ G_0^{\lambda|3} - \boldsymbol{\theta}_{3,0}^\top \boldsymbol{\phi}_0 \right] \boldsymbol{\phi}_0, \\ \boldsymbol{\theta}_{3,2} &= \boldsymbol{\theta}_{3,1} + \alpha_1 \left[ G_1^{\lambda|3} - \boldsymbol{\theta}_{3,1}^\top \boldsymbol{\phi}_1 \right] \boldsymbol{\phi}_1, \\ \boldsymbol{\theta}_{3,3} &= \boldsymbol{\theta}_{3,2} + \alpha_2 \left[ G_2^{\lambda|3} - \boldsymbol{\theta}_{3,2}^\top \boldsymbol{\phi}_2 \right] \boldsymbol{\phi}_2, \end{aligned}$$

More generally, the weight vector  $\boldsymbol{\theta}_{t,k}$  is incrementally defined by the equation:

$$\boldsymbol{\theta}_{t,k} = \boldsymbol{\theta}_{t,k-1} + \alpha_{k-1} \left[ G_{k-1}^{\lambda|t} - \boldsymbol{\theta}_{t,k-1}^\top \boldsymbol{\phi}_{k-1} \right] \boldsymbol{\phi}_{k-1} \quad (9)$$

for  $0 < k \leq t$ , with  $\boldsymbol{\theta}_{t,0} = \boldsymbol{\theta}_{init}$ .

Note that  $\boldsymbol{\theta}_{i,j}$  for  $j < i$  are temporary ‘helper’ vectors. In the context of the true online forward view,  $\boldsymbol{\theta}_t$  refers to  $\boldsymbol{\theta}_{t,t}$ . For example, the value estimate of state  $s$  at time step  $t$  is:

$$\hat{v}_t(s) = \boldsymbol{\theta}_t^\top \boldsymbol{\phi}(s) = \boldsymbol{\theta}_{t,t}^\top \boldsymbol{\phi}(s).$$

We call the algorithm that computes  $\boldsymbol{\theta}_{t,t}$  at each time step  $t$  the *truncated  $\lambda$ -return algorithm*. Note that for  $\lambda = 0$  the following holds:

$$G_t^{0|t'} = G_{t,\boldsymbol{\theta}_t}^{(1)} = R_{t+1} + \gamma \boldsymbol{\theta}_t^\top \boldsymbol{\phi}_{t+1}.$$

Hence, for  $\lambda = 0$  the truncated  $\lambda$ -return algorithm is equivalent to TD(0). For  $\lambda = 1$ , at the end of an episode ( $t' = T$ ) the following holds:

$$G_t^{1|T} = G_{t,\boldsymbol{\theta}_{T-1}}^{(T-t)} = G_t.$$

Hence, for  $\lambda = 1$  the values at the end of an episode computed by the truncated  $\lambda$ -return algorithm are equal to those computed by (an online, step-size version of) Monte-Carlo.

The truncated  $\lambda$ -return algorithm is an expensive method, and requires storage of all observed states and rewards. In the next section, we introduce a new online TD( $\lambda$ ) version that implements the strict-online forward view efficiently using eligibility traces.

<sup>2</sup>For ease of exposition, we focus on the linear case. However, our approach can be extended to the non-linear case in a straightforward way.

## 4. True Online TD( $\lambda$ )

This section presents the online TD( $\lambda$ ) method that matches the new online forward view exactly. First the algorithm itself is presented, then its equivalence to the new online forward view is proven. The section finishes with empirical results demonstrating the benefits of the algorithm.

### 4.1. The Algorithm

True online TD( $\lambda$ ) forms the backward view of the truncated  $\lambda$ -return algorithm. Like traditional TD( $\lambda$ ), it updates feature weights proportional to a decaying eligibility trace.

The three update equations that specify the algorithm are as follows:

$$\delta_t = R_{t+1} + \gamma \boldsymbol{\theta}_t^\top \boldsymbol{\phi}_{t+1} - \boldsymbol{\theta}_{t-1}^\top \boldsymbol{\phi}_t \quad (10)$$

$$\mathbf{e}_t = \gamma \lambda \mathbf{e}_{t-1} + \alpha_t \boldsymbol{\phi}_t - \alpha_t \gamma \lambda [\mathbf{e}_{t-1}^\top \boldsymbol{\phi}_t] \boldsymbol{\phi}_t \quad (11)$$

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t + \delta_t \mathbf{e}_t + \alpha_t [\boldsymbol{\theta}_{t-1}^\top \boldsymbol{\phi}_t - \boldsymbol{\theta}_t^\top \boldsymbol{\phi}_t] \boldsymbol{\phi}_t \quad (12)$$

Comparing these three equations with equations (3), (4) and (5) — the update equations for traditional TD( $\lambda$ ) — reveals that the true online TD( $\lambda$ ) updates are different in the following ways:

- $\delta_t$ :  $\boldsymbol{\theta}_{t-1}^\top \boldsymbol{\phi}_t$  is used instead of  $\boldsymbol{\theta}_t^\top \boldsymbol{\phi}_t$ .
- $\mathbf{e}_t$ : the term  $-\alpha_t \gamma \lambda [\mathbf{e}_{t-1}^\top \boldsymbol{\phi}_t] \boldsymbol{\phi}_t$  is added.
- $\boldsymbol{\theta}_{t+1}$ : the term  $\alpha_t [\boldsymbol{\theta}_{t-1}^\top \boldsymbol{\phi}_t - \boldsymbol{\theta}_t^\top \boldsymbol{\phi}_t]$  is added.

Algorithm 2 shows pseudo-code that implements the true online TD( $\lambda$ ) equations. For  $\lambda = 0$ , the algorithm is equivalent to (regular) TD(0).

---

#### Algorithm 2 true online TD( $\lambda$ )

---

```

initialize  $\boldsymbol{\theta}$  arbitrarily
loop {over episodes}
  initialize  $\mathbf{e} = \mathbf{0}$ 
  initialize  $S$ 
   $\hat{v}_S \leftarrow \boldsymbol{\theta}^\top \boldsymbol{\phi}(S)$ 
  repeat {for each step in the episode}
    generate reward  $R$  and next state  $S'$  for  $S$ 
     $\hat{v}_{S'} \leftarrow \boldsymbol{\theta}^\top \boldsymbol{\phi}(S')$ 
     $\delta \leftarrow R + \gamma \hat{v}_{S'} - \hat{v}_S$ 
     $\mathbf{e} \leftarrow \gamma \lambda \mathbf{e} + \alpha [1 - \gamma \lambda \mathbf{e}^\top \boldsymbol{\phi}(S)] \boldsymbol{\phi}(S)$ 
     $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \delta \mathbf{e} + \alpha [\hat{v}_S - \boldsymbol{\theta}^\top \boldsymbol{\phi}(S)] \boldsymbol{\phi}(S)$ 
     $\hat{v}_S \leftarrow \hat{v}_{S'}$ 
     $S \leftarrow S'$ 
  until  $S$  is terminal
end loop
    
```

---

## 4.2. Equivalence to the Online Forward View

In this section, we prove that the values computed by true online TD( $\lambda$ ) are exactly the same as those computed by the truncated  $\lambda$ -return algorithm. We start by deriving two lemmas.

**Lemma 1**  $G_t^{\lambda|t'}$ , used by the truncated  $\lambda$ -return algorithm, is related to  $\delta_t$ , used by true online TD( $\lambda$ ) ((10)), as follows:

$$G_t^{\lambda|t'+1} - G_t^{\lambda|t'} = (\gamma\lambda)^{t'-t} \delta_{t'}.$$

**Proof** The truncated  $\lambda$ -return algorithm uses  $G_t^{\lambda|t'}$  ((8)) with  $\theta_{\tau(n)} = \theta_{t+n-1}$ , resulting in:

$$G_t^{\lambda|t'} = (1-\lambda) \sum_{n=1}^{t'-t-1} \lambda^{n-1} G_{t,\theta_{t+n-1}}^{(n)} + \lambda^{t'-t-1} G_{t,\theta_{t'-1}}^{(t'-t)}.$$

From this it follows that  $G_t^{\lambda|t'+1}$  is:

$$G_t^{\lambda|t'+1} = (1-\lambda) \sum_{n=1}^{t'-t} \lambda^{n-1} G_{t,\theta_{t+n-1}}^{(n)} + \lambda^{t'-t} G_{t,\theta_{t'}}^{(t'+1-t)}.$$

Subtracting  $G_t^{\lambda|t'}$  from  $G_t^{\lambda|t'+1}$  results in

$$G_t^{\lambda|t'+1} - G_t^{\lambda|t'} = \lambda^{t'-t} [G_{t,\theta_{t'}}^{(t'+1-t)} - G_{t,\theta_{t'-1}}^{(t'-t)}]. \quad (13)$$

Similarly, using (6), it can be shown that subtracting  $G_{t,\theta_{t'-1}}^{(t'-t)}$  from  $G_{t,\theta_{t'}}^{(t'+1-t)}$  yields:

$$\begin{aligned} G_{t,\theta_{t'}}^{(t'+1-t)} - G_{t,\theta_{t'-1}}^{(t'-t)} &= \gamma^{t'-t} [R_{t'+1} + \gamma \theta_{t'}^\top \phi_{t'+1} - \\ &\quad \theta_{t'-1}^\top \phi_{t'}] \\ &= \gamma^{t'-t} \delta_{t'}. \end{aligned}$$

Substituting this result in (13) yields the lemma.  $\blacksquare$

Using Lemma 1, the following lemma can be derived.

**Lemma 2** The following equation holds:

$$\theta_{t+1,t} - \theta_{t,t} = \gamma\lambda\delta_t e_{t-1},$$

with  $\delta_t$  and  $e_{t-1}$  defined by (10) and (11), respectively.

**Proof** The proof structure is a proof by induction. First we prove that if for some  $k$ ,  $1 \leq k < t$ , the following equation holds:

$$\theta_{t+1,k} - \theta_{t,k} = (\gamma\lambda)^{t+1-k} \delta_t e_{k-1}, \quad (14)$$

then it also holds for  $k+1$ . We then show that it holds for  $k=1$ . Hence, by induction, it also holds for  $k=t$ , proving the lemma.

Subtracting  $\theta_{t,k+1}$  from  $\theta_{t+1,k+1}$  using (9) yields:

$$\begin{aligned} \theta_{t+1,k+1} - \theta_{t,k+1} &= \theta_{t+1,k} - \theta_{t,k} \\ &\quad + \alpha_k [G_k^{\lambda|t+1} - G_k^{\lambda|t} \\ &\quad - (\theta_{t+1,k} - \theta_{t,k})^\top \phi_k] \phi_k \end{aligned}$$

If (14) holds for  $k$ , this can be rewritten as (using Lemma 1):

$$\begin{aligned} \theta_{t+1,k+1} - \theta_{t,k+1} &= (\gamma\lambda)^{t+1-k} \delta_t e_{k-1} \\ &\quad + \alpha_k [(\gamma\lambda)^{t-k} \delta_t \\ &\quad - (\gamma\lambda)^{t+1-k} \delta_t (e_{k-1})^\top \phi_k] \phi_k \\ &= (\gamma\lambda)^{t-k} \delta_t [\gamma\lambda e_{k-1} + \alpha_k \phi_k \\ &\quad - \alpha_k \gamma\lambda (e_{k-1}^\top \phi_k) \phi_k] \\ &= (\gamma\lambda)^{t-k} \delta_t \theta_k \end{aligned}$$

This proves that if (14) holds for  $k$ , it holds for  $k+1$ .

Computing  $\theta_{t+1,1} - \theta_{t,1}$  using (9) yields:

$$\begin{aligned} \theta_{t+1,1} - \theta_{t,1} &= \theta_{t+1,0} - \theta_{t,0} \\ &\quad + \alpha_k [G_0^{\lambda|t+1} - G_0^{\lambda|t} \\ &\quad - (\theta_{t+1,0} - \theta_{t,0})^\top \phi_0] \phi_0 \end{aligned}$$

Using  $\theta_{t+1,0} = \theta_{t,0} = \theta_{init}$  and Lemma 1 this reduces to:

$$\theta_{t+1,1} - \theta_{t,1} = \alpha_k (\gamma\lambda)^t \delta_t \phi_0$$

This demonstrates that (14) holds for  $k=1$ .<sup>3</sup> By induction, it follows that (14) also holds for  $k=t$ , proving the lemma.  $\blacksquare$

We are now ready to prove our main theorem.

**Theorem 1**  $\theta_{t,t}$ , computed by the truncated  $\lambda$ -return algorithm, is the same as  $\theta_t$ , computed by true online TD( $\lambda$ ), for all time steps  $t$ .

**Proof** Using (9),  $\theta_{t+1,t+1}$  can be written as:

$$\theta_{t+1,t+1} = \theta_{t+1,t} + \alpha_t [G_t^{\lambda|t+1} - \theta_{t+1,t}^\top \phi_t] \phi_t. \quad (15)$$

Rewriting Lemma 2 as:

$$\theta_{t+1,t} = \theta_{t,t} + \gamma\lambda\delta_t e_{t-1},$$

<sup>3</sup> Technically, initialization of  $e_t$  occurs for  $t=-1$ , that is,  $e_{-1} = \mathbf{0}$ . Using (11) it follows that  $e_0 = \alpha_0 \phi_0$ .

and substituting this in (15) yields:

$$\begin{aligned}
 \theta_{t+1,t+1} &= \theta_{t,t} + \gamma\lambda\delta_t e_{t-1} + \alpha_t [R_{t+1} + \theta_{t,t}^\top \phi_{t+1}] \phi_t \\
 &\quad - \alpha_t [(\theta_{t,t} + \gamma\lambda\delta_t e_{t-1})^\top \phi_t] \phi_t \\
 &= \theta_{t,t} + \gamma\lambda\delta_t e_{t-1} + \alpha_t [\delta_t + \theta_{t-1,t-1}^\top \phi_t] \phi_t \\
 &\quad - \alpha_t [\theta_{t,t}^\top \phi_t + \gamma\lambda\delta_t e_{t-1}^\top \phi_t] \phi_t \\
 &= \theta_{t,t} + \gamma\lambda\delta_t e_{t-1} + \alpha_t \delta_t \phi_t \\
 &\quad - \alpha_t \gamma\lambda\delta_t [e_{t-1}^\top \phi_t] \phi_t \\
 &\quad + \alpha_t [\theta_{t-1,t-1}^\top \phi_t - \theta_{t,t}^\top \phi_t] \phi_t \\
 &= \theta_{t,t} + \delta_t e_t + \alpha_t [\theta_{t-1,t-1}^\top \phi_t - \theta_{t,t}^\top \phi_t] \phi_t
 \end{aligned}$$

This final equation shows that the weight update from one time step to the other for the truncated  $\lambda$ -return is the same as the weight update for true online TD( $\lambda$ ) ((12)). ■

### 4.3. Empirical Results

We compare the performance of true online TD( $\lambda$ ) with traditional TD( $\lambda$ ) with accumulating and replacing traces on a random-walk task. The random-walk task is shown in Figure 1 for  $N = 6$  ( $N$  being the total number of states, including the terminal state). In our experiment we use  $N = 11$ . The transition probability in the direction of the terminal state,  $p$ , is set to 0.9. Initial  $\theta$  is  $\mathbf{0}$  and  $\gamma = 0.99$ .

We used linear function approximation with two types of features, resulting in two different tasks. The feature values of these two tasks are shown in Table 1 (for  $N = 6$ ). In task 1, there are (at most) 3 non-zero features for each state. In task 2, the number of non-zero features is between 1 and  $N - 1$ . For the terminal state all features have value 0. The  $L^2$ -norm of  $\phi(s)$  is 1 for all non-terminal states.

For the considered features of both tasks, the true state values can be represented exactly with linear function approximation. We specify the performance of a method as the root-mean-squared (RMS) error of the value estimates of all non-terminal states at the end of an episode with respect to their true values (which are analytically determined), averaged over the first 10 episodes and 100 independent runs.

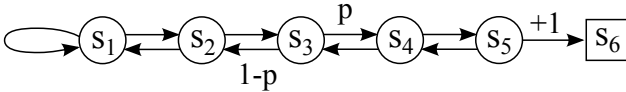


Figure 1. Random-walk for  $N = 6$  states (experiments used  $N = 11$ ). Transition probability to the right is  $p$ ; transition probability to the left is  $1 - p$ . All rewards are 0, except transition to the terminal state ( $s_6$ ), which results in a reward of +1. The initial state is  $s_1$ .

Figure 2 shows the performance on both tasks for  $\alpha$  from 0 to 1.5 with steps of 0.01 and  $\lambda$  from 0 to 0.9 with steps of 0.1 and from 0.9 to 1.0 with steps of 0.025. Figure 3

Table 1. Feature values for random-walk task 1 and task 2 for  $N = 6$ .

		$s_1$	$s_2$	$s_3$	$s_4$	$s_5$	$s_6$
Task 1	$\phi_1$	1	$1/\sqrt{2}$	$1/\sqrt{3}$	0	0	0
	$\phi_2$	0	$1/\sqrt{2}$	$1/\sqrt{3}$	$1/\sqrt{3}$	0	0
	$\phi_3$	0	0	$1/\sqrt{3}$	$1/\sqrt{3}$	$1/\sqrt{3}$	0
	$\phi_4$	0	0	0	$1/\sqrt{3}$	$1/\sqrt{3}$	0
	$\phi_5$	0	0	0	0	$1/\sqrt{3}$	0
Task 2	$\phi_1$	1	$1/\sqrt{2}$	$1/\sqrt{3}$	$1/\sqrt{4}$	$1/\sqrt{5}$	0
	$\phi_2$	0	$1/\sqrt{2}$	$1/\sqrt{3}$	$1/\sqrt{4}$	$1/\sqrt{5}$	0
	$\phi_3$	0	0	$1/\sqrt{3}$	$1/\sqrt{4}$	$1/\sqrt{5}$	0
	$\phi_4$	0	0	0	$1/\sqrt{4}$	$1/\sqrt{5}$	0
	$\phi_5$	0	0	0	0	$1/\sqrt{5}$	0

shows the performance for the different  $\lambda$  values at the best  $\alpha$  value.

While the return has an upper bound of 1, accumulating traces get errors above 1 on both tasks, indicating divergence of values. While replacing traces does not result in divergence of values for the considered step-sizes, it has no effect in the second task. True online TD( $\lambda$ ) is the only method that achieves a performance benefit (with respect to TD(0)) on both tasks.

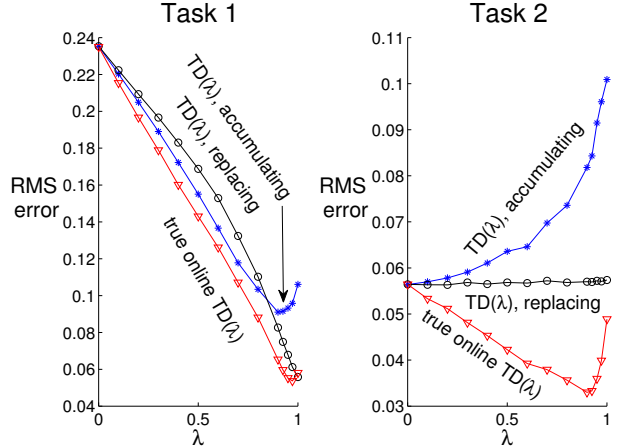


Figure 3. RMS error of state values at the end of each episode, averaged over the first 10 episodes, as well as 100 independent runs, for different values of  $\lambda$  at the best value of  $\alpha$ .

## 5. Control

The true online TD( $\lambda$ ) algorithm (Algorithm 2) can be easily modified for control. Simply using a feature vector consisting of state-action features (i.e., using  $\phi(s, a)$  instead of  $\phi(s)$ ) changes the algorithm to a true online Sarsa( $\lambda$ ) algorithm.

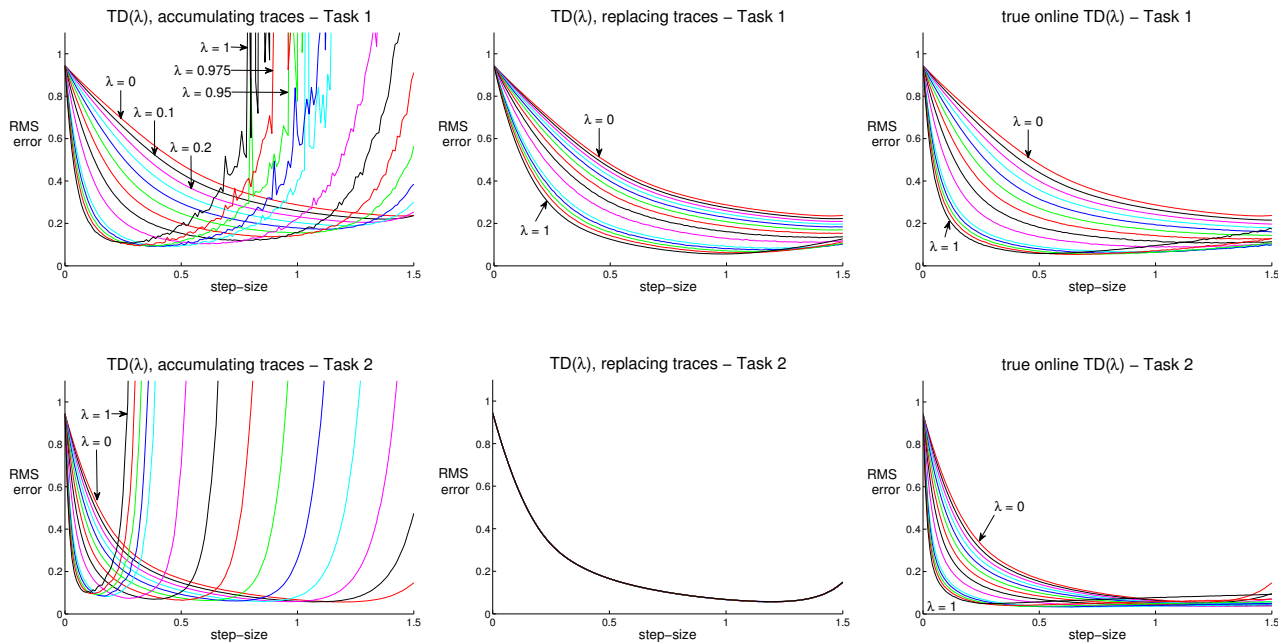


Figure 2. RMS error of state values at the end of each episode, averaged over the first 10 episodes, as well as 100 independent runs, for different values of  $\alpha$  and  $\lambda$ .

Figure 4 compares true online Sarsa( $\lambda$ ) with the traditional Sarsa( $\lambda$ ) implementation on the standard mountain car task (Sutton & Barto, 1998) using 10 tilings of each  $10 \times 10$  tiles. Results are plotted for  $\lambda = 0.9$  and  $\alpha = \alpha_0/10$ , for  $\alpha_0$  from 0.2 to 2.0 with steps of 0.2. Clearing/no clearing refers to whether the trace values of non-selected actions are set to 0 (clearing) or not (no clearing), in case of replacing traces. The results suggest that the true online principle is also effective in a control setting.

## 6. Conclusion

We presented for the first time an online version of the forward view which forms the theoretical and intuitive foundation for the TD( $\lambda$ ) algorithm. In addition, we have presented a new variant of TD( $\lambda$ ), with the same computational complexity as the classical algorithm, which we call true online TD( $\lambda$ ). We proved that true online TD( $\lambda$ ) matches the new online forward view exactly, in contrast to classical online TD( $\lambda$ ), which only approximates its forward view. In addition, we demonstrated empirically that true online TD( $\lambda$ ) outperforms conventional TD( $\lambda$ ) on three benchmark problems. It seems, by adhering more truly to the original goal of TD( $\lambda$ )—matching an intuitively clear forward view even in the online case—that we have found a new algorithm that simply improves on TD( $\lambda$ ).

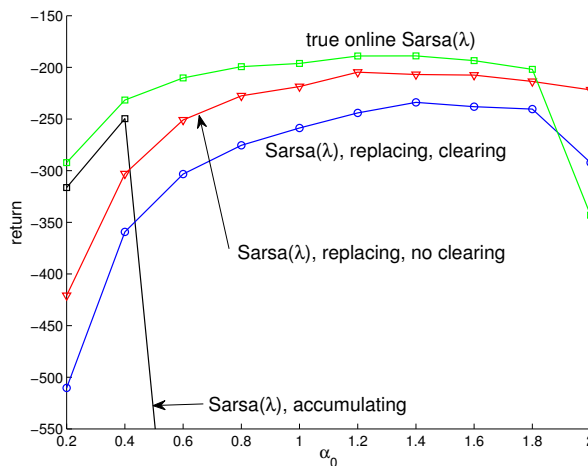


Figure 4. Average return over first 20 episodes on mountain car task for  $\lambda = 0.9$  and different  $\alpha_0$ . Results are averaged over 100 independent runs.

## Acknowledgements

The authors thank Hado van Hasselt and Rupam Mahmood for extensive discussions leading to the refinement of these ideas. This work was supported by grants from Alberta Innovates – Technology Futures and the National Science and Engineering Research Council of Canada.



## References

- Barnard, E. (1993). Temporal-difference methods and Markov models. *IEEE Transactions on Systems, Man, and Cybernetics* 23(2):357–365.
- Kaelbling, L. P., Littman, M. L., and Moore, A. P. (1996). Reinforcement Learning: A Survey. *Journal of Artificial Intelligence Research* 4:237–285.
- Maei, H., Sutton, R. S. (2010). GQ( $\lambda$ ): A general gradient algorithm for temporal-difference prediction learning with eligibility traces. In: *Proceedings of the Third Conference on Artificial General Intelligence*, pp. 91–96.
- Schapire, R. E., and Warmuth, M. K. (1996). On the worst-case analysis of temporal-difference learning algorithms. *Machine Learning* 22(1/2/3):95–121.
- Singh, S. P., and Sutton, R. S. (1996). Reinforcement learning with replacing eligibility traces. *Machine Learning* 22(1/2/3):123–158.
- Sutton, R. S. (1988). Learning to predict by the methods of temporal differences. *Machine Learning* 3(1):9–44.
- Sutton, R. S., and Barto, A. G. (1998). *Reinforcement Learning: An Introduction*. MIT Press, Cambridge, Massachusetts.
- Sutton, R. S., Maei, H. R., Precup, D., Bhatnagar, S., Silver, D., Szepesvári, Cs. and Wiewiora, E. (2009). Fast gradient-descent methods for temporal-difference learning with linear function approximation. In *Proceedings of the 26th International Conference on Machine Learning*, pp. 993–1000.
- Sutton, R. S., Mahmood, A. R., Precup, D. (in preparation). A new Q( $\lambda$ ).
- Szepesvári, Cs. (2010). *Algorithms for Reinforcement Learning*. Morgan and Claypool Press.