



# Technical Update: Least-Squares Temporal Difference Learning

JUSTIN A. BOYAN

jab@itasoftware.com (<http://www.boyan.com/justin/>)

ITA Software, 141 Portland Street, Cambridge, MA 02139, USA

**Editor:** Satinder Singh

**Abstract.**  $TD(\lambda)$  is a popular family of algorithms for approximate policy evaluation in large MDPs.  $TD(\lambda)$  works by incrementally updating the value function after each observed transition. It has two major drawbacks: it may make inefficient use of data, and it requires the user to manually tune a stepsize schedule for good performance. For the case of linear value function approximations and  $\lambda = 0$ , the Least-Squares TD (LSTD) algorithm of Bradtke and Barto (1996, *Machine learning*, 22:1–3, 33–57) eliminates all stepsize parameters and improves data efficiency.

This paper updates Bradtke and Barto’s work in three significant ways. First, it presents a simpler derivation of the LSTD algorithm. Second, it generalizes from  $\lambda = 0$  to arbitrary values of  $\lambda$ ; at the extreme of  $\lambda = 1$ , the resulting new algorithm is shown to be a practical, incremental formulation of supervised linear regression. Third, it presents a novel and intuitive interpretation of LSTD as a *model-based* reinforcement learning technique.

**Keywords:** reinforcement learning, temporal difference learning, value function approximation, linear least-squares methods

## 1. Background

We address the problem of approximating the value function  $V^\pi$  of a fixed policy  $\pi$  in a large Markov decision process (Bertsekas & Tsitsiklis, 1996; Sutton & Barto, 1998). This is an important subproblem of several algorithms for sequential decision making, including optimistic policy iteration (Bertsekas & Tsitsiklis, 1996) and STAGE (Boyan & Moore, 1998).  $V^\pi(x)$  simply predicts the expected long-term sum of future rewards obtained when the process starts in state  $x$  and follows policy  $\pi$  until termination. For simplicity we will assume that  $\pi$  is proper (guaranteed to terminate), which implies that  $V^\pi$  is well-defined without the use of a discount factor; however, the algorithms presented here extend straightforwardly to the discounted case.

For small Markov chains whose transition probabilities are all explicitly known, computing  $V^\pi$  is a trivial matter of solving a system of linear equations. However, in many practical applications, the transition probabilities of the chain are available only implicitly—either in the form of a simulation model or in the form of an agent’s actual experience executing  $\pi$  in its environment. In either case, we must compute  $V^\pi$  or an approximation thereof (denoted  $\tilde{V}^\pi$ ) solely from a collection of trajectories sampled from the chain. This is where the  $TD(\lambda)$  family of algorithms applies.

$TD(\lambda)$  was introduced in Sutton (1988); excellent summaries may now be found in several books (Bertsekas & Tsitsiklis, 1996; Sutton & Barto, 1998). For each state on each

observed trajectory,  $\text{TD}(\lambda)$  incrementally adjusts the coefficients of  $\tilde{V}^\pi$  toward new target values. The target values depend on the parameter  $\lambda \in [0, 1]$ . At  $\lambda = 1$ , the target at each visited state  $x_t$  is the *Monte-Carlo return*, i.e., the actual observed sum of future rewards  $R_t + R_{t+1} + \dots + R_{\text{END}}$ . This is an unbiased sample of  $V^\pi(x_t)$ , but may have significant variance since it depends on a long stochastic sequence of rewards. At the other extreme,  $\lambda = 0$ , the target value is set by a sampled one-step lookahead:  $R_t + \tilde{V}^\pi(x_{t+1})$ . This value has lower variance—the only random component is a single state transition—but is biased by the potential inaccuracy of the lookahead estimate of  $V^\pi$ . The parameter  $\lambda$  trades off between bias and variance. Empirically, several studies have found intermediate values of  $\lambda$  to perform best (Sutton, 1988; Sutton & Barto, 1998).

$\text{TD}(\lambda)$  has been shown to converge to a good approximation of  $V^\pi$  when *linear architectures* are used, assuming a suitable decreasing schedule of stepsizes for the incremental weight updates (Tsitsiklis & Van Roy, 1997). Linear architectures—which include lookup tables, state aggregation methods, CMACs, radial basis function networks with fixed bases, and multi-dimensional polynomial regression—approximate  $V^\pi(x)$  by first mapping the state  $x$  to a feature vector  $\phi(x) \in \mathfrak{R}^K$ , and then computing a linear combination of those features,  $\phi(x)^\top \beta$ . Figure 1 gives a convenient form of  $\text{TD}(\lambda)$  that exploits this representation. Table 1 summarizes the notation used throughout this paper.

On each transition, the algorithm computes the scalar one-step TD error  $R_t + (\phi(x_{t+1}) - \phi(x_t))^\top \beta$ , and apportions that error among all state features according to their respective *eligibilities*  $\mathbf{z}_t$ . The eligibility vector may be seen as an algebraic trick by which  $\text{TD}(\lambda)$  propagates rewards backward over the current trajectory without having to remember the trajectory explicitly. Each feature's eligibility at time  $t$  depends on the trajectory's history and on  $\lambda$ :  $\mathbf{z}_t = \sum_{i=t_0}^t \lambda^{t-i} \phi(x_i)$ , where  $t_0$  is the time at which the current trajectory started. In

---

**TD( $\lambda$ ) for approximate policy evaluation:**

- Given:*
- a *simulation model* for a proper policy  $\pi$  in MDP  $X$ ;
  - a *featurizer*  $\phi : X \rightarrow \mathfrak{R}^K$  mapping states to feature vectors,  $\phi(\text{END}) \stackrel{\text{def}}{=} \mathbf{0}$ ;
  - a parameter  $\lambda \in [0, 1]$ ; and
  - a sequence of *stepsizes*  $\alpha_1, \alpha_2, \dots$  for incremental coefficient updating.

*Output:* a coefficient vector  $\beta$  for which  $V^\pi(x) \approx \beta \cdot \phi(x)$ .

Set  $\beta := \mathbf{0}$  (or an arbitrary initial estimate),  $t := 0$ .

for  $n := 1, 2, \dots$  do: {

Set  $\delta := 0$ .

Choose a start state  $x_t \in X$ .

Set  $\mathbf{z}_t := \phi(x_t)$ .

while  $x_t \neq \text{END}$ , do: {

Simulate one step of the process, producing a reward  $R_t$  and next state  $x_{t+1}$ .

Set  $\delta := \delta + \mathbf{z}_t (R_t + (\phi(x_{t+1}) - \phi(x_t))^\top \beta)$ .

Set  $\mathbf{z}_{t+1} := \lambda \mathbf{z}_t + \phi(x_{t+1})$ .

Set  $t := t + 1$ .

}

Set  $\beta := \beta + \alpha_n \delta$ .

}

---

Figure 1. Ordinary  $\text{TD}(\lambda)$  for linearly approximating the undiscounted value function of a fixed proper policy.

Table 1. A glossary of notation used in this paper.

Symbol	Type	Meaning
$X$	set	State space for a Markov Decision Problem (MDP)
$R_t$	scalar	The one-step reward received at time $t$
$V^\pi$	$: X \rightarrow \mathfrak{R}$	The true value function of a given policy $\pi$ in the MDP
$\tilde{V}^\pi$	$: X \rightarrow \mathfrak{R}$	A linear approximation to $V^\pi$ given by $\tilde{V}^\pi(x) = \phi(x)^\top \beta$
$K$	$\in \mathbb{N}$	Dimensionality of the feature vector representing each state
$\phi$	$: X \rightarrow \mathfrak{R}^K$	The mapping from state to real-valued feature vector
$\beta$	$\in \mathfrak{R}^K$	The vector of coefficients that specifies $\tilde{V}^\pi$
$\lambda$	scalar $\in [0,1]$	Bias-variance parameter for temporal-difference learning
$\alpha_i$	scalar $\in (0,1)$	Stepsize for incremental coefficient updating in TD( $\lambda$ )
$\mathbf{z}_t$	$\in \mathfrak{R}^K$	“Eligibility vectors” used to trace a trajectory’s history
$\beta_\lambda$	$\in \mathfrak{R}^K$	The coefficients to which TD( $\lambda$ ) converges (dependent on $\lambda$ )
$\mathbf{d}$	$\in \mathfrak{R}^K$	Vector used in convergence analysis of TD( $\lambda$ ) (see Eq. (1))
$\mathbf{C}$	$K \times K$ matrix	Matrix used in convergence analysis of TD( $\lambda$ ) (see Eq. (1))
$\mathbf{b}$	$\in \mathfrak{R}^K$	Vector produced by LSTD( $\lambda$ ) for direct solution of $\beta$ (see Eq. (2))
$\mathbf{A}$	$K \times K$ matrix	Matrix produced by LSTD( $\lambda$ ) for direct solution of $\beta$ (see Eq. (2))

the case of TD(0), only the current state’s features are eligible to be updated, so  $\mathbf{z}_t = \phi(x_t)$ ; whereas in TD(1), the features of all states seen so far on the current trajectory are eligible, so  $\mathbf{z}_t = \sum_{i=t_0}^t \phi(x_i)$ .

To what weights does TD( $\lambda$ ) converge? Examining the update rule for  $\delta$  in figure 1, it is not difficult to see that the coefficient changes made by TD( $\lambda$ ) after an observed trajectory  $(x_0, x_1, \dots, x_L, \text{END})$  have the form  $\beta := \beta + \alpha_n(\mathbf{d} + \mathbf{C}\beta + \omega)$ , where

$$\mathbf{d} = \mathbb{E} \left\{ \sum_{i=0}^L \mathbf{z}_i R_i \right\}; \quad \mathbf{C} = \mathbb{E} \left\{ \sum_{i=0}^L \mathbf{z}_i (\phi(x_{i+1}) - \phi(x_i))^\top \right\}; \quad (1)$$

and  $\omega =$  zero-mean noise. The expectations are taken with respect to the distribution of trajectories through the Markov chain. It is shown in Bertsekas and Tsitsiklis (1996, §6.3.4) that  $\mathbf{C}$  is negative definite and that the noise  $\omega$  has sufficiently small variance, which together with the stepsize conditions mentioned above, imply that  $\beta$  converges to a fixed point  $\beta_\lambda$  satisfying  $\mathbf{d} + \mathbf{C}\beta_\lambda = \mathbf{0}$ . In effect, TD( $\lambda$ ) solves this system of equations by performing stochastic gradient descent on a potential function  $\|\beta - \beta_\lambda\|^2$ . It never explicitly represents  $\mathbf{d}$  or  $\mathbf{C}$ . The changes to  $\beta$  depend only on the most recent trajectory, and after those changes are made, the trajectory and its rewards are simply forgotten. This approach, while requiring little computation per iteration, wastes data and may require sampling many trajectories to reach convergence.

One technique for using data more efficiently is “experience replay” (Lin, 1993) explicitly remember all trajectories ever seen, and whenever asked to produce an updated set of coefficients, perform repeated passes of TD( $\lambda$ ) over all the saved trajectories until

convergence. This technique is similar to the batch training methods commonly used to train neural networks. However, in the case of linear function approximators, there is another way.

## 2. The least-squares TD( $\lambda$ ) algorithm

The Least-Squares TD( $\lambda$ ) algorithm, or LSTD( $\lambda$ ), converges to the same coefficients  $\beta_\lambda$  that TD( $\lambda$ ) does. However, instead of performing gradient descent, LSTD( $\lambda$ ) builds explicit estimates of the  $\mathbf{C}$  matrix and  $\mathbf{d}$  vector (actually, estimates of a constant multiple of  $\mathbf{C}$  and  $\mathbf{d}$ ), and then solves  $\mathbf{d} + \mathbf{C}\beta_\lambda = \mathbf{0}$  directly. The actual data structures that LSTD( $\lambda$ ) builds from experience are the matrix  $\mathbf{A}$  (of dimension  $K \times K$ , where  $K$  is the number of features) and the vector  $\mathbf{b}$  (of dimension  $K$ ):

$$\mathbf{b} = \sum_{i=0}^t \mathbf{z}_i R_i \quad \mathbf{A} = \sum_{i=0}^t \mathbf{z}_i (\phi(x_i) - \phi(x_{i+1}))^\top \quad (2)$$

After  $n$  independent trajectories have been observed,  $\mathbf{b}$  is an unbiased estimate of  $n\mathbf{d}$ , and  $\mathbf{A}$  is an unbiased estimate of  $-n\mathbf{C}$ . Thus,  $\beta_\lambda$  can be estimated as  $\mathbf{A}^{-1}\mathbf{b}$ . As is standard in least-squares algorithms, Singular Value Decomposition is used to invert  $\mathbf{A}$  robustly (Press et al., 1992). The complete LSTD( $\lambda$ ) algorithm is specified in figure 2.

When  $\lambda = 0$ , LSTD(0) reduces precisely to Bradtke and Barto's LSTD algorithm, which they derived using a more complex approach based on regression with instrumental variables (Bradtke & Barto, 1996). At the other extreme, when  $\lambda = 1$ , LSTD(1) produces the same  $\mathbf{A}$  and  $\mathbf{b}$  that would be produced by supervised linear regression on training pairs

---

### LSTD( $\lambda$ ) for approximate policy evaluation:

*Given:* a simulation model, featurizer, and  $\lambda$  as in ordinary TD( $\lambda$ ).

(No stepsize schedules or initial estimates of  $\beta$  are necessary.)

*Output:* a coefficient vector  $\beta$  for which  $V^\pi(x) \approx \beta \cdot \phi(x)$ .

Set  $\mathbf{A} := \mathbf{0}$ ,  $\mathbf{b} := \mathbf{0}$ ,  $t := 0$ .

**for**  $n := 1, 2, \dots$  **do:** {

Choose a start state  $x_t \in X$ .

Set  $\mathbf{z}_t := \phi(x_t)$ .

**while**  $x_t \neq \text{END}$ , **do:** {

Simulate one step of the chain, producing a reward  $R_t$  and next state  $x_{t+1}$ .

Set  $\mathbf{A} := \mathbf{A} + \mathbf{z}_t (\phi(x_t) - \phi(x_{t+1}))^\top$ . /\* outer product \*/

Set  $\mathbf{b} := \mathbf{b} + \mathbf{z}_t R_t$ .

Set  $\mathbf{z}_{t+1} := \lambda \mathbf{z}_t + \phi(x_{t+1})$ .

Set  $t := t + 1$ .

}

Whenever updated coefficients are desired: Set  $\beta := \mathbf{A}^{-1}\mathbf{b}$ . /\* Use SVD. \*/

}

---

Figure 2. A least-squares version of TD( $\lambda$ ) (compare figure 1). Note that  $\mathbf{A}$  has dimension  $K \times K$ , and  $\mathbf{b}$ ,  $\beta$ ,  $\mathbf{z}$ , and  $\phi(x)$  all have dimension  $K \times 1$ .

of {state features  $\mapsto$  observed Monte-Carlo returns}; the proof is given in the appendix. Thanks to the algebraic trick of the eligibility vectors, LSTD(1) builds the regression matrices *fully incrementally*—without having to store the trajectory while waiting to observe the eventual outcome. When trajectories through the chain are long, this provides significant memory savings over linear regression.

The computation per timestep required to update  $\mathbf{A}$  and  $\mathbf{b}$  is the same as least-squares linear regression:  $O(K^2)$ , where  $K$  is the number of features. LSTD( $\lambda$ ) must also perform a matrix inversion at a cost of  $O(K^3)$  whenever  $\beta$ 's coefficients are needed—typically, once per complete trajectory. (If updated coefficients are required more frequently, then the  $O(K^3)$  cost can be avoided by recursive-least-squares (Bradtke & Barto, 1996) or Kalman-filtering techniques [(Bertsekas & Tsitsiklis, 1996), §3.2.2], which update  $\beta$  on each timestep at a cost of only  $O(K^2)$ .) LSTD( $\lambda$ ) performs more computation per observation than incremental TD( $\lambda$ ), which updates the coefficients using only  $O(K)$  computation per timestep. However, as pointed out in Bradtke and Barto (1996), LSTD( $\lambda$ ) offers several significant advantages:

- Least-squares algorithms “extract more information from each additional observation” (Bradtke & Barto, 1996) and would thus be expected to converge with fewer training samples.
- TD( $\lambda$ )’s convergence can be slowed dramatically by a poor choice of the stepsize parameters  $\alpha_n$ . LSTD( $\lambda$ ) eliminates these parameters.
- TD( $\lambda$ )’s performance is sensitive to the initial estimate for  $\beta_\lambda$ . LSTD( $\lambda$ ) does not rely on an arbitrary initial estimate.
- TD( $\lambda$ ) is also sensitive to the ranges of the individual features. LSTD( $\lambda$ ) is not.

Section 4 below presents experimental results comparing the data efficiency of gradient-based and least-squares-based TD learning.

### 3. LSTD( $\lambda$ ) as model-based learning

Surprisingly, the move from an incremental, gradient-based update rule to a direct, least-squares-based update rule for TD( $\lambda$ ) turns out to be mathematically equivalent to a move from a model-free to a model-based reinforcement learning algorithm. This equivalence provides interesting new intuitions about the space of temporal difference learning algorithms.

To begin, let us restrict our attention to the case of a small discrete state space  $X$ , over which  $V^\pi$  can be represented and learned exactly by a lookup table. A classical model-based algorithm for learning  $V^\pi$  from simulated trajectory data would proceed as follows:

1. From the state transitions and rewards observed so far, build in memory an *empirical model* of the Markov chain. The sufficient statistics of this model are
  - a vector  $\mathbf{n}$  recording the number of times each state has been visited;
  - a matrix  $\mathbf{T}$  recording the observed state-transition counts:  $\mathbf{T}_{ij}$  = how many times  $x_j$  was seen to directly follow  $x_i$ ; and
  - a vector  $\mathbf{s}$  recording, for each state, the sum of all one-step rewards observed on transitions leaving that state.

2. Whenever a new estimate of the value function  $V^\pi$  is desired, solve the linear system of Bellman equations corresponding to the current empirical model. Writing  $\mathbf{N} = \mathbf{diag}(\mathbf{n})$ , the solution vector of  $V^\pi$  values is given by

$$\mathbf{v} = (\mathbf{N} - \mathbf{T})^{-1}\mathbf{s}. \quad (3)$$

This model-based technique contrasts with  $\text{TD}(\lambda)$ , a model-free approach to the same problem.  $\text{TD}(\lambda)$  does not maintain any statistics on observed transitions and rewards; it simply updates the components of  $\mathbf{v}$  directly. In the limit, assuming a lookup-table representation, both converge to the optimal  $V^\pi$ . The advantage of  $\text{TD}(\lambda)$  is its low computational burden per step; the advantage of the classical model-based method is that it makes the most of the available training data. The relative advantages of model-based and model-free reinforcement learning methods have been investigated in, e.g., Sutton (1990), Moore and Atkeson (1993) and Atkeson and Santamaria (1997).

Where does  $\text{LSTD}(\lambda)$  fit in? In fact, for the case of  $\lambda = 0$ , it precisely duplicates the classical model-based method sketched above. The assumed lookup-table representation for  $\tilde{V}^\pi$  means that we have one independent feature per state: the feature vector  $\phi$  corresponding to state 1 is  $(1, 0, 0, \dots, 0)$ ; corresponding to state 2 is  $(0, 1, 0, \dots, 0)$ ; etc. Referring to the algorithm of figure 2, we see that  $\text{LSTD}(0)$  performs the following operations upon each observed transition:

$$\mathbf{b} := \mathbf{b} + \phi(x_t)R_t \quad \mathbf{A} := \mathbf{A} + \phi(x_t)(\phi(x_t) - \phi(x_{t+1}))^\top \quad (4)$$

Clearly, the role of  $\mathbf{b}$  is to sum all the rewards observed at each state, exactly as the vector  $\mathbf{s}$  does in the classical technique.  $\mathbf{A}$ , meanwhile, accumulates the statistics  $(\mathbf{N} - \mathbf{T})$ . To see this, note that the outer product in Eq. (4) is a matrix consisting of an entry of  $+1$  on the single diagonal element corresponding to state  $x_t$ ; an entry of  $-1$  on the element in row  $x_t$ , column  $x_{t+1}$ ; and all the rest zeroes. Summing one such sparse matrix for each observed transition gives  $\mathbf{A} \equiv \mathbf{N} - \mathbf{T}$ . Finally,  $\text{LSTD}(0)$  performs the inversion  $\beta := \mathbf{A}^{-1}\mathbf{b} = (\mathbf{N} - \mathbf{T})^{-1}\mathbf{s}$ , giving the same solution as in Eq. (3).

Thus, when  $\lambda = 0$ , the  $\mathbf{A}$  and  $\mathbf{b}$  matrices built by  $\text{LSTD}(\lambda)$  effectively record a model of all the observed transitions. What about when  $\lambda > 0$ ? Again,  $\mathbf{A}$  and  $\mathbf{b}$  record the sufficient statistics of an empirical Markov model—but in this case, the model being captured is one whose single-step transition probabilities directly encode the multi-step  $\text{TD}(\lambda)$  backup operations. That is, the model links each state  $x$  to all the downstream states that follow  $x$  on any trajectory, and records how much influence each has on estimating  $\tilde{V}^\pi(x)$  according to  $\text{TD}(\lambda)$ . In the case of  $\lambda = 0$ , the  $\text{TD}(\lambda)$  backups correspond to the one-step transitions, resulting in the equivalence described above. The opposite extreme, the case of  $\lambda = 1$ , is also interesting: the empirical Markov model corresponding to  $\text{TD}(1)$ 's backups is the chain in which each state  $x$  leads directly to absorption, and  $\beta$  then simply computes the average Monte-Carlo return at each state. For general  $\lambda$ , we produce the statistics of the “simple beta-model” of multi-scale reinforcement learning (Sutton, 1995). In short, if we assume a lookup-table representation for the function  $\tilde{V}^\pi$ , we can view the  $\text{LSTD}(\lambda)$  algorithm as performing these two steps:

1. It implicitly uses the observed simulation data to build a Markov chain. This chain compactly models all the backups that  $\text{TD}(\lambda)$  would perform on the data.
2. It solves the chain by performing a matrix inversion.

The lookup-table representation for  $\tilde{V}^\pi$  is intractable in practical problems; in practice,  $\text{LSTD}(\lambda)$  operates on states only via their (linearly dependent) feature representations  $\phi(x)$ . In this case,  $\text{LSTD}(\lambda)$  implicitly builds a *compressed* version of the empirical model's transition matrix  $\mathbf{N} - \mathbf{T}$  and summed-reward vector  $\mathbf{s}$ :

$$\mathbf{b} = \Phi^\top \mathbf{s} \quad \mathbf{A} = \Phi^\top (\mathbf{N} - \mathbf{T}) \Phi \quad (5)$$

where  $\Phi$  is the  $|X| \times K$  matrix representation of the function  $\phi: X \rightarrow \mathfrak{R}^K$ . From the compressed empirical model,  $\text{LSTD}(\lambda)$  computes these coefficients for  $\tilde{V}^\pi$ :

$$\beta_\lambda = \mathbf{A}^{-1} \mathbf{b} = (\Phi^\top (\mathbf{N} - \mathbf{T}) \Phi)^{-1} (\Phi^\top \mathbf{s}). \quad (6)$$

Ideally, these coefficients  $\beta_\lambda$  would be equivalent to the *empirical optimal* coefficients  $\beta_\lambda^*$ . The empirical optimal coefficients are those that would be found by building the full uncompressed empirical model (represented by  $\mathbf{N} - \mathbf{T}$  and  $\mathbf{s}$ ), using a lookup table to solve for that model's value function ( $\mathbf{v} = (\mathbf{N} - \mathbf{T})^{-1} \mathbf{s}$ ), and then performing a least-squares linear fit from the state features  $\Phi$  to the lookup-table value function:

$$\beta_\lambda^* \stackrel{\text{def}}{=} (\Phi^\top \Phi)^{-1} (\Phi^\top \mathbf{v}) = (\Phi^\top \Phi)^{-1} \Phi^\top (\mathbf{N} - \mathbf{T})^{-1} \mathbf{s}. \quad (7)$$

It can be shown that Eqs. (6) and (7) are indeed equivalent for the case of  $\lambda = 1$ , because that setting of  $\lambda$  implies that  $\mathbf{T} = \mathbf{0}$ , thus  $(\mathbf{N} - \mathbf{T})^{-1}$  is diagonal and commutes. However, for the case of  $\lambda < 1$ , solving the compressed empirical model does not in general produce the optimal least-squares fit to the solution of the uncompressed model (Tsitsiklis & Van Roy, 1997).

#### 4. Experimental comparison of $\text{TD}(\lambda)$ and $\text{LSTD}(\lambda)$

This section reports experimental results comparing  $\text{TD}(\lambda)$  and  $\text{LSTD}(\lambda)$  on two domains: a simple illustrative Markov chain and the large-scale game of backgammon. We begin with the simple chain, pictured in figure 3. It consists of 13 states, each of which is represented by four state features as shown. We seek to represent the value function compactly as a linear function of these four state features. In fact, this domain's optimal  $V^\pi$  function is exactly linear in these features: the optimal coefficients  $\beta_\lambda^*$  are  $(-24, -16, -8, 0)$ .

Because  $V^\pi$  is exactly linear in this domain,  $\text{LSTD}(\lambda)$  is guaranteed to converge with probability 1 to the optimal  $\beta_\lambda^*$  for any  $\lambda$ .  $\text{TD}(\lambda)$  is also guaranteed to converge to  $\beta_\lambda^*$ , under the additional condition that an appropriate schedule of stepsizes is chosen. The following three criteria on the schedule  $(\alpha_n)$  are sufficient:  $\alpha_n \geq 0 \forall n$ ;  $\sum_{n=1}^{\infty} \alpha_n = \infty$ ; and

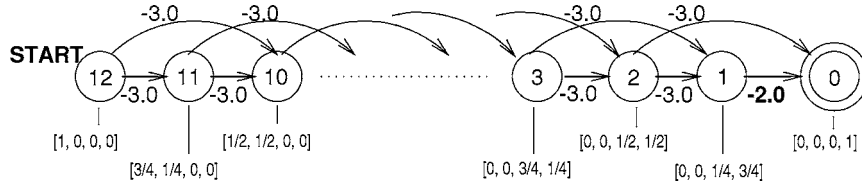


Figure 3. A 13-state Markov chain. In states 2–12, each outgoing arc is taken with probability 0.5. For value function approximation, each state is represented by four features, as follows: the representations for states 12, 8, 4, and 0 are, respectively, [1, 0, 0, 0], [0, 1, 0, 0], [0, 0, 1, 0], and [0, 0, 0, 1]; and the representations for the other states are obtained by linearly interpolating between these.

$\sum_{n=1}^{\infty} \alpha_n^2 < \infty$ . Our experiments use schedules that satisfy these criteria, having the form:

$$\alpha_n \stackrel{\text{def}}{=} a_0 \frac{n_0 + 1}{n_0 + n} \quad n = 1, 2, \dots \tag{8}$$

The parameter  $a_0$  determines the initial stepsize, and  $n_0$  determines how gradually the stepsize decreases over time. Each  $\text{TD}(\lambda)$  experiment was run with six different stepsize schedules, corresponding to the six combinations of  $a_0 \in \{0.1, 0.01\}$  and  $n_0 \in \{10^2, 10^3, 10^6\}$ . These settings produce a range of learning rates that is typical for applications of gradient descent.  $\text{TD}(\lambda)$  also requires an initial setting for the coefficients  $\beta$ ; we initialized  $\beta$  to  $\mathbf{0}$  on each experiment.

Comparative results are given in figures 4 and 5. Figure 4 focuses on the case of  $\lambda = 0.4$ , comparing the learning curve for  $\text{LSTD}(\lambda)$  against those of all six schedules of  $\text{TD}(\lambda)$ . Each point plotted represents the average of 10 trials. The plot shows clearly that for

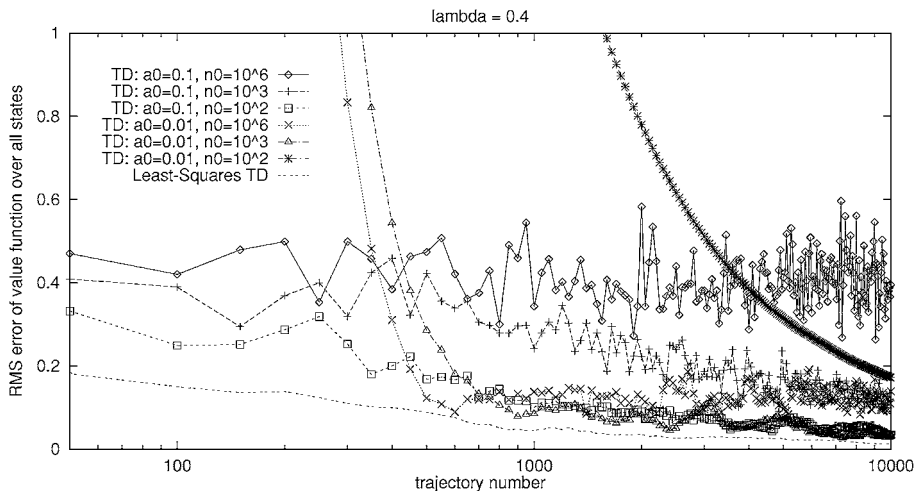


Figure 4. Performance of  $\text{TD}(0.4)$  and  $\text{LSTD}(0.4)$  on the sample domain. Note the log scale on the  $x$ -axis. All points plotted represent the average of 10 trials.



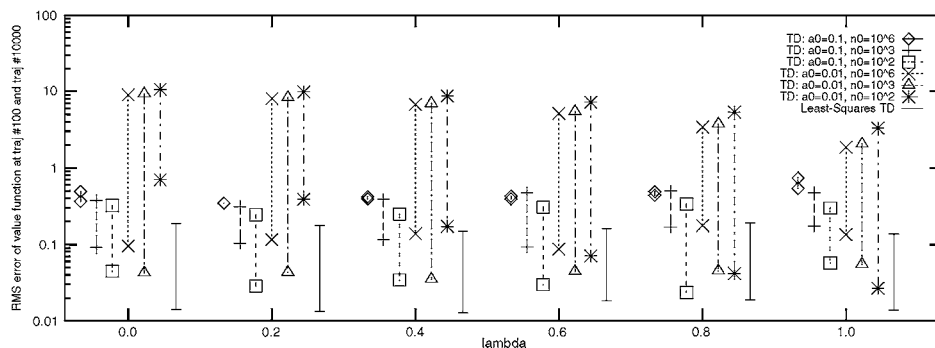


Figure 5. Summary of results at six settings of  $\lambda$ . At each setting, seven algorithms are compared:  $\text{TD}(\lambda)$  (with six different stepsize schedules) and  $\text{LSTD}(\lambda)$ . The plotted segment shows the mean RMS value function approximation error after 100 trajectories (top of segment) and 10,000 trajectories (bottom of segment). Note the log scale on the y-axis.  $\text{LSTD}(\lambda)$  is best in all cases.

$\lambda = 0.4$ ,  $\text{LSTD}(\lambda)$  learns a good approximation to  $V^\pi$  in fewer trials than any of the  $\text{TD}(\lambda)$  experiments, and performs better asymptotically as well.

Figure 5 graphically summarizes six learning-curve plots similar to figure 4, corresponding to varying  $\lambda$  over the settings  $\{0, 0.2, 0.4, 0.6, 0.8, 1.0\}$ . The results may be summarized as follows:

- Across all values of  $\lambda$ ,  $\text{LSTD}(\lambda)$  learns a good approximation to  $V^\pi$  in fewer trials than any of the  $\text{TD}(\lambda)$  experiments, and performs better asymptotically as well.
- The performance of  $\text{TD}(\lambda)$  depends critically on the stepsize schedule chosen.  $\text{LSTD}(\lambda)$  has no tunable parameters other than  $\lambda$  itself.
- Varying  $\lambda$  has a relatively small effect on  $\text{LSTD}(\lambda)$ 's performance.

#### 4.1. Large-scale example: Backgammon policy evaluation

We now turn to a larger-scale test domain, familiar in the reinforcement-learning literature: the game of backgammon. In this larger domain, we compare  $\text{TD}(\lambda)$  with  $\text{LSTD}(\lambda)$  in terms of overall computational efficiency as well as data efficiency.

We generate trajectories by opposing two fixed, pre-trained backgammon policies: “Fiona,” a neural network patterned after TD-Gammon (Tesauro, 1994), and “pubeval,” a benchmark backgammon evaluator.<sup>1</sup> We sidestep the two-player aspect of backgammon by defining a Markov chain consisting of only those states in which it is Fiona’s turn to play. That is, each single transition covers two backgammon moves, one by Fiona and one by pubeval. The value function we seek to learn specifies, for any backgammon position, the expected winnings of Fiona over pubeval, starting from that position with Fiona to roll the dice and move first.

In the small-scale domain of the previous section, both  $\text{TD}(\lambda)$  and  $\text{LSTD}(\lambda)$  were known to converge to the optimal value function  $V^\pi$ , because  $V^\pi$  was exactly linear in that domain’s

features. In backgammon, by contrast, we cannot guarantee convergence to the optimal value function—nor do we even have that function available for measuring learning performance. Instead, we measure performance against a “gold standard” value function approximation  $\tilde{G}(x)$ , built by training LSTD(1) on a set of 100,000 sample games (amounting to nearly 3 million positions). The learning architecture for  $\tilde{G}(x)$  was quadratic regression over 15 backgammon-specific high-level features.<sup>2</sup> To measure the performance of a new value function approximation  $\tilde{V}^\pi$ , then, we compute the sampled RMS difference between  $\tilde{V}^\pi$  and  $\tilde{G}$ :

$$\text{Error}(\tilde{V}^\pi) \stackrel{\text{def}}{=} \left( \frac{1}{|\mathcal{S}|} \sum_{x \in \mathcal{S}} (\tilde{V}^\pi(x) - \tilde{G}(x))^2 \right)^{1/2}$$

over an independent test set  $\mathcal{S}$  of 10,000 backgammon positions.

Our experiments compared TD( $\lambda$ ) with LSTD( $\lambda$ ) using both linear and quadratic approximations over the 15 backgammon features noted above. The coefficient vectors corresponding to these approximations comprise  $K = 16$  and  $K = 136$  parameters, respectively. The parameter  $\lambda$  was set arbitrarily to 0.9. For TD( $\lambda$ ), four different learning rate schedules were used:  $a_0 \in \{10^{-2}, 10^{-3}, 10^{-4}, 10^{-5}\}$ ,  $n_0 = 1000$ . The coefficients  $\beta$  were always initialized to  $\mathbf{0}$ . Performance measurements were taken after every 50 games. Each algorithm was trained on the same 9 sets of 10,000 games each, and the average error curves over these 9 runs are plotted.

Results are shown in figure 6. The upper row of the figure shows the average error curves for the linear case (at left) and quadratic case (at right) as a function of the number of sample trajectories observed. As in the smaller-scale results presented earlier, we can see that LSTD( $\lambda$ ) brings down the approximation error using far less sample data than TD( $\lambda$ ). We also see again that TD( $\lambda$ ) is quite sensitive to the learning rate schedule chosen; in fact, in the quadratic case, the highest learning rate of  $a_0 = 0.01$  produced error curves greater than the bounds of the plot.

The lower row of figure 6 plots the same error curves, but with respect to elapsed time rather than the number of trajectories observed. Experiments were run on a 296 MHz Ultrasparc II. The linear regression results (at left) show that when  $\beta$  has relatively few parameters, the time penalty for using LSTD( $\lambda$ ) is slight. With quadratic regression (at right), by contrast, the large matrix inversions performed by LSTD( $\lambda$ ) cause it to run markedly more slowly than TD( $\lambda$ ) (note the semilog scale). Nevertheless, because of its better data efficiency, the RMS error of LSTD( $\lambda$ ) surpasses that of the best TD( $\lambda$ ) run after only about 10 seconds of computation.

In interpreting these timing results, two important points should be considered. First, LSTD( $\lambda$ ) only performed its coefficient update  $\beta := \mathbf{A}^{-1} \mathbf{b}$  once every 50 trajectories, when needed to collect performance statistics. Practical applications may require updating the  $\beta$  coefficients more or less frequently, with a corresponding increase or decrease in computational requirements. Second, the timings reported here measure only the time spent in the algorithmic calculations of TD( $\lambda$ ) or LSTD( $\lambda$ ); the time spent in trajectory generation is excluded. In many practical domains, including backgammon, the time spent in generating sample state transitions in the Markov chain may exceed the time required for

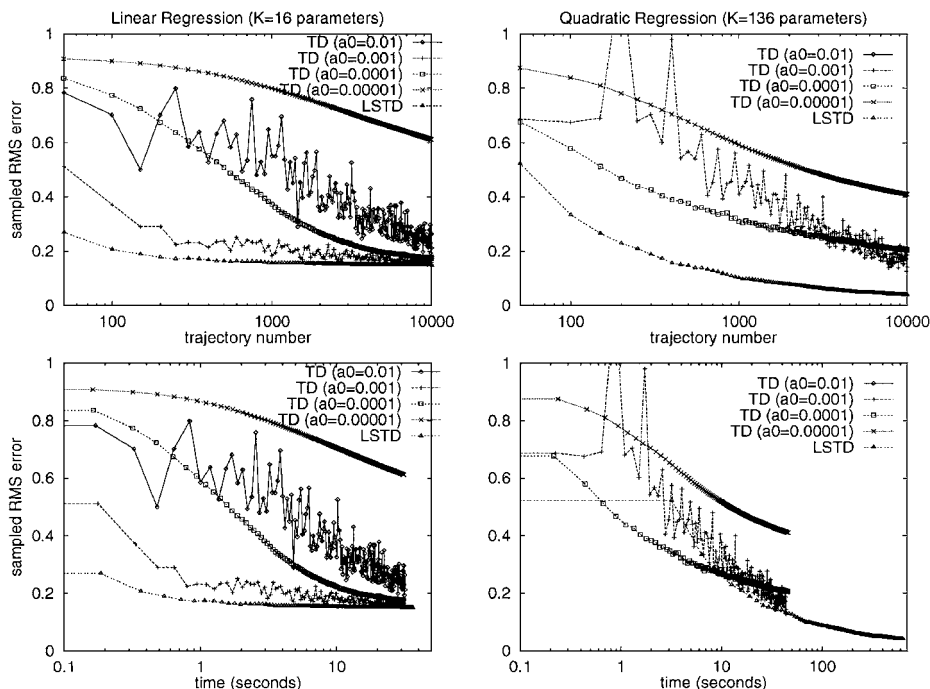


Figure 6. Results of TD(0.9) and LSTD(0.9) on backgammon policy evaluation. The upper two graphs plot learning curves with respect to the number of trajectories sampled, for linear and quadratic models, respectively. The lower two graphs plot the identical learning curves, but with respect to the amount of elapsed wall-clock time used by the learning algorithms (excluding time spent generating moves). All points plotted represent the average of 9 trials.

temporal-difference learning. When data is expensive to generate, maximizing the data efficiency of learning as LSTD( $\lambda$ ) does will reduce overall running times despite the greater computation spent in learning.

## 5. Conclusions and future work

We have argued, both from the experimental results above and from the deep connection to model-based reinforcement learning presented in Section 3, that the least-squares formulation of TD learning makes better use of simulation data than TD( $\lambda$ ). As a practical matter, the decision of when to prefer one algorithm over the other will depend on the application domain. If a domain has many features and simulation data is available cheaply, then incremental methods such as TD( $\lambda$ ) may have better real-time performance than least-squares methods (Sutton, 1992). On the other hand, some reinforcement learning applications have been successful with small numbers of features (e.g., Singh & Bertsekas, 1997; Boyan & Moore, 1998), and in these situations LSTD( $\lambda$ ) should be superior.

LSTD( $\lambda$ ) has been applied successfully in the context of STAGE, a reinforcement learning algorithm for combinatorial optimization (Boyan, 1998). An exciting possibility for future work is to apply LSTD( $\lambda$ ) in the context of approximation algorithms for general Markov decision problems. LSTD( $\lambda$ ) provides an alternative to TD( $\lambda$ ) for the inner loop of optimistic policy iteration (Bertsekas & Tsitsiklis, 1996), and should enable good control policies to be discovered with fewer trial simulations.

## Appendix

### *Equivalence of LSTD(1) and linear regression*

We show here that when  $\lambda = 1$ , the LSTD( $\lambda$ ) algorithm of figure 2 produces an approximate value function which is equivalent to that which would be produced by standard, non-incremental, least-squares linear regression. To be precise, assume we are given a sample trajectory  $(x_0, x_1, \dots, x_L, \text{END})$  of a Markov chain, with corresponding feature vectors  $\phi_t$  and one-step rewards  $R_t$  on each step. From this trajectory, a supervised-learning system learning from Monte-Carlo returns would generate the following training pairs:

$$\begin{aligned}\phi_0 &\mapsto R_0 + R_1 + \dots + R_L \\ \phi_1 &\mapsto R_1 + \dots + R_L \\ &\vdots \\ \phi_L &\mapsto R_L\end{aligned}$$

Performing standard least-squares linear regression on the above training set would produce the coefficients  $\beta = \mathbf{A}^{-1}\mathbf{b}$ , with  $\mathbf{A}$  and  $\mathbf{b}$  computed as follows:

$$\mathbf{A}_{\text{LR}} = \sum_{i=0}^L \phi_i \phi_i^\top \quad \mathbf{b}_{\text{LR}} = \sum_{i=0}^L \phi_i y_i \quad \text{where } y_i = \sum_{j=i}^L R_j$$

We now show that, thanks to the algebraic trick of the eligibility vectors  $\mathbf{z}_t$ , LSTD(1) builds the equivalent  $\mathbf{A}$  and  $\mathbf{b}$  fully incrementally—without having to store the trajectory while waiting to observe the eventual outcome  $y_i$ .

**Proof:** With simple algebraic manipulations, the sums built by LSTD(1)'s  $\mathbf{A}$  and  $\mathbf{b}$  telescope neatly into  $\mathbf{A}_{\text{LR}}$  and  $\mathbf{b}_{\text{LR}}$ , as follows:

$$\begin{aligned}\mathbf{A} &= \sum_{i=0}^L \mathbf{z}_i (\phi_i - \phi_{i+1})^\top \\ &= \sum_{i=0}^L \left( \sum_{j=0}^i \phi_j \right) (\phi_i - \phi_{i+1})^\top \quad (\text{by definition of } \mathbf{z}_t)\end{aligned}$$

$$\begin{aligned}
&= \sum_{i=0}^L \sum_{j=0}^i \phi_j \phi_i^\top - \sum_{i=0}^L \sum_{j=0}^i \phi_j \phi_{i+1}^\top \\
&= \left( \phi_0 \phi_0^\top + \sum_{i=1}^L \sum_{j=0}^i \phi_j \phi_i^\top \right) - \left( \sum_{k=1}^{L+1} \sum_{j=0}^{k-1} \phi_j \phi_k^\top \right) && \text{(substituting } k \equiv i + 1 \text{)} \\
&= \phi_0 \phi_0^\top + \sum_{i=1}^L \sum_{j=0}^i \phi_j \phi_i^\top - \sum_{k=1}^L \sum_{j=0}^{k-1} \phi_j \phi_k^\top && \text{(since } \phi_{L+1} \stackrel{\text{def}}{=} \mathbf{0} \text{)} \\
&= \phi_0 \phi_0^\top + \sum_{i=1}^L \left( \sum_{j=0}^i \phi_j \phi_i^\top - \sum_{j=0}^{i-1} \phi_j \phi_i^\top \right) && \text{(substituting } i \equiv k \text{)} \\
&= \sum_{i=0}^L \phi_i \phi_i^\top \\
&= \mathbf{A}_{\text{LR}}, \text{ as desired;}
\end{aligned}$$

and

$$\begin{aligned}
\mathbf{b} &= \sum_{i=0}^L \mathbf{z}_i R_i \\
&= \sum_{i=0}^L \left( \sum_{j=0}^i \phi_j \right) R_i && \text{(by definition of } \mathbf{z}_i \text{)} \\
&= \sum_{i=0}^L \sum_{j=0}^L \mathbf{1}(j \leq i) \phi_j R_i && \text{(where } \mathbf{1}(\text{True}) \stackrel{\text{def}}{=} 1, \mathbf{1}(\text{False}) \stackrel{\text{def}}{=} 0 \text{)} \\
&= \sum_{j=0}^L \sum_{i=0}^L \mathbf{1}(j \leq i) \phi_j R_i \\
&= \sum_{j=0}^L \phi_j \sum_{i=j}^L R_i \\
&= \mathbf{b}_{\text{LR}}, \text{ as desired.}
\end{aligned}$$

These reductions prove that the contributions to  $\mathbf{A}$  and  $\mathbf{b}$  by any single trajectory are identical in LSTD(1) and least-squares linear regression. In both algorithms, contributions from multiple trajectories are simply summed into the matrices. Thus, LSTD(1) and linear regression compute the same statistics and, ultimately, the same coefficients for the approximated value function.  $\square$

### Acknowledgments

Thanks to Andrew Moore, Jeff Schneider, Geoff Gordon, Jeremy Frank, Hamid Berenji, and the anonymous reviewers for many helpful comments. Thanks to Marc Ringuette for

helping with the design of the backgammon evaluator. This work was supported by a NASA GSRP fellowship while the author was at Carnegie Mellon University.

## Notes

1. The Fiona network, designed by Marc Ringuette and myself, consists of 200 input units, 20 sigmoidal hidden units, and 3 linear output units. Pubeval is a linear network designed by Gerry Tesauro, available from <http://www.cs.cmu.edu/afs/cs/project/connect/code/tesauro/>.
2. The 15 backgammon features used for training were measures of the race status, amount of contact, prime strength, inner board strength, checkers on bar, and immediate hit probability. Please contact the author for details.

## References

- Atkeson, C. G., & Santamaria, J. C. (1997). A comparison of direct and model-based reinforcement learning. In *International Conference on Robotics and Automation*.
- Bertsekas, D., & Tsitsiklis, J. (1996). *Neuro-dynamic programming*. Belmont, MA: Athena Scientific.
- Boyan, J. A. (1998). Learning evaluation functions for global optimization. Ph.D. Thesis, Carnegie Mellon University.
- Boyan, J. A., & Moore, A. W. (1998) Learning evaluation functions for global optimization and Boolean satisfiability. In *Proceedings of the Fifteenth National Conference on Artificial Intelligence (AAAI)*.
- Bradtke, S. J., & Barto, A. G. (1996). Linear least-squares algorithms for temporal difference learning. *Machine Learning*, 22:1-3, 33-57.
- Lin, L.-J. (1993). Reinforcement learning for robots using neural networks. Ph.D. Thesis, Carnegie Mellon University.
- Moore, A. W., & Atkeson, C. G. (1993). Prioritized sweeping: Reinforcement learning with less data and less time. *Machine Learning*, 13, 103-130.
- Press, W. H., Teukolsky, S. A., Vetterling, W. T., & Flannery, B. P. (1992). *Numerical recipes in C: The art of scientific computing*. (2nd ed.), Cambridge: Cambridge University Press.
- Singh, S., & Bertsekas, D. (1997). Reinforcement learning for dynamic channel allocation in cellular telephone systems. In M. C. Mozer, M. I. Jordan, & T. Petsche (Eds.), *NIPS-9* (p. 974). Cambridge, MA: The MIT Press.
- Sutton, R. S. (1988). Learning to predict by the methods of temporal differences. *Machine Learning*, 3.
- Sutton, R. S. (1990). Integrated architectures for learning, planning, and reacting based on approximating dynamic programming. In *Proceedings of the Seventh International Conference on Machine Learning*. San Mateo, CA: Morgan Kaufmann.
- Sutton, R. S. (1992). Gain adaptation beats least squares. In *Proceedings of the 7th Yale Workshop on Adaptive and Learning Systems* (pp. 161-166).
- Sutton, R. S. (1995). TD models: Modeling the world at a mixture of time scales. In *Machine Learning: Proceedings of the 12th International Conference* (pp. 531-539). San Mateo, CA: Morgan Kaufmann.
- Sutton, R. S., & Barto, A. G. (1998). *Reinforcement learning: An introduction*. Cambridge, MA: MIT Press.
- Tesauro, G. (1994). TD-Gammon, a self-teaching backgammon program, achieves master-level play. *Neural Computation*, 6:2, 215-219.
- Tsitsiklis, J. N., & Van Roy, B. (1997). An analysis of temporal-difference learning with function approximation. *IEEE Trans. Auto. Control*, 42:5, 674-690.

Received March 16, 1999

Revised November 3, 1999

Accepted November 3, 1999

Final manuscript November 3, 1999