

Temporal Difference Learning with Neural Networks - Study of the Leakage Propagation Problem

Hugo Penedones*, Damien Vincent*, Hartmut Maennel,
Sylvain Gelly, Timothy Mann, André Barreto

DeepMind and Google Brain - London, UK and Zurich, Switzerland

{hugopen, damienv, hartmutm, sylvaingelly, timothymann, andrebarreto}@google.com

July 10, 2018

Abstract

Temporal-Difference learning (TD) [Sutton, 1988] with function approximation can converge to solutions that are worse than those obtained by Monte-Carlo regression, even in the simple case of on-policy evaluation. To increase our understanding of the problem, we investigate the issue of approximation errors in areas of sharp discontinuities of the value function being further propagated by bootstrap updates. We show empirical evidence of this *leakage propagation*, and show analytically that it must occur, in a simple Markov chain, when function approximation errors are present. For reversible policies, the result can be interpreted as the tension between two terms of the loss function that TD minimises, as recently described by [Ollivier, 2018]. We show that the upper bounds from [Tsitsiklis and Van Roy, 1997] hold, but they do not imply that leakage propagation occurs and under what conditions. Finally, we test whether the problem could be mitigated with a better state representation, and whether it can be learned in an unsupervised manner, without rewards or privileged information.

1 Introduction

Reinforcement Learning [Sutton and Barto, 1998, 2018] is a framework for studying sequential decision making

*equal contribution

processes. It deals with the setup where an agent interacts with an environment and at each time step it sees the current *state* (or just an *observation*) and takes one *action*, which will determine the immediate *reward* R_t it will get in the next state.

Typically, the agent is interested in maximising the *expected value* of the return G_t , which is often defined as the discounted sum of all rewards in the trajectory.

$$G_t \doteq R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots$$

A policy $\pi(a|s)$ is a probability distribution over actions, conditioned on the current state. A *state-value function*, v^π , is a function mapping from states to real numbers, defined as the expected return, when the agent starts in a given state s and then behaves according to the policy π .

$$v_\pi(s) \doteq \mathbb{E}_\pi[G_t | S_t = s]$$

The environment is formally represented as a Markov Decision Process (MDP), which is a 5-tuple $(\mathcal{S}, \mathcal{A}, P, r, \gamma)$ where \mathcal{S} is the set of states, \mathcal{A} the set of actions, $P(s'|s, a)$ the probability that action a in state s will lead to state s' , $r(s, a, s')$ is the reward function, and $\gamma \in [0, 1]$ is the discount factor.

1.1 On-Policy Evaluation

Estimating the state-value function for the policy that collected the data is the problem known as *on-policy evaluation*, in contrast with the (harder) problem of *off-policy evaluation*, in which one tries to estimate the value function of a different policy [Precup et al., 2001, Dann et al., 2014]. Note that in the *evaluation* or *prediction* problem, one does *not* attempt to find a better or optimal policy, which would be the role of *policy improvement* or *optimisation*. Instead, the goal is just to learn the value function of the current policy as accurately as possible. In this article we focus on the setup of *offline* or *batch* learning [Lange et al., 2012, Riedmiller, 2005], in which the dataset is pre-collected, rather than a setup where data is being collected *online* simultaneously with the learning process.

Assuming that our value function approximation \hat{v}_π is parameterised by \mathbf{w} (e.g. all neural network weights), and v_π is the ground truth value function, the goal would be to minimise the Mean-Squared Value Error (MSVE):

$$MSVE(\mathbf{w}) \doteq \|v_\pi(s) - \hat{v}_\pi(s, \mathbf{w})\|_\mu^2, \quad (1)$$

where μ is the stationary distribution over states, induced by the policy π and the environment dynamics.

Monte-Carlo regression (MC) and Temporal Difference Learning (TD) are the two main algorithms to estimate value functions. The main difference between them is easily visible in their update rules.

Monte-Carlo regression update:

$$\mathbf{w} \leftarrow \mathbf{w} + \epsilon[G_t - \hat{v}(S_t, \mathbf{w})]\nabla\hat{v}(S_t, \mathbf{w})$$

TD learning update:

$$\mathbf{w} \leftarrow \mathbf{w} + \epsilon[R_t + \gamma\hat{v}(S_{t+1}, \mathbf{w}) - \hat{v}(S_t, \mathbf{w})]\nabla\hat{v}(S_t, \mathbf{w})$$

where ϵ is the learning rate.

MC uses the actual return G_t , available at the end of the trajectory. On the other hand, TD approximates it with the immediate reward plus the discounted estimate of the next state, $R_t + \gamma\hat{v}(S_{t+1}, \mathbf{w})$, something that is known as *bootstrapping*.

1.2 The leakage propagation problem

In the tabular case, where the value function is approximated by a table with one entry per state, Temporal Dif-

ference (TD) learning converges to the true value function, see [Sutton and Barto, 2018] section 6.3. Not only does it converge to the true value function in the limit, like Monte Carlo regression (MC), it does so at a faster rate. Sutton and Barto [2018] explain that “Batch Monte Carlo methods always find the estimates that minimise mean-squared error on the training set, whereas batch TD always finds the estimates that would be exactly correct for the maximum-likelihood model of the Markov process”.

However, when combining TD with function approximation, such as the case of neural networks, this no longer holds. Parametric function approximators are typically forced to make approximation trade-offs. In some areas of the state space, there might be sharp discontinuities in the (true) value function, but a neural network function approximator might make a poor fit in those areas, due to shortage of data, or due to limited capacity. Such poor approximations can degrade the accuracy of the value function estimation, in both Monte-Carlo Regression and Temporal-Difference Learning, however with TD Learning the problem is aggravated by further propagating those bad estimates to nearby regions of the state space, after each bootstrap update. We call this problem *leakage propagation*¹.

It is well known that even in the limit the hypothesis that minimises the squared TD error may be worse than the MC solution by a factor of $\frac{1}{1-\gamma}$, where $\gamma \in [0, 1)$ is the discount factor of the MDP [Tsitsiklis and Van Roy, 1997]. A footnote in the same paper even mentions that a much tighter bound, with factor $\frac{1}{\sqrt{1-\gamma^2}}$, can be obtained.

These results are upper bounds on the error of the TD fixed point, but we would like to understand better how and under which conditions these errors actually occur. We expand on the relationship with this previous work on section 2.3.

¹Do not confuse the sharp discontinuities in *state space* with the sharp discontinuities of the value function when plotted as a function of *time*. Discontinuities in time are much more natural, and often happen after the collection of a big reward, in an environment where the next reward is far into the future. If the observations immediately before and after the collection of the reward are very different, this wouldn't cause a leakage propagation problem.

1.3 Collecting experimental evidence

Our first step is to demonstrate the occurrence of *leakage propagation* experimentally. For this, we define a toy environment in which we can expect sharp discontinuities in the value function to happen. In our setup, an agent navigates a two-dimensional space taking fixed-length steps towards some chosen direction. The state is represented by the real-valued (x, y) coordinates. There might be walls, which the agent can not cross, and one or more circle-shaped areas that give a positive reward.

The task is to estimate the state-value function of a policy in these environments as accurately as possible. Our training dataset is made of 100 independent trajectories, using a random policy initialised at random starting locations. We set the trajectory length limit to 2000 steps, but in the experiments we also show what would happen if the steps limit was lower.

For evaluation, we compare the learned value function, with a very precise estimate of the ground truth value function. The ground truth estimate was computed by simulating, in each point of an uniform grid, 1000 trajectories and averaging their Monte-Carlo returns. See Figures (1, 2 and 3) for visualisations of the layouts and the ground truth value functions.

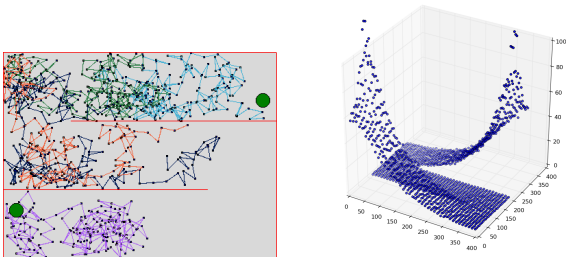


Figure 1: Random policy trajectories in a S-shaped layout with two reward zones (map 1), and its ground truth value function.

As we can see from figures (4 and 5), MC incurs some approximation error around the walls, but the problem is largely contained in those narrow regions, however in the case of TD, this leakage is propagated much further to other regions.

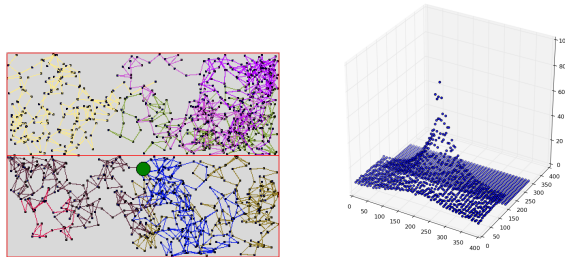


Figure 2: Random policy trajectories in a layout with two separate rooms (map 2), having a reward zone only in the lower room, and the ground truth value function.

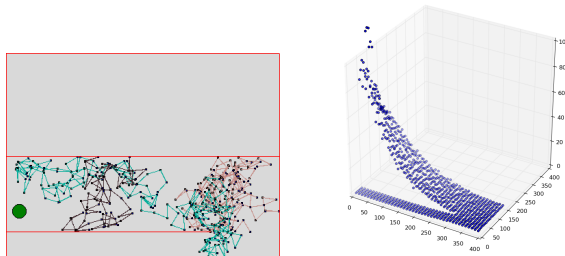


Figure 3: Left: Random policy trajectories in U-shape layout (map 3). The upper room was not used to generate data (only for padding). Right: the associated ground truth value function.

2 A more formal understanding

Now that we have an intuition and demonstrated the problem in practice, we want to understand it more formally. A recent paper by Ollivier [2018] describing the loss function that approximate TD minimises can shed some light. The result in that paper is restricted to *reversible* policies, where the probability of a transition and the reverse transition are related by $\mu(s)P(s, s') = \mu(s')P(s', s) \forall s, s'$, but it happens to match our navigation task with random policies.

Ollivier shows that if the transition probability P is reversible with respect to its stationary distribution μ , then TD Learning is doing gradient descent over the loss:

$$\gamma \|\hat{v} - v\|_{Dir}^2 + (1 - \gamma) \|\hat{v} - v\|_{\mu}^2 \quad (2)$$

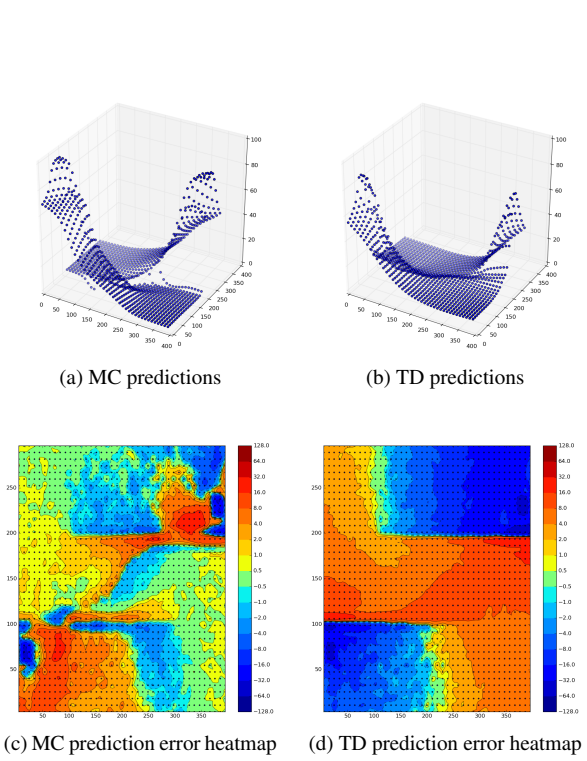


Figure 4: Predictions of MC (left) and TD (right) and error heatmaps compared to ground truth on the S-shaped environment (map 1). Notice that when using TD learning, the middle corridor in the environment, which should have near-zero value, is now being over-estimated, due to leakage propagation from across both walls.

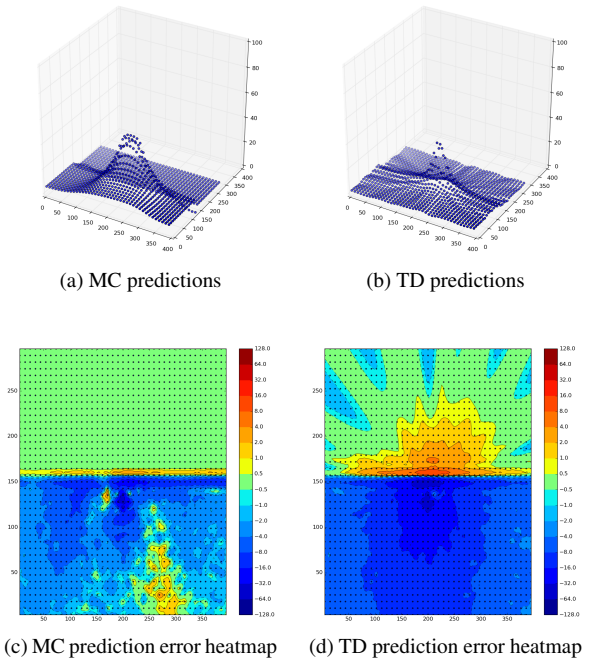


Figure 5: Predictions of MC (left) and TD (right) and error heat-maps compared to ground truth on the environment with two separate rooms (map 2). Notice that when using TD learning, the upper part of the environment, which should have value zero, is now being over-estimated, due to leakage propagation from across the wall. The side that contains the reward zone, also seems under-estimated, suggesting a possible “negative” leakage from the zero area. In MC the leakage effect is much smaller and it is restricted to regions immediately close to the wall.

where γ is the discount factor of the MDP and the L_2 norm weighted by the stationary distribution, is just:

$$\|\hat{v} - v\|_{\mu}^2 \doteq \sum_s \mu(s) [\hat{v}(s) - v(s)]^2$$

The Dirichlet norm $\|\hat{v} - v\|_{Dir}^2$ is defined as:

$$\frac{1}{2} \sum_{s,s'} \mu(s) P(s, s') [(\hat{v}(s') - v(s')) - (\hat{v}(s) - v(s))]^2$$

How does this relate to the leakage propagation problem? Well, note that the second term of the loss in equation (2) optimises for the mean square value error that we care about, but the first term optimises for a different norm. In particular, the Dirichlet norm cares about the *differences* in values between consecutive states in our approximation, to be the same as the differences in the true value function. Or equivalently, that the approximation error is the same in the current and next state. This implies that the Dirichlet norm does not distinguish among functions that are shifted by an additive constant. More importantly, as we will see in an example below, under some circumstances the two terms might be under some tension and a solution will compromise between the two, depending on the value of the discount factor γ .

2.1 Constraints on sharpness and Leakage Propagation

As an example, imagine now a scenario in which the state space is the set of non-negative integers \mathbb{N} , all transitions go from one integer $s \in \mathbb{N}$ to either $s + 1$ or $s - 1$, except for $s = 0$, in which case the only possible transition goes to $s = 1$. Assume further that all rewards are 0, so the true value function is 0 everywhere, and our policy chooses for $s > 0$ the transition $s \rightarrow s + 1$ with probability p and $s \rightarrow s - 1$ with probability $q = 1 - p$.

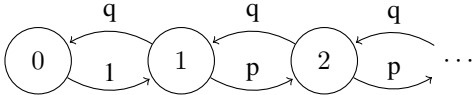


Figure 6: Simple infinite Markov chain with transition probabilities p and q .

We will assume that $0 < p < 0.5$, since for $p \geq 0.5$ we would not have a stationary distribution, the process would “wander off to infinity”. For $p < 0.5$ the stationary distribution μ on $S = \mathbb{N}$ has probabilities proportional to

$$q, 1, \frac{p}{q}, \frac{p^2}{q^2}, \frac{p^3}{q^3}, \dots \quad (3)$$

We assume also that the class of permissible approximations has as (only) condition that $v(0) = \alpha$ for some $\alpha > 0$. This is a very idealised case of a wall like in Figure (2) - imagine there are also states for integers $s < 0$, and they have positive value functions, but they cannot be reached from $s \geq 0$, and our function class somehow fixes the value at $s = 0$ to $v(0) = \alpha$ when the correct values at $s < 0$ are given.

Assuming we start with some estimate for our value function $v : \mathbb{N} \rightarrow \mathbb{R}$, what is the solution to which TD will converge? Without restrictions, there are two possibilities, however we will assume that we start from a bounded estimate (it would also be enough to assume $v \in L^2(\mathbb{N}, \mu)$, which is a weaker condition).

We will first analyse this directly, and later see how this can also be interpreted as minimum of (2). For compactness we abbreviate $v(s)$ as v_s .

The average TD update at a state $s > 0$ with learning rate ϵ would be

$$v_s \leftarrow v_s + \epsilon(p\gamma \cdot v_{s+1} + q\gamma \cdot v_{s-1} - v_s) \quad (4)$$

By rewriting the right hand side as

$$(1 - \epsilon) \cdot v_s + \epsilon \cdot \gamma(p \cdot v_{s+1} + q \cdot v_{s-1})$$

we see that if v is bounded by $|v(s)| < B$ for some $B \in \mathbb{R}$, the same is true for the updated value function. So the updates can only converge to a bounded solution (and it is also true that the $L_2(\mathbb{N}, \mu)$ norm does not increase, which we would use if we only assume we started from a $v \in L_2(\mathbb{N}, \mu)$.)

The update (4) leaves v_s for the state $s > 0$ invariant iff

$$v_{s+1} - \frac{1}{p\gamma} v_s + \frac{q}{p} \cdot v_{s-1} = 0 \quad (5)$$

Finding the roots of the characteristic equation:

$$r^2 - \frac{1}{p\gamma} r + \frac{q}{p} = 0 \quad (6)$$

we get the two solutions:

$$r_{1,2} = \frac{1 \pm \sqrt{1 - 4pq\gamma^2}}{2p\gamma}$$

which are real and positive for any $\gamma < 1$ since $4pq < 1$. Let r_1 be the smaller one of them. Since $\gamma = (p+q) \cdot \gamma < 1$, we have $q\gamma < 1 - p\gamma$ and

$$1 - 4pq\gamma^2 > 1 - 4q\gamma(1 - p\gamma) = (1 - 2p\gamma)^2$$

from which it follows that

$$1 - \sqrt{1 - 4pq\gamma^2} < 2p\gamma$$

and hence $r_1 < 1$. This computation also shows that

$$r_1 \rightarrow 1 \quad \text{for } \gamma \rightarrow 1 \quad (7)$$

Since (6) shows that $r_1 \cdot r_2 = q/p$, we have

$$0 < r_1 < 1 < q/p < r_2$$

We can write the general solution for the recurrence as:

$$v_s = c_1 r_1^s + c_2 r_2^s$$

Since v is bounded (in fact, $v \in L^2(\mathbb{N}, \mu)$ is enough), we must have $c_2 = 0$. On the other hand $s = 0$ gives $c_1 + c_2 = \alpha$, and therefore:

$$v_s = \alpha r_1^s \quad \text{for } s = 0, 1, 2, \dots$$

In other words, in the optimal solution, the value function v that TD finds decays exponentially, as states move away from the constrained state, see Figure (7).

Because of (7), the closer the discount factor γ is to 1, the slower the decay will be, which means more leakage propagation. This makes intuitive sense, because it means we are bootstrapping more heavily. If $\gamma \rightarrow 1$ in fact, we get that $v_s \rightarrow \alpha$, which is a totally flat, grossly overestimating region. When $\gamma \rightarrow 0$, $r_1 \approx q\gamma \rightarrow 0$ and therefore $v_s \rightarrow 0$ for $s > 0$, which means that there is no leakage propagation and the mean-squared error is minimal. This corresponds to the non-interesting situation of just predicting immediate rewards, however.

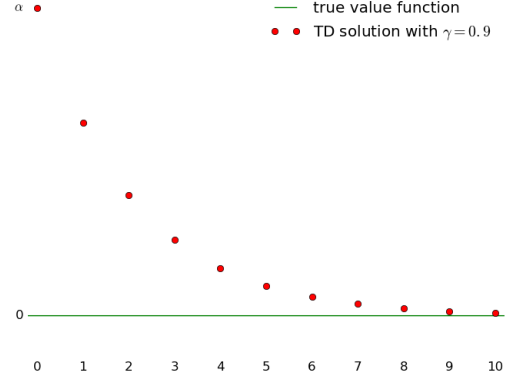


Figure 7: Plot of the analytic solution for TD in a simple discrete 1-d problem. Assuming the true value function (green) is zero, but we are forced to make an approximation mistake in state $s = 0$, by setting its value to $\alpha > 0$. After that mistake is incurred, TD will not make a sharp correction to the value function, but rather decay it exponentially towards zero (red). This formally demonstrates the leakage propagation effect with approximate TD.

2.2 Relation to Ollivier [2018]

We now relate this to (2). As we can see directly from (3), our simple example satisfies the reversibility condition, as $\mu(s_t)P(s_t, s_{t+1}) = c \frac{p^t}{q^{t-1}} = \mu(s_{t+1})P(s_{t+1}, s_t)$, where c is the constant needed to get the stationary distribution probabilities to sum to 1. Since the true value function in our example is $v_t = 0$ for all t and noting that $p\mu(s_{t-1}) = q\mu(s_t)$ for $t > 1$, we sum all over transitions to get the Dirichlet term, $\|\tilde{v} - v\|_{Dir}^2$, of the loss function (2):

$$\mu(0)\alpha^2 + \sum_{t=1}^{\infty} \frac{1}{2} [p\mu(t)(\tilde{v}_{t+1} - \tilde{v}_t)^2 + q\mu(t)(\tilde{v}_t - \tilde{v}_{t-1})^2]$$

while the L_2 term $\|\tilde{v} - v\|_{\mu}^2$ is just $\sum_{t=0}^{\infty} \mu(s_t) \tilde{v}_t^2$. If we solve for $\frac{\partial}{\partial \tilde{v}_i} (\gamma \|\tilde{v} - v\|_{Dir}^2 + (1 - \gamma) \|\tilde{v} - v\|_{\mu}^2) = 0$ to find its minima, we recover the recurrence previously mentioned in (5).

In the paper Ollivier [2018] a finite state space was assumed, but the computation of the gradient of (2) is a

purely algebraic transformation, so it is still valid for our case if everything converges.

Now we can understand the leakage as a consequence of the tension between the two components of the loss: according to the L_2 norm one would set all $v_s = 0$ for $s \neq 0$, achieving the minimal first part of the cost $\gamma \cdot |v|_\mu^2 = \gamma \cdot \alpha^2$. However, according to the second part $(1-\gamma) \cdot |v|_{Dir}^2$, the best would be to actually set all $v_s = \alpha$, assuring that the estimates are “as flat” as the true value function in that region, and incurring zero cost for the Dirichlet norm part. Given that the full loss is a mixture of the two convex loss functions which have no minimum in common, neither of the minima of these two components can be a minimum for the mixture, so in particular, leakage must occur.

2.3 Relation to Tsitsiklis and Van Roy [1997]

We can look at this example also from the point of view of Tsitsiklis and Van Roy [1997]. This paper defines operators $T^{(0)}$ and Π on $L^2(S, \mu)$, where μ is the stable distribution on S , Π is the orthogonal projection on the function class, and $T^{(0)}$ is the average TD(0) update (see below).

Our case does not directly satisfy the assumptions in Tsitsiklis and Van Roy [1997], since our function space is not a vector space, but an affine space: For

$$a \doteq (\alpha, 0, 0, \dots)$$

we can write it as $a + W$ with the vector space

$$W \doteq \{(v_0, v_1, \dots) \in L^2(\mathbb{N}, \mu) \mid v_0 = 0\}$$

We will consider both the orthogonal projection Π onto $a + W$ and the orthogonal projection Π_0 onto W :

$$\begin{aligned} \Pi_0(v_0, v_1, \dots) &\doteq (0, v_1, v_2, \dots) \\ \Pi(v_0, v_1, \dots) &\doteq (\alpha, v_1, v_2, \dots) \end{aligned}$$

Since we do not consider other $T^{(\lambda)}$, we will just write T for the operator $T^{(0)}$ on $L^2(S, \mu)$, it is given by

$$Tv(s) \doteq \begin{cases} \gamma \cdot v(1) & \text{for } s = 0 \\ \gamma \cdot (p \cdot v(s+1) + q \cdot v(s-1)) & \text{for } s > 0 \end{cases}$$

So our update (4) leaves \tilde{v} invariant if and only if it is a fixed point of ΠT , as in Tsitsiklis and Van Roy [1997].

On the other hand the ideal solution $v^* = (0, 0, \dots)$ is a fixed point of T .

To get the “ $1/(1-\gamma)$ bound” of this paper, we adjust their computation to our situation:

$$\begin{aligned} |\tilde{v} - v^*| &\leq |\tilde{v} - \Pi v^*| + |\Pi v^* - v^*| & (8) \\ &= |\Pi T \tilde{v} - \Pi T v^*| + |\Pi v^* - v^*| \\ &= |\Pi_0 T(\tilde{v} - v^*)| + |\Pi v^* - v^*| \end{aligned}$$

The proof of Lemma 1 (p. 680) in Tsitsiklis and Van Roy [1997] then shows that $|Tv| \leq \gamma|v|$ for all $v \in L^2(S, \mu)$, and since obviously $|\Pi_0 v| \leq |v|$, also $\Pi_0 T$ is a contraction with factor $\leq \gamma$, from which we get

$$|\tilde{v} - v^*| \leq \frac{1}{1-\gamma} |\Pi v^* - v^*| \quad (9)$$

In fact, we can replace (8) by the equality

$$|\tilde{v} - v^*|^2 = |\tilde{v} - \Pi v^*|^2 + |\Pi v^* - v^*|^2$$

because $\tilde{v} - \Pi v^* \in W$ and $\Pi v^* - v^*$ is orthogonal to W . This gives (in otherwise the same way) the sharper bound

$$|\tilde{v} - v^*| \leq \frac{1}{\sqrt{1-\gamma^2}} |\Pi v^* - v^*| \quad (10)$$

which is mentioned in the footnote of Tsitsiklis and Van Roy [1997].

The above proofs used the facts

- $Tv^* = v^*$
- $\Pi T \tilde{v} = \tilde{v}$
- $\Pi_0 T$ is contracting with factor γ

about T, \tilde{v}, v^* . To see whether the factor $1/\sqrt{1-\gamma^2}$ that occurred in this proof is related to our “leakage observation”, consider another operator / estimate pair

$$\begin{aligned} T'v(s) &\doteq \gamma \cdot v(s) \\ \tilde{v}' &\doteq (\alpha, 0, 0, 0, \dots) \end{aligned}$$

which has the same properties (even though \tilde{v}' is without “leakage”). The proof gives the same bound, but in fact we even have

$$|\tilde{v}' - v^*| = |\Pi v^* - v^*| \quad (11)$$

The reason that $1/\sqrt{1-\gamma^2}$ appeared in (10) was that we only know that $\Pi_0 T$ is contracting with factor γ , so without knowledge about \tilde{v} we only can say that

$$|\Pi_0 T(\tilde{v} - v^*)| \leq \gamma \cdot |\tilde{v} - v^*|$$

from which (10) followed. However, to prove that the constant in (10) cannot be lowered from $1/\sqrt{1-\gamma^2}$ to 1 (as in (11)), we would need to give a lower bound to $|\Pi_0 T(\tilde{v} - v^*)|$. In the T', \tilde{v}' example “without leakage” this expression is 0; if we do have leakage (i.e. values $\tilde{v}_s > 0$ for $s > 0$) in \tilde{v} , then this expression is > 0 . So in a sense the possibility of leakage is indeed responsible for a factor > 1 in (10) in this example, but whether we have leakage or not does not seem to follow from the general arguments in the above proof.

3 Better state representations

We saw that the existence of sharp discontinuities in the state-value function causes TD to miss-behave in combination with function approximation. In this section, we investigate what would happen if we had a better state representation and whether we can hope to learn it from the data, without privileged information and without using the reward signal.

We start by noting that if two states in our trajectories have very similar input features, but it takes a large number of steps to connect them, following the policy, we do not want the function approximator to extrapolate from one to the other, by continuity or smoothness. Instead, we want these two states to be pushed apart, in some initial embedding space, so that the naïve extrapolation does not happen by default. However, if two states occur just after each other in the same trajectory, we assume it is safer to keep them close to each other in the embedding space. Even if they have a big intermediate reward in between, there is already enough incentive from the TD error to differentiate between the two.

In order to compare different state representations, hand-designed or learned, we start by defining a two stage architecture, composed of an embedding network followed by a value network.

3.1 Embedding and value networks

Our baseline neural network is a Multi-Layer Perceptron (MLP) which takes as input the state representation s and outputs a single scalar $\hat{v}(s, \mathbf{w})$. To experiment with different representations, keeping the model capacity constant, we split the MLP into two sets of layers: the embedding layers, with weights \mathbf{w}_e and the value layers, with weights \mathbf{w}_v . The final output is therefore a composition of the two functions: $\hat{v}(s, \mathbf{w}) \equiv g(f(s, \mathbf{w}_e), \mathbf{w}_v)$.

The training procedure is now divided into two stages: first learn the embedding network f and freeze its weights \mathbf{w}_e ; and *then* learn the top network layers g to approximate the value function.

In the next subsections we suggest some possible directions to create such representations. We start by illustrating the opportunity gap with a hand-designed transformation and then try out two unsupervised ways to learn the intermediate representation from data.

3.2 Hand-designed Oracle Embedding

As a simple sanity check, we hand-craft a transformation that separates nearby points across the walls perfectly, assuming privileged knowledge about the map layout. If this does not work, there is little hope for learning such representation without supervision.

Figure (8) shows the transformation of the input space that we defined for the S-shaped map layout. As we can see in the results from Figure (9), the solution found by TD becomes indeed much sharper and the MSVE is greatly reduced.

This is encouraging, but we have assumed privileged information to create the transformation. What if we could do it in an unsupervised manner, purely from the existing trajectories data?

3.3 Time-proximity Embeddings

There have been several recent works [Sermanet et al., 2017, Savinov et al., 2018, Aytar et al., 2018] around the idea of predicting whether two states or observations, are temporally close to each other, based on a dataset of trajectories generated by a policy. To test whether these could be used to induce a good intermediate representation, we have implemented such a temporal proximity

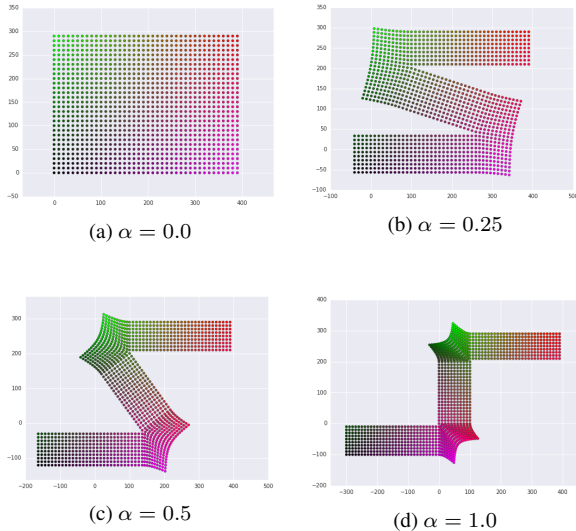


Figure 8: Hand-designed “oracle” embedding with different degrees of separation across the labyrinth walls for the S-shaped layout (map 1). This assumes privileged knowledge about the exact location of the walls. The amount of separation of the independent “corridors” in this layout is controlled by varying the α parameter from 0 to 1.

classifier. We defined multiple bins for different ranges of (discounted) time intervals into the future and trained it with self-supervised information from the trajectories. More specifically, the input of the classifier is a pair of states from the same episode Δt steps apart, where Δt is sampled according to a geometric distribution of discount γ (the same as in the MDP), sampling more often frames close in time. The K bins are created such that they are equally likely, i.e. the upper bound of bin # k is then given by:

$$T_k = \frac{\ln(1-p)}{\ln \gamma} \quad \text{with} \quad p = \frac{k}{K}$$

The last bin is reserved for states that are very far apart in time. To increase the diversity in training, we also sample pairs of states coming each from a different episode. Two states sampled this way are assumed to be far apart in time with high probability, thus we label them as belonging to the last bin # K . The estimated distribution of the time bins for a pair of states s_1 and s_2 is given by

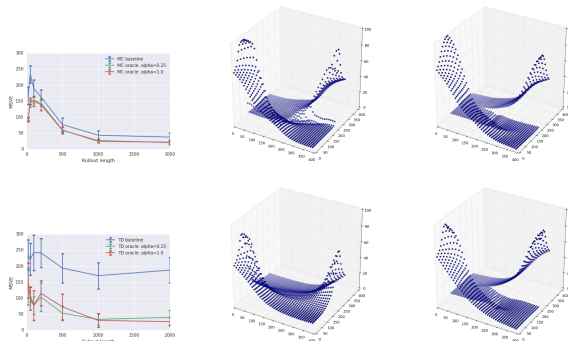


Figure 9: Left: MSVE with and without the oracle embedding. Estimated value function without the oracle (middle) and with the oracle (right). Top: MC, Bottom: TD. Note the sharp discontinuities in the walls regions, absence of leakage, and lower MSVE, especially in the TD case.

$g_{\mathbf{w}_c}(f(s_1, \mathbf{w}_e), f(s_2, \mathbf{w}_e))$ where both f and g are functions parameterised with a neural network with weights \mathbf{w}_e and \mathbf{w}_c , respectively. After training is complete we only keep $f(s, \mathbf{w}_e)$, which defines the embedding function.

The results in Figure (11) show that by taking the intermediate layer representation from the time proximity classifier, and then performing TD to estimate the value function on top of it, the overall Mean Square Value Error tends to be lower. Naturally, the improvement is not as drastic as with the Oracle transformation, but it achieves a significant gain without privileged knowledge. We now investigate an alternative way of learning an intermediate representation, relying on a well-established concept from the RL literature.

3.4 Successor Features

The successor features concept was introduced in Kulkaarni et al. [2016] and Barreto et al. [2017] as an extension of the successor representation [Dayan, 1993] for a continuous state space. The successor representation $\Psi(s, s')$ between two states s and s' is defined in the tabular case as the expected discounted time to reach s' from s when

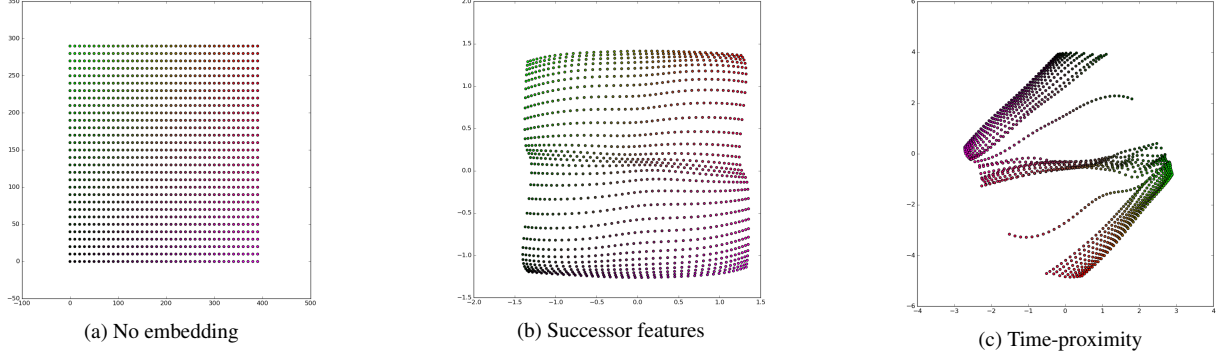


Figure 10: Warping of the original 2-dimensional space on the S-shaped layout (map 1), using the embedding networks trained with the successor features or the time-proximity classifier.

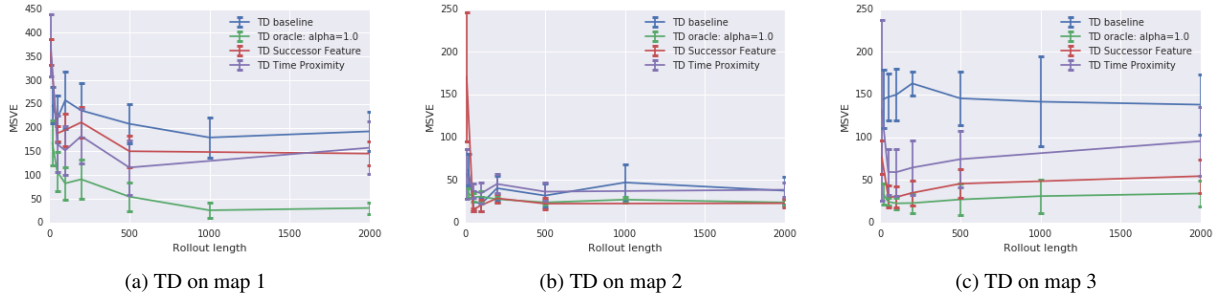


Figure 11: MSVE of value function estimates learnt on top of the time proximity and the successor feature embeddings, with varying amounts of training data available (changing trajectory lengths). The hand-designed oracle embedding curves are also shown, for comparison.

following policy π :

$$\Psi_{\pi, s'}(s) = \mathbb{E}_{\pi} \left[\sum_t \gamma^t \mathbb{1}_{S_t = s'} | S_0 = s \right]$$

The successor representation is suitable as an intermediate embedding, as the value function can be obtained directly from this embedding through a linear operation:

$$v_{\pi}(s) = \sum_{s' \in S} \Psi_{\pi, s'}(s) r(s')$$

In the continuous setting, the infinite family of functions $\mathbb{1}_{s'}$ is not suitable anymore, all the expectations would be 0. The successor features use a finite set of functions ϕ_k instead:

$$\Psi_{\pi, k}(s) = \mathbb{E}_{\pi} \left[\sum_t \gamma^t \phi_k(S_t) | S_0 = s \right].$$

In the labyrinth toy environment, we use 2 features: $\phi_0(s) = x$ and $\phi_1(s) = y$. The successor features now define the average discounted position of the agent in the

labyrinth. An example is given by Figure (12): one can notice that an initial position that is close to the wall can only get further away from the wall along the trajectories. Hence the successor features of two points in opposite sides of a wall will be significantly different.

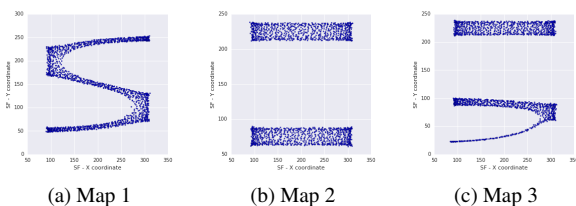


Figure 12: Successor Features targets with $\gamma = 0.99$

As we can see from Figure (11), the successor representation for states allows TD to find a state-value function with significantly lower Mean-Squared Value Error. In our setup it is *crucial* that the successor features are first learned with Monte Carlo targets and *then* used by Temporal Difference to learn the value function. If we would learn the successor features with TD, we would immediately face the leakage propagation problem which we intend to avoid.

4 Experimental details

Below we show the choices of hyper-parameters related to the environments, data generation, neural network architecture and the training procedure.

Data	
Environments size (W x H)	400 x 300
Environment discount factor	0.99
Environment reward areas radius	10
Environment reward	+30
Agent number of actions	360 (angles)
Agent step size	20
Trajectory max length	2000
Number of training trajectories	100

Multi-layer Perceptron architecture

Embedding layers sizes	[20, 20, 20, 2]
Value layers sizes	[30, 30, 1]
Non-linearities	tanh

Training hyper-parameters

Mini-batch size	32
Optimiser algorithm	Adam
Learning rate	0.001
Beta1	0.9
Beta2	0.999
Epsilon	1e-08
Embedding training steps	40000
Value training steps	40000

Note that we fixed the output of the embedding network to always be 2-dimensional, so that we could easily visualise it. This is also the case for the hand-designed embedding. To generate the error bars in our plots, we run each configuration with 20 different random seeds.

5 Related Work

Reinforcement Learning books by Sutton and Barto [1998, 2018] and Szepesvári [2010], cover the general ideas behind on-policy evaluation via TD and MC, with and without function approximation. In the following sub-sections we point to more specific research articles in the literature.

5.1 Temporal Difference Learning

TD Learning with function approximation has been analysed by Tsitsiklis and Van Roy [1997], with convergence results for the on-policy linear case, and showcasing divergence examples with off-policy data. Baird [1995] introduces convergent algorithms that minimise the Bellman residual directly, but they can be slow. More concerning, if the environment is stochastic, they would require two independent transitions from the same state. This is often unrealistic, most notably when dealing with very large or continuous state spaces. Baird [1995] also describes how a value function with some parameter sharing, can propagate information “the wrong direction”, by not respecting temporal causality. Singh et al. [1995] deals with aggregation of states, in a soft-clustering setup,

which does not cover other function approximators of interest, like neural networks.

Sutton et al. [2009b,a] introduce two modified algorithms for TD with linear function approximation that provably converge in the off-policy setting. Adam and White [2016] provide a comparison of linear TD methods and recommendations on which variant to pick depending on the problem constraints. Bhatnagar et al. [2009] introduces a version of TD learning that converges under the off-policy setting with a large class of smooth function approximators, such as neural networks. However, the solutions found might have poor Mean-squared Value Error (MSVE).

In a recent paper, Ollivier [2018] shows that for the particular case of *reversible policies*, TD does gradient descent of a loss function that includes a term with the Dirichlet norm $\|\hat{v} - v\|_{Dir}^2$. This formulation provides a more natural way to think about the fixed point solution for TD, for reversible policies, as a tension between two competing objectives.

Great part of the RL literature focus on online learning, however, there is also a set of works devoted to the offline, or *batch* case. In this setup, all the data is assumed to be pre-collected and available at training time. In analogy with the well established least-squares regression, Bradtke and Barto [1996] introduces the first linear least-squares algorithms for temporal-difference learning (LSTD). This has the advantage of making optimal use of all the training data available (as opposed to some “forgetting” in online gradient methods), however it has higher memory resources requirements, and it is limited to linear function approximation. Boyan [1999] provides a simpler derivation of LSTD and extends that work by generalising from $\lambda = 1$ to arbitrary lambda values.

Mann et al. [2016] deals with a setup very similar to ours: accurate on-policy evaluation from a dataset of trajectories. They suggest tuning the λ parameter of TD(λ) with leave-one-trajectory-out cross-validation, in order to minimise the final MSVE. However, the suggested algorithm can only do it efficiently for linear function approximation.

5.2 Successor Representations

The concept of *Successor Representation* was originally introduced by [Dayan, 1993]. That influential paper, al-

ready sets the main intuition that we pursue in the second half of this work: “*For static tasks, generalisation is typically sought by awarding similar representations to states that are nearby in some space. This concept extends to tasks involving predictions over time, except that adjacency is defined in terms of similarity of the future course of the behaviour of a dynamic system.*”. In our work though, we focus on the setup of continuous state spaces and neural networks trained with gradient methods, rather than finite state Markov Decision Processes (MDP) and linear function approximators. We also do not assume knowledge of environment dynamics and/or the policy specification, restricting ourselves to a dataset of trajectories, as in Batch RL. Dayan [1993], who experimented with a control task, hinted that some of the problems of having a policy-dependent representation would not be an issue if one would care only about pure estimation tasks. It turns out this is our problem of interest, as we just want to be more accurate at evaluating the current policy, rather than finding a better one.

Extensions to continuous state spaces, named *Deep Successor Representation* or *Successor Features* were introduced by [Kulkarni et al., 2016] and [Barreto et al., 2017], respectively. [Kulkarni et al., 2016] gives a constructive way of learning these representations, by using auxiliary losses that shape the representation to keep the necessary information to predict the immediate reward and next observation. Their aim is to solve control/navigation tasks and the automatic extraction of sub-goals, often from raw pixels of the observations. Barreto et al. [2017] focuses more on using successor features for transfer learning, in setups where the reward function changes but the environment dynamics remains the same.

All these works highlight the importance of the factorisation of the value function into a dot product of the successor features and a weight vector, and mention that the representation can be learned by any RL algorithm (e.g. Monte-Carlo or TD itself). In our work, though, we put less emphasis on the factorisation and only use successor features as a way to distinguish deceptively similar states. In addition, we use it as a mechanism to avoid some of the problems of using TD with neural networks. In this setup, we recommend *against* using TD to learn the Successor Features. We recommend first learning the successor features with Monte Carlo regression and *then* use them to learn the value function on top with Temporal Difference,

so as to avoid the problems of leakage propagation.

Proto-value functions by Mahadevan [2005], Mahadevan and Maggioni [2007] also try to tackle the problem of having a better representation for learning value functions. The original paper provided inspiration for our representation learning approaches: "... *learned not from rewards, but instead from analysing the topology of the state space*", and the experimental scenarios we test: "*states close in Euclidean distance may be far apart on the manifold (e.g. two states on opposite sides of a wall)*". It has been shown that proto-value functions and the successor representation in discrete state spaces are fundamentally the same thing, and that *successor features* are a way of extending the concept to continuous state spaces. We therefore focused more on the use of successor features.

5.3 Deep RL and the deadly triad

Sutton and Barto [2018] explain that three "deadly" components are needed to have divergence: *function approximation*, *bootstrapping* and *off-policy* data. Great progress has been made at improving the stability of Deep Reinforcement Learning with techniques like Target Networks and Experience Replays [Mnih et al., 2015, Schaul et al., 2015], but these focus more on the setup of *off-policy* learning for control tasks.

In our work, we do not have the problem of off-policy data, therefore are protected from divergence. Our focus is not on whether TD converges or not, but what solution it converges to, and how good it is. We studied a particular pathology, that we called leakage propagation, under a somewhat restricted scenario of random walks in navigation tasks. However, we believe it illustrates well the problems that should also be visible in the broader setup.

6 Conclusions and Future Work

We described the *leakage propagation* problem that can happen in on-policy evaluation with Temporal-Difference learning and neural networks, and tried to deepen our understanding of it. We showed visual empirical evidence to develop intuition about the problem, in environments with sharp discontinuities in seemingly nearby states. We then gave a formal understanding of the phenomenon in

a simple reversible MDP, by analytically finding the minimum of the Dirichlet and the Euclidean norms mixture. We showed that privileged-knowledge representations can mitigate the problem and bring significant estimation accuracy gains; we then suggested ways of using unsupervised learning to get partway there. The additional unsupervised losses are simple to implement, and make use of topological information about the trajectories, that neither Monte-Carlo nor Temporal-Difference learning explicitly exploit. These heuristics worked reasonably well in the simple two-dimensional navigation environments. However, we make no claims that these are directly extensible to more complex problems with high-dimensional inputs, as one would have to ensure that no important information is discarded in the process of learning the embedding. This requires further research. Another line of future research would be to devise a criterion to decide, per state, whether doing a TD bootstrap update is "trustworthy" or whether it would be better to rely on longer n-step updates. This would be an interesting alternative to the representation learning approach.

Acknowledgements

We would like to thank Cosmin Paduraru, Will Dabney, Tom Schaul and Hado van Hasselt for helpful and insightful discussions.

References

- Adam Adam and Martha White. Investigating practical linear temporal difference learning. In *Proceedings of the 2016 International Conference on Autonomous Agents & Multiagent Systems*, pages 494–502. International Foundation for Autonomous Agents and Multiagent Systems, 2016.
- Yusuf Aytar, Tobias Pfaff, David Budden, Tom Le Paine, Ziyu Wang, and Nando de Freitas. Playing hard exploration games by watching youtube. *arXiv preprint arXiv:1805.11592*, 2018.
- Leemon Baird. Residual algorithms: Reinforcement learning with function approximation. In *Machine Learning Proceedings 1995*, pages 30–37. Elsevier, 1995.

- André Barreto, Will Dabney, Rémi Munos, Jonathan J Hunt, Tom Schaul, Hado P van Hasselt, and David Silver. Successor features for transfer in reinforcement learning. In *Advances in neural information processing systems*, pages 4055–4065, 2017.
- Shalabh Bhatnagar, Doina Precup, David Silver, Richard S Sutton, Hamid R Maei, and Csaba Szepesvári. Convergent temporal-difference learning with arbitrary smooth function approximation. In *Advances in Neural Information Processing Systems*, pages 1204–1212, 2009.
- Justin A Boyan. Least-squares temporal difference learning. In *ICML*, pages 49–56, 1999.
- Steven J Bradtko and Andrew G Barto. Linear least-squares algorithms for temporal difference learning. *Machine learning*, 22(1-3):33–57, 1996.
- Christoph Dann, Gerhard Neumann, and Jan Peters. Policy evaluation with temporal differences: A survey and comparison. *The Journal of Machine Learning Research*, 15(1):809–883, 2014.
- Peter Dayan. Improving generalization for temporal difference learning: The successor representation. *Neural Computation*, 5(4):613–624, 1993.
- Tejas D Kulkarni, Ardavan Saeedi, Simanta Gautam, and Samuel J Gershman. Deep successor reinforcement learning. *arXiv preprint arXiv:1606.02396*, 2016.
- Sascha Lange, Thomas Gabel, and Martin Riedmiller. Batch reinforcement learning. In *Reinforcement learning*, pages 45–73. Springer, Berlin, Heidelberg, 2012.
- Sridhar Mahadevan. Proto-value functions: Developmental reinforcement learning. In *Proceedings of the 22nd international conference on Machine learning*, pages 553–560. ACM, 2005.
- Sridhar Mahadevan and Mauro Maggioni. Proto-value functions: A laplacian framework for learning representation and control in markov decision processes. *Journal of Machine Learning Research*, 8(Oct):2169–2231, 2007.
- Timothy A Mann, Hugo Penedones, Shie Mannor, and Todd Hester. Adaptive lambda least-squares temporal difference learning. *arXiv preprint arXiv:1612.09465*, 2016.
- Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. Human-level control through deep reinforcement learning. *Nature*, 518(7540):529, 2015.
- Yann Ollivier. Approximate temporal difference learning is a gradient descent for reversible policies. *arXiv preprint arXiv:1805.00869*, 2018.
- Doina Precup, Richard S Sutton, and Sanjoy Dasgupta. Off-policy temporal-difference learning with function approximation. In *ICML*, pages 417–424, 2001.
- Martin Riedmiller. Neural fitted q iteration—first experiences with a data efficient neural reinforcement learning method. In *European Conference on Machine Learning*, pages 317–328. Springer, Berlin, Heidelberg, 2005.
- Nikolay Savinov, Alexey Dosovitskiy, and Vladlen Koltun. Semi-parametric topological memory for navigation. *arXiv preprint arXiv:1803.00653*, 2018.
- Tom Schaul, John Quan, Ioannis Antonoglou, and David Silver. Prioritized experience replay. *arXiv preprint arXiv:1511.05952*, 2015.
- Pierre Sermanet, Corey Lynch, Yevgen Chebotar, Jasmine Hsu, Eric Jang, Stefan Schaal, and Sergey Levine. Time-contrastive networks: Self-supervised learning from video. *arXiv preprint arXiv:1704.06888*, 2017.
- Satinder P Singh, Tommi Jaakkola, and Michael I Jordan. Reinforcement learning with soft state aggregation. In *Advances in neural information processing systems*, pages 361–368, 1995.
- Richard S Sutton. Learning to predict by the methods of temporal differences. *Machine learning*, 3(1):9–44, 1988.

Richard S Sutton and Andrew G Barto. *Reinforcement Learning: An Introduction*. MIT press Cambridge, 1998.

Richard S Sutton and Andrew G Barto. *Reinforcement Learning: An Introduction - second edition*. MIT press Cambridge, 2018.

Richard S Sutton, Hamid R Maei, and Csaba Szepesvári. A convergent $o(n)$ temporal-difference algorithm for off-policy learning with linear function approximation. In *Advances in neural information processing systems*, pages 1609–1616, 2009a.

Richard S Sutton, Hamid Reza Maei, Doina Precup, Shalabh Bhatnagar, David Silver, Csaba Szepesvári, and Eric Wiewiora. Fast gradient-descent methods for temporal-difference learning with linear function approximation. In *Proceedings of the 26th Annual International Conference on Machine Learning*, pages 993–1000. ACM, 2009b.

Csaba Szepesvári. Algorithms for reinforcement learning. *Synthesis lectures on artificial intelligence and machine learning*, 4(1):1–103, 2010.

John N Tsitsiklis and Benjamin Van Roy. Analysis of temporal-difference learning with function approximation. In *Advances in neural information processing systems*, pages 1075–1081, 1997.