

ON TEMPORAL DIFFERENCE ALGORITHMS FOR CONTINUOUS SYSTEMS

Alexandre Donzé

Verimag Laboratory

2 avenue de Vignate, 38100 Gières, France

Keywords: Optimal Control of Continuous Systems, Dynamic Programming, Optimal Value function, Reinforcement Learning, Temporal Differences Algorithms.

Abstract: This article proposes a general, intuitive and rigorous framework for designing temporal differences algorithms to solve optimal control problems in continuous time and space. Within this framework, we derive a version of the classical TD(λ) algorithm as well as a new TD algorithm which is similar, but designed to be more accurate and to converge as fast as TD(λ) for the best values of λ without the burden of finding these values.

1 INTRODUCTION

Dynamic programming (DP) is a powerful technique to compute optimal feedback controllers which was developed at the early stages of the field of optimal control (Bellman, 1957). It has been used since then in a wide variety of domains ranging from control of finite state machines, Markov Decision Process (MDP), to optimal control of continuous and hybrid systems. For all these problems, its goal is to compute a so-called *value function* (VF) mapping the states of the system to numbers that reflect the optimal cost to reach a certain objective from these states. Once this function computed, the optimal feedback control can be easily deduced from its values. The main drawback of DP lies in its computational cost. Computing the VF over the whole state space of the system can be prohibitively expensive if the number of states is large, particularly if it is infinite, e.g. in the case of continuous problems. This led some to develop *relaxed* DP techniques, e.g. (Rantzer, 2005).

In the reinforcement learning (RL) community, the temporal difference (TD) family of algorithms gained a particular interest after several spectacular experimental success, such as TD-Gammon (Tesauro, 1995), a player of backgammon that outperformed professional human players. These algorithms work by “learning” the VF while observing simulated trajectories. This family of algorithms is most typically designed for discrete MDPs, but various attempts were also made to adapt it to optimal control

of continuous systems, see (Sutton and Barto, 1998) for a survey. Recently, one of these works gave very promising and impressive results in the control of a continuous mechanical system with a large number of degrees of freedom and input variables (respectively 14 and 6) (Coulom, 2002). This motivated our interest in studying the applicability of TD algorithms to continuous optimal control problems. The main problem of this study is the lack of theoretical results such as convergence proofs, uniqueness of solutions and so on. For TD algorithms, results are generally available in their original context of MDPs and are not trivial to adapt to the continuous case. On the optimal control side, the most advanced theoretical framework available is that of functional analysis applied to the numerical solutions of Hamilton Jacobi Bellman (HJB) equations. To the best of our knowledge, the most complete and rigorous analysis of RL algorithm in this context is (Munos, 2000), but it does not include TD algorithms. Algorithms under study are numerical schemes usually used to solve partial differential equations: finite differences, finite elements, etc. Apart from their ease of implementation, one potential advantage of TD methods over other numerical schemes to compute the solution of HJB equations is that as it is based on simulation of real trajectories, it takes into account the structure of the reach set of the system. Hence, the approximation of the VF will then naturally converge faster in the ‘interesting’ zones inside the reach set. If the optimality is not the prior objective, the correct exploitation of this feature can

allow to design more rapidly a sub-optimal controller that meets the desired specifications.

This article attempts to provide a clarified framework for TD algorithms applied to continuous problems. Our hope is that by clearly identifying and isolating different problems arising in the design of TD algorithms, it will help to develop new and better TD algorithms, and to analyze their behavior more rigorously. Our first result in this direction is the derivation of a version of the classical TD(λ) algorithm along with a new TD algorithm, we called TD(\emptyset), designed to be more accurate and to converge as fast as TD(λ) for the best values of λ without the burden of finding these values. The works closest to this one are (Doya, 1996), (Doya, 2000) and (Coulom, 2002). Apart from the TD(\emptyset) algorithm, a contribution w.r.t. these works is a slightly higher level of abstraction, which allows an easier intuitive interpretation of the behaviour of TD algorithms and the role of their parameters, in particular the λ parameter. It is this interpretation of the role of λ that led to the design of TD(\emptyset) algorithm.

The paper is organized as follows. Section 2 recalls the principles of DP, first applied to the simple case of deterministic discrete transition systems, then adapted to continuous problems. Section 3 introduces our continuous framework for TD algorithms, the TD(λ) algorithm instantiated in this framework and our variant TD(\emptyset). Section 4 present some experimental results obtained on two test problems.

2 DYNAMIC PROGRAMMING

2.1 DP For Discrete Systems

We first consider in this section a purely discrete deterministic system with a state set \mathcal{X} , an input or action set \mathcal{U} and a transition map $\rightarrow: \mathcal{X} \times \mathcal{U} \mapsto \mathcal{X}$. For each pair (x, u) of state-action, a bounded, non-negative cost value $c(x, u)$ reflects the price of taking action u from state x . Given a state x_0 and an infinite input sequence $u = (u_n)$, we can thus define the cost-to-go or *value function* of the trajectory $x_0 \xrightarrow{u_0} x_1 \xrightarrow{u_1} \dots$ as:

$$V^u(x_0) = \sum_{n=0}^{\infty} \gamma^n c(x_n, u_n) \quad (1)$$

where γ is a discount factor lying strictly between 0 and 1, which prevents the infinite sum V^u from diverging. Since $V^u(x_0)$ represents the cost of applying the input sequence u from x_0 , the optimal control problem consists in finding a sequence u^* that minimizes this cost, i.e. $V^{u^*}(x_0) \leq V^u(x_0), \forall u \in \mathcal{U}^{\mathbb{N}}$. Note that $V^{u^*}(x_0)$ does not depend on u^* and thus will be noted $V^*(x_0)$. Instead of searching an optimal

sequence, DP aims at computing directly the optimal value function $V^*(x)$ for all $x \in \mathcal{X}$. It relies on *Bellman equation* satisfied by V^* which is easily deduced from (1):

$$V^*(x) = \min_{u, x \xrightarrow{u} x'} c(x, u) + \gamma V^*(x') \quad (2)$$

Once V^* has been computed, an optimal sequence u^* is obtained by solving the right hand side of Bellman equation, which gives the state feedback controller:

$$u^*(x) = \arg \min_{u, x \xrightarrow{u} x'} c(x, u) + \gamma V^*(x') \quad (3)$$

Since (2) is a fix point equation, V^* can be computed using fix point iterations. This gives the following *value iteration algorithm*. Convergence of this algo-

Algorithm 1 Value Iteration

- 1: **Init** $V^0, i \leftarrow 0$
 - 2: **repeat**
 - 3: **for all** $x \in \mathcal{X}$ **do**
 - 4: $V^{i+1}(x) \leftarrow \min_{u, x \xrightarrow{u} x'} c(x, u) + \gamma V^i(x')$
 - 5: **end for**
 - 6: $i \leftarrow i + 1$
 - 7: **until** V has converged
-

rithm is guaranteed by the discount factor γ which makes the iteration a contraction. Indeed, it is easy to see that

$$\|V_{i+1} - V_i\|_{\infty} \leq \gamma \|V_i - V_{i-1}\|_{\infty}$$

2.2 DP in Continuous Time and Space

In this section, we adapt the previous algorithm to the continuous case. Let \mathcal{X} and \mathcal{U} now be bounded subsets of \mathbb{R}^n and \mathbb{R}^m respectively, and $f: \mathcal{X} \times \mathcal{U} \mapsto \mathbb{R}^n$ the dynamics of the system, so that for all $t \geq 0$,

$$\dot{x}(t) = f(x(t), u(t)) \quad (4)$$

Cost and cost-to-go functions have their continuous counter parts: for each x and u , $c(x, u)$ is a non-negative scalar bounded by a constant $\bar{c} > 0$, and for any initial state x_0 and any input function $u(\cdot)$,

$$V^{u(\cdot)}(x_0) = \int_0^{\infty} e^{-s\gamma t} c(x(t), u(t)) dt$$

Where $x(t)$ and $u(t)$ satisfy (4) with $x(0) = x_0$. The optimal value function, still not depending on $u(\cdot)$, is such that $\forall x \in \mathcal{X}$,

$$V^*(x) = \min_{u(\cdot) \in \mathcal{U}^{\mathbb{R}^+}} \int_0^{\infty} e^{-s\gamma t} c(x(t), u(t)) dt$$

A straightforward way to adapt Algorithm 1 in this continuous context is first to discretize time and space and then to find an equivalent of the Bellman equation (2). Let us fix a time step Δt and a regular grid \mathcal{X}_ϵ of resolution ϵ covering \mathcal{X} . We evaluate V^* on each point of the grid and interpolate for any x in \mathcal{X} outside \mathcal{X}_ϵ (see figure 1).

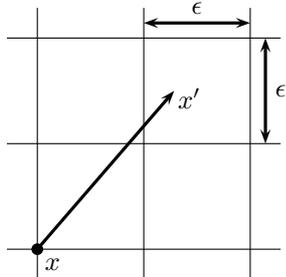


Figure 1: Simple grid of resolution ϵ . V^i is known at the grid points, e.g. at x , whereas $V^i(x')$ has to be interpolated

In order to get a transition function, one can classically perform an Euler integration of (4):

$$x \xrightarrow{u} x' \Leftrightarrow x' = x + \Delta x \quad (5)$$

where $\Delta x = f(x, u)\Delta t$

Moreover we can write :

$$\begin{aligned} V^u(x_0) &= \int_0^{\Delta t} e^{-s_\gamma t} c(x(t), u(t)) dt \\ &\quad + \int_{\Delta t}^{\infty} e^{-s_\gamma t} c(x(t), u(t)) dt \\ &\simeq c(x(0), u(0))\Delta t + e^{-s_\gamma \Delta t} V^u(x(0) + \Delta x) \\ &\simeq c(x_0, u_0)\Delta t + \gamma V^u(x_0 + \Delta x) \end{aligned} \quad (6)$$

Where we note $e^{-s_\gamma \Delta t} = \gamma$ then V^* satisfies

$$V^*(x) \simeq \min_{u, x \xrightarrow{u} x'} c(x, u)\Delta t + \gamma V^*(x') \quad (7)$$

which provides a decent equivalent of the discrete Bellman equation. From here, nothing prevents us from applying Algorithm 1 using (7) as a fix point iteration. The complete algorithm is given below.

Algorithm 2 Continuous Value Iteration

- 1: **Init** $V_\epsilon^0, i \leftarrow 0$
 - 2: **repeat**
 - 3: **for** all $x \in \mathcal{X}$ **do**
 - 4: $V_\epsilon^{i+1}(x) \leftarrow \min_{u, x \xrightarrow{u} x'} c(x, u)\Delta t + \gamma V_\epsilon^i(x')$
 - 5: **end for**
 - 6: $i \leftarrow i + 1$
 - 7: **until** V_ϵ has converged
-

Again, Algorithm 2 is guaranteed to converge for the same reasons as Algorithm 1 does, and, moreover,

in (Munos, 2000) it is proven that under certain regularity conditions, as ϵ and Δt tend toward zero, V_ϵ tends toward the exact optimal VF of the continuous problem.

2.3 Discussion

Limitations As presented above, the adaptation of discrete value iteration algorithm to the continuous case may seem simple. Unfortunately, in practice it may fail for several reasons: first, as it is a direct application of DP, it has the limitation that Bellman called originally the *curse of dimensionality* (Bellman, 1957) which expresses the fact that complexity of these algorithms is exponential in the dimension of the system. This makes them impractical for problems in high dimensions. Another limitation of Algorithm 2 that is more specific to the continuous case is the simplicity of the underlying numerical methods. Explicit Euler integration used in (5) is the simplest method to solve a differential equations and is known to have severe limitations. In particular, it is of order one and is sensitive to stiffness. Also, in (6), the rectangle method is used to approximate integral on interval $[0, \Delta t]$ which is also a method of order one. A lot can be done to benefit from more clever numerical schemes to solve (4) or to better approximate the integral in (6) using more advanced methods of numerical analysis. For instance, backward Euler integration is preferable to explicit Euler integration (Doya, 2000) and a cost-less trick to improve approximation (6) consists in, instead of considering that $e^{-s_\gamma t} c(x, u)$ is constant on $[0, \Delta t]$, only considering that $c(x, u)$ is constant and integrating formally the exponential, which gives, if $\gamma = e^{-s_\gamma \Delta t}$,

$$V^*(x) \simeq \min_{u, x \xrightarrow{u} x'} c(x, u) \frac{(1 - \gamma)}{s_\gamma} + \gamma V^*(x') \quad (8)$$

which is a more precise continuous equivalent of the discrete Bellman equation than (7) that can be used in Algorithm 2 (Coulom, 2002).

Function approximators The major reason for the curse of dimensionality is the discretization of the state space into a simple grid, an object whose size grows exponentially with the number of dimensions. Among classical alternative solutions are variable resolution discretization (Munos and Moore, 1999), sparse coarse coding (Sutton, 1996), or neural networks (Coulom, 2002), (Tesauro, 1995). It is worth recalling that all these techniques, including simple grids, approximate a function defined over (a bounded subset of) \mathcal{R}^n using a finite number of parameters $\{\omega_i, i \in \mathbb{N}\}$. In the case of simple grids, these parameters are associated to the exponentially many

grid points. In the case of neural networks, the parameters are the weights of the “neurons” in the network, whose number is dimension independent. Although methods using sophisticated function approximators like neural networks proved to be potentially very powerful in some practical cases, such as for swimmers in (Coulom, 2002), they offer in general no guarantee of convergence or optimality. In fact, function approximators are double-edged swords. On one hand, they provide *generalization*, that is, while the VF is updated at a point x , it is also updated in some neighborhood of x , or more precisely in all points affected by the change of the parameters of the approximator; the bad side is *interference* which is in fact an excessive generalization, e.g. if an update in x affects the VF in y in such a way that it destroys the benefits of a previous update in y . Understanding and controlling the effects of generalization, in particular in the case of non linear function approximators, is a very deep and complicated numerical analysis problem.

Hamilton-Jacobi-Bellman (HJB) Equation Here we present the connection with the HJB Equation. In fact, (7) and (8) can be seen as finite differences schemes approximating the HJB equation, that we can retrieve from those equations by dividing each term by Δt and making Δt going to zero, which leads to

$$\min_{u \in \mathcal{U}} \left(c(x, u) + \frac{\partial V^*}{\partial x} \cdot f(x, u) - s_\gamma V^*(x) \right) = 0 \quad (9)$$

From this definition of the HJB equation, we define the *Hamiltonian*:

$$\mathcal{H}(x) = \min_{u \in \mathcal{U}} \left(c(x, u) + \frac{\partial V^i}{\partial x} \cdot f(x, u) - s_\gamma V^i(x) \right) \quad (10)$$

Algorithms for computing V^* thus often try to minimize \mathcal{H} in order to find a solution to (9). The problem is that there might be an infinite number of *generalized* solutions of (9) i.e. functions that satisfy $\mathcal{H}(x) = 0$ *almost everywhere*, while being very far from the true value function V^* (Munos, 2000). Thus, if no guarantee is given by the algorithm other than the minimization of $\mathcal{H}(x)$, then no one can tell whether the computed solution is optimal or even near to the optimal VF.

(Munos, 2000) uses the theory of viscosity solutions to prove that Algorithm 2 actually converges toward the true value function but this result is limited to the case of grids.

3 TEMPORAL DIFFERENCE ALGORITHMS

DP algorithms presented so far compute the next estimation $V^{i+1}(x)$ based on $V^i(x)$ and $V^i(x')$ where

x' is in the spatial neighborhood of x . The slightly different point of view taken by temporal difference (TD) algorithms is that instead of considering that x' is the *neighbor in space* of x they consider it as its *neighbor in time* along a trajectory. This way, with the same algorithm, $V^{i+1}(x)$ would be updated from $V^i(x(0))$ and $V^i(x(\Delta t))$, where $x = x(0)$ and $x' = x(\Delta t)$ are then viewed in the context of a whole trajectory $(x(t))_{t>0}$. Then, pursuing in this spirit, it is clear that if $V^i(x(\Delta t))$ holds interesting information for the update of $V^i(x(0))$, then it is also the case for $V^i(x(t))$ for all $t > 0$. Henceforth, TD algorithms update their current approximation of the VF thanks to data extracted from complete trajectories.

3.1 General Framework

In section 2, we discussed the fact that DP algorithms were bound to different issues and problematics: function approximation, numerical analysis etc. This was the case for value iteration algorithms and their adaptation to continuous time and space, and this is also the case for TD algorithms. In this subsection, we propose a general, high level framework common to most of TD algorithms that allow us to separate those different issues, assuming that some are solved and focusing on the others. This high-level framework is represented by Algorithm 3.

Algorithm 3 Generic TD algorithm

```

1: Init  $V^0, i \leftarrow 0$ 
2: repeat
3:   Choose  $x_0 \in \mathcal{X}$ 
4:   Compute  $(x(t))_{t \in [0, T]}$  using  $u(\cdot)$ , starting from  $x_0$ 
5:   Compute update trace  $e(\cdot)$  from  $x(\cdot)$  and  $u(\cdot)$ 
6:   for all  $t \in \mathbb{R}^+$  do
7:      $V^{i+1}(x(t)) \leftarrow V^i(x(t)) + e(t)$ 
8:   end for
9:    $i \leftarrow i + 1$ 
10: until Stopping condition is true
    
```

To implement an instance of this generic algorithm, we need the non-trivial following elements:

- **A simulator** that can generate a trajectory of duration $T > 0$ $(x(t))_{t \in [0, T]}$ given an initial state x_0 and an input function u .
- **A function approximator** to be used to represent any function V from \mathcal{X} to \mathbb{R}^+ and allowing us to simply write $V(x)$ for any x in \mathcal{X} e.g. lines 7. A question arising here is, as already discussed in previous sections, whether to use simple grids, linear function approximators, or non linear function approximators such as neural networks etc.

- **An update function** which, given a function V , a state x and a quantity e , alters V^i into V^{i+1} so that $V^{i+1}(x)$ is equal to or near $V^i(x) + e$, which is written line 7 as an assignment statement. When parametrized function approximators are used, updating can be done e.g. by gradient descent on the parameters (Doya, 2000).
- **A policy**, or controller, which returns a control input at each time t , needed by the simulator to compute trajectories. In the case of *optimistic policy iteration* (Tsitsiklis, 2002), u is the feedback control obtained from V^i using the arg min of left hand side of (10). It is said to be *optimistic* because if V^i is the optimal VF, then this gives an optimal control input.
- **An update trace generator** used to extract information from the computed trajectory $x(\cdot)$ and to store it into the update trace $e(\cdot)$ (line 5). In the field of RL, $e(\cdot)$ is classically called the *eligibility trace* (Sutton and Barto, 1998).

Variants of TD algorithms differ in the way this last point is implemented, as discussed in the following sections.

3.2 Continuous TD(λ)

3.2.1 Formal algorithm

The objective of TD algorithms is to build an improved estimation of the VF based on the data of computed trajectories. The main additional information provided by a trajectory $x(\cdot)$ is the values of the cost function along its course. If we combine these values with the current estimation V^i , we can then build a new estimation for a given horizon $\tau > 0$, which is:

$$\tilde{V}_\tau(x_0) = \int_0^\tau e^{-s\gamma t} c(x, u) dt + e^{-s\gamma\tau} V^i(x(\tau)) \quad (11)$$

In other words, the estimation $\tilde{V}_\tau(x_0)$ relies on a portion of the trajectory of duration τ and on the old estimation $V^i(x(\tau))$ at the end of this portion. We remark that if V^i is already the optimal value function V^* , then $\tilde{V}_\tau = V^*$ for all $\tau > 0$ and this equality is, in fact, a generalization of the Bellman equation.

A question to answer is then, how to choose horizon τ ? In the case of TD(λ), the algorithm constructs a new estimation \tilde{V}_λ which is a combination of the $\tilde{V}_\tau(x_0)$ for all $\tau > 0$. In our continuous framework, we define for $s_\lambda > 0$:

$$\tilde{V}_\lambda^i(x_0) = \int_0^\infty s_\lambda e^{-s_\lambda\tau} \tilde{V}_\tau^i(x_0) d\tau \quad (12)$$

Several remarks about this definition are:

- Equation (12) means that \tilde{V}_λ^i is constructed from all $\tilde{V}_\tau^i, \tau > 0$, each of them contributing with an exponentially decreasing amplitude represented by the term $e^{-s_\lambda\tau}$. In other words, the further \tilde{V}_τ^i looks into the trajectory, the less it contributes to \tilde{V}_λ^i .
- This definition is sound in the sense that it can be seen as a convex combination of $\tilde{V}_\tau^i, \tau > 0$. In effect, it is true that $\int_0^\infty s_\lambda e^{-s_\lambda\tau} d\tau = 1$. Thus, if each \tilde{V}_τ^i is a sound estimation of V^* , then \tilde{V}_λ^i is a sound estimation of V^* .
- This definition is consistent with the original definition of TD(λ) algorithm (Sutton and Barto, 1998). In fact, by choosing a fixed time step $\Delta\tau$, assuming that \tilde{V}_τ^i is constant on $[\tau, \tau + \Delta\tau]$ and summing $s_\lambda e^{-s_\lambda\tau}$ on this interval, we get:

$$\tilde{V}_\lambda^i(x_0) \simeq (1 - \lambda) \sum_{k=0}^\infty \lambda^k \tilde{V}_{k\Delta\tau}^i(x_0) \quad (13)$$

$$\text{where } \lambda = e^{-s_\lambda\Delta\tau}$$

which is the usual discrete TD(λ) estimation built upon the k-step TD estimation (Tsitsiklis, 2002)

$$\begin{aligned} \tilde{V}_{k\Delta\tau}^i(x_0) &= \int_0^{k\Delta\tau} e^{-s\gamma t} c(x, u) dt + e^{-s\gamma k\Delta\tau} V^i(x(k\Delta\tau)) \\ &\simeq \sum_{j=0}^{k-1} \gamma^j c(x(j\Delta\tau), u(j\Delta\tau)) + \\ &\quad e^{-s\gamma k\Delta\tau} V^i(x(k\Delta\tau)) \\ &= \sum_{j=0}^{k-1} \gamma^j c(x_j, u_j) + \gamma^k V^i(x_k) \end{aligned} \quad (14)$$

where $\gamma = e^{-s\gamma\Delta\tau}$

3.2.2 Implementation

In practice, we can only compute trajectories with a finite duration T . As a consequence, \tilde{V}_τ^i can only be computed for $\tau \leq T$. In this context, (12) is replaced by

$$\begin{aligned} \tilde{V}_\lambda^i(x_0) &= \int_0^T s_\lambda e^{-s_\lambda\tau} \tilde{V}_\tau^i(x_0) d\tau + e^{-s_\lambda T} \tilde{V}_T^i(x_0) \\ &\simeq (1 - \lambda) \sum_{k=0}^{N-1} \lambda^k \tilde{V}_{k\Delta\tau}^i(x_0) + \lambda^N \tilde{V}_T^i(x_0) \end{aligned} \quad (15)$$

where $T = N\Delta\tau$

Combining (15) and (14), one can show that:

$$V^i(x_0) - \tilde{V}_\lambda^i(x_0) \simeq \sum_{j=0}^N \lambda^j \gamma^j \delta_j \quad (16)$$

where

$$\delta_j = c(x_j, u_j) + \gamma V^i(x_{j+1}) - V^i(x_j) \quad (17)$$

is what is usually referred to as the *temporal difference error* (Tsitsiklis, 2002). This expression provides an easy implementation given in Algorithm 4

which refines line 5 of Algorithm 3. There, we assume that the trajectory has already been computed with a fixed time step $\Delta t = \Delta\tau$, that is at times $t_j = j\Delta t, 0 \leq j \leq N + 1$.

Algorithm 4 TD(λ): Update traces

```

1: for  $j = 0$  to  $N$  do
2:    $\delta_j \leftarrow c(x_j, u_j) + \gamma V^i(x_{j+1}) - V^i(x_j)$ 
3: end for
4:  $e(t_N) \leftarrow \delta_N$ 
5: for  $k = 1$  to  $N$  do
6:    $e(t_{N-k}) \leftarrow e(t_k) + (\lambda\gamma)e(t_{N-k+1})$ 
7: end for
    
```

Note that the computation complexity of the update trace is in $O(N)$ where $N + 2$ is the number of computed points in the trajectory x . Thus the overall complexity of the algorithms depends only on the number of trajectories to be computed in order to obtain a good approximation of V^* .

3.2.3 Qualitative interpretation of TD(λ)

In the previous sections, we started by presenting TD(λ) as an algorithm that compute a new estimation of the VF using a trajectory and older estimations. This allowed us to provide a continuous formulation of this algorithm and an intuition of why it should converge towards the true VF. Then, using a fixed time step and numerical approximations for implementation purposes, we derived equation (16) and Algorithm (4). These provide another intuition of how this algorithm behaves.

The TD error δ_j (17) can be seen as a local error in $x(t_j)$ (in fact, it is an order one approximation of the Hamiltonian $\mathcal{H}(x(t_j))$). Thus, (16) means that the local error in $x(t_j)$ affects the global error estimated by TD(λ) in $x(t_0)$ with a shortness factor equal to $(\gamma\lambda)^j$. The values of λ ranges from 0 to 1 (in the continuous formulation, s_λ ranges from 0 to ∞). When $\lambda = 0$, then only local errors are considered, as in value iteration algorithms. When $\lambda = 1$ then errors along the trajectory are fully reported to x_0 . Intermediate values of λ are known to provide better results than these extreme values. But how to choose the best value for λ remains an open question in the general case. Our intuition, backed by experiments, tends to show that higher values of λ often produce larger updates, resulting in a faster convergence, at least at the beginning of the process, but also often return a coarser approximation of the VF, when it does not simply diverge. On the other hand, smaller values of λ result in a slower convergence but toward a more precise approximation. In the next section, we use this intuition to design a variant of TD(λ) that combines the qualities of high and low values of λ .

3.3 TD(\emptyset)

3.3.1 Idea

The new TD algorithm that we propose is based on the intuition about local and global updates presented in the previous section. Global updates are those performed by TD(1) whereas local updates are those used by TD(0). The idea is that global updates should only be used if they are “relevant”. In other cases, local updates should be performed. To decide whether the global update is “relevant” or not, we use a monotonicity argument: from a trajectory $x(\cdot)$, we compute an over-approximation $\bar{V}(x(t))$ of $V^*(x(t))$, along with the TD error $\delta(x(t))$. Then, if $\bar{V}(x(t))$ is less than the current estimation of the value function $V^i(x(t))$, it is chosen as a new estimation to be used for the next update. In the other case, $V^i(x(t)) + \delta(x(t))$ is used instead.

Let us first remark that since c is bounded and $s_\lambda > 0$, then $V^*(x) \leq V_{\max}, \forall x \in \mathcal{X}$, where

$$V_{\max} = \int_0^\infty e^{-s_\lambda t} \bar{c} dt = \frac{\bar{c}}{s_\lambda} \quad (18)$$

This upper bound of V^* represents the cost-to-go of an hypothetical forever worse trajectory, that is, a trajectory for which at every moment, the pair state input (x, u) has the worse cost \bar{c} . Thus, V_{\max} could be chosen as a trivial over-approximation of $V^*(x(t))$. In this case, our algorithm would be equivalent to TD(0). But if we assume that we compute a trajectory $x(\cdot)$ on the interval $[0, T]$, then a better over-approximation can be obtained:

$$\bar{V}(x_0) = \int_0^T e^{-s_\lambda t} c(x, u) dt + e^{-s_\lambda T} V_{\max} \quad (19)$$

It is easy to see that (19) is indeed an over-approximation of $V^*(x(0))$: it represents the cost of a trajectory that would begin as the computed trajectory $x(\cdot)$ on $[0, T]$, which is at best optimal on this finite interval, and then from T to ∞ it behaves as the ever worse trajectory. Thus, $\bar{V}(x_0) \geq V^*(x_0)$.

3.3.2 Continuous Implementation of TD(\emptyset)

In section 3.2.2 we fixed a time step Δt and gave a numerical scheme to compute estimation \tilde{V}_λ . This was useful in particular to make the connection with discrete TD(λ). In this section, we give a continuous implementation of TD(\emptyset) by showing that the computation of \bar{V} can be coupled with that of the trajectory $x(\cdot)$ in the solving of a unique dynamical system.

Let $x_V(t) = \int_0^t e^{-s_\lambda r} c(x, u) dr$. Then,

$$\bar{V}(x_0) = x_V(T) + e^{-s_\lambda T} V_{\max}$$

and more generally, for all $t \in [0, T]$,

$$\bar{V}(x(t)) = e^{s_\gamma t}(x_V(T) - x_V(t)) + e^{-s_\gamma(T-t)}V_{\max}$$

Where x_V can be computed together with x by solving the problem:

$$\begin{cases} \dot{x}(t) &= f(x(t), u(t)) \\ \dot{x}_V(t) &= e^{-s_\gamma t}c(x(t), u(t)) \\ x(0) = x_0 &, \quad x_V(0) = 0 \end{cases} \quad (20)$$

From there, we can give Algorithm 5.

Algorithm 5 TD(\emptyset)

- 1: **Init** V^0 with V_{\max} , $i \leftarrow 0$
 - 2: **repeat**
 - 3: Choose $x_0 \in \mathcal{X}$
 - 4: Compute $(x(t))_{t \in [0, T]}$ and $(x_V(t))_{t \in [0, T]}$ by solving (20)
 - 5: Compute $\bar{V}(x(\cdot))$ and $\delta(x(\cdot))$
 - 6: **for all** $t \in [0, T]$ **do**
 - 7: $V^{i+1}(x(t)) \leftarrow \min(\bar{V}(x(t)), V^i(x(t)) + \delta(x(t)))$
 - 8: **end for**
 - 9: **until** Stopping condition is **true**
-

Several remarks are worth mentioning:

- Initializing V^0 to V_{\max} imposes a certain monotonicity with respect to i . This monotonicity is not strict since when local updates are made, nothing prevents $\delta(x(t))$ from being positive, but during the first trajectories at least, as long as they pass through unexplored states, \bar{V}_T will be automatically better than the pessimistic initial value. Note also that if V^* is known at some special states (e.g. at stationary points), convergence can be fastened by initializing V^0 to these values.
- Computing x and x_V (and hence \bar{V}) together in the same ordinary differential equation (ODE) facilitates the use of variable time step size integration, the choice of which can be left to a specialized efficient ODE solver, and thus permits a better control of the numerical error at this level.

4 EXPERIMENTAL RESULTS

To compare the performances of TD(\emptyset) with TD(λ) for various values of λ , we implemented and applied Algorithms 4 and 5 to deterministic, continuous versions of two classical problems in the field of RL:

The continuous walker: a one dimensional problem in which a robot must exit the zone $[-1, 1]$ as quickly as possible. The system equations are given by

$$\dot{x} = \begin{cases} u & \text{if } -1 < x < 1 \\ \min(0, u) & \text{if } x = 1 \\ \max(0, u) & \text{if } x = -1 \end{cases} \quad (21)$$

The input u represents the speed of the walker in either direction. It is bounded in absolute value by 1. Thus, in the context of optimal control, we can restrict to a binary decision problem where $u = 1$ or $u = -1$. The cost function is:

$$c(x) = \begin{cases} 1 & \text{if } -1 < x < 1 \\ 0 & \text{if } x = -1 \text{ or } x = 1 \end{cases} \quad (22)$$

The swing-up of a pendulum with limited torque.

The goal is to drive a pendulum to the vertical position. The dynamics of the system are described by a second order non linear differential equation:

$$\begin{aligned} \ddot{\theta} &= -\mu\dot{\theta} + g \sin \theta + u \\ -10 < \dot{\theta} < 10, & \quad -3 < u < 3 \end{aligned} \quad (23)$$

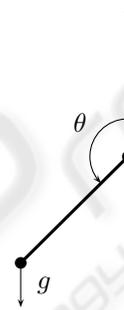


Figure 2: Pendulum

The cost function is $c(\theta) = (1 - \cos(\theta))$. It is thus minimal, equal to 0 when the pendulum is up ($\theta = 0$) and maximal, equal to 2, when the pendulum is down. The variable μ is a friction parameter and g the gravitational constant.

For both of these problems, we used regular grids and linear interpolation to represent values of V^i . We call a *sweep* a set of trajectories for which the set of initial states cover all points of the grid. For both problems, we then applied TD(λ) and TD(\emptyset) to such set of trajectories and after each sweep, we measured the 2-norm of the Hamiltonian error (noted $\|\delta\|_2$) over the state space. This allowed to observe the evolution of this error depending on the number of sweeps of the state space performed.

The results are given in figure 3 and 4. An interesting thing to observe is that for the simple problem of the walker, the value of λ that worked best in our experiments was 0, whereas for the pendulum, this value was more around 0.7. Not surprisingly, in the first case, we see that the error curve of TD(\emptyset) closely follows that of TD(0), even performing slightly better. For the second problem, TD(\emptyset) seems to perform better than TD(λ) for any value of λ . This leads to the very encouraging observation that TD(\emptyset) seems to behave at least as well as TD(λ) for the best value of λ , independently of this value and of the problem. Of course, more experiments are needed to validate this conclusion.

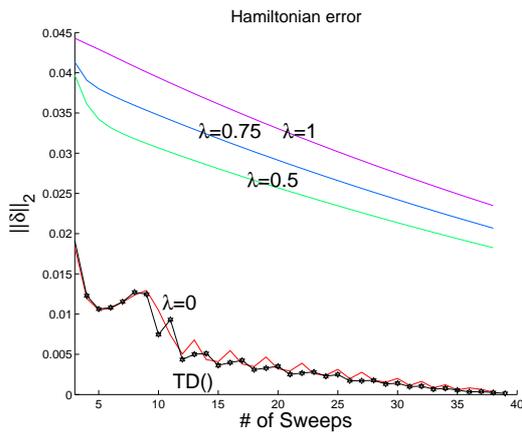


Figure 3: Hamiltonian error for walker

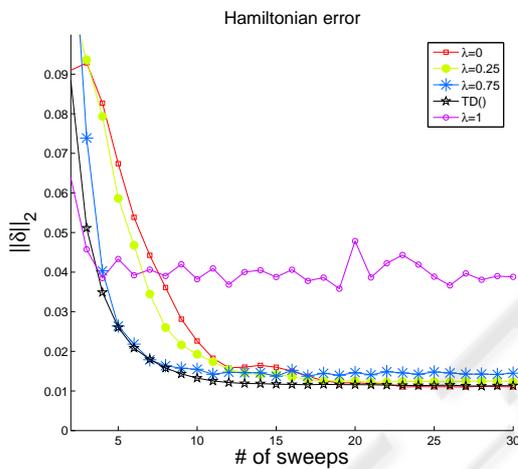


Figure 4: Hamiltonian error for pendul

5 CONCLUSIONS

We have presented a framework for the design of temporal differences algorithms for optimal control problems in continuous time and space. This framework is attractive as it clearly separates different issues related to the use of such algorithms, namely:

1. How to choose and compute trajectories - which initial state, and which choice of u - in order to ensure that the VF will be accurately and quickly computed in the interesting part of the state space?
2. How to represent the VF efficiently in order to break the curse of dimensionality?
3. How to update the VF based on simulated trajectories?

After describing the context of DP and our framework, we focused on question 3 and presented a continuous version of TD(λ) which was shown to be consistent after appropriate discretization with the origi-

nal discrete version of the algorithm. We then discussed its efficiency with respect to the parameter λ , and presented a variant, TD(\emptyset), that we found to be experimentally as efficient as TD(λ) for the best values of λ . Future work will include improvement of this algorithm and its application to higher dimensional problems, and consequently an investigation of the above mentioned Questions 1 and 2, a necessary step toward scalability.

ACKNOWLEDGMENTS

I am grateful to Oded Maler, Thao Dang and Eugene Asarin for helpful discussions, comments and suggestions. This work was partially supported by the European Community projects IST-2001-33520 CC (Control and Computation).

REFERENCES

Bellman, R. (1957). *Dynamic Programming*. Princeton University Press, Princeton, New Jersey.

Coulom, R. (2002). *Reinforcement Learning Using Neural Networks, with Applications to Motor Control*. PhD thesis, Institut National Polytechnique de Grenoble.

Doya, K. (1996). Temporal difference learning in continuous time and space. In Touretzky, D. S., Mozer, M. C., and Hasselmo, M. E., editors, *Advances in Neural Information Processing Systems*, volume 8, pages 1073–1079. The MIT Press.

Doya, K. (2000). Reinforcement learning in continuous time and space. *Neural Computation*, 12(1):219–245.

Munos, R. (2000). A study of reinforcement learning in the continuous case by the means of viscosity solutions. *Machine Learning*, 40(3):265–299.

Munos, R. and Moore, A. (1999). Variable resolution discretization for high-accuracy solutions of optimal control problems. In *International Joint Conference on Artificial Intelligence*.

Rantzer, A. (2005). On relaxed dynamic programming in switching systems. *IEE Proceedings special issue on Hybrid Systems*. Invited paper, to appear.

Sutton, R. S. (1996). Generalization in reinforcement learning: Successful examples using sparse coarse coding. In *Advances in Neural Information Processing Systems* 8, pages 1038–1044. MIT Press.

Sutton, R. S. and Barto, A. G. (1998). *Reinforcement Learning: An Introduction*. MIT Press, Cambridge, MA.

Tesauro, G. (1995). Temporal difference learning and TD-Gammon. *Communications of the ACM*, 38(3):58–68.

Tsitsiklis, J. N. (2002). On the convergence of optimistic policy iteration. *Journal of Machine Learning Research*, 3:59–72.