

TD(λ) Converges with Probability 1

PETER DAYAN
TERRENCE J. SEJNOWSKI
CNL, The Salk Institute, P.O. Box 85800, San Diego, CA 92186-5800

(dayan@helmholtz.sdsc.edu)
(tsejnowski@uscd.edu)

Editor: Richard Sutton

Abstract. The methods of temporal differences (Samuel, 1959; Sutton, 1984, 1988) allow an agent to learn accurate predictions of stationary stochastic future outcomes. The learning is effectively stochastic approximation based on samples extracted from the process generating the agent's future.

Sutton (1988) proved that for a special case of temporal differences, the expected values of the predictions converge to their correct values, as larger samples are taken, and Dayan (1992) extended his proof to the general case. This article proves the stronger result that the predictions of a slightly modified form of temporal difference learning converge with probability one, and shows how to quantify the rate of convergence.

Keywords. reinforcement learning, temporal differences, Q-learning

1. Introduction

Temporal difference (TD) learning is a way of extracting information from observations of sequential stochastic processes so as to improve predictions of future outcomes. Its key insight is that estimates from successive states should be self-consistent—for instance, the prediction made at one state about a terminal outcome should be related to the prediction from the next state, since this transition is obviously one of the ways of getting to the termination point. The TD algorithm investigated here was invented by Sutton (1988), and uses the difference between such predictions to drive modifications to the parameters that generate them. In fact, Sutton defined a whole class of such TD algorithms, TD(λ), which look at these differences further and further ahead in time, weighted exponentially less according to their distance by the parameter λ .

TD(λ) algorithms have wide application, from modeling classical conditioning in animal learning (Sutton & Barto, 1987) to generating a prize-winning backgammon-playing program (Tesauro, 1992). However, the theory underlying the algorithms is not so well developed. Sutton (1988) proved the first theorem, which demonstrated under certain conditions that the mean TD estimates of the terminal reward or return from a particular form of absorbing Markov process converge to the appropriate values. Using some insightful analysis by Watkins (1989) about how TD methods for control relate to dynamic programming (Ross, 1983), Dayan (1992) generalized Sutton's theorem to cover the case where the differences between the predictions from many successive states are taken into account.

Unfortunately, mere convergence of the mean is a very weak criterion. For instance, consider a sequence \mathcal{J}_n of independent, identically distributed, random variables with a finite expectation H . Trivially, $\lim_{n \rightarrow \infty} \mathcal{E}[\mathcal{J}_n] = H$, but in no useful sense do the \mathcal{J}_n converge. Convergence of the mean is not the same as convergence *in* mean, which would

require that $\lim_{n \rightarrow \infty} \mathcal{E}[|\mathcal{J}C_n - H|] = 0$. The latter does not hold for the simple example. Convergence with probability one is one of the more desirable forms of stochastic convergence, and this article concentrates on proving that it holds for TD.

That the simplest form of TD converges with probability one in a very special case was shown by Dayan (1992). This pointed out the equivalence of TD and Watkins' (1989) Q-learning stochastic control learning algorithm, in the case where at no stage is there any choice of possible action. Watkins (1989) and Watkins et al. (1992) proved convergence of Q-learning with probability one, and this assurance therefore extends to TD.

The present article applies some general theory by Kushner and Clark (1978), which was originally directed at the Robbins-Monro procedure (Robbins & Monro, 1951) to prove that TD(λ) with a linear representation converges with probability one. Kuan and White (1990, 1991) have applied the same theory to the stochastic convergence of static and dynamic backpropagation; they provide a more easily approachable introduction to the methods and also more directly applicable convergence conditions.

2. Definition of TD(λ)

We will treat the same Markov estimation problem that Sutton (1988) used, and will generally adopt his notation. Consider the case of an absorbing Markov chain with stochastic terminal returns, defined by sets and values:

\mathcal{J}		Terminal states
\mathcal{N}		Non-terminal states
$q_{ij} \in [0, 1]$	$i \in \mathcal{N}, j \in \mathcal{N}$	Transition probabilities between non-terminal states
$s_{ij} \in [0, 1]$	$i \in \mathcal{N}, j \in \mathcal{J}$	Transition probabilities to terminal states
$\mathbf{x}_i \in \mathcal{R}^c$	$i \in \mathcal{N}$	Vectors representing non-terminal states
v_j	$j \in \mathcal{J}$	Expected terminal return from state j
μ_i	$i \in \mathcal{N}$	Probabilities of starting at state i
		where $\sum_{i \in \mathcal{N}} \mu_i = 1$

The estimation system is fed complete sequences $\mathbf{x}_{i_1}, \mathbf{x}_{i_2}, \dots, \mathbf{x}_{i_m}$ of observation vectors, together with their scalar terminal return v . It has to generate for every non-terminal state $i \in \mathcal{N}$ a prediction of the expected value $\mathcal{E}[v|i]$ for starting from that state. If the transition matrix of the Markov chain were completely known, these predictions could be computed as

$$\bar{e}_i^* \equiv \mathcal{E}[v|i] = \sum_{j \in \mathcal{J}} s_{ij} v_j + \sum_{j \in \mathcal{N}} q_{ij} \sum_{k \in \mathcal{J}} s_{jk} v_k + \sum_{j \in \mathcal{N}} q_{ij} \sum_{k \in \mathcal{N}} q_{jk} \sum_{l \in \mathcal{J}} s_{kl} v_l + \dots \quad (1)$$

where \bar{e}^* is the vector of correct predictions.

Again, following Sutton, let $[M]_{ab}$ denote the ab^{th} entry of any matrix M , $[\mathbf{u}]_a$ denote the a^{th} component of any vector \mathbf{u} , Q denote the square matrix with components $[Q]_{ab} = q_{ab}$, $a, b \in \mathcal{N}$, and \mathbf{h} denote the vector whose components are $[\mathbf{h}]_a = \sum_{b \in \mathcal{J}} s_{ab} v_b$, for $a \in \mathcal{N}$. Then from equation (1)

$$\mathcal{E}[\nu | i] = \left[\sum_{k=0}^{\infty} Q^k \mathbf{h} \right]_i = [(\mathbf{I} - Q)^{-1} \mathbf{h}]_i. \quad (2)$$

As Sutton showed, the existence of the limit in this equation follows from the fact that Q is the transition matrix for the non-terminal states of an absorbing Markov chain, which, with probability one, will ultimately terminate.

During the learning phase, linear TD(λ) generates successive vectors $\mathbf{w}_1, \mathbf{w}_2, \dots$, changing \mathbf{w} after each complete observation sequence. Define $V_n^\lambda(i) = \mathbf{w}_n \cdot \mathbf{x}_i$ as the prediction of the terminal return starting from state i , at stage n in learning. Then, during one such sequence, $V_n^\lambda(i_t)$ are the intermediate predictions of these terminal returns, and, abusing notation somewhat, define also $V_n^\lambda(i_{m+1}) = \nu$, the observed terminal return.¹ TD(λ) changes \mathbf{w} according to

$$\mathbf{w}_{n+1} = \mathbf{w}_n + \sum_{i=1}^m \left\{ \alpha_{n+1} [V_n^\lambda(i_{t+1}) - V_n^\lambda(i_t)] \sum_{k=1}^i \lambda^{i-k} \nabla_{\mathbf{w}_n} V_n^\lambda(i_k) \right\} \quad (3)$$

where α_{n+1} is the learning rate for the n^{th} trial.

Sutton proved the following theorem for $\lambda = 0$; Dayan (1992) extended it to the case of general $0 < \lambda < 1$.

Theorem T. For any absorbing Markov chain, for any distribution of starting probabilities μ_i such that there are no inaccessible states, for any outcome distributions with finite expected values ν_j , and for any linearly independent set of observation vectors $\{\mathbf{x}_i | i \in \mathcal{X}\}$, there exists an $\epsilon > 0$ such that, if $\alpha_n = \alpha$ where $0 < \alpha < \epsilon$ and for any initial weight vector, the predictions of linear TD(λ) (with weight updates after each sequence) converge in expected value to the ideal predictions of equation (2); that is, if \mathbf{w}_n denotes the weight vector after n sequences have been experienced, then

$$\lim_{n \rightarrow \infty} \mathcal{E}[\mathbf{w}_n, \mathbf{x}_i] = \mathcal{E}[\nu | i] = [(\mathbf{I} - Q)^{-1} \mathbf{h}]_i, \quad \forall i \in \mathcal{X}. \quad (4)$$

In proving this for $\lambda < 1$, it turns out that there is a linear update matrix for the mean estimates of terminal return from each state after n sequences have been observed. If \mathbf{w}_n are the actual weights after sequence n , and if $\bar{\mathbf{w}}_{n+1}$ are the expected weights after the next sequence, then

$$X^T \bar{\mathbf{w}}_{n+1} = X^T \mathbf{w}_n - \alpha_{n+1} X^T X D [\mathbf{I} - (1 - \lambda) Q (\mathbf{I} - \lambda Q)^{-1}] (X^T \mathbf{w}_n - \bar{\mathbf{e}}^*), \quad (5)$$

where X is the matrix whose columns are the vectors representing the non-terminal states; $[X]_{ab} = [\mathbf{x}_a]_b$. Furthermore, the mean estimates converge appropriately because, if the vectors representing the states are independent (i.e., if X is full rank), then $-X^T X D [\mathbf{I} - (1 - \lambda) Q (\mathbf{I} - \lambda Q)^{-1}]$ has a full set of eigenvalues, each of whose real parts are negative. It turns out that these, together with some other conditions that are mainly guaranteed by the finiteness of the Markov chain, are just what justify the use of the powerful stochastic convergence proof methods of Kushner and Clark (1978). The next two sections show how.

3. Convergence proof

On pages 21–24 and 26–27 of their book, Kushner and Clark (1978) consider the following problem (changing the notation to fit with the above). $z_i \in \mathfrak{R}^c$ are random variables, $k : \mathfrak{R}^c \rightarrow \mathfrak{R}^c$ is a stochastic function whose mean for every $z \in \mathfrak{R}^c$ is $\bar{k}(z)$, $0 < \alpha_n < 1$ is a sequence of real numbers such that

$$\sum_{i=1}^{\infty} \alpha_n = \infty \text{ but } \sum_{i=1}^{\infty} \alpha_n^2 < \infty \quad (6)$$

and

$$z_{n+1} = z_n + \alpha_{n+1} \bar{k}(z_n) + \alpha_{n+1} d_n, \text{ and} \quad (7)$$

$$d_n = [k(z_n) - \bar{k}(z_n)] \quad (8)$$

where d_n acts as a zero mean noise process.

Define a discrete ‘time’ as $t_n = \sum_{i=1}^n \alpha_i$ and a piecewise linear interpolation of the sequence $\{z_n\}$ as

$$z^{\mathcal{Z}}(t) = \frac{t_{n+1} - t}{\alpha_n} z_n + \frac{t - t_n}{\alpha_n} z_{n+1} \quad (9)$$

and functions that are left shifts of this as $z^n(t) \equiv z^{\mathcal{Z}}(t + t_n)$. Then, as a consequence of theorem 2.3.1 of Kushner and Clark (1978), if the variance of the noise d_n is bounded and $\{z_n\}$ is bounded with probability 1, then, also with probability 1, the sequence of functions $\{z^n(t)\}$ has a convergent subsequence that tends to some function $z(t)$ that satisfies the equation

$$\frac{dz}{dt} = \bar{k}(z) \quad (10)$$

Kushner and Clark (1978) also show that $z_n \rightarrow z_0$ as $n \rightarrow \infty$ if there is a particular constant solution z_0 of this differential equation that is asymptotically stable in the following way: for any bounded solution $z(t)$ and any $\epsilon > 0$, there is a $\delta > 0$ such that $|z(t) - z_0| < \epsilon$ for $t \geq 0$ if $|z(0) - z_0| < \delta$, and $z(t) \rightarrow z_0$ as $t \rightarrow \infty$.

Consider this in the context of TD(λ). Define $z_n \equiv X^T w_n$ as the vector containing all the predictions at the n^{th} trial. Then in the appropriate form of the update equation (7), the first two terms on the right-hand side, which take care of the mean, are just equation (5). Therefore, the associated differential equation equivalent to expression (10) is

$$\frac{dz}{dt} = \bar{k}(z) = -X^T X D [I - (1 - \lambda) Q (I - \lambda Q)^{-1}] (z - \bar{e}^*). \quad (11)$$

As mentioned above, from Sutton (1988) and Dayan (1992), if the vectors representing the states are independent (i.e., if X is full rank), the negated growth matrix in equation (11) has a full set of eigenvalues, all of whose real parts are negative. Therefore, the differential equation is asymptotically stable about \bar{e}^* in the above manner.

So, if the variance of the noise were bounded, and $\{z_n\}$ were bounded with probability 1, then, from the above, we would know that $z_n \rightarrow \bar{e}^*$ (equivalently, $w_n \rightarrow w^*$ where $X^T w^* = \bar{e}^*$ as X is full rank) as $n \rightarrow \infty$, with probability 1, i.e., TD(λ) would converge to the right answer with probability 1.

Although far weaker conditions would probably suffice (Kushner, 1984), it is adequate to bound k using a projection technique. Choose a large real bound \mathcal{B} . If z_{n+1} , updated according to equation (7), would lie outside the hypercube defined as $\{-\mathcal{B}, \mathcal{B}\}^c$, it is projected orthogonally back along the offending axes so that it lies on the surface. This makes both $\{z_n\}$ and the variance $\mathcal{V}[k(z)]$ bounded and so, since also the Markov chain is absorbing, Kushner and Clark's (1978) theorem 2.3.1 holds for TD(λ), provided that $\bar{e}^* \in \{-\mathcal{B}, \mathcal{B}\}^c$.

4. Rate of convergence

Kushner and Clark (1978) go on to show how to quantify the rate of convergence of stochastic algorithms such as the Robbins–Monro procedure. In the notation of this article, they consider setting $\alpha_n = \gamma(n+1)^{-r}$, for some $1 \geq r > 1/2$, define $F_n = z_n - \bar{e}^*$ and $u_n = (n+1)^s F_n$, and deem the rate of convergence the largest $s \in (0, 1)$ “for which the asymptotic part of $\{u_n\}$ makes sense as a non-degenerate but ‘stable’ process” (Kushner & Clark, 1978, p. 233). r determines how fast learning can progress—the larger it is, the quicker.

If we define H as the Jacobian matrix of $\bar{k}(\cdot)$ at \bar{e}^* , i.e.,

$$H = -X^T X D [I - (1 - \lambda)Q(I - \lambda Q)^{-1}], \quad (12)$$

then theorem 7.3.1 (Kushner & Clark, 1978, p. 245) proves the following theorem (ignoring conditions that are clearly true in the present case):

Theorem P. If

1. $\exists \bar{e}^*$ such that $z_n \rightarrow \bar{e}^*$ as $n \rightarrow \infty$, with probability 1, and $\bar{k}(\bar{e}^*) = 0$,
2. Either:
 - a) $r = 1$, $s = 1/2$, and $\bar{H} \equiv \gamma H + sI$ has the real part of all its eigenvalues strictly less than 0, or
 - b) $r < 1$, $s = r/2$, and $\bar{H} \equiv \gamma H$ has the real part of all its eigenvalues strictly less than 0
3. There is a matrix M such that $\mathcal{E}[\{k(z_n) - \bar{k}(z_n)\} \{k(z_n) - \bar{k}(z_n)\}^T] \rightarrow M$ as $n \rightarrow \infty$.
4. There are real $\delta > 0$ and $\omega < \infty$ such that $\mathcal{E}[|k(z_n) - \bar{k}(z_n)|^{2+\delta}] \leq \omega$ for all n

then u_n converges in distribution to a normally distributed random variable.

Conditions 1 and 2b are guaranteed by the previous theorem (condition 2a is stronger and may hold in certain particular cases), and conditions 3 and 4 hold since the Markov chain is absorbing and z_n converges with probability 1. Therefore, TD(λ) converges at least as fast as $r/2$, where r can be chosen to be 1 if condition 2a holds.

For case 2a, the asymptotic variance of the normal distribution of \mathbf{u}_n is related to the covariance matrix M , which in turn is related to the variance of the update operator. Decreasing λ should decrease this variance (Watkins, 1989).² However, decreasing λ can also increase the bias, and this slows optimal convergence. This trade-off between bias and variance (Watkins, 1989; Geman, Bienenstock, & Doursat, 1992) is characteristic of such algorithms.

5. Discussion

The application of powerful stochastic approximation algorithms in this article is neither the most general possible nor the most elegant. Kushner and Clark's (1978) work has been further extended (e.g., Kushner, 1984; Benveniste et al., 1990) and convergence with probability 1 could be proved under less restrictive conditions. Jaakkola, Jordan and Singh (personal communication) have a more refined and comprehensive proof that generalizes the Watkins et al. (1992) result that Q-learning converges with probability 1.

Nevertheless, we have shown that TD(λ) converges with probability one, under the standard stochastic convergence constraints on the learning rate given in equation (6) and the other stated conditions. The maximal rate of convergence of this algorithm is determined by the eigenvalues of the update matrix. This gives for TD(λ) a similar assurance that other approximation algorithms enjoy.

Although these theorems provide mathematical assurance for the convergence of TD(λ), the actual rate of convergence can often be too slow for real-world problems, especially for state spaces of high dimensionality. Developing good representations of states is of critical importance in achieving good performance with this as well as other classes of reinforcement learning algorithms.

Acknowledgments

We are most grateful to Chung Ming Kuan and two anonymous reviewers for making detailed comments on this article and improving the proof, and to Rich Sutton and Halbert White for starting us on the course that led to the article. Thanks also to Andy Barto, Vijaykumar Gullapalli, Satinder Singh, and Chris Watkins for invaluable discussions. Support was from SERC and the Howard Hughes Medical Institute.

Notes

1. Sutton used P_t^n for $V_n^\lambda(t_i)$.
2. For case 2b, M becomes irrelevant.

References

- Benveniste, A., Métivier, M., & Priouret, P. (1990). *Adaptive algorithms and stochastic approximation*. Berlin: Springer-Verlag.
- Dayan, P. (1992). The convergence of TD(λ) for general λ . *Machine Learning*, 8, 341-362.

- Geman, S., Bienenstock, E., & Doursat, R. (1991). Neural networks and the bias/variance dilemma. *Neural Computation*, 4, 1-58.
- Kuan, C.M., & White, H. (1990). *Recursive m-estimation, non-linear regression and neural network learning with dependent observations* (discussion paper). Department of Economics, University of California at San Diego.
- Kuan, C.M., & White, H. (1991). *Strong convergence of recursive m-estimators for models with dynamic latent variables* (discussion paper 91-05). Department of Economics, University of California at San Diego.
- Kushner, H.J. (1984). *Approximation and weak convergence methods for random processes, with applications to stochastic systems theory*. Cambridge, MA: MIT Press.
- Kushner, H.J., & Clark, D. (1978). *Stochastic approximation methods for constrained and unconstrained systems*. Berlin: Springer-Verlag.
- Robbins, H., & Monro, S. (1951). A stochastic approximation method. *Annals of Mathematical Statistics*, 22, 400-407.
- Ross, S. (1983). *Introduction to stochastic dynamic programming*. New York: Academic Press.
- Samuel, A.L. (1959). Some studies in machine learning using the game of checkers. *IBM Journal of Research and Development*, 3, 311-229.
- Sutton, R.S. (1984). *Temporal credit assignment in reinforcement learning*. Ph.D. thesis, Department of Computer Science, University of Massachusetts, Amherst, MA.
- Sutton, R.S. (1988). Learning to predict by the methods of temporal difference. *Machine Learning*, 3, 9-44.
- Sutton, R.S., & Barto, A.G. (1987). A temporal-difference model of classical conditioning. GTE Laboratories Report TR87-509-2. Waltham, MA.
- Tesauro, G. (1992). Practical issues in temporal difference learning. *Machine Learning*, 8, 257-278.
- Watkins, C.J.C.H. (1989). *Learning from delayed rewards*. Ph.D. thesis, King's College, University of Cambridge, England.
- Watkins, C.J.C.H., & Dayan, P. (1992). Q-learning. *Machine Learning*, 8, 279-292.

Received October 8, 1992

Accepted January 8, 1993

Final Manuscript March 17, 1993