

Stabilizing Temporal Difference Learning via Implicit Stochastic Recursion

Hwanwoo Kim^{*1}, Panos Toulis^{†2}, and Eric Laber^{‡1}

¹Department of Statistical Science, Duke University

²Booth School of Business, University of Chicago

Abstract

Temporal difference (TD) learning is a foundational algorithm in reinforcement learning (RL). For nearly forty years, TD learning has served as a workhorse for applied RL as well as a building block for more complex and specialized algorithms. However, despite its widespread use, TD procedures are generally sensitive to step size specification. A poor choice of step size can dramatically increase variance and slow convergence in both on-policy and off-policy evaluation tasks. In practice, researchers use trial and error to identify stable step sizes, but these approaches tend to be ad hoc and inefficient. As an alternative, we propose implicit TD algorithms that reformulate TD updates into fixed point equations. Such updates are more stable and less sensitive to step size without sacrificing computational efficiency. Moreover, we derive asymptotic convergence guarantees and finite-time error bounds for our proposed implicit TD algorithms, which include implicit TD(0), TD(λ), and TD with gradient correction (TDC). Our results show that implicit TD algorithms are applicable to a much broader range of step sizes, and thus provide a robust and versatile framework for policy evaluation and value approximation in modern RL tasks. We demonstrate these benefits empirically through extensive numerical examples spanning both on-policy and off-policy tasks.

KEYWORDS: reinforcement learning, temporal difference learning, stochastic approximation, policy evaluation, implicit recursion

1 Introduction

Temporal difference (TD) learning, originally introduced by Sutton [34], is a cornerstone of reinforcement learning (RL). Combining the strengths of Monte Carlo methods and dynamic programming, TD learning enables incremental updates using temporally correlated data, making it

^{*}hwanwoo.kim@duke.edu

[†]panos.toulis@chicagobooth.edu

[‡]eric.laber@duke.edu

both simple and efficient for policy evaluation. This foundational algorithm underpins many modern RL techniques and has been applied successfully in a wide range of domains, including robotics, finance, and neuro-imaging, where accurate value prediction is critical for evaluation and control [28, 29, 32]. In real-world scenarios, Markov decision processes (MDPs) often operate in large state spaces, making exact value estimation computationally infeasible. A common approach to address this issue is to apply TD learning with linear function approximation. This approach makes TD learning a practical and scalable solution even for high-dimensional problems [3, 44].

Since the seminal work by Tsitsiklis and Roy [44] on the asymptotic convergence of TD algorithms with linear function approximation, numerous theoretical analyses have been conducted under a wide range of assumptions and settings [4, 11, 27, 30, 33]. For instance, Dalal et al. [11] conducted a finite-time error analysis under the assumption of i.i.d. streaming data. Bhandari et al. [4] extended this work to Markovian data by incorporating a projection step to analyze the mean path of TD iterates. Srikant and Ying [33] further derived finite-time error bounds for TD algorithms with Markovian data without requiring a projection step; their approach relied on novel refinements of stochastic approximation methods via Lyapunov function-based stability analysis. More recently, Mitra [27] established finite-time error bounds under Markovian data assumption through an elegant induction approach.

In off-policy learning problems, the data are generated under one (behavior) policy but evaluation or improvement is desired under a different (target) policy. Off-policy TD methods must correct for this distributional mismatch, typically via importance sampling ratios or projected Bellman gradients. To that end, numerous algorithms have been proposed [15, 25, 36, 37, 38] with provable asymptotic convergence guarantees under linear approximation and mild mixing conditions on the underlying Markov process. Among the long list of off-policy evaluation algorithms, temporal difference learning with gradient correction (TDC) algorithm has shown to demonstrate superior empirical performance [13, 15, 25]. Rigorous theoretical studies on TDC include finite-time error bounds under i.i.d. data [12] and under Markovian data with a projection step [46].

While TD algorithms are celebrated for their efficiency, they remain sensitive to the choice of step size in both on-policy and off-policy regimes. From a practitioner’s perspective, larger step sizes can accelerate convergence but at increased risk of numerical instability or divergence [10, 11, 39]; conversely, smaller step sizes, e.g., chosen to satisfy conservative rates conditions, guarantee stability but at the cost of slow progress. Adaptive step size mechanisms—such as those proposed by Dabney and Barto [10], which adjust rates based on temporal-error signals, or the state-dependent rules of Hutter and Legg [18]—offer some relief, but typically rely on heuristics, incur extra computational burden, and lack comprehensive theoretical support. From theoretician’s perspective, existing finite-time error bounds for both TD and TDC impose restrictive conditions on the choice of step size [4, 33, 46], again highlighting the issue of step size calibration. Thus, there remains a pressing demand for numerically stable and computationally efficient adaptive schemes with provable convergence guarantees.

Implicit stochastic recursions, as exemplified by implicit stochastic gradient descent [SGD; 40,

41, 42], provide a promising framework for improving stability in TD learning. Implicit SGD reformulates standard gradient-based recursion into a fixed point equation, where the updated parameters are constrained by both the current and new values. This formulation introduces a natural stabilizing effect, reducing sensitivity to step size and preventing divergence. Unlike explicit update methods, which directly apply gradient steps, implicit SGD imposes data-adaptive stabilization in gradient updates to control large deviations, ensuring robustness while maintaining computational simplicity. Utilizing this key idea behind the implicit SGD, we provide a principled approach to resolve step size instabilities for both on-policy and off-policy TD algorithms.

1.1 Contributions

We extend and formalize the idea of implicit recursions in TD learning, which was introduced for $TD(\lambda)$ in an unpublished manuscript by Tamar et al. [39]. We propose implicit $TD(0)$, implicit TDC, and projected implicit TD and TDC algorithms, thus creating an encompassing framework for implicit TD update rules. In implicit TD learning, the standard TD recursion is reformulated into a fixed point equation, which brings the stabilizing effects of implicit updates and thereby reduces sensitivity to the choice of step size.

Our work substantially extends Tamar et al. [39], which offers only a preliminary stability analysis of implicit $TD(\lambda)$ under a restrictive zero-reward assumption. In contrast, we provide rigorous theoretical justification for the improved numerical stability of implicit TD algorithms without relying on such unrealistic assumptions. This analysis includes asymptotic convergence guarantees for implicit $TD(0)$ and $TD(\lambda)$ algorithms, as well as finite-time error bounds for their projected variants. Moreover, we establish finite-time error bounds of the implicit version of the projected TDC. We show that, in many problems, these bounds hold independently of the step size. Furthermore, we demonstrate that the proposed implicit TD algorithms offer substantial improvements in stability and robustness while retaining the computational efficiency of standard TD methods.

Our contributions in this paper can thus be summarized as follows:

- development of implicit $TD(0)$, $TD(\lambda)$, and TDC algorithms with/ without a projection step;
- building connections between implicit and standard TD algorithms to demonstrate that implicit updates can be made with virtually no additional computational cost (in Lemma 3.1 & Lemma 3.2);
- asymptotic convergence guarantees for implicit $TD(0)$ and $TD(\lambda)$ algorithms with a decreasing step size sequence (in Theorem 4.7);
- finite-time error bounds for projected implicit $TD(0)$ and $TD(\lambda)$ with a constant step size (in Theorem 4.10 & Theorem 4.12);
- asymptotic convergence of projected implicit $TD(0)$ and $TD(\lambda)$ with a decreasing step size sequence (in Theorem 4.14 & Theorem 4.15);

- finite-time error bounds for projected implicit TDC algorithms with both constant and decreasing step size sequences (in Theorem 4.18 & Theorem 4.21);
- substantial relaxation on the requirement for step size in establishing finite-time error bounds for TD(0), TD(λ), and TDC algorithms;
- empirical demonstration of superior numerical stability of the proposed implicit TD(0) and TD(λ) in synthetic random walk, Markov reward process environments as well as continuous domain control problems;
- demonstration of implicit TDC’s substantially improved numerical stability and value function approximation over TDC in the celebrated Baird’s counterexample [1].

In Section 2, we provide the mathematical framework for TD algorithms with linear function approximation and demonstrate their instability with respect to the choice of step size. In Section 3, we formulate implicit TD(0), TD(λ), and TDC algorithms both with and without projection. In Section 4, we present theoretical justifications for proposed implicit TD(0), TD(λ), and TDC algorithms. We present both asymptotic convergence results and finite-time error bounds with constant and decreasing step size schedules. In Section 5, we demonstrate the superior numerical stability of implicit TD algorithms over standard TD algorithms in a range of environments. Finally, in Section 6, we provide a summary and concluding remarks.

2 Background

2.1 Value function

We consider a discrete-time Markov decision process with finite state space \mathcal{X} , finite action space \mathcal{A} , target policy $\pi_* : \mathcal{X} \rightarrow \mathcal{A}$, transition kernel $P(x'|a, x)$ for $x, x' \in \mathcal{X}$, $a \in \mathcal{A}$, discount factor $\gamma \in (0, 1)$, and bounded reward function $r : \mathcal{X} \times \mathcal{A} \rightarrow \mathbb{R}_{\geq 0}$. In addition, we assume there is a fixed and known feature mapping $\phi : \mathcal{X} \rightarrow \mathbb{R}^d$. Let x_n denote the state at time n , $r_n := r(x_n, a_n)$ the reward, and $\phi_n := \phi(x_n)$ the feature mapping. The primary object of interest is the value function

$$V(x) = \mathbb{E}_{\pi_*} \left(\sum_{n=1}^{\infty} \gamma^n r_n \middle| x_1 = x \right),$$

where expectation is over the sequence of states generated according to the time-homogeneous transition kernel $P_{\pi_*}(x'|x) = \sum_{a \in \mathcal{A}} P(x'|a, x)\pi_*(a|x)$. The goal of on-policy evaluation is to approximate the value function under the assumption that the observed data $(x_n)_{n \in \mathbb{N}}$ are generated by the transition kernel P_{π_*} induced by the target policy.¹ However, in some settings, collecting data under π_* , can be impossible or impractical. In such cases, off-policy evaluation is used to estimate the value function using data generated under a different, known behavioral policy π_b .

¹Since any Markov reward process can be viewed as a Markov decision process under a fixed, deterministic policy, approximating its value function is a special case of the on-policy evaluation problem.

For both policy regimes, we will assume the Markov chain $(x_n)_{n \in \mathbb{N}}$ admits a unique steady-state distribution μ_π , corresponding to the policy π (either π_* or π_b) that governs the observed data dynamics.

When the state-space, \mathcal{X} , is high-dimensional, it is generally infeasible to compute V exactly. In such cases, one must posit additional structure on the value function. As is commonly done in practice, we consider linear function approximation in which it is assumed that for some weight vector $w_* \in \mathbb{R}^d$, the value function satisfies

$$V(x) \approx V_{w_*}(x) = \phi(x)^T w_*.$$

The problem of estimating V thus reduces to estimating w_* . Define $\Phi = \left[\phi(x)^T \right]_{x \in \mathcal{X}}$, and $V_{w_*} = \Phi w_*$. Throughout, we assume Φ is of full-column rank. Such an assumption is natural, as otherwise, we can attain the same quality of approximation after removing some of the features.

2.2 Temporal difference learning

TD(0) and TD(λ) algorithms [34, 35] constitute a widely used class of stochastic approximation methods for estimating the value function V from accumulating data. Under the linear approximation, these algorithms provide a recursive estimator of w_* . For $n \in \mathbb{N}$, recall that $\phi_n = \phi(x_n)$, $r_n = r(x_n, a_n)$, and $\phi_{n+1} = \phi(x_{n+1})$. The TD(0) update rule is given by

$$\begin{aligned} w_{n+1} &= w_n + \alpha_n \delta_n \phi_n, \\ \delta_n &:= r_n + \gamma \phi_{n+1}^T w_n - \phi_n^T w_n, \end{aligned} \tag{1}$$

where α_n is the step size for the n^{th} iteration, and δ_n is referred to as the TD error. The update rule for the TD(λ) algorithm, parametrized by $\lambda \in [0, 1]$, is given by

$$\begin{aligned} w_{n+1} &= w_n + \alpha_n \delta_n e_n, \\ \delta_n &:= r_n + \gamma \phi_{n+1}^T w_n + (\lambda \gamma) e_{n-1}^T w_n - e_n^T w_n, \\ e_n &:= \phi_n + (\lambda \gamma) e_{n-1}, \quad e_0 = 0, \end{aligned} \tag{2}$$

where e_n is known as the eligibility trace, which contains information on all previously visited states. Note that the TD(λ) algorithm subsumes TD(0) and the Monte Carlo evaluation (TD(1)) as special cases. In prior numerical experiments, TD(λ) with $\lambda \in (0, 1)$ has shown superior performance over TD(0) and the Monte Carlo in approximating the value function [35].

While TD(λ) with $\lambda \in [0, 1)$ is effective in many on-policy evaluation tasks, it can become numerically unstable in off-policy settings [1, 35]. To address this, several extensions have been proposed to improve performance [36, 37, 38]. Here, we focus on the TDC algorithm with linear approximation [37]. To this end, let us define $\rho_n = \rho(a_n, x_n) := \pi_*(a_n|x_n)/\pi_b(a_n|x_n)$ which denotes the ratio of target and behavioral policy, also known as an importance weight. With $\phi_n = \phi(x_n)$, $r_n = r(x_n, a_n)$, $\phi_{n+1} = \phi(x_{n+1})$ and $\delta_n = r_n + \gamma \phi_{n+1}^T w_n - \phi_n^T w_n$, the TDC update is

given by

$$w_{n+1} = w_n + \alpha_n \rho_n \delta_n \phi_n - \alpha_n \rho_n \gamma \phi_{n+1} \phi_n^T u_n, \quad (3)$$

$$u_{n+1} = u_n + \beta_n \rho_n \delta_n \phi_n - \beta_n \rho_n \phi_n \phi_n^T u_n, \quad (4)$$

where $(\alpha_n)_{n \in \mathbb{N}}$ and $(\beta_n)_{n \in \mathbb{N}}$ are non-negative step size sequences. The primary iterates $(w_n)_{n \in \mathbb{N}}$ parameterize the value function, while the auxiliary iterates $(u_n)_{n \in \mathbb{N}}$ serve to modify the direction of TD updates. The term $\alpha_n \rho_n \gamma \phi_{n+1} \phi_n^T u_n$ in the primary iterate update is known as the gradient correction term. It is common to assume $\alpha_n \ll \beta_n$, and hence, the auxiliary iterates update on a faster time-scale than that of the target parameter iterates.

To facilitate the analysis of temporal difference learning, Bhandari et al. [4] further incorporated an additional projection step to ensure iterates $(w_n)_{n \in \mathbb{N}}$ fall into an ℓ_2 -ball of radius R . In addition to the recursive updates in (1) and (2), their update includes the projection step

$$\Pi_R(w) = \underset{w': \|w'\| \leq R}{\operatorname{argmin}} \|w - w'\| = \begin{cases} R w / \|w\| & \text{if } \|w\| > R \\ w & \text{otherwise,} \end{cases}$$

where the projection radius R is chosen to be sufficiently large to guarantee $\|w_*\| \leq R$. Such a projection step not only serves as a way to facilitate finite-time error analysis for TD(0) and TD(λ) [4], but also prevents potential divergent behavior. In the same spirit, the projected TDC algorithm, which incorporates following projection steps to (3) and (4)

$$\Pi_{R_w}(w) = \begin{cases} R_w w / \|w\| & \text{if } \|w\| > R_w \\ w & \text{otherwise} \end{cases} \quad \text{and} \quad \Pi_{R_u}(u) = \begin{cases} R_u u / \|u\| & \text{if } \|u\| > R_u \\ u & \text{otherwise} \end{cases}$$

has been studied in depth, and finite-time error bounds of TDC were established [46].

2.3 Stochastic approximation

The aforementioned TD algorithms fall into a broader class of linear stochastic approximation methods [2, 22, 31, 33], characterized by updates of the form

$$w_{n+1} = w_n + \alpha_n (b_n - A_n w_n), \quad n \in \mathbb{N},$$

where α_n is the step size for the n^{th} iteration, and (b_n, A_n) are random quantities. Under suitable technical assumptions on α_n, b_n and A_n , various types of convergence of the stochastic approximation algorithms have been established [2, 6, 24, 31]. Of particular relevance to our setting are convergence results when randomness in (b_n, A_n) is induced by the underlying time-homogeneous Markov chain $(x_n)_{n \in \mathbb{N}}$, which is assumed to mix at a geometric rate. Let \mathbb{E}_∞ denote expectation with respect to the steady-state distribution of $(x_n)_{n \in \mathbb{N}}$. Define $b = \mathbb{E}_\infty(b_n)$ and $A = \mathbb{E}_\infty(A_n)$. The so-called Robbins-Monro condition on the step size, i.e., $\sum_{n=1}^\infty \alpha_n = \infty$ and $\sum_{n=1}^\infty \alpha_n^2 < \infty$,

combined with suitable assumptions on A and b guarantees convergence of iterates w_n to w_* , where w_* is a solution of the equation $Aw = b$ [e.g., see 2, 3, 44].

Rewriting the TD(0) and TD(λ) updates as above, it can be shown that both algorithms fall into the class of linear stochastic approximation algorithms. A range of approaches utilizing existing convergence results for stochastic approximation methods [3, 44], mean-path analysis [4], Lyapunov-function based analysis [33] and mathematical induction [27] have established asymptotic and finite error bounds of TD(0) and TD(λ) iterates, respectively, to the solution w_* satisfying

$$\begin{aligned}\mathbb{E}_\infty(\phi_n \phi_n^T - \gamma \phi_n \phi_{n+1}^T)w_* &= \mathbb{E}_\infty(r_n \phi_n), \\ \mathbb{E}_\infty(e_{-\infty:n} \phi_n^T - \gamma e_{-\infty:n} \phi_{n+1}^T)w_* &= \mathbb{E}_\infty(r_n e_{-\infty:n}),\end{aligned}$$

where $e_{-\infty:n} = \sum_{k=-\infty}^n (\lambda\gamma)^{n-k} \phi_k$ is the steady-state eligibility trace.

The previously discussed linear stochastic approximation framework naturally extends to two-time scale linear stochastic approximation, characterized by coupled iterative updates operating at distinct timescales. A subclass of these algorithms is of the form:

$$\begin{aligned}w_{n+1} &= w_n + \alpha_n (b_n + A_n w_n + B_n u_n), \quad n \in \mathbb{N} \\ u_{n+1} &= u_n + \beta_n (b_n + A_n w_n + C_n u_n), \quad n \in \mathbb{N}\end{aligned}$$

where (A_n, B_n, C_n, b_n) are random quantities driven by the underlying Markov processes, and the sequences $(\alpha_n)_{n \in \mathbb{N}}$ and $(\beta_n)_{n \in \mathbb{N}}$ are positive step sizes satisfying: $\alpha_n \ll \beta_n$ and $\alpha_n/\beta_n \rightarrow 0$. This separation of scales ensures that the auxiliary iterates u_n evolve faster than the primary iterates w_n . Under additional assumptions and mixing conditions on the underlying Markov chain $(x_n)_{n \in \mathbb{N}}$, these iterates converge to a solution of coupled equilibrium equations involving expectations under the steady-state distribution [5, 6, 20, 21].

In fact, the TDC algorithm precisely fits this two-time scale framework. The primary parameter w_n and the auxiliary parameter u_n are updated concurrently but at differing step sizes α_n and β_n , respectively, with $\alpha_n \ll \beta_n$. The auxiliary iterates u_n are introduced to estimate a correction term that aligns the update direction of the primary iterates w_n with the gradient of a suitable objective function. This adjustment allows the TDC algorithm to become gradient-based, differentiating it from the standard TD method, which does not correspond to gradient-based updates. Convergence analyses of TDC explicitly leverage the two-time scale stochastic approximation theory, demonstrating that the iterates asymptotically approach the equilibrium solutions of a coupled linear system [36, 37, 46]. The two-time scale structure captures the interplay between primary and auxiliary iterates inherent to the TDC algorithm, providing rigorous convergence results and finite-time error bounds [12, 17, 46] within a unified theoretical framework.

2.4 Numerical instability

Despite the widespread use of TD algorithms, their sensitivity to step size selection presents a persistent practical challenge. While larger step sizes can speed up convergence, the updates may become unstable and cause divergence [10, 11, 39]. Conversely, using smaller step sizes improves numerical stability but can significantly slow the learning process. The primary issue stems from the recursive nature of TD methods, where updates are based on estimates that rely on prior updates, causing errors to propagate and potentially compound over time. Analogous to standard TD algorithms, it has been extensively documented that the choice of step size sequences for the TDC algorithm also requires careful calibration and restricts the usage of large step sizes, which may be inefficient [13, 15].

To demonstrate the numerical instability caused by an inappropriate choice of step size, consider the value function approximation within a simple 11-state random walk environment, as well as the celebrated Baird’s example [1]. One hundred independent experiments were conducted for the random walk environment and Baird’s counterexample. The average performance of each method and the true value function are depicted as lines, whereas the shaded bands represent variability corresponding to one standard deviation above and below the mean. Detailed descriptions of both environments are provided in Section 5. In the random walk environment, our goal is to approximate the true value function, depicted by the red dotted line in Figure 1, using cosine and sine basis functions. As shown in the left subfigure of Figure 1, the estimate obtained from the standard TD(0), depicted by the blue line, closely matches the true value function when using a small constant step size ($\alpha_n = 0.05, \forall n \in \mathbb{N}$). However, with a larger constant step size ($\alpha_n = 1.5, \forall n \in \mathbb{N}$), the approximated value function obtained using the standard TD(0) method results in values drastically deviating from the true value function, as illustrated in the right subfigure of Figure 1. For the Baird’s example, Figure 2 demonstrates substantial deterioration in the weight parameter estimates produced by the TDC algorithm as the constant step sizes change from $(\alpha_n, \beta_n, \forall n \in \mathbb{N}) = (0.005, 0.05)$ to $(\alpha_n, \beta_n, \forall n \in \mathbb{N}) = (0.01, 0.1)$, highlighting the sensitivity of the standard TDC algorithm to step size selection.

To address the numerical instability of TD algorithms—such as those shown in Figures 1—a variety of strategies have been proposed in the literature. Hutter and Legg [18] introduced adaptive step size schedules based on discounted state visitation counters. While this method demonstrated improved stability in several settings, it can diverge in continuous domains. Mahmood et al. [26] proposed an alternative adaptive step size scheme that requires tuning meta-parameters governing the decay rate. However, to avoid divergent behavior, they recommended using the step size schedule with sufficiently small initial values. This approach thus suffers from the same instability issues as standard TD methods, as it remains sensitive to the choice of step size. The Alpha-Bound algorithm [10] introduced an adaptive bound on the effective step size and demonstrated improved stability over prior approaches. Nonetheless, it requires storing vector-valued quantities across all past TD iterations and can still exhibit high variance and divergence when initialized with a large step size. Importantly, none of the aforementioned step size adaptation methods provide theoretical

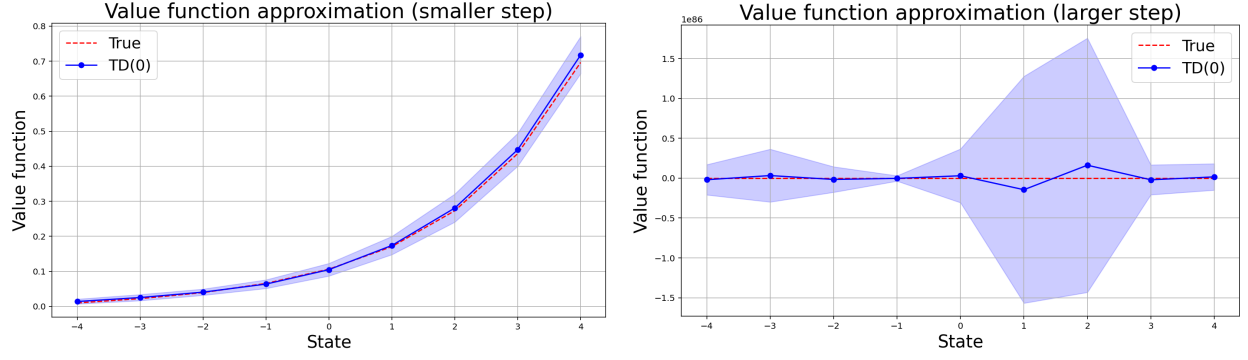


Figure 1: Left: Value function approximation in the random walk environment with a constant step size $\alpha_n = 0.05$. The estimated value function (solid blue line) by TD(0) closely matches the true value function (red dotted line), and the pointwise one standard deviation bands (shaded region) remain narrow, indicating stable performance. Right: Value function approximation in the random walk environment with a constant step size $\alpha_n = 1.5$. The standard TD(0) estimates diverge substantially from the true value function, and the pointwise one standard deviation bands inflate to extreme magnitudes, reflecting loss of numerical stability under a moderately large step size.

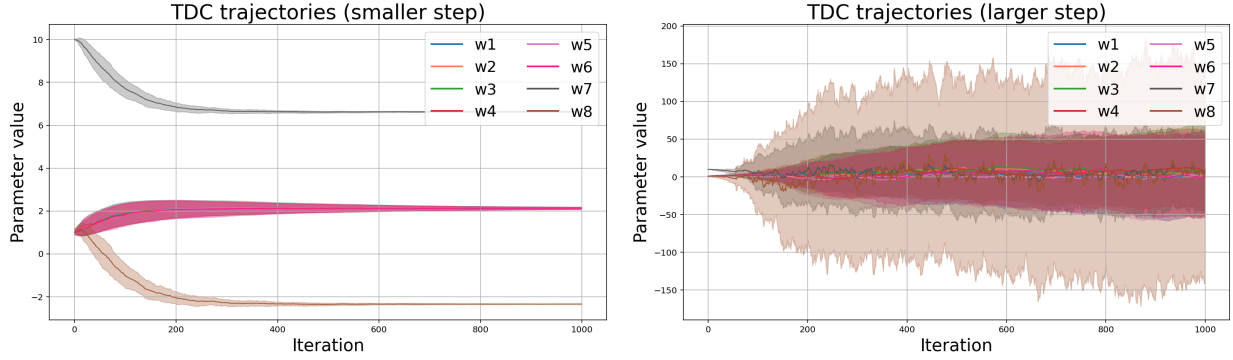


Figure 2: Left: Estimated weight parameter trajectories in Baird's counterexample with constant step sizes $(\alpha_n, \beta_n) = (0.005, 0.05)$. The TDC estimates (solid lines) converge toward a fixed point, and the pointwise one standard deviation bands (shaded region) shrink, indicating stable performance. Right: Estimated weight parameter trajectories in Baird's example with constant step sizes $(\alpha_n, \beta_n) = (0.01, 0.1)$. The TDC estimates demonstrate pronounced oscillations, and pointwise one standard deviation bands are substantially wider, reflecting TDC's sensitivity to step size selection.

guarantees or error bounds on the convergence of TD iterates. For a comprehensive discussion of other existing deterministic and stochastic step size strategies, we refer the reader to a review by George and Powell [14].

3 Implicit temporal difference learning

In this section, we introduce implicit TD algorithms, which are designed to alleviate the numerical instability discussed in Section 2.4. The key idea behind implicit updates is in rewriting the recursive update as a fixed point equation, where the future iterate appears on both the left-

and right-hand sides of the update rule. To provide a concrete example, consider the standard stochastic gradient descent (SGD) algorithm [7, 8] applied to an objective function f ,

$$w_{n+1} = w_n + \alpha_n \nabla f(w_n; \xi_n), \quad n \geq 1, \quad (5)$$

where α_n is a non-negative step size and ξ_n represents the random component involved in computing the n^{th} stochastic gradient. An implicit version of the aforementioned SGD algorithm, given by

$$w_{n+1}^{\text{im}} = w_n^{\text{im}} + \alpha_n \nabla f(w_{n+1}^{\text{im}}; \xi_n), \quad n \geq 1, \quad (6)$$

was proposed and analyzed in Toulis et al. [40], Toulis and Airolidi [41, 42], and Chee et al. [9]. Note that the highlighted term in (6) indicates that the gradient is evaluated at the future iterate, resulting in a fixed point equation. For a range of objective functions f , this equation admits a simple update rule [40]. Moreover, such an implicit update has shown to provide significant stability improvements over the standard SGD algorithm given in (5). Motivated by this central idea of implicit SGD, our goal in the following two subsections is to develop implicit variants of TD algorithms.

3.1 Implicit TD(0)/TD(λ) algorithms

In this subsection, inspired by the principles of implicit SGD, we reformulate the TD update rules as fixed-point equations. Recall that $\phi_n = \phi(x_n)$, $r_n = r(x_n, a_n)$, and $\phi_{n+1} = \phi(x_{n+1})$. Based on this formulation, we propose the following implicit TD(0) algorithm

$$\begin{aligned} w_{n+1}^{\text{im}} &= w_n^{\text{im}} + \alpha_n \delta_n^{\text{im}} \phi_n, \\ \delta_n^{\text{im}} &:= r_n + \gamma \phi_{n+1}^\top w_n^{\text{im}} - \phi_n^\top w_{n+1}^{\text{im}}, \end{aligned} \quad (7)$$

and the implicit TD(λ) algorithm [39]

$$\begin{aligned} w_{n+1}^{\text{im}} &= w_n^{\text{im}} + \alpha_n \delta_n^{\text{im}} e_n, \\ \delta_n^{\text{im}} &:= r_n + \gamma \phi_{n+1}^\top w_n^{\text{im}} + \lambda \gamma e_{n-1}^\top w_n^{\text{im}} - e_n^\top w_{n+1}^{\text{im}}, \\ e_n &:= \phi_n + (\lambda \gamma) e_{n-1}, \quad e_0 = 0, \end{aligned} \quad (8)$$

where $(\alpha_n)_{n \in \mathbb{N}}$ is a sequence of non-negative step sizes. Combining the future iterate value w_{n+1}^{im} from both sides, (7) can be rewritten as

$$(I + \alpha_n \phi_n \phi_n^\top) w_{n+1}^{\text{im}} = w_n^{\text{im}} + \alpha_n (r_n + \gamma \phi_{n+1}^\top w_n^{\text{im}}) \phi_n.$$

Analogously, equation (8) can be reexpressed as

$$(I + \alpha_n e_n e_n^\top) w_{n+1}^{\text{im}} = w_n^{\text{im}} + \alpha_n (r_n + \gamma \phi_{n+1}^\top w_n^{\text{im}} + \lambda \gamma e_{n-1}^\top w_n^{\text{im}}) e_n.$$

Using the Sherman-Morrison-Woodbury formula, we have

$$(I + \alpha_n \phi_n \phi_n^T)^{-1} = I - \frac{\alpha_n}{1 + \alpha_n \|\phi_n\|^2} \phi_n \phi_n^T \quad \text{and} \quad (I + \alpha_n e_n e_n^T)^{-1} = I - \frac{\alpha_n}{1 + \alpha_n \|e_n\|^2} e_n e_n^T.$$

These expressions provide insight into why implicit TD algorithms are more stable than standard TD. First, note that the norms of the update matrices shown above are both less than or equal to one, resulting in a stable update from w_n^{im} to w_{n+1}^{im} . In each iteration, implicit algorithms utilize both feature and eligibility trace information to impose adaptive shrinkage on the running iterates. In Algorithm 1 shown below, we present a concrete description of the implicit TD algorithms, with and without the projection step.

Algorithm 1 Implicit TD(0)/TD(λ)

Input: initial guess w_1^{im} , initial state x_1 , step size $(\alpha_n)_{n \in \mathbb{N}}$, eligibility weight parameter λ (for TD(λ)), projection radius $R > 0$ (for projected version)

For $n = 1, \dots, N$, **do**:

1. Obtain values of the reward r_n and next state x_{n+1} .
2. Compute the temporal difference error:

$$\delta_n^{\text{im}} = r_n + \gamma \phi_{n+1}^T w_n^{\text{im}} - \phi_n^T w_n^{\text{im}}$$

3. For TD(0), update:

$$w_{n+1}^{\text{im}} = w_n^{\text{im}} + \frac{\alpha_n}{1 + \alpha_n \|\phi_n\|^2} \delta_n^{\text{im}} \phi_n$$

For TD(λ), update:

$$w_{n+1}^{\text{im}} = w_n^{\text{im}} + \frac{\alpha_n}{1 + \alpha_n \|e_n\|^2} \delta_n^{\text{im}} e_n, \\ e_n = \phi_n + (\lambda \gamma) e_{n-1}, \text{ with } e_0 = 0$$

4. (For projected implicit TD) If $\|w_{n+1}^{\text{im}}\| > R$:

$$w_{n+1}^{\text{im}} = \frac{R}{\|w_{n+1}^{\text{im}}\|} w_{n+1}^{\text{im}}$$

Output: final estimate w_{N+1}^{im} .

We begin our analysis of implicit TD algorithms by establishing a connection to standard TD updates. This result is presented in Lemma 3.1 that follows.

Lemma 3.1. *An implicit update of TD(0) given in (7) can be written as*

$$w_{n+1}^{\text{im}} = w_n^{\text{im}} + \tilde{\alpha}_n \left(r_n + \gamma \phi_{n+1}^\top w_n^{\text{im}} - \phi_n^\top w_n^{\text{im}} \right) \phi_n, \quad (9)$$

where $\tilde{\alpha}_n = \frac{\alpha_n}{1 + \alpha_n \|\phi_n\|^2}$. Similarly, the implicit TD(λ) given in (8) can be expressed as

$$w_{n+1}^{\text{im}} = w_n^{\text{im}} + \tilde{\alpha}_n \left(r_n + \gamma \phi_{n+1}^\top w_n^{\text{im}} - \phi_n^\top w_n^{\text{im}} \right) e_n, \quad (10)$$

where $\tilde{\alpha}_n = \frac{\alpha_n}{1 + \alpha_n \|e_n\|^2}$.

From Lemma 3.1, we see that implicit TD(0) and TD(λ) algorithms, respectively, move along the same direction as the standard TD(0) and TD(λ). Unlike the standard TD algorithms, in implicit TD algorithms, an additional source of shrinkage in running iterates is provided through adaptive step sizes $(\tilde{\alpha}_n)_{n \in \mathbb{N}}$, which scale inversely proportional to the norm of the feature or eligibility trace. Lemma 3.1 highlights that implicit updates can be made without much additional computational cost, as the implicit TD(0) and TD(λ) algorithms amount to using random step sizes $(\tilde{\alpha}_n)_{n \in \mathbb{N}}$. In combination with a projection step discussed in Section 2.2, we introduce projected implicit TD algorithms, which can further enhance numerical stability.

3.2 Implicit TDC algorithm

In the same spirit as the implicit TD(0) and TD(λ) algorithms, here we introduce an implicit version of the TDC algorithm, which we refer to as the implicit TDC algorithm. Recall that $\phi_n = \phi(x_n)$, $r_n = r(x_n, a_n)$, $\phi_{n+1} = \phi(x_{n+1})$, and $\rho_n = \pi_*(a_n|x_n)/\pi_b(a_n|x_n)$. We propose the following implicit version of the aforementioned TDC algorithm:

$$w_{n+1}^{\text{im}} = w_n^{\text{im}} + \alpha_n \rho_n (r_n \phi_n + \gamma \phi_n \phi_{n+1}^T w_n^{\text{im}} - \gamma \phi_{n+1} \phi_n^T u_n^{\text{im}}) - \alpha_n \rho_n \phi_n \phi_n^T w_{n+1}^{\text{im}}, \quad (11)$$

$$u_{n+1}^{\text{im}} = u_n^{\text{im}} + \beta_n \rho_n (r_n \phi_n + \gamma \phi_n \phi_{n+1}^T w_n^{\text{im}} - \gamma \phi_{n+1} \phi_n^T u_n^{\text{im}}) - \beta_n \rho_n \phi_n \phi_n^T u_{n+1}^{\text{im}}. \quad (12)$$

In these expressions, the sequences $(\alpha_n)_{n \in \mathbb{N}}$ and $(\beta_n)_{n \in \mathbb{N}}$ are non-negative step sizes. To gain insight into the numerical stability of the implicit TDC update, observe that the implicit TDC update for the primary parameter w^{im} can be rewritten as follows

$$\begin{aligned} (I + \alpha_n \rho_n \phi_n \phi_n^T) w_{n+1}^{\text{im}} &= w_n^{\text{im}} + \alpha_n \rho_n (r_n \phi_n + \gamma \phi_n \phi_{n+1}^T w_n^{\text{im}} - \gamma \phi_{n+1} \phi_n^T u_n^{\text{im}}) \\ \Leftrightarrow w_{n+1}^{\text{im}} &= (I + \alpha_n \rho_n \phi_n \phi_n^T)^{-1} \{w_n^{\text{im}} + \alpha_n \rho_n (r_n \phi_n + \gamma \phi_n \phi_{n+1}^T w_n^{\text{im}} - \gamma \phi_{n+1} \phi_n^T u_n^{\text{im}})\} \\ \Leftrightarrow w_{n+1}^{\text{im}} &= (I - \alpha'_n \rho_n \phi_n \phi_n^T) \{w_n^{\text{im}} + \alpha_n \rho_n (r_n \phi_n + \gamma \phi_n \phi_{n+1}^T w_n^{\text{im}} - \gamma \phi_{n+1} \phi_n^T u_n^{\text{im}})\} \end{aligned}$$

where $\alpha'_n = \frac{\alpha_n}{1 + \alpha_n \rho_n \|\phi_n\|^2}$. Similarly, the implicit TDC update for the auxiliary parameter u^{im} can be rewritten as follows

$$\begin{aligned} (I + \beta_n \rho_n \phi_n \phi_n^T) u_{n+1}^{\text{im}} &= u_n^{\text{im}} + \beta_n \rho_n \delta_n^{\text{im}} \phi_n \Leftrightarrow u_{n+1}^{\text{im}} = (I + \beta_n \rho_n \phi_n \phi_n^T)^{-1} (u_n^{\text{im}} + \beta_n \rho_n \delta_n^{\text{im}} \phi_n) \\ &\Leftrightarrow u_{n+1}^{\text{im}} = (I - \beta'_n \rho_n \phi_n \phi_n^T) (u_n^{\text{im}} + \beta_n \rho_n \delta_n^{\text{im}} \phi_n) \end{aligned}$$

where $\delta_n^{\text{im}} = r_n + \gamma \phi_{n+1}^T w_n^{\text{im}} - \phi_n^T u_n^{\text{im}}$ and $\beta'_n = \frac{\beta_n}{1 + \beta_n \rho_n \|\phi_n\|^2}$. It turns out that the implicit TDC admits a succinct update expression, requiring the same order of computational cost as the standard TDC algorithm. A complete characterization of the implicit TDC update is provided in Lemma 3.2.

Lemma 3.2. *Implicit TDC algorithm given in (11) and (12) can be written as*

$$w_{n+1}^{im} = w_n^{im} + \alpha'_n \rho_n \delta_n^{im} \phi_n - \alpha_n \rho_n \gamma (\phi_n^T u_n^{im}) \{ \phi_{n+1} - \alpha'_n \rho_n (\phi_n^T \phi_{n+1}) \phi_n \}, \quad (13)$$

$$u_{n+1}^{im} = u_n^{im} + \beta'_n \rho_n \delta_n^{im} \phi_n - \beta'_n \rho_n \phi_n \phi_n^T u_n^{im}, \quad (14)$$

where $\alpha'_n = \frac{\alpha_n}{1 + \alpha_n \rho_n \|\phi_n\|^2}$, $\beta'_n = \frac{\beta_n}{1 + \beta_n \rho_n \|\phi_n\|^2}$ and $\delta_n^{im} = r_n + \gamma \phi_{n+1}^T w_n^{im} - \phi_n^T u_n^{im}$.

Compared to the standard TDC updates given in (3) and (4), Lemma 3.2 reveals that the implicit TDC algorithm closely resembles the standard TDC algorithm, but with adjusted step sizes and a modified correction term. Specifically, in the primary parameter update of the implicit TDC algorithm, α'_n serves as a data-adaptive version of the step size α_n . Similarly, β'_n replaces the original step size β_n in the auxiliary parameter update. Regarding the gradient correction term, the implicit TDC algorithm adjusts the TD update in the direction of

$$-\alpha_n \rho_n \gamma (\phi_n^T u_n^{im}) \{ \phi_{n+1} - \alpha'_n \rho_n (\phi_n^T \phi_{n+1}) \phi_n \},$$

in contrast to the standard TDC algorithm's correction term

$$-\alpha_n \rho_n \gamma (\phi_n^T u_n) \phi_{n+1}.$$

Roughly speaking, the standard TDC algorithm leverages the full information contained in ϕ_{n+1} , whereas the implicit TDC algorithm effectively filters out the component of ϕ_{n+1} that is aligned with ϕ_n . This reduces the correlation between consecutive gradient correction terms, which can enhance numerical stability. In Algorithm 2 shown below, we present a concrete description of the implicit TDC algorithm that was introduced in this section.

4 Theoretical analysis

In this section, we provide a theoretical analysis of our proposed implicit TD algorithms. We begin by listing out assumptions and definitions that will be used throughout this section. Unless explicitly noted otherwise, $\|\cdot\|$ denotes the Euclidean norm for vectors and the corresponding induced norm for matrices. The first assumption we introduce imposes restrictions on the data generating process.

Assumption 4.1. *[Aperiodicity and irreducibility of Markov chain] The Markov chain $(x_n)_{n \in \mathbb{N}}$ is aperiodic and irreducible with a unique steady-state distribution μ_π with $\mu_\pi(x) > 0$ for all $x \in \mathcal{X}$. In the on-policy evaluation setting, we assume that $\mu_\pi = \mu_{\pi_*}$ for some target policy μ_{π_*} . In the off-policy evaluation setting, we assume $\mu_\pi = \mu_{\pi_b}$, where π_b is the behavioral policy used to generate the data.*

Note that Assumption 4.1, together with the finiteness of the state space, implies that the Markov chain $(x_n)_{n \in \mathbb{N}}$ mixes at a uniform geometric rate [23], i.e., $(x_n)_{n \in \mathbb{N}}$ is uniformly ergodic. That is,

Algorithm 2 Implicit TDC

Input: initial guess $w_1^{\text{im}}, u_1^{\text{im}}$, initial state x_1 , step size $(\alpha_n)_{n \in \mathbb{N}}$, step size $(\beta_n)_{n \in \mathbb{N}}$, projection radius $R_w, R_u \in \mathbb{R}_{>0}$ (for projected version)

For $n = 1, \dots, N$, **do**:

1. Obtain values of the reward r_n and next state x_{n+1}
2. Compute the temporal difference error:

$$\delta_n^{\text{im}} = r_n + \gamma \phi_{n+1}^T w_n^{\text{im}} - \phi_n^T w_n^{\text{im}}$$

3. Update:

$$\begin{aligned} w_{n+1}^{\text{im}} &= w_n^{\text{im}} + \alpha'_n \rho_n \delta_n^{\text{im}} \phi_n - \alpha_n \rho_n \gamma (\phi_n^T w_n^{\text{im}}) \{ \phi_{n+1} - \alpha'_n \rho_n (\phi_n^T \phi_{n+1}) \phi_n \} \\ u_{n+1}^{\text{im}} &= u_n^{\text{im}} + \beta'_n \rho_n \delta_n^{\text{im}} \phi_n - \beta'_n \rho_n \phi_n \phi_n^T u_n^{\text{im}} \end{aligned}$$

$$\text{with } \alpha'_n = \frac{\alpha_n}{1 + \alpha_n \rho_n \|\phi_n\|^2} \text{ and } \beta'_n = \frac{\beta_n}{1 + \beta_n \rho_n \|\phi_n\|^2}$$

4. For projected implicit TDC:

if $\|w_{n+1}^{\text{im}}\| > R_w$:

$$w_{n+1}^{\text{im}} = \frac{R_w}{\|w_{n+1}^{\text{im}}\|} w_{n+1}^{\text{im}}$$

if $\|u_{n+1}^{\text{im}}\| > R_u$:

$$u_{n+1}^{\text{im}} = \frac{R_u}{\|u_{n+1}^{\text{im}}\|} u_{n+1}^{\text{im}}$$

Output: final estimate w_{N+1}^{im} .

there exist constants $m > 0$ and $\rho \in (0, 1)$ such that

$$\sup_{x \in \mathcal{X}} d_{\text{TV}} \{ \mathbb{P}(x_n \in \cdot \mid x_1 = x), \mu_\pi \} \leq m \rho^n \quad \forall n \in \mathbb{N}, \quad (15)$$

where $d_{\text{TV}}(P, Q)$ denotes the total-variation distance between probability measures P and Q . Here, the initial distribution of x_1 is the steady-state distribution μ_π , i.e., (x_1, x_2, \dots) is a stationary sequence. We next list out some assumptions on the environment and feature mapping used for approximating the value function.

Assumption 4.2. *[Bounded reward] There exists $r_{\max} > 0$, such that $\|r_n\| \leq r_{\max}$ with probability one, for all $n \in \mathbb{N}$.*

Assumption 4.3. *[Normalized features] We assume that $\|\phi_n\| \leq 1$ with probability one, for all $n \in \mathbb{N}$.*

Assumption 4.4. *[Full rank] Define $\Phi = [\phi(x)^T]_{x \in \mathcal{X}}$ as the full state matrix where the k^{th} row corresponds to ϕ evaluated at the k^{th} state in \mathcal{X} . We assume that Φ is full rank.*

Assumptions 4.1, 4.2, 4.3 and 4.4 are widely accepted in the literature [3, 4, 33, 44]. They are considered to be mild as they encompass many real world RL environments. In particular, Assumption 4.3 and Assumption 4.4 can be satisfied by removing redundant features and normalizing. In our

theory, the combined role of Assumption 4.1 and Assumption 4.4 is to preclude irregularities in the long-term behavior of the TD algorithm since, under these assumptions, the steady-state feature covariance matrix,

$$\Sigma = \Phi^T D \Phi = \sum_{x \in \mathcal{X}} \mu_\pi(x) \phi(x) \phi(x)^T,$$

is positive definite, where we set $D := \text{diag}\{\pi(x)\}_{x \in \mathcal{X}}$. We will denote the minimum eigenvalue of Σ as λ_{\min} . Moreover, thanks to Assumption 4.3, we have that $\lambda_{\min} \in (0, 1)$. Lastly, for the statement of the finite-time error bounds, we introduce the mixing time of the Markov chain $(x_n)_{n \in \mathbb{N}}$ which appears in the bounds we establish.

Definition 4.5 (Mixing time). *Given a threshold $\epsilon > 0$, constants $\rho \in (0, 1)$ and $m \in (0, \infty)$, the mixing time of the uniformly ergodic Markov chain $(x_n)_{n \in \mathbb{N}}$ is defined as*

$$\tau_\epsilon = \min\{n \in \mathbb{N} \mid m\rho^n \leq \epsilon\}.$$

For the $\text{TD}(\lambda)$ algorithm, a modified definition of mixing time, which reflects the geometric weighting of the eligibility trace will be used in the finite-time error bound expression. A formal definition is given below.

Definition 4.6 (Modified mixing time). *Given a trace-decay parameter $\lambda \in (0, 1)$, a discount factor $\gamma \in (0, 1)$, and a threshold $\epsilon > 0$, the modified mixing time of the uniformly ergodic Markov chain $(x_n)_{n \in \mathbb{N}}$ is defined as*

$$\tau_{\lambda, \epsilon} = \max\{\tau_\epsilon, \tau_\epsilon^\lambda\}, \quad \text{where } \tau_\epsilon^\lambda := \min\{n \in \mathbb{N} \mid (\lambda\gamma)^n \leq \epsilon\}.$$

To understand how these quantities behave as ϵ decreases, consider the case where $\epsilon = O(1/t^s)$ for some $s > 0$. Under this condition, it can be shown that both τ_ϵ and $\tau_{\lambda, \epsilon}$ grow at a rate of $O(\log t)$.

4.1 Asymptotic analysis for implicit TD without projection

Under the aforementioned assumptions, we can now establish the mean square convergence of the implicit $\text{TD}(0)$ and $\text{TD}(\lambda)$ algorithms.

Theorem 4.7 (Asymptotic convergence of implicit TD). *Under Assumptions 4.1-4.4, the implicit $\text{TD}(0)$ or $\text{TD}(\lambda)$ with a step size $\alpha_n = cn^{-s}$, for some constant $c > 0$ and $s \in (0.5, 1]$,*

$$\lim_{n \rightarrow \infty} \mathbb{E}\{\|w_n^{\text{im}} - w_*\|^2\} = 0.$$

The main challenge in proving convergence of the implicit algorithms is that, unlike standard TD algorithms, where the deterministic step sizes satisfy the Robbins-Monro condition, i.e., $\sum_{n=1}^{\infty} \alpha_n = \infty$, $\sum_{n=1}^{\infty} \alpha_n^2 < \infty$, the effective step sizes $(\tilde{\alpha}_n)_{n \in \mathbb{N}}$ for implicit algorithms are random as discussed in Lemma 3.1. To this end, we first establish the upper and lower bounds of the random step size $\tilde{\alpha}_n$ in terms of the deterministic step size α_n . Extending the approach taken in Srikant and Ying

[33], whose results were developed for the deterministic step size, we establish mean square error bounds of implicit TD algorithms for a sufficiently large time n using Lyapunov function based error analysis. Taking the limit of such bounds, we obtain the asymptotic convergence of implicit TD algorithms.

Remark 4.8. *Just like in standard TD algorithms [27, 33], for a sufficiently small constant step size $\alpha_n = \alpha, \forall n \in \mathbb{N}$, it is possible to establish finite-time error bounds for implicit TD algorithms. While the theoretical guarantee with the constant step size only holds for a sufficiently small α , implicit TD algorithms demonstrate improved numerical stability in comparison to standard TD algorithms over a wide range of α values, which we will confirm empirically in Section 5.*

4.2 Finite-time analysis of implicit TD with projection

To justify the robustness of implicit TD algorithms, we establish finite-time analyses of implicit TD algorithms with an additional projection step. The benefit of adding a projection step is in obtaining an upper bound of the TD update, i.e., $\delta_n \phi_n$ or $\delta_n e_n$. Since the projection step guarantees that all running iterates w_n^{im} lie inside the ball of radius $R > 0$, we get the following upper bounds for the TD updates.

Proposition 4.9. *[Lemma 6, 17 of Bhandari et al. [4]] Given any projection radius $R > 0$, for $w \in \{u : \|u\| \leq R\}$, we have*

$$\begin{aligned}\|\delta_n \phi_n\| &= \|(r_n + \gamma \phi_{n+1}^T w - \phi_n^T w) \phi_n\| \leq G := r_{\max} + 2R \\ \|\delta_n e_n\| &= \|(r_n + \gamma \phi_{n+1}^T w - \phi_n^T w) e_n\| \leq B := \frac{r_{\max} + 2R}{1 - \lambda \gamma},\end{aligned}$$

for all $n \in \mathbb{N}$.

Bhandari et al. [4] used these bounds to control the magnitudes of the stochastic updates at each iteration, ensuring that the deviation of the projected TD update from the mean-path TD update remains uniformly bounded.²

We use Proposition 4.9 to derive finite-time error bounds and asymptotic convergence for implicit TD algorithms. Our analysis extends the proof strategy of Bhandari et al. [4], who use the bounds in Proposition 4.9 to ensure that the deviation of the projected TD update from the mean-path TD update remains uniformly bounded.

Theorem 4.10 (Finite-time analysis for projected implicit TD(0)). *Suppose that Assumptions 4.1-4.4 hold with a constant step size $\alpha = \alpha_1 = \dots = \alpha_N$. Suppose also that $2\alpha(1 - \gamma)\lambda_{\min} < 1 + \alpha$. Then, the projected implicit TD(0) iterates with $R \geq \|w_*\|$ satisfy*

$$\mathbb{E} \left\{ \|w_{N+1}^{\text{im}} - w_*\|^2 \right\} \leq e^{-\frac{2\alpha(1-\gamma)\lambda_{\min}}{1+\alpha} N} \|w_1^{\text{im}} - w_*\|^2 + \frac{\alpha(1+\alpha)G^2(9+12\tau_\alpha)}{2(1-\gamma)\lambda_{\min}}.$$

²While Proposition 4.9 holds for any $R > 0$, convergence to the optimal weight parameter w_* requires that $R > \|w_*\|$. For a specific choice of R that satisfies this condition, we refer the reader to Bhandari et al. [4]. In practice, one can set $R > 0$ large enough just to prevent possible divergent behavior of TD iterates.

Remark 4.11. Under the assumptions of Theorem 4.10, we have $\lambda_{\min} \in (0, 1)$ and hence the condition $2\alpha(1 - \gamma)\lambda_{\min} < 1 + \alpha$ is met when $\gamma \in [0.5, 1)$. As such, the above finite-time bound can hold regardless of the step size choice. In comparison, note that the bound for the projected $TD(0)$ obtained in [4] requires $2\alpha(1 - \gamma)\lambda_{\min} < 1$, which does not hold for a moderately large step size. This requirement highlights the standard $TD(0)$ algorithm's potential sensitivity to the choice of step size. In contrast, the implicit TD algorithms can exhibit greater robustness across a wider range of constant step size values.

Next, we provide a finite-time error bound for the implicit $TD(\lambda)$ algorithm.

Theorem 4.12 (Finite-time analysis for projected implicit $TD(\lambda)$). *Suppose that Assumptions 4.1-4.4 hold with a constant step size $\alpha = \alpha_1 = \dots = \alpha_N$. Suppose also that $2\alpha(1 - \lambda\gamma)^2(1 - \kappa)\lambda_{\min} < 1 + \alpha$ where $\kappa = \frac{\gamma(1-\lambda)}{1-\lambda\gamma}$. Then, the projected implicit $TD(\lambda)$ iterates with $R \geq \|w_*\|$ satisfy*

$$\mathbb{E} \left\{ \|w_{N+1}^{im} - w_*\|^2 \right\} \leq e^{-\frac{2\alpha(1-\lambda\gamma)^2(1-\kappa)\lambda_{\min}}{1+\alpha}N} \|w_1^{im} - w_*\|^2 + \frac{(1 + \alpha) \{ \alpha B^2(24\tau_{\lambda,\alpha} + 15) + 2B^2 \}}{2(1 - \lambda\gamma)^2(1 - \kappa)\lambda_{\min}}.$$

Remark 4.13. Note that $(1 - \lambda\gamma)^2(1 - \kappa) = (1 - \lambda\gamma)(1 - \gamma)$. Hence, for $\gamma \in [0.5, 1)$, just like in the case of the projected implicit $TD(0)$, the above finite-time error bound holds regardless of the constant step size. Thanks to the additional factor of $(1 - \lambda\gamma)$, the result applies to a broader class of problems, indicating enhanced numerical stability over projected implicit $TD(0)$. In particular, for $\lambda \geq \frac{1}{2\gamma}$, the bound holds regardless of the choice of step size.

The theoretical results shown above are under a constant step size regime, where the running iterates w_N^{im} do not necessarily converge to w_* . With a decreasing step size sequence, we can establish the following asymptotic convergence results for both the implicit $TD(0)$ and $TD(\lambda)$ algorithms.

Theorem 4.14 (Asymptotic convergence of projected implicit $TD(0)$). *Suppose that Assumptions 4.1-4.4 hold. For $\alpha_1 > 0$ and $N > \tau_{\alpha_N}$, with a step size sequence $\alpha_n = \frac{\alpha_1}{\alpha_1 \lambda_{\min}(1-\gamma)(n-1)+1}$, the projected implicit $TD(0)$ iterates with $R \geq \|w_*\|$ achieves*

$$\mathbb{E} \left\{ \|w_{N+1}^{im} - w_*\|^2 \right\} = \tilde{O}(1/N),$$

where \tilde{O} is big- O suppressing logarithmic factors. In particular,

$$\mathbb{E} \left\{ \|w_{N+1}^{im} - w_*\|^2 \right\} \rightarrow 0 \quad \text{as } N \rightarrow \infty.$$

Theorem 4.15 (Asymptotic convergence of projected implicit $TD(\lambda)$). *Suppose that Assumptions 4.1-4.4 hold. For $\alpha_1 > 0, \kappa = \frac{\gamma(1-\lambda)}{1-\lambda\gamma}$ and $N > 2\tau_{\alpha_N}$, with a step size sequence $\alpha_n = \frac{\alpha_1}{\alpha_1 \lambda_{\min}(1-\kappa)(n-1)+1}$, the projected implicit $TD(0)$ iterates with $R \geq \|w_*\|$ achieves*

$$\mathbb{E} \left\{ \|w_{N+1}^{im} - w_*\|^2 \right\} = \tilde{O}(1/N),$$

where \tilde{O} is big- O suppressing logarithmic factors. In particular,

$$\mathbb{E} \left\{ \|w_{N+1}^{\text{im}} - w_*\|^2 \right\} \rightarrow 0 \quad \text{as } N \rightarrow \infty.$$

4.3 Finite-time analysis of implicit TDC with projection

In this subsection, we establish finite-time error bounds for the proposed projected implicit TDC algorithm with both decreasing step size schedules and constant step sizes. Recall that in the off-policy evaluation setting, the goal is to approximate the value function under target policy with data generated from the behavioral policy. To this end, we need to restrict the class of behavioral policy to a reasonable set of policies to ensure optimal value function approximation is identifiable. Widely accepted assumptions for such a guarantee are listed below [45, 46].

Assumption 4.16. *[Importance weights] There exist $\rho_{\max} \in (0, \infty)$ such that, for all $n \in \mathbb{N}$, $\rho(a_n, x_n) = \pi_*(a_n|x_n)/\pi_b(a_n|x_n) \leq \rho_{\max}$.*

Assumption 4.16 is a mild condition to guarantee that the support of behavioral policy is as large as the target policy. Furthermore, with a slight abuse of notation, in the context of off-policy evaluation³, let us define

$$A := \mathbb{E}_{\mu_{\pi_b}} \left[\rho(x, a) \phi(x) \{ \gamma \phi(x') - \phi(x) \}^\top \right], \quad C := -\mathbb{E}_{\mu_{\pi_b}} \left[\rho(x, a) \phi(x) \phi(x)^\top \right],$$

as well as $b := \mathbb{E}_{\mu_{\pi_b}} [\rho(x, a) r(x, a) \phi(x)]$. It can be shown that an optimal linear function approximation with respect to the mean-square projected Bellman error is obtained when $w_* = A^{-1}b$ [37, 46]. We make the following assumption on matrices A and C to guarantee the existence of unique optimal linear function approximation representation of the target value function.

Assumption 4.17. *[Problem solvability] The matrix A and C are nonsingular. We denote the minimum absolute eigenvalue of the matrix C to be $\lambda_c > 0$. Furthermore, there exist λ_u and λ_w such that $\lambda_{\max}(2C) \leq \lambda_u < 0$ and $\lambda_{\max}(2A^\top C^{-1}A) \leq \lambda_w < 0$.*

Recall that in the implicit TDC algorithm, there are two sequences of iterates: w_n^{im} , which parameterizes the value function of interest, and u_n^{im} , which serves as an auxiliary variable to compute the gradient correction term for the primary iterate. To facilitate the error analysis, we introduce the tracking error vector $v_n := u_n^{\text{im}} - u_n^*$, where $u_n^* := -C^{-1}(b + Aw_n^{\text{im}})$ denotes the stationary point of the ODE: $u' = b + Aw^{\text{im}} + Cu$. In short, for a fixed value of w_n^{im} , u_n^* represents the point to which the auxiliary iterates u_n^{im} would converge. The tracking error v_n thus quantifies the deviation of the auxiliary iterates from their instantaneous stationary point, providing a handle to assess how much of the overall error in the primary iterate w_n^{im} can be attributed to imperfect tracking by the auxiliary sequence. We first establish finite-time error bounds of the implicit TDC algorithm with a decreasing sequence of step sizes $(\alpha_n)_{n \in \mathbb{N}}$ and $(\beta_n)_{n \in \mathbb{N}}$.

³In the context of on-policy evaluation, definitions on the page 6 implies $A := \mathbb{E}_\infty [\phi(x) \{ \gamma \phi(x') - \phi(x) \}^\top]$ and $b := \mathbb{E}_\infty [r(x, a) \phi(x)]$

Theorem 4.18 (Finite-time analysis for implicit TDC with decreasing step sizes). *Given Assumptions 4.1, 4.2, 4.3, 4.16 and 4.17, suppose $\alpha_n = \frac{c_\alpha}{n^\sigma}$, $\beta_n = \frac{c_\beta}{n^\nu}$, with $0 < \nu < \sigma < 1$, $c_\alpha(|\lambda_w| - \rho_{\max}) < 1$ and $c_\beta(|\lambda_u| - \rho_{\max}) < 1$. Then for any $\epsilon \in (0, \sigma - \nu]$, $\epsilon' \in (0, 0.5]$, the projected implicit TDC with $R_w \geq \|w_*\|$ and $R_u \geq 2\rho_{\max}(\gamma + 1)R_w/\lambda_c$ yields*

$$\begin{aligned}\mathbb{E}\|w_n^{im} - w_*\|^2 &= O\left(e^{\frac{-|\lambda_w|c_\alpha n^{1-\sigma}}{(1+c_\alpha\rho_{\max})(1-\sigma)}}\right) + O\left(\frac{\log n}{n^\sigma}\right) + O\left(\frac{\log n}{n^\nu} + h(\sigma, \nu)\right)^{1-\epsilon'} \\ \mathbb{E}\|v_n\|^2 &= O\left(\frac{\log n}{n^\nu}\right) + O(h(\sigma, \nu)), \quad h(\sigma, \nu) = \begin{cases} \frac{1}{n^\nu}, & \sigma > 1.5\nu, \\ \frac{1}{n^{2(\sigma-\nu)-\epsilon}}, & \nu < \sigma \leq 1.5\nu. \end{cases}\end{aligned}$$

Remark 4.19. *The key result in Theorem 4.18 is that implicit TDC offers greater flexibility than standard TDC in choosing a step size schedule when there exists a large discrepancy between the target and behavioral policy. For example, the condition on the step size, i.e., $c_\alpha(|\lambda_w| - \rho_{\max}) < 1$, holds for any $c_\alpha > 0$ if $|\lambda_w| - \rho_{\max} \leq 0$. Moreover, even if $|\lambda_w| - \rho_{\max} > 0$, the condition $c_\alpha < 1/(|\lambda_w| - \rho_{\max})$ permits much wider range of initial step sizes in comparison to the requirement $c_\alpha < 1/|\lambda_w|$ for the standard TDC [46]. The same logic also applies to c_β .*

Remark 4.20. *Note that the step size condition is automatically satisfied for all c_α and c_β when $\rho_{\max} \geq \max\{|\lambda_w|, |\lambda_u|\}$. This implies that, unlike standard TDC, the implicit TDC algorithm permits a more flexible choice of step sizes, particularly when the discrepancy between the target and behavioral policies becomes large. At the same time, the increased difficulty of learning the value function under a behavioral policy that is far from the target policy is captured in the leading term of the error bound for the primary iterate. In particular, the decaying rate of the leading term slows as the gap between the behavioral and target policies increases.*

Theorem 4.21 (Finite-time analysis for implicit TDC with a constant step size). *Given Assumptions 4.1, 4.2, 4.3, 4.16 and 4.17, suppose $\alpha_n = c_\alpha$, $\beta_n = c_\beta$, with $c_\alpha(|\lambda_w| - \rho_{\max}) < 1$ and $c_\beta(|\lambda_u| - \rho_{\max}) < 1$. Then, for all $n \in \mathbb{N}$, the projected implicit TDC with $R_w \geq \|w_*\|$ and $R_u \geq 2\rho_{\max}(\gamma + 1)R_w/\lambda_c$ yields,*

$$\begin{aligned}\mathbb{E}\|v_{n+1}\|^2 &\leq (1 - \underline{\beta}|\lambda_u|)^n \|v_1\|^2 + C_v \\ \mathbb{E}\|w_{n+1}^{im} - w_*\|^2 &\leq (1 - \underline{\alpha}|\lambda_w|)^n \|w_1 - w_*\|^2 + C_w\end{aligned}$$

where $\underline{\alpha} = \frac{c_\alpha}{1+c_\alpha\rho_{\max}}$, $\underline{\beta} = \frac{c_\beta}{1+c_\beta\rho_{\max}}$ and

$$\begin{aligned}C_v &= O(\max\{c_\beta\tau_{c_\beta}, c_\beta^2\tau_{c_\beta}\}) + O(\max\{c_\alpha, c_\alpha^2\}) + O(\max\{c_\alpha/c_\beta, c_\alpha^2/c_\beta\}) \\ C_w &= O(\max\{c_\alpha, c_\alpha^4\}) + O(\sqrt{C_v} + c_\alpha\sqrt{C_v}) + O(\max\{c_\alpha, c_\alpha^3\}\tau_{c_\alpha}).\end{aligned}$$

Remark 4.22. *Similar to the case with diminishing step sizes, implicit TDC substantially relaxes the restrictions on the choice of constant step sizes c_α and c_β in comparison to standard TDC [46], which requires $c_\alpha < 1/|\lambda_w|$ and $c_\beta < 1/|\lambda_u|$. Under a constant step-size schedule, the algorithm*

converges into a neighborhood of the true solution w_* . Larger values of c_α and c_β accelerate the rate at which w_n approaches the neighborhood but also enlarge its radius, settling farther from w_* . Conversely, smaller step sizes shrink the size of the neighborhood at the expense of slower convergence. The finite-time bounds therefore demonstrate a clear trade-off: one must balance the speed of convergence against the size of the neighborhood of convergence. In the limiting regime $c_\alpha \rightarrow 0$ and $c_\beta \rightarrow 0$ with $c_\alpha/c_\beta \rightarrow 0$, the neighborhood radius vanishes and $w_n^{im} \rightarrow w_*$ as $n \rightarrow \infty$.

5 Numerical experiments

5.1 Random walk with absorbing states

In this subsection, we consider a one-dimensional random walk environment with 11 integer-valued states arranged on a real line, with zero at the center. The two endpoints—the leftmost and rightmost states—are absorbing and thus omitted from the value function approximation. The reward is zero for all states except for the rightmost state, where the reward is one. A total number of 100 independent experiments were run with a discount factor $\gamma = 0.9$ and a projection radius $R = 10$. We employ TD(0), implicit TD(0), projected TD(0), and projected implicit TD(0), and show their average performance as well as one standard deviation bands as shades in Figure 3. In all experiments, we use a sequence of constant step sizes between 0.05 and 1.5.

Based on the top left plot in Figure 3, we observe that as the step size increases, the average value approximation error increases for all four algorithms. We also observe that both implicit TD(0) and projected implicit TD(0) had a smaller increase in value approximation error compared to TD(0) and projected TD(0). For a small step size $\alpha = 0.05$, all four algorithms provided accurate value function approximation as shown in the top right plot in Figure 3. However, for a moderately large step size ($\alpha = 1.5$), both TD(0) and projected TD(0) suffered from numerical instability, yielding poor value function approximation results. This can be seen in the bottom right plot in Figure 3 as well as Figure 1. Unlike TD(0) or projected TD(0), the implicit procedures remained numerically stable, which can be seen in the bottom left plot in Figure 3. We also employed standard and implicit TD(0.5) algorithms, and observed qualitatively identical results. These results confirm our theoretical results, particularly in Theorem 4.10 and 4.12. See the Supplementary Material for details.

5.2 Reward process with 100-dimensional states

Motivated by Zhang et al. [47], we construct a synthetic Markov reward process with 100 states whose transition probability matrix is generated at random. For each state, we sample 99 independent Uniform(0, 1) samples, sort them, and take successive differences—treating 0 and 1 as boundary points—to form a valid transition probability distribution. Repeating this procedure for all 100 states and stacking the resulting distributions row-wise yields the full transition probability matrix P . Rewards are assigned by drawing one Uniform(0, 1) sample per state and collecting these into the reward vector r . We set the discount factor to $\gamma = 0.9$.

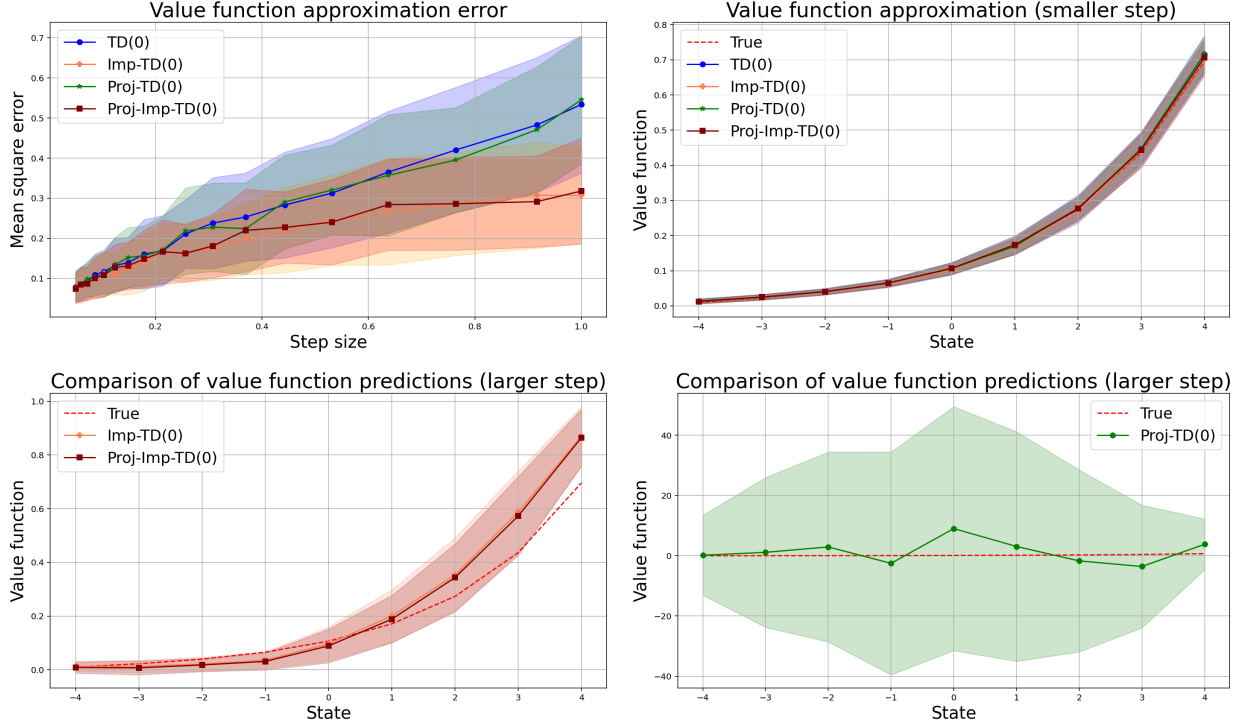


Figure 3: All figures pertain to the random walk environment. **Top left:** Value function approximation error versus constant step size over the interval $[0.05, 1]$. Implicit TD(0) exhibits a more gradual increase in value approximation error as the step size grows, reflecting its enhanced robustness to large step sizes. **Top right:** Value function approximation with $\alpha_n = 0.05$. Both standard and implicit TD(0) algorithms accurately recover the true value function, with tight confidence bands. **Bottom left:** Value function approximation with $\alpha_n = 1.5$ using implicit TD(0). Unlike the standard TD(0), with a moderately large constant step size, implicit TD(0) algorithms remain numerically stable. **Bottom right:** Value function approximation with $\alpha_n = 1.5$ using projected TD(0). Even with projection, the standard TD(0) algorithm exhibits pronounced instability, reflected in the enlarged confidence band.

In this example, the true value function can be analytically computed and is given by $v_* = (I - \gamma P)^{-1}r$. Our job is to approximate the true value function via Φw , where each row of $\Phi \in \mathbb{R}^{100 \times 20}$ represents a normalized random binary feature. The oracle parameter w_* was obtained by solving $\min_w \|\Phi w - v_*\|$. Both standard and implicit TD algorithms were run for $N = 10^5$ iterations with $\lambda \in \{0, 0.5\}$ under the decaying step size schedule $\alpha_n = 300/n$. We set a vacuously large projection radius $R = 5000$ and conducted 20 independent experiments. Figure 4 depicts the mean estimation error, with shaded bands indicating one standard deviation. Figures 4 and Table 5a present parameter estimation results for standard versus implicit TD(0) and TD(λ) algorithms.

In Table 5a, we see that the mean estimation error for standard TD(0) is 5.356 (std 3.279), while for implicit TD(0) it is 0.117 (std 0.044), a reduction of roughly 98%. Figure 4 (left) further shows that, within the first 50 iterations, standard TD(0) trajectory deviates from the true parameter, whereas the implicit TD(0) algorithm immediately reduces the estimation error. After 10^5 itera-

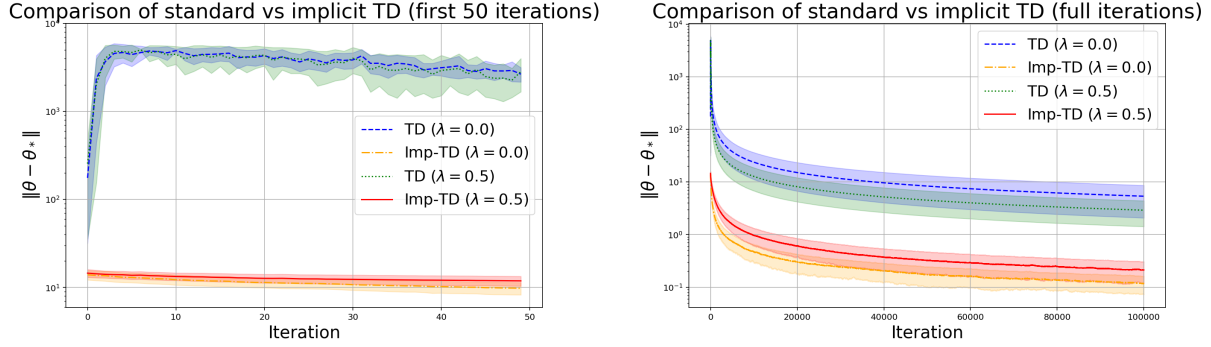


Figure 4: Parameter estimation error in a synthetic 100-state Markov reward process, comparing standard TD and implicit TD for $\lambda = 0$ and $\lambda = 0.5$. Step size was set to $\alpha_n = 300/n$. **Left:** Over the first 50 iterations, standard TD exhibits pronounced error amplification before slowly decaying, whereas implicit TD yields an immediate error reduction. **Right:** Over 10^5 iterations, implicit TD consistently converges toward the optimal parameter with superior accuracy, whereas standard TD remains hindered by its initial error amplification.

tions (Figure 4, right), standard TD(0) plateaus at a high error, but implicit TD(0) has already reached near-zero error. This comparison extends to TD(0.5) variants as well. From Table 5a, we see that standard TD(0.5) achieves mean error 2.906 (std 1.484), while implicit TD(0.5) attains mean 0.212 (std 0.094). Although standard TD(0.5) roughly halved the estimation error relative to TD(0), implicit TD(0.5) algorithms nonetheless outperformed the standard methods by an order of magnitude and exhibited smaller variance across independent runs. Implicit TD methods consistently improved numerical stability, allowing the use of large step sizes for fast early learning, and produced both lower bias and lower variance in the final parameter estimates, for both TD(0) and TD(0.5).

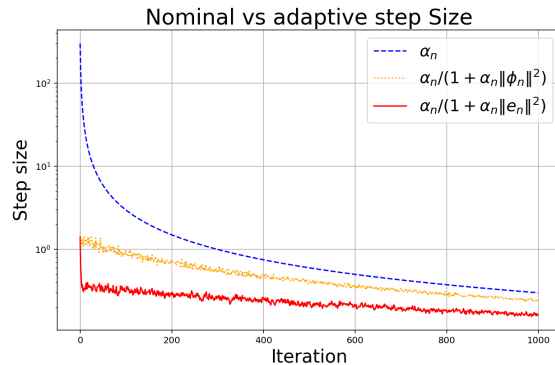
In Figure 5b, we provide a plot of decreasing step size $\alpha_n = 300/n$ versus effective step sizes for implicit TD(0): $\alpha_n/(1 + \alpha_n\|\phi_n\|^2)$ and implicit TD(0.5): $\alpha_n/(1 + \alpha_n\|e_n\|^2)$. The figure shows that all three step size schedules decrease to zero. In the meantime, the effective step sizes for the implicit algorithms are not necessarily monotonic, as they depend on the random quantities ϕ_n and e_n . Such an adaptive step size prevents numerical instability by appropriately scaling down drastic temporal difference updates.

5.3 Policy evaluation for continuous domain control

In this subsection, we test the robustness of implicit updates in classical control tasks. We considered both Acrobot and Mountain Car environments available through **Gymnasium** library in Python [43]. The Acrobot environment consists of a two link pendulum system with two joints, where only the joint between the two links is actuated. The episode begins with both links hanging downward, and the objective is to swing the end of the lower link upward to reach a specified target height. The agent receives a reward of -1 at each time step until the goal is achieved, which ends the episode with a reward of 0. In the Mountain Car environment, a car is positioned between two

Method	λ	Mean \pm Std
Standard TD	0.0	5.356 ± 3.279
Implicit TD	0.0	0.117 ± 0.044
Standard TD	0.5	2.906 ± 1.484
Implicit TD	0.5	0.212 ± 0.094

(a) Final average parameter estimation error



(b) Nominal step size vs effective step size

Figure 5: **Left:** Average final parameter estimation errors and standard deviation for standard and implicit TD algorithms in a synthetic 100-state Markov reward process. Implicit TD yields substantially improved parameter estimation error when used with a large initial step size. **Right:** Nominal step size $\alpha_n = 300/n$ versus effective step size trajectories $\alpha_n/(1 + \alpha_n \|\phi_n\|^2)$ for TD(0) and $\alpha_n/(1 + \alpha_n \|e_n\|^2)$ for TD(0.5). Although the effective step sizes in the implicit algorithms exhibit non-monotonicity, they eventually converge to zero.

hills, where the goal is to reach the top of the right hill. The car’s engine is underpowered, so the agent must build momentum by driving back and forth. The state includes position and velocity; actions apply force left, right, or none. Each step incurs a reward of -1 , and the episode ends upon reaching the goal.

We applied both the standard and implicit TD(0) algorithms to state-action value function approximation in the Acrobot and Mountain Car environments. In each case, the state-action value function was approximated by radial basis features $\phi_n \in \mathbb{R}^{100}$, and we measured performance by the empirical root mean square temporal difference error (RMSTDE) computed over 1000 input values. We used a decaying step size schedule $\alpha_n = \alpha_1/n$, $\alpha_1 \in \{1.0, 10.0\}$ with a radius $R = 100$ for Acrobot and $\alpha_1 \in \{1.0, 5.0\}$ with $R = 1000$ for Mountain Car. A total of 20 independent experiments were conducted. Figure 6 presents the mean RMSTDE for both environments, with shaded regions covered by one standard deviation bands.

The results for the Acrobot environment are shown in Figure 6 (left) and Table 1 (left). With an initial step size of $\alpha_1 = 1.0$, standard TD(0) attained a mean RMSTDE of 0.546 (std 0.167), whereas implicit TD(0) yielded a higher mean RMSTDE of 0.655 but with markedly lower variability (std 0.062). When α_1 was increased to 10.0, standard TD(0) performed much worse than implicit TD(0), achieving a mean RMSTDE of 2.585 (std 2.308) against a mean RMSTDE of 0.428 (std 0.099). This demonstrates that the implicit procedure remains stable and even benefits from larger step sizes, while the standard TD procedure suffers from a large initial step size and greater run-to-run variance.

The results for the Mountain Car environment are shown in Figure 6 (right) and Table 1 (right). In this environment, the advantage of implicit updates under an aggressive step size is more evident. With $\alpha_1 = 1.0$, both methods performed similarly (standard: mean RMSTDE of 0.379 with std

0.088; implicit: mean RMSTDE of 0.324 with std 0.043). But with $\alpha_1 = 5.0$, standard update drastically deteriorated (mean RMSTDE of 19.827 with std 10.395), whereas implicit version obtained an improved error (mean RMSTDE of 0.162 with std 0.042). These results demonstrate that implicit algorithms retain the ease of implementation of standard methods while substantially enhancing numerical stability even in continuous domain control problems.

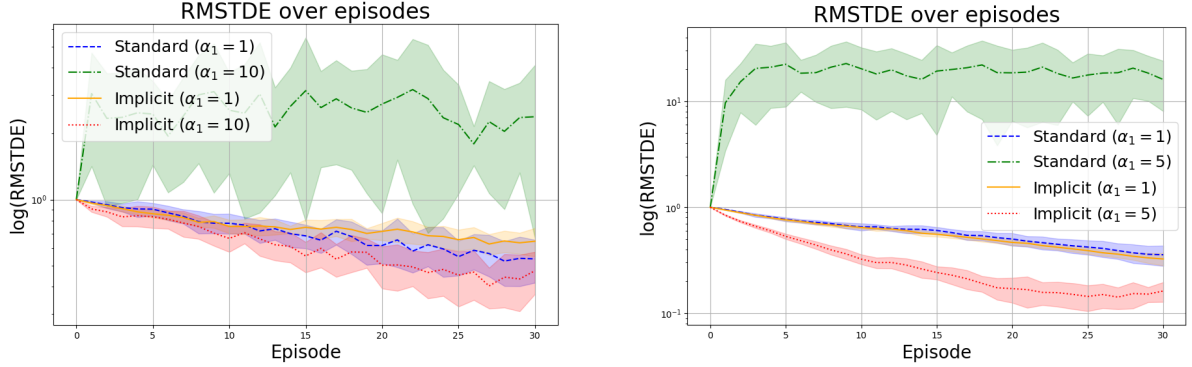


Figure 6: Average root mean square temporal difference error (RMSTDE) versus episode index for standard and implicit TD(0). Shaded bands indicate pointwise one standard deviation over 20 independent runs. **Left:** In the Acrobot environment with a large initial step size, implicit TD(0) delivers accelerated RMSTDE decay and reduced variance, whereas standard TD(0) exhibits poor convergence and amplified variance. **Right:** In the Mountain Car environment, implicit TD(0) showcases rapid RMSTDE reduction under a large initial step size, reflecting superior numerical stability relative to standard TD(0).

Acrobot (RMSTDE)		
Method	α_1	Mean \pm Std
Standard	1.0	0.546 \pm 0.167
Standard	10.0	2.585 \pm 2.308
Implicit	1.0	0.655 \pm 0.062
Implicit	10.0	0.428 \pm 0.099

Mountain Car (RMSTDE)		
Method	α_1	Mean \pm Std
Standard	1.0	0.379 \pm 0.088
Standard	5.0	19.827 \pm 10.395
Implicit	1.0	0.324 \pm 0.043
Implicit	5.0	0.162 \pm 0.042

Table 1: Final root mean square temporal difference error (RMSTDE) for standard and implicit TD(0) on the Acrobot (left) and Mountain Car (right) environments. In the Acrobot environment, implicit TD(0) matches standard TD(0) at $\alpha_1 = 1$, while at $\alpha_1 = 10$ it substantially reduces both error magnitude and variance compared to standard TD(0). In the Mountain Car environment, implicit TD(0) offers modest improvement at $\alpha_1 = 1$ and, at $\alpha_1 = 5$, prevents the severe RMSTDE explosion exhibited by standard TD(0), thereby demonstrating superior numerical stability.

5.4 Baird’s counterexample

In this last subsection, we consider a celebrated off-policy evaluation problem, known as the Baird’s counterexample. This is a classical benchmark problem in reinforcement learning, specif-

ically constructed to expose instability and convergence issues in off-policy TD algorithms when combined with linear function approximation. Originally introduced by Baird et al. [1], this example is notable for its simplicity yet significant theoretical and practical implications. The environment consists of seven states, with one center state connected directly to six peripheral states. The behavioral policy used in Baird’s example uniformly selects one of the six peripheral states with equal probability ($1/6$), while the target policy deterministically transitions to the center state, creating a distinct discrepancy between the two policies. This deliberate mismatch poses substantial difficulties for algorithms relying on off-policy updates. Linear function approximation is employed in this counterexample, characterized by eight distinct features designed to create an inherently challenging setting for standard TD methods. Every peripheral state has its own unique feature, and there is one extra feature that remains active in every state, including the center. This particular choice of feature representation leads to nontrivial correlation among the features, further complicating the convergence and stability of existing TD based algorithms.

We performed 100 independent experiments and report the mean and standard deviation of the outcomes. The results depicted in Figures 7 and 8 demonstrate the critical role that step size selection plays in the performance of TDC and implicit TDC algorithms when applied to Baird’s example. In the constant step size setting, a smaller step size ($\alpha_1 = 0.005, \beta_1 = 0.05$) yields stable and convergent behavior for both TDC and implicit TDC. While implicit TDC demonstrated lower final errors in root mean square value error (RMSVE), TDC obtained a smaller root mean square projected Bellman error (RMSPBE). However, when a larger constant step size is used ($\alpha_1 = 0.025, \beta_1 = 0.25$), standard TDC diverges, exhibiting extremely large errors, i.e., RMSPBE value of 41.754 (std 21.570) and RMSVE value of 71.425 (std 33.891). In contrast, implicit TDC remains stable and achieves substantially smaller errors, i.e., RMSPBE value of 0.284 (std 0.202), RMSVE value of 0.521 (std 0.176).

Under decreasing step sizes, where $\alpha_n = \alpha_1/n^{(99/100)}$ and $\beta_n = \beta_1/n^{(2/3)}$, aforementioned patterns continue to highlight the implicit TDC’s numerical robustness. Standard TDC suffered from numerical instability when the initial step sizes are chosen to be large (e.g., $\alpha_1 = 1.0, \beta_1 = 10.0$), resulting in error amplification, indicated from RMSPBE value of 8.219 (std 3.514) and RMSVE value of 59.899 (std 33.652). In stark contrast, implicit TDC obtains significantly improved errors, i.e., RMSPBE value of 0.851 (std 0.423) and RMSVE value of 1.861 (std 1.139) with large initial step sizes. Crucially, implicit updates maintain stability even with large initial step sizes, avoiding the numerical instabilities seen in standard TDC.

Figure 8 further demonstrates improved numerical stability of implicit TDC by presenting trajectories of the parameter updates. While the standard TDC trajectories oscillate significantly, especially with larger step sizes, implicit TDC trajectories quickly stabilize and remain bounded. These empirical results validate our theoretical findings in Theorems 4.18 and 4.21, underline the critical importance of selecting appropriate step sizes in off-policy algorithms, and strongly advocate for the robustness of the implicit TDC approach in off-policy evaluation examples such as Baird’s counterexample.

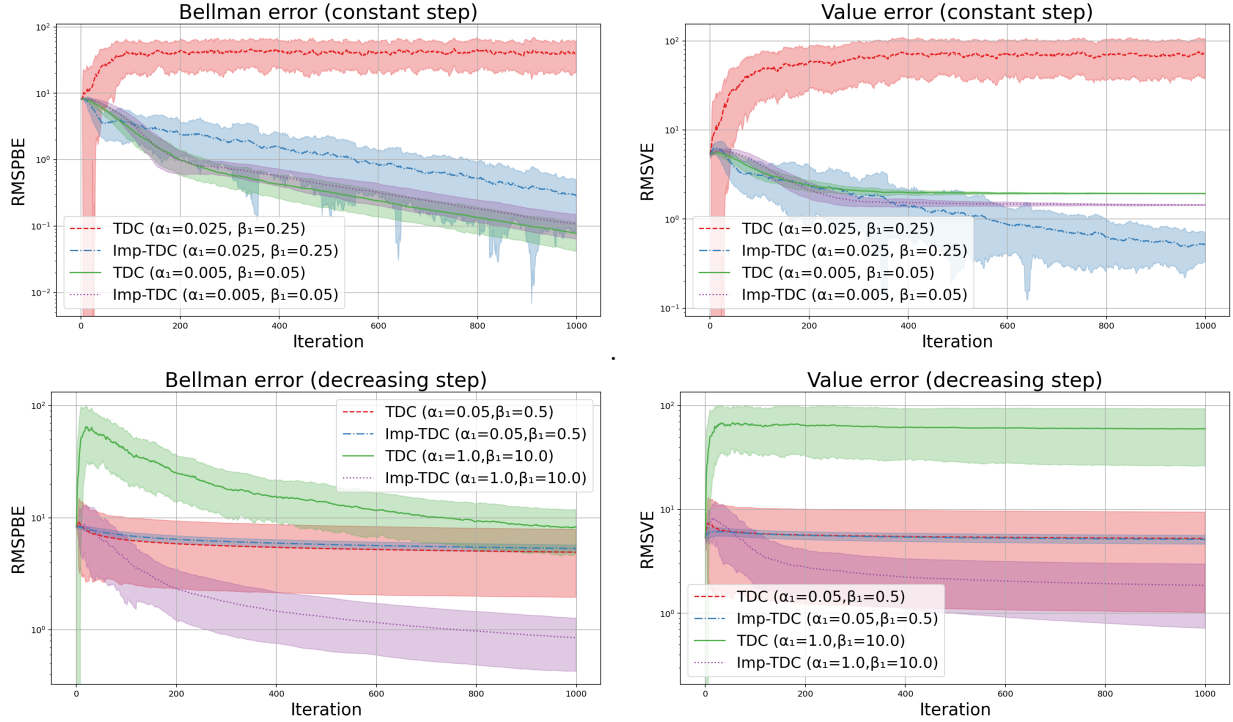


Figure 7: All figures pertain to Baird’s counterexample. **Top left (Bellman error, constant step):** Implicit TDC achieves rapid projected Bellman error decay under both constant step size configurations, whereas standard TDC amplifies error and plateaus at a large step size configuration. **Top right (value error, constant step):** Implicit TDC consistently yields lower value error with a large step size configuration, unlike standard TDC, which fails to reduce error and underperforms even with a small step size configuration. **Bottom left (Bellman error, decreasing step):** Implicit TDC significantly reduces projected Bellman error with a large initial step size, reaching far lower errors than standard TDC. **Bottom right (value error, decreasing step):** Implicit TDC maintains consistent value error decline, while standard TDC exhibits elevated, plateauing error under a large initial step size.

6 Conclusion

This paper introduces implicit version of TD algorithms, ranging from TD(0) and TD(λ) for on-policy evaluation to TDC for off-policy evaluation. Combined with a feature approximation framework, we extend the classical TD algorithm to address the critical challenge of step size sensitivity. By reformulating TD updates as fixed point equations, we show that implicit TD enhances robustness in algorithm convergence. Our theoretical contributions include results on mean square convergence and finite-time error bounds of the projected implicit TD algorithms. The proposed algorithms are computationally efficient and scalable, making them well-suited for high-dimensional state spaces, as illustrated in several empirical applications. Looking ahead, we believe that an interesting area for future research is the application of implicit TD algorithms to policy learning and multi-agent RL. Furthermore, implicit algorithms could be beneficial in actor-critic algorithms, where sensitivity to step size specification remains an issue.

Step size	Method	α_1	β_1	RMSPBE	RMSVE
Constant	TDC	0.005	0.05	0.078 ± 0.037	1.938 ± 0.010
	Imp-TDC	0.005	0.05	0.107 ± 0.043	1.445 ± 0.021
	TDC	0.025	0.25	41.754 ± 21.570	71.425 ± 33.891
	Imp-TDC	0.025	0.25	0.284 ± 0.202	0.521 ± 0.176
Decreasing	TDC	0.05	0.50	4.926 ± 2.968	5.270 ± 4.241
	Imp-TDC	0.05	0.50	5.314 ± 0.418	5.164 ± 0.505
	TDC	1.00	10.00	8.219 ± 3.514	59.899 ± 33.652
	Imp-TDC	1.00	10.00	0.851 ± 0.423	1.861 ± 1.139

Table 2: Final root mean square projected Bellman error (RMSPBE) and root mean square value error (RMSVE) for standard and implicit TDC on Baird’s counterexample under constant and decreasing step size schedules. Implicit TDC matches standard TDC at small step size configurations; however, under an aggressive constant schedule ($\alpha_1 = 0.025, \beta_1 = 0.25$) and a decreasing schedule with large initial values ($\alpha_1 = 1.0, \beta_1 = 10.0$), it suppresses errors to near zero, whereas standard TDC exhibits severe amplification.

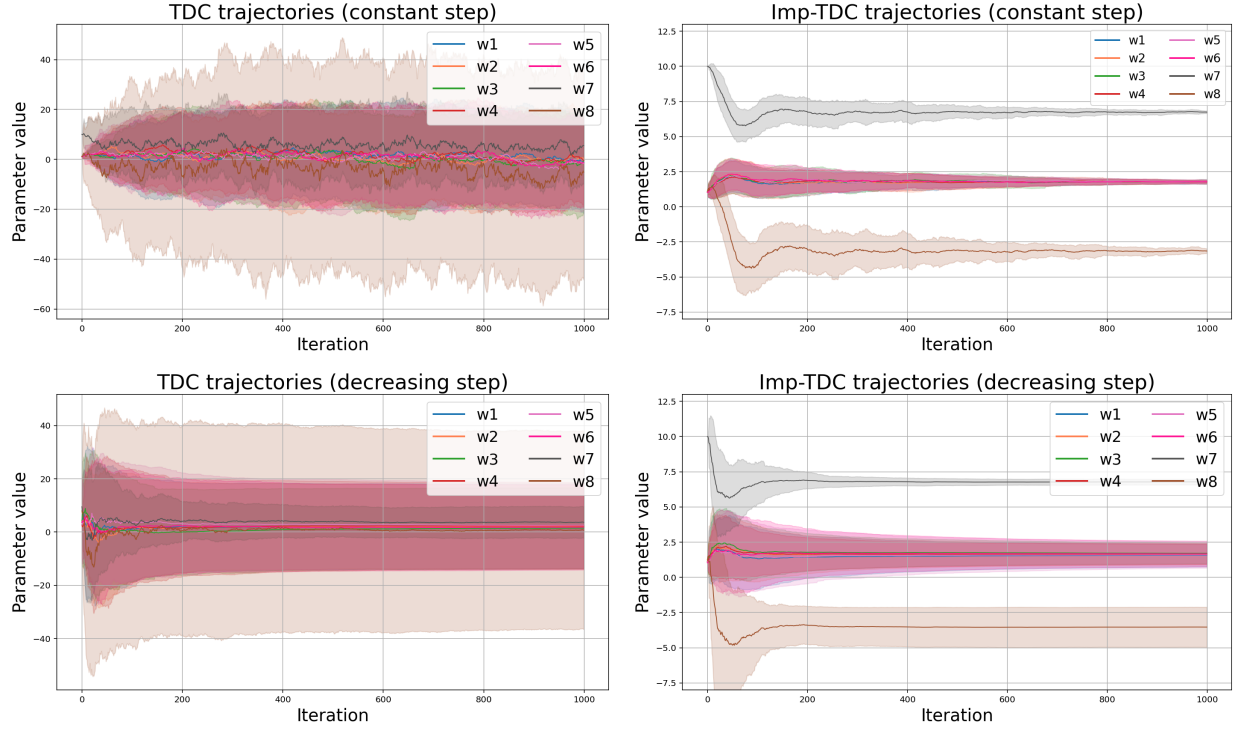


Figure 8: All figures pertain to Baird’s counterexample. Trajectories of the estimated weight parameters for standard TDC (left column) and implicit TDC (right column) on Baird’s counterexample under constant (top row) and decreasing (bottom row) step size schedules. **Top left (standard, constant step):** Large oscillations of the estimated weight parameter trajectories reflect loss of numerical stability in standard TDC under a moderately large constant step size configuration. **Top right (implicit, constant step):** Implicit TDC rapidly drives all weights toward a fixed point, demonstrating superior numerical stability under a moderately large constant step size configuration. **Bottom left (standard, decreasing step):** Standard TDC exhibits large run-to-run variance under a decaying step size schedule with a large initial value. **Bottom right (implicit, decreasing step):** Implicit TDC yields fast convergence with substantially small run-to-run variance under a decaying step size schedule with a large initial value, underscoring its robustness to the choice of step size.

A Proofs of preliminary results

We list and establish foundational lemmas on eligibility trace and implicit update, which will be heavily used in establishing asymptotic convergence as well as finite-time error bounds. Unless explicitly stated, $\|\cdot\|$ implies the Euclidean norm for vectors and their induced norm for matrices.

Lemma A.1. *Given a trace decaying parameter $\lambda \in (0, 1)$ and a discount factor $\gamma \in (0, 1)$, $\|e_n\| \leq \frac{1}{1-\lambda\gamma}$, for all $n \in \mathbb{N}$.*

Proof. Recall that $e_n = \sum_{i=1}^n (\lambda\gamma)^{n-i} \phi_i$. Using the triangle inequality with normalized features, we have

$$\|e_n\| \leq \sum_{i=1}^n (\lambda\gamma)^{n-i} \leq \sum_{i=0}^{\infty} (\lambda\gamma)^i = \frac{1}{1-\lambda\gamma}$$

□

We now provide a proof for Lemma 3.1, which establishes a connection between standard TD algorithms and their implicit counterpart.

Proof of Lemma 3.1. Rearranging terms for the implicit TD(0) update, we have

$$(I + \alpha_n \phi_n \phi_n^T) w_{n+1}^{\text{im}} = w_n^{\text{im}} + \alpha_n (r_n + \gamma \phi_{n+1}^\top w_n^{\text{im}}) \phi_n$$

Multiplying the inverse of $(I + \alpha_n \phi_n \phi_n^T)$ both sides, we get

$$\begin{aligned} w_{n+1}^{\text{im}} &= (I + \alpha_n \phi_n \phi_n^T)^{-1} \left\{ w_n^{\text{im}} + \alpha_n (r_n + \gamma \phi_{n+1}^\top w_n^{\text{im}}) \phi_n \right\} \\ &= \left(I - \frac{\alpha_n}{1 + \alpha_n \|\phi_n\|^2} \phi_n \phi_n^T \right) \left\{ w_n^{\text{im}} + \alpha_n (r_n + \gamma \phi_{n+1}^\top w_n^{\text{im}}) \phi_n \right\}. \end{aligned}$$

where the second equality follows from the Sherman-Morrison-Woodbury identity. Expanding terms out, we have

$$\begin{aligned} w_{n+1}^{\text{im}} &= w_n^{\text{im}} + \alpha_n r_n \phi_n + \alpha_n \gamma \phi_{n+1}^\top w_n^{\text{im}} \phi_n - \frac{\alpha_n}{1 + \alpha_n \|\phi_n\|^2} \phi_n^T w_n^{\text{im}} \phi_n - \frac{\alpha_n^2 r_n \|\phi_n\|^2}{1 + \alpha_n \|\phi_n\|^2} \phi_n - \frac{\alpha_n^2 \gamma \|\phi_n\|^2 \phi_{n+1}^\top w_n^{\text{im}}}{1 + \alpha_n \|\phi_n\|^2} \phi_n \\ &= w_n^{\text{im}} + \alpha_n r_n \left(1 - \frac{\alpha_n \|\phi_n\|^2}{1 + \alpha_n \|\phi_n\|^2} \right) \phi_n + \alpha_n \gamma \phi_{n+1}^\top w_n^{\text{im}} \left(1 - \frac{\alpha_n \|\phi_n\|^2}{1 + \alpha_n \|\phi_n\|^2} \right) \phi_n - \frac{\alpha_n}{1 + \alpha_n \|\phi_n\|^2} \phi_n^T w_n^{\text{im}} \phi_n \\ &= w_n^{\text{im}} + \frac{\alpha_n}{1 + \alpha_n \|\phi_n\|^2} (r_n + \gamma \phi_{n+1}^\top w_n^{\text{im}} - \phi_n^T w_n^{\text{im}}) \phi_n, \end{aligned}$$

where, in the second equality, we collected terms of common factors and obtained a succinct expression in the third equality. Analogously, for the implicit TD(λ) algorithm, we have

$$(I + \alpha_n e_n e_n^T) w_{n+1}^{\text{im}} = w_n^{\text{im}} + \alpha_n (r_n + \gamma \phi_{n+1}^\top w_n^{\text{im}} + \lambda \gamma e_{n-1}^T w_n^{\text{im}}) e_n.$$

Multiplying by inverse of $(I + \alpha_n e_n e_n^T)$, we get

$$w_{n+1}^{\text{im}} = (I + \alpha_n e_n e_n^T)^{-1} \left\{ w_n^{\text{im}} + \alpha_n (r_n + \gamma \phi_{n+1}^\top w_n^{\text{im}} + \lambda \gamma e_{n-1}^T w_n^{\text{im}}) e_n \right\}$$

Using the Sherman-Morrison-Woodbury identity, we get

$$w_{n+1}^{\text{im}} = \left(I - \frac{\alpha_n}{1 + \alpha_n \|e_n\|^2} e_n e_n^T \right) \left\{ w_n^{\text{im}} + \alpha_n (r_n + \gamma \phi_{n+1}^\top w_n^{\text{im}} + \lambda \gamma e_{n-1}^T w_n^{\text{im}}) e_n \right\}.$$

Expanding terms and collecting terms, we have

$$\begin{aligned} w_{n+1}^{\text{im}} &= w_n^{\text{im}} + \alpha_n r_n e_n + \alpha_n \gamma \phi_{n+1}^\top w_n^{\text{im}} e_n + \alpha_n \lambda \gamma e_{n-1}^T w_n^{\text{im}} e_n \\ &\quad - \frac{\alpha_n}{1 + \alpha_n \|e_n\|^2} e_n^T w_n^{\text{im}} e_n - \frac{\alpha_n^2 r_n \|e_n\|^2}{1 + \alpha_n \|e_n\|^2} e_n - \frac{\alpha_n^2 \gamma \|e_n\|^2 \phi_{n+1}^\top w_n^{\text{im}}}{1 + \alpha_n \|e_n\|^2} e_n - \frac{\alpha_n^2 \lambda \gamma \|e_n\|^2 e_{n-1}^T w_n^{\text{im}}}{1 + \alpha_n \|e_n\|^2} e_n \\ &= w_n^{\text{im}} + \left(\alpha_n r_n e_n - \frac{\alpha_n^2 r_n \|e_n\|^2}{1 + \alpha_n \|e_n\|^2} e_n \right) + \left(\alpha_n \gamma \phi_{n+1}^\top w_n^{\text{im}} e_n - \frac{\alpha_n^2 \gamma \|e_n\|^2 \phi_{n+1}^\top w_n^{\text{im}}}{1 + \alpha_n \|e_n\|^2} e_n \right) \\ &\quad + \left(\alpha_n \lambda \gamma e_{n-1}^T w_n^{\text{im}} e_n - \frac{\alpha_n^2 \lambda \gamma \|e_n\|^2 e_{n-1}^T w_n^{\text{im}}}{1 + \alpha_n \|e_n\|^2} e_n \right) - \frac{\alpha_n}{1 + \alpha_n \|e_n\|^2} e_n^T w_n^{\text{im}} e_n \\ &= w_n^{\text{im}} + \alpha_n r_n \left(1 - \frac{\alpha_n \|e_n\|^2}{1 + \alpha_n \|e_n\|^2} \right) e_n + \alpha_n \gamma \phi_{n+1}^\top w_n^{\text{im}} \left(1 - \frac{\alpha_n \|e_n\|^2}{1 + \alpha_n \|e_n\|^2} \right) e_n \\ &\quad + \alpha_n \lambda \gamma e_{n-1}^T w_n^{\text{im}} \left(1 - \frac{\alpha_n \|e_n\|^2}{1 + \alpha_n \|e_n\|^2} \right) e_n - \frac{\alpha_n}{1 + \alpha_n \|e_n\|^2} e_n^T w_n^{\text{im}} e_n \\ &= w_n^{\text{im}} + \frac{\alpha_n}{1 + \alpha_n \|e_n\|^2} (r_n + \gamma \phi_{n+1}^\top w_n^{\text{im}} + \lambda \gamma e_{n-1}^T w_n^{\text{im}} - e_n^T w_n^{\text{im}}) e_n. \end{aligned}$$

□ Next, we provide deterministic upper and lower bounds of the random step size $\tilde{\alpha}_n$.

Lemma A.2. *Given a positive, deterministic non-increasing sequence $(\alpha_n)_{n \in \mathbb{N}}$, the sequence $(\tilde{\alpha}_n)_{n \in \mathbb{N}}$ given by*

$$\tilde{\alpha}_n = \begin{cases} \frac{\alpha_n}{1 + \alpha_n \|\phi_n\|^2} & \text{for } TD(0) \\ \frac{\alpha_n}{1 + \alpha_n \|e_n\|^2} & \text{for } TD(\lambda) \end{cases}$$

respectively satisfy

$$\begin{aligned} \frac{\alpha_n}{1 + \alpha_n} &\leq \tilde{\alpha}_n \leq \alpha_n \quad \text{for } TD(0), \\ \frac{(1 - \lambda \gamma)^2 \alpha_n}{(1 - \lambda \gamma)^2 + \alpha_n} &\leq \tilde{\alpha}_n \leq \alpha_n \quad \text{for } TD(\lambda), \end{aligned}$$

with probability one.

Proof. Since $1 + \alpha_n \|\phi_n\|^2 \geq 1$, we have $\tilde{\alpha}_n \leq \alpha_n$ for $TD(0)$. Analogously $1 + \alpha_n \|e_n\|^2 \geq 1$ implies $\tilde{\alpha}_n \leq \alpha_n$ for $TD(\lambda)$. To prove the lower bounds, note that $\frac{1}{1 + \alpha_n \|\phi_n\|^2} \geq \frac{1}{1 + \alpha_n}$ and $\frac{1}{1 + \alpha_n \|e_n\|^2} \geq \frac{(1 - \lambda \gamma)^2}{(1 - \lambda \gamma)^2 + \alpha_n}$, where the first identity is due to $\|\phi_n\| \leq 1$ and the second identity follows from Lemma

A.1. Therefore, we get

$$\begin{aligned}\tilde{\alpha}_n &\geq \frac{\alpha_n}{1 + \alpha_n} \quad \text{for TD}(0), \\ \tilde{\alpha}_n &\geq \frac{(1 - \lambda\gamma)^2 \alpha_n}{(1 - \lambda\gamma)^2 + \alpha_n} \quad \text{for TD}(\lambda),\end{aligned}$$

with probability one. \square

We next provide a proof of Lemma 3.2 establishing the relationship between the standard TDC and the implicit TDC algorithm.

Proof of Lemma 3.2. The implicit TDC update for the target parameter w^{im} can be rewritten as

$$\begin{aligned}(I + \alpha_n \rho_n \phi_n \phi_n^T) w_{n+1}^{\text{im}} &= w_n^{\text{im}} + \alpha_n \rho_n (r_n \phi_n + \gamma \phi_n \phi_{n+1}^T w_n^{\text{im}} - \gamma \phi_{n+1} \phi_n^T u_n^{\text{im}}) \\ \Leftrightarrow w_{n+1}^{\text{im}} &= (I + \alpha_n \rho_n \phi_n \phi_n^T)^{-1} \{w_n^{\text{im}} + \alpha_n \rho_n (r_n \phi_n + \gamma \phi_n \phi_{n+1}^T w_n^{\text{im}} - \gamma \phi_{n+1} \phi_n^T u_n^{\text{im}})\} \\ \Leftrightarrow w_{n+1}^{\text{im}} &= (I - \alpha'_n \rho_n \phi_n \phi_n^T) \{w_n^{\text{im}} + \alpha_n \rho_n (r_n \phi_n + \gamma \phi_n \phi_{n+1}^T w_n^{\text{im}} - \gamma \phi_{n+1} \phi_n^T u_n^{\text{im}})\} \quad (16)\end{aligned}$$

for $\alpha'_n = \frac{\alpha_n}{1 + \alpha_n \rho_n \|\phi_n\|^2}$. Expanding the right hand side of (16), we have

$$\begin{aligned}w_{n+1}^{\text{im}} &= w_n^{\text{im}} + \alpha_n \rho_n r_n \phi_n + \alpha_n \rho_n \gamma (\phi_{n+1}^T w_n^{\text{im}}) \phi_n - \alpha_n \rho_n \gamma (\phi_n^T u_n^{\text{im}}) \phi_{n+1} \\ &\quad - \alpha'_n \rho_n \alpha_n \rho_n r_n \|\phi_n\|^2 \phi_n - \alpha'_n \rho_n \alpha_n \rho_n \gamma \|\phi_n\|^2 (\phi_{n+1}^T w_n^{\text{im}}) \phi_n + \alpha'_n \rho_n \alpha_n \rho_n \gamma (\phi_n^T \phi_{n+1}) (\phi_n^T u_n^{\text{im}}) \phi_n \\ &\quad - \alpha'_n \rho_n (\phi_n^T w_n^{\text{im}}) \phi_n \\ &= w_n^{\text{im}} + \alpha_n \rho_n r_n (1 - \alpha'_n \rho_n \|\phi_n\|^2) \phi_n + \alpha_n \rho_n \gamma (1 - \alpha'_n \rho_n \|\phi_n\|^2) (\phi_{n+1}^T w_n^{\text{im}}) \phi_n - \alpha'_n \rho_n (\phi_n^T w_n^{\text{im}}) \phi_n \\ &\quad - \alpha_n \rho_n \gamma (\phi_n^T u_n^{\text{im}}) \phi_{n+1} + \alpha'_n \rho_n \alpha_n \rho_n \gamma (\phi_n^T \phi_{n+1}) (\phi_n^T u_n^{\text{im}}) \phi_n \\ &= w_n^{\text{im}} + \alpha'_n \rho_n \delta_n^{\text{im}} \phi_n - \alpha_n \rho_n \gamma (\phi_n^T u_n^{\text{im}}) \{\phi_{n+1} - \alpha'_n \rho_n (\phi_n^T \phi_{n+1}) \phi_n\}\end{aligned}$$

where $\delta_n^{\text{im}} \phi_n := r_n \phi_n + \gamma \phi_n \phi_{n+1}^T w_n^{\text{im}} - \phi_n \phi_n^T w_n^{\text{im}}$. An implicit TDC update for auxiliary parameter u^{im} can be analogously derived. Consider

$$\begin{aligned}(I + \beta_n \rho_n \phi_n \phi_n^T) u_{n+1}^{\text{im}} &= u_n^{\text{im}} + \beta_n \rho_n \delta_n^{\text{im}} \phi_n \Leftrightarrow u_{n+1}^{\text{im}} = (I + \beta_n \rho_n \phi_n \phi_n^T)^{-1} (u_n^{\text{im}} + \beta_n \rho_n \delta_n^{\text{im}} \phi_n) \\ \Leftrightarrow u_{n+1}^{\text{im}} &= (I - \beta'_n \rho_n \phi_n \phi_n^T) (u_n^{\text{im}} + \beta_n \rho_n \delta_n^{\text{im}} \phi_n) \quad (17)\end{aligned}$$

for $\beta'_n = \frac{\beta_n}{1 + \beta_n \rho_n \|\phi_n\|^2}$. Re-expressing the right hand side of (17), we obtain

$$\begin{aligned}u_{n+1}^{\text{im}} &= u_n^{\text{im}} + \beta_n \rho_n \delta_n^{\text{im}} \phi_n - \beta'_n \rho_n \phi_n^T u_n^{\text{im}} \phi_n - \beta'_n \rho_n \beta_n \rho_n \|\phi_n\|^2 \delta_n^{\text{im}} \phi_n \\ &= u_n^{\text{im}} - \beta'_n \rho_n \phi_n^T u_n^{\text{im}} \phi_n + \beta_n \rho_n \delta_n^{\text{im}} \phi_n (1 - \beta'_n \rho_n \|\phi_n\|^2) \\ &= u_n^{\text{im}} + \beta'_n \rho_n \delta_n^{\text{im}} \phi_n - \beta'_n \rho_n \phi_n \phi_n^T u_n^{\text{im}}\end{aligned}$$

\square

B Theoretical analysis for implicit TD(0) and TD(λ)

We only deal with time-homogeneous Markov processes whose steady-state distribution is well-defined. To simplify our presentation, for the TD(0) algorithm, let us define

$$\begin{aligned} S_n(w) &:= r_n \phi_n + \gamma \phi_n \phi_{n+1}^T w - \phi_n \phi_n^T w = b_n + A_n w, \\ S(w) &:= \mathbb{E}_\infty \{r_n \phi_n\} + \mathbb{E}_\infty \{\gamma \phi_n \phi_{n+1}^T\} w - \mathbb{E}_\infty \{\phi_n \phi_n^T\} w = b + A w, \end{aligned}$$

where $A_n = \gamma \phi_n \phi_{n+1}^T - \phi_n \phi_n^T$, $A = \mathbb{E}_\infty \{A_n\}$, $b_n = r_n \phi_n$, $b = \mathbb{E}_\infty \{b_n\}$. Here \mathbb{E}_∞ is the expectation with respect to the steady-state distribution of the Markov process $(x_n)_{n \in \mathbb{N}}$. Similarly, for the TD(λ) algorithm,

$$\begin{aligned} S_n(w) &:= r_n e_n + \gamma e_n \phi_{n+1}^T w - e_n \phi_n^T w = b_n + A_n w, \\ S(w) &:= \mathbb{E}_\infty \{r_n e_{-\infty:n}\} + \mathbb{E}_\infty \{\gamma e_{-\infty:n} \phi_{n+1}^T\} w - \mathbb{E}_\infty \{e_{-\infty:n} \phi_n^T\} w = b + A w, \end{aligned}$$

where $e_{-\infty:n} := \sum_{k=0}^{\infty} (\lambda \gamma)^k \phi_{n-k}$ represents the steady-space eligibility trace and $A_n = \gamma e_n \phi_{n+1}^T - e_n \phi_n^T$, $A = \mathbb{E}_\infty \{\gamma e_{-\infty:n} \phi_{n+1}^T\} - \mathbb{E}_\infty \{e_{-\infty:n} \phi_n^T\} = \lim_{n \rightarrow \infty} \mathbb{E} \{A_n\}$, $b_n = r_n e_n$ and $b = \mathbb{E}_\infty \{r_n e_{-\infty:n}\} = \lim_{n \rightarrow \infty} \mathbb{E} \{b_n\}$. Tsitsiklis and Roy [44] has shown that the limit point of TD algorithms, denoted by w_* solves the equation $S(w) = 0$. Before establishing the asymptotic convergence of implicit TD(0) and TD(λ), we first provide bounds between A_n and A as well as b_n and b utilizing the geometric mixing condition induced by assumption 4.1.

Lemma B.1. *From Assumption 4.1, for every $n, \tau \geq 0$, $n \geq \tau$, there exists some $\tilde{\rho} \in [0, 1)$ and a constant \tilde{m} , such that*

- $\|\mathbb{E} \{A_n | X_{n-\tau} = x\} - A\| \leq \tilde{m} \tilde{\rho}^\tau$
- $\|\mathbb{E} \{b_n | X_{n-\tau} = x\} - b\| \leq \tilde{m} \tilde{\rho}^\tau$.

Proof. Due to time-homogeneity of transition probabilities, the statement is equivalent to Lemma 6.7 in [3]. \square

Let us define a mixing time for A_n and b_n like we did for the underlying Markov process.

Definition B.2. *Given a threshold $\epsilon > 0$, the mixing time for A_n and b_n is given by*

$$\tilde{\tau}_\epsilon = \min\{n \in \mathbb{N} \mid \tilde{m} \tilde{\rho}^n \leq \epsilon\}.$$

B.1 Asymptotic convergence analysis for implicit TD(0) and TD(λ)

We closely follow the approach taken in [33] with a few modifications made to accommodate the data-adaptive step size of implicit TD algorithms. For the analysis of implicit algorithms, we focus on the step sizes $(\alpha_n)_{n \in \mathbb{N}}$ satisfying the following condition: 1) $(\alpha_n)_{n \in \mathbb{N}}$ is a non-increasing sequence and 2) there exists $n^* > 0$ and $\kappa \geq 1$ such that for any $n \geq n^*$, we have $n - \tilde{\tau}_{\alpha_n} > 0$,

$\alpha_{n-\tilde{\tau}_{\alpha_n}} \tilde{\tau}_{\alpha_n} \leq \frac{1}{4c_\lambda}$, $c_\lambda := \frac{2}{1-\lambda\gamma} \geq 1$ and $\alpha_{n-\tilde{\tau}_{\alpha_n}} \leq \kappa\alpha_n$. Notice the step size sequence $\alpha_n = cn^{-s}$, for some $c > 0, s \in (0.5, 1]$ satisfy these conditions. From Assumption 4.1 and Lemma B.1, we have $\tilde{\tau}_{\alpha_n} = O(\log n)$. Therefore, we know $n - \tilde{\tau}_{\alpha_n} \rightarrow \infty$ and $\tilde{\tau}_{\alpha_n}/(n - \tilde{\tau}_{\alpha_n})^s \rightarrow 0$. Furthermore, we have $\alpha_{n-\tilde{\tau}_{\alpha_n}}/\alpha_n = \{n/(n - \tilde{\tau}_{\alpha_n})\}^s$, which converges to 1 as $n \rightarrow \infty$. Hence, for large $n \in \mathbb{N}$, there must exist $\kappa \geq 1$ satisfying the above condition. We begin listing preliminary results needed to prove the asymptotic convergence results. To simplify notations, we use $\theta_n := w_* - w_n^{\text{im}}$. We first introduce upper bounds for the norm of the TD update direction.

Lemma B.3. *For all $n \in \mathbb{N}$,*

$$\|A_n\| \leq c_\lambda := \frac{2}{1-\lambda\gamma},$$

for both $TD(0)$ and $TD(\lambda)$. Furthermore, for all $n \in \mathbb{N}$,

$$\|A_n w_* + b_n\| \leq S_{\max} := \frac{2\|w_*\| + r_{\max}}{1-\lambda\gamma},$$

with probability one.

Proof. Notice that

$$\|A_n\| = \begin{cases} \|\gamma\phi_n\phi_{n+1}^T - \phi_n\phi_n^T\| \leq (\gamma+1) & \text{for } TD(0), \\ \|\gamma e_n\phi_{n+1}^T - e_n\phi_n^T\| \leq \frac{\gamma+1}{1-\lambda\gamma} & \text{for } TD(\lambda), \end{cases}$$

which can be deduced from the normalized features assumption and Lemma A.1 with the triangle inequality. The first statement is the direct consequence of the facts $\gamma < 1$ and $\frac{1}{1-\lambda\gamma} > 1$. In a similar vein, recall that

$$\|A_n w_* + b_n\| = \begin{cases} \|\gamma\phi_n\phi_{n+1}^T w_* - \phi_n\phi_n^T w_* + r_n\phi_n\| \leq (\gamma+1)\|w_*\| + r_{\max} & \text{for } TD(0), \\ \|\gamma e_n\phi_{n+1}^T w_* - e_n\phi_n^T w_* + r_n e_n\| \leq \frac{(\gamma+1)\|w_*\| + r_{\max}}{1-\lambda\gamma} & \text{for } TD(\lambda), \end{cases}$$

which follow from the normalized features, bounded reward assumptions, and Lemma A.1 with the triangle inequality. Since $\gamma < 1$ and $\frac{1}{1-\lambda\gamma} > 1$, we get the second statement. \square

Lemma B.4. *Let $n \geq n^*$ with $\ell = n - \tilde{\tau}_{\alpha_n}$. The following statements hold*

1. $\|\theta_n - \theta_\ell\| \leq 2c_\lambda\alpha_\ell\tilde{\tau}_{\alpha_n}(\|\theta_\ell\| + S_{\max}),$
2. $\|\theta_n - \theta_\ell\| \leq 4c_\lambda\alpha_\ell\tilde{\tau}_{\alpha_n}(\|\theta_n\| + S_{\max}),$
3. $\|\theta_n - \theta_\ell\|^2 \leq 32c_\lambda^2\alpha_\ell^2\tilde{\tau}_{\alpha_n}^2(\|\theta_n\|^2 + S_{\max}^2) \leq 8c_\lambda\alpha_\ell\tilde{\tau}_{\alpha_n}(\|\theta_n\|^2 + S_{\max}^2).$

with probability one.

Proof. Statement 1: We begin proving the first statement. For $\ell < t \leq n$, note that

$$\begin{aligned}
\theta_t &:= w_t^{\text{im}} - w_* \\
&= w_{t-1}^{\text{im}} - w_* + \tilde{\alpha}_{t-1}(A_{t-1}w_{t-1}^{\text{im}} + b_{t-1}) \\
&= w_{t-1}^{\text{im}} - w_* + \tilde{\alpha}_{t-1}A_{t-1}(w_{t-1}^{\text{im}} - w_*) + \tilde{\alpha}_{t-1}(A_{t-1}w_* + b_{t-1}) \\
&= \theta_{t-1} + \tilde{\alpha}_{t-1}(A_{t-1}\theta_{t-1} + A_{t-1}w_* + b_{t-1}),
\end{aligned}$$

where in the second line, we use the definition of w_t^{im} , and in the third line, we add and subtract $\tilde{\alpha}_{t-1}A_{t-1}w_*$. The last line is due to the definition of θ_{t-1} . Therefore, we have

$$\begin{aligned}
\|\theta_t - \theta_{t-1}\| &= \|\tilde{\alpha}_{t-1}(A_{t-1}\theta_{t-1} + A_{t-1}w_* + b_{t-1})\| \\
&\leq \alpha_{t-1} \|A_{t-1}\theta_{t-1} + A_{t-1}w_* + b_{t-1}\| \\
&\leq \alpha_{t-1}(c_\lambda \|\theta_{t-1}\| + S_{\max}),
\end{aligned} \tag{18}$$

where the first inequality follows from Lemma A.2 and in the second inequality, we used Lemma B.3 with the triangle inequality. Using the reverse triangle inequality, we get

$$\begin{aligned}
\|\theta_t\| &\leq (1 + c_\lambda \alpha_{t-1})\|\theta_{t-1}\| + \alpha_{t-1}S_{\max} \\
&\leq (1 + c_\lambda \alpha_{t-1}) \cdots (1 + c_\lambda \alpha_\ell) \|\theta_\ell\| + (1 + c_\lambda \alpha_{t-1}) \cdots (1 + c_\lambda \alpha_{\ell+1}) \alpha_\ell S_{\max} \\
&\quad + \cdots + (1 + c_\lambda \alpha_{t-1}) \alpha_{t-2} S_{\max} + \alpha_{t-1} S_{\max},
\end{aligned} \tag{19}$$

and the second inequality follows from recursive applications of (19). Thanks to the non-increasingness of $(\alpha_n)_{n \in \mathbb{N}}$, we know $(1 + c_\lambda \alpha_k) \leq 1 + c_\lambda \alpha_\ell$, $\alpha_k \leq \alpha_\ell$ for all $k \leq \ell$, which give us

$$\begin{aligned}
\|\theta_t\| &\leq (1 + c_\lambda \alpha_\ell)^{t-\ell} \|\theta_\ell\| + (1 + c_\lambda \alpha_\ell)^{t-\ell-1} \alpha_\ell S_{\max} + (1 + c_\lambda \alpha_\ell)^{t-\ell-2} \alpha_\ell S_{\max} \\
&\quad + \cdots + (1 + c_\lambda \alpha_\ell) \alpha_\ell S_{\max} + \alpha_\ell S_{\max} \\
&= (1 + c_\lambda \alpha_\ell)^{t-\ell} \|\theta_\ell\| + \left\{ \frac{(1 + c_\lambda \alpha_\ell)^{t-\ell} - 1}{c_\lambda} \right\} S_{\max} \\
&\leq (1 + c_\lambda \alpha_\ell)^{\tilde{\tau}_{\alpha_n}} \|\theta_\ell\| + \left\{ \frac{(1 + c_\lambda \alpha_\ell)^{\tilde{\tau}_{\alpha_n}} - 1}{c_\lambda} \right\} S_{\max},
\end{aligned} \tag{20}$$

where the last inequality is due to $t - \ell \leq n - \ell = \tilde{\tau}_{\alpha_n}$. Recall from the choice of step size, we know $\alpha_\ell \tilde{\tau}_{\alpha_n} \leq \frac{1}{4c_\lambda}$, which gives us $c_\lambda \alpha_\ell \leq \frac{1}{4\tilde{\tau}_{\alpha_n}} \leq \frac{\log 2}{\tilde{\tau}_{\alpha_n} - 1}$. Furthermore, for $x \leq \frac{\log 2}{\tilde{\tau}_{\alpha_n} - 1}$, one can show that $(1 + x)^{\tilde{\tau}_{\alpha_n}} \leq 1 + 2x\tilde{\tau}_{\alpha_n}$. Therefore, we have $(1 + c_\lambda \alpha_\ell)^{\tilde{\tau}_{\alpha_n}} \leq 1 + 2c_\lambda \alpha_\ell \tilde{\tau}_{\alpha_n}$. Plugging this upper bound back in (20), we get

$$\|\theta_t\| \leq (1 + 2c_\lambda \alpha_\ell \tilde{\tau}_{\alpha_n}) \|\theta_\ell\| + 2\alpha_\ell \tilde{\tau}_{\alpha_n} S_{\max} \leq 2\|\theta_\ell\| + 2\alpha_\ell \tilde{\tau}_{\alpha_n} S_{\max}, \tag{21}$$

where the last inequality follows from the fact that $c_\lambda \alpha_\ell \leq \frac{1}{4\tilde{\tau}_{\alpha_n}}$.

We now obtain the upper bound of $\|\theta_n - \theta_\ell\|$. Notice that

$$\|\theta_n - \theta_\ell\| \leq \sum_{t=\ell}^{n-1} \|\theta_{t+1} - \theta_t\| \leq \sum_{t=\ell}^{n-1} \alpha_t (c_\lambda \|\theta_t\| + S_{\max}) \leq c_\lambda \alpha_\ell \left\{ \sum_{t=\ell}^{n-1} \|\theta_t\| \right\} + \alpha_\ell (n - \ell) S_{\max},$$

where the first inequality follows from the triangle inequality, the second inequality is due to (18), and the third inequality is thanks to the non-increasingness of the sequence step size sequence. Plugging the bound we obtained in (21), we get

$$\begin{aligned} \|\theta_n - \theta_\ell\| &\leq c_\lambda \alpha_\ell \tilde{\tau}_{\alpha_n} (2\|\theta_\ell\| + 2\alpha_\ell \tilde{\tau}_{\alpha_n} S_{\max}) + \alpha_\ell \tilde{\tau}_{\alpha_n} S_{\max} \\ &= 2c_\lambda \alpha_\ell \tilde{\tau}_{\alpha_n} \|\theta_\ell\| + 2c_\lambda \alpha_\ell^2 \tilde{\tau}_{\alpha_n}^2 S_{\max} + \alpha_\ell \tilde{\tau}_{\alpha_n} S_{\max} \\ &\leq 2c_\lambda \alpha_\ell \tilde{\tau}_{\alpha_n} \|\theta_\ell\| + c_\lambda \alpha_\ell \tilde{\tau}_{\alpha_n} S_{\max} + c_\lambda \alpha_\ell \tilde{\tau}_{\alpha_n} S_{\max} \\ &= 2c_\lambda \alpha_\ell \tilde{\tau}_{\alpha_n} \|\theta_\ell\| + 2c_\lambda \alpha_\ell \tilde{\tau}_{\alpha_n} S_{\max}, \end{aligned} \tag{22}$$

where the second inequality is due to positivity of $\alpha_\ell \tilde{\tau}_{\alpha_n} S_{\max}$ with $2\alpha_\ell \tilde{\tau}_{\alpha_n} \leq 1$ and $c_\lambda \geq 1$.

Statement 2: From the triangle inequality, we know $\|\theta_\ell\| \leq \|\theta_n - \theta_\ell\| + \|\theta_n\|$. Plugging this to (22), we get

$$\|\theta_n - \theta_\ell\| \leq 2c_\lambda \alpha_\ell \tilde{\tau}_{\alpha_n} \|\theta_n - \theta_\ell\| + 2c_\lambda \alpha_\ell \tilde{\tau}_{\alpha_n} \|\theta_n\| + 2c_\lambda \alpha_\ell \tilde{\tau}_{\alpha_n} S_{\max}.$$

With the fact $\alpha_\ell \tilde{\tau}_{\alpha_n} \leq \frac{1}{4c_\lambda}$, we get

$$\|\theta_n - \theta_\ell\| \leq \frac{1}{2} \|\theta_n - \theta_\ell\| + 2c_\lambda \alpha_\ell \tilde{\tau}_{\alpha_n} \|\theta_n\| + 2c_\lambda \alpha_\ell \tilde{\tau}_{\alpha_n} S_{\max}.$$

Subtracting $\frac{1}{2} \|\theta_n - \theta_\ell\|$ from both sides and multiplying by two, we get

$$\|\theta_n - \theta_\ell\| \leq 4c_\lambda \alpha_\ell \tilde{\tau}_{\alpha_n} \|\theta_n\| + 4c_\lambda \alpha_\ell \tilde{\tau}_{\alpha_n} S_{\max}. \tag{23}$$

Statement 3: Applying $(a + b)^2 \leq 2a^2 + 2b^2$ to (23), we have

$$\|\theta_n - \theta_\ell\|^2 \leq 32c_\lambda^2 \alpha_\ell^2 \tilde{\tau}_{\alpha_n}^2 \|\theta_n\|^2 + 32c_\lambda^2 \alpha_\ell^2 \tilde{\tau}_{\alpha_n}^2 S_{\max}^2 \leq 8c_\lambda \alpha_\ell \tilde{\tau}_{\alpha_n} \|\theta_n\|^2 + 8c_\lambda \alpha_\ell \tilde{\tau}_{\alpha_n} S_{\max}^2,$$

where the last inequality follows from the fact $\alpha_\ell \tilde{\tau}_{\alpha_n} \leq \frac{1}{4c_\lambda}$. □

Lemma B.5. For $n \geq n^*$, $\ell = n - \tilde{\tau}_{\alpha_n}$ with $A = \begin{cases} \mathbb{E}_\infty \{ \gamma \phi_n \phi_{n+1}^T - \phi_n \phi_n^T \} & \text{for } TD(0) \\ \mathbb{E}_\infty \{ \gamma e_n \phi_{n+1}^T - e_n \phi_n^T \} & \text{for } TD(\lambda) \end{cases}$

$$\left| \mathbb{E} \left\{ \theta_n^T (\theta_{n+1} - \theta_n - \tilde{\alpha}_n A \theta_n) \middle| \theta_\ell, x_\ell \right\} \right| \leq c_1 \alpha_n^2 \tilde{\tau}_{\alpha_n} \mathbb{E} \{ \|\theta_n\|^2 | \theta_\ell, x_\ell \} + c_2 \alpha_n^2 \tilde{\tau}_{\alpha_n},$$

for some constants $c_1, c_2 > 0$.

Proof. Recall that

$$\begin{aligned}
\theta_{n+1} &= w_{n+1}^{\text{im}} - w_* \\
&= w_n^{\text{im}} - w_* + \tilde{\alpha}_n(A_n w_n^{\text{im}} + b_n) \\
&= w_n^{\text{im}} - w_* + \tilde{\alpha}_n A_n (w_n^{\text{im}} - w_*) + \tilde{\alpha}_n (A_n w_* + b_n) \\
&= \theta_n + \tilde{\alpha}_n (A_n \theta_n + A_n w_* + b_n),
\end{aligned}$$

where in the first and last equality, we used the definition of θ_n , and the second equality is due to the definition of w_{n+1}^{im} . The third equality follows from adding and subtracting $\tilde{\alpha}_n A_n w_*$ and the last equality is due to the definition of θ_n . Then, we have

$$\begin{aligned}
\mathbb{E} \left\{ \theta_n^T (\theta_{n+1} - \theta_n - \tilde{\alpha}_n A \theta_n) \middle| \theta_\ell, x_\ell \right\} &= \mathbb{E} \left\{ \tilde{\alpha}_n \theta_n^T (A_n \theta_n + A_n w_* + b_n - A \theta_n) \middle| \theta_\ell, x_\ell \right\} \\
&= \mathbb{E} \left\{ \tilde{\alpha}_n \theta_n^T (A_n w_* + b_n) \middle| \theta_\ell, x_\ell \right\} + \mathbb{E} \left\{ \tilde{\alpha}_n \theta_n^T (A_n - A) \theta_n \middle| \theta_\ell, x_\ell \right\}.
\end{aligned} \tag{24}$$

We will now provide an upper bound of each term in (24).

Step 1: Let us first consider the leading term in (24). Recall that $\frac{\alpha_n}{1+\alpha_n} < \tilde{\alpha}_n \leq \alpha_n$ holds almost surely for TD(0). Since

$$\begin{aligned}
\mathbb{E} \left\{ \tilde{\alpha}_n \theta_n^T (A_n w_* + b_n) \middle| \theta_\ell, x_\ell \right\} &\leq \max \left[\frac{\alpha_n}{1+\alpha_n} \mathbb{E} \left\{ \theta_n^T (A_n w_* + b_n) \middle| \theta_\ell, x_\ell \right\}, \alpha_n \mathbb{E} \left\{ \theta_n^T (A_n w_* + b_n) \middle| \theta_\ell, x_\ell \right\} \right], \\
\mathbb{E} \left\{ \tilde{\alpha}_n \theta_n^T (A_n w_* + b_n) \middle| \theta_\ell, x_\ell \right\} &\geq \min \left[\frac{\alpha_n}{1+\alpha_n} \mathbb{E} \left\{ \theta_n^T (A_n w_* + b_n) \middle| \theta_\ell, x_\ell \right\}, \alpha_n \mathbb{E} \left\{ \theta_n^T (A_n w_* + b_n) \middle| \theta_\ell, x_\ell \right\} \right],
\end{aligned}$$

we know

$$\left| \mathbb{E} \left\{ \tilde{\alpha}_n \theta_n^T (A_n w_* + b_n) \middle| \theta_\ell, x_\ell \right\} \right| \leq \alpha_n \left| \mathbb{E} \left\{ \theta_n^T (A_n w_* + b_n) \middle| \theta_\ell, x_\ell \right\} \right|.$$

The same holds for TD(λ) almost surely, with $\frac{\alpha_n}{1+\alpha_n}$ replaced by $\frac{(1-\lambda\gamma)\alpha_n}{(1-\lambda\gamma)^2+\alpha_n}$. Therefore, for both TD(0) and TD(λ), we get

$$\begin{aligned}
\left| \mathbb{E} \left\{ \tilde{\alpha}_n \theta_n^T (A_n w_* + b_n) \middle| \theta_\ell, x_\ell \right\} \right| &\leq \alpha_n \left| \mathbb{E} \left\{ \theta_n^T (A_n w_* + b_n) \middle| \theta_\ell, x_\ell \right\} \right| \\
&= \alpha_n \left| \mathbb{E} \left\{ \theta_\ell^T (A_n w_* + b_n) \middle| \theta_\ell, x_\ell \right\} + \mathbb{E} \left\{ (\theta_n - \theta_\ell)^T (A_n w_* + b_n) \middle| \theta_\ell, x_\ell \right\} \right| \\
&\stackrel{(i)}{\leq} \alpha_n \left| \theta_\ell^T \mathbb{E} \left\{ (A_n w_* + b_n) \middle| \theta_\ell, x_\ell \right\} \right| + \alpha_n \mathbb{E} \left\{ \|\theta_n - \theta_\ell\| \|A_n w_* + b_n\| \middle| \theta_\ell, x_\ell \right\} \\
&\stackrel{(ii)}{\leq} \alpha_n \|\theta_\ell\| \left\| \mathbb{E} \left\{ (A_n w_* + b_n) \middle| \theta_\ell, x_\ell \right\} \right\| + \alpha_n \mathbb{E} \left\{ \|\theta_n - \theta_\ell\| \middle| \theta_\ell, x_\ell \right\} S_{\max},
\end{aligned} \tag{25}$$

where (i) follows from the linearity of expectation with the Cauchy-Schwarz and triangle inequality,

(ii) from the Cauchy-Schwarz inequality with the fact $\|A_n w_* + b_n\| \leq S_{\max}$. Furthermore, note that

$$\begin{aligned} \left\| \mathbb{E} \left\{ (A_n w_* + b_n) \middle| \theta_\ell, x_\ell \right\} \right\| &= \left\| \mathbb{E} \left\{ (A_n w_* + b_n) \middle| \theta_\ell, x_\ell \right\} - (A w_* + b) \right\| \\ &\leq \left\| \mathbb{E} \left\{ A_n \middle| \theta_\ell, x_\ell \right\} - A \right\| \|w_*\| + \left\| \mathbb{E} \left\{ b_n \middle| \theta_\ell, x_\ell \right\} - b \right\| \\ &\leq \alpha_n (\|w_*\| + 1), \end{aligned} \quad (26)$$

where in the first inequality, we used the fact $A w_* + b = 0$, the second inequality follows from the triangle inequality, and for the last inequality, we used Lemma B.1. Plugging (26) into (25) and invoking Lemma B.4, we get

$$\begin{aligned} \left| \mathbb{E} \left\{ \tilde{\alpha}_n \theta_n^T (A_n w_* + b_n) \middle| \theta_\ell, x_\ell \right\} \right| &\leq \alpha_n^2 (\|w_*\| + 1) \|\theta_\ell\| + 2c_\lambda \alpha_n \alpha_\ell \tilde{\tau}_{\alpha_n} (\|\theta_\ell\| + S_{\max}) S_{\max} \\ &\leq \alpha_\ell^2 (\|w_*\| + 1) \|\theta_\ell\| + 2c_\lambda \alpha_\ell^2 \tilde{\tau}_{\alpha_n} (\|\theta_\ell\| + S_{\max}) S_{\max} \\ &= \alpha_\ell^2 c_{w_*} \|\theta_\ell\| + 2c_\lambda \alpha_\ell^2 \tilde{\tau}_{\alpha_n} (\|\theta_\ell\| + S_{\max}) S_{\max} \end{aligned} \quad (27)$$

where the second inequality follows from the fact that $\alpha_n \leq \alpha_\ell$ since $n \leq \ell$ and the last equality follows from the definition $c_{w_*} := \|w_*\| + 1$. Note that by definition $c_{w_*} \leq S_{\max} + 1$, where $S_{\max} = \frac{2\|w_*\| + r_{\max}}{1 - \lambda\gamma}$.

Step 2: Next we bound the second term, which can be re-expressed as

$$\mathbb{E} \left\{ \tilde{\alpha}_n \theta_n^T (A_n - A) \theta_n \middle| \theta_\ell, x_\ell \right\} = \mathbb{E} \left\{ \tilde{\alpha}_n \theta_\ell^T (A_n - A) \theta_\ell \middle| \theta_\ell, x_\ell \right\} \quad (28)$$

$$+ \mathbb{E} \left\{ \tilde{\alpha}_n (\theta_n - \theta_\ell)^T (A_n - A) (\theta_n - \theta_\ell) \middle| \theta_\ell, x_\ell \right\} \quad (29)$$

$$+ \mathbb{E} \left\{ \tilde{\alpha}_n (\theta_n - \theta_\ell)^T (A_n - A) \theta_\ell \middle| \theta_\ell, x_\ell \right\} \quad (30)$$

$$+ \mathbb{E} \left\{ \tilde{\alpha}_n \theta_\ell^T (A_n - A) (\theta_n - \theta_\ell) \middle| \theta_\ell, x_\ell \right\}. \quad (31)$$

To get a bound for the term in (28), recall that, for TD(0),

$$\begin{aligned} \mathbb{E} \left\{ \tilde{\alpha}_n \theta_\ell^T (A_n - A) \theta_\ell \middle| \theta_\ell, x_\ell \right\} &\leq \max \left[\alpha_n \mathbb{E} \left\{ \theta_\ell^T (A_n - A) \theta_\ell \middle| \theta_\ell, x_\ell \right\}, \frac{\alpha_n}{1 + \alpha_n} \mathbb{E} \left\{ \theta_\ell^T (A_n - A) \theta_\ell \middle| \theta_\ell, x_\ell \right\} \right] \\ \mathbb{E} \left\{ \tilde{\alpha}_n \theta_\ell^T (A_n - A) \theta_\ell \middle| \theta_\ell, x_\ell \right\} &\geq \min \left[\alpha_n \mathbb{E} \left\{ \theta_\ell^T (A_n - A) \theta_\ell \middle| \theta_\ell, x_\ell \right\}, \frac{\alpha_n}{1 + \alpha_n} \mathbb{E} \left\{ \theta_\ell^T (A_n - A) \theta_\ell \middle| \theta_\ell, x_\ell \right\} \right] \end{aligned}$$

from which we have

$$\left| \mathbb{E} \left\{ \tilde{\alpha}_n \theta_\ell^T (A_n - A) \theta_\ell \middle| \theta_\ell, x_\ell \right\} \right| \leq \alpha_n \left| \mathbb{E} \left\{ \theta_\ell^T (A_n - A) \theta_\ell \middle| \theta_\ell, x_\ell \right\} \right|.$$

Again, the result holds for TD(λ) by the same argument with $\frac{\alpha_n}{1 + \alpha_n}$ replaced by $\frac{(1 - \lambda\gamma)^2 \alpha_n}{(1 - \lambda\gamma)^2 + \alpha_n}$. Applying the Cauchy-Schwarz inequality with Lemma B.1, we get

$$\left| \mathbb{E} \left\{ \tilde{\alpha}_n \theta_\ell^T (A_n - A) \theta_\ell \middle| \theta_\ell, x_\ell \right\} \right| \leq \alpha_n \|\theta_\ell\|^2 \|\mathbb{E}[A_n | x_\ell] - A\| \leq \alpha_n^2 \|\theta_\ell\|^2. \quad (32)$$

From the Cauchy-Schwarz inequality and triangle inequality, we get the bound for the second term in (29), given by

$$\begin{aligned} \left| \mathbb{E} \left\{ \tilde{\alpha}_n (\theta_n - \theta_\ell)^T (A_n - A) (\theta_n - \theta_\ell) \middle| \theta_\ell, x_\ell \right\} \right| &\leq \alpha_n \mathbb{E} \left\{ \|\theta_n - \theta_\ell\|^2 (\|A_n\| + \|A\|) \middle| \theta_\ell, x_\ell \right\} \\ &\leq 2c_\lambda \alpha_n \mathbb{E} \left\{ \|\theta_n - \theta_\ell\|^2 \middle| \theta_\ell, x_\ell \right\}, \end{aligned} \quad (33)$$

where in the second inequality, we have used the fact that both $\|A\|$, $\|A_n\|$ are bounded by c_λ . Finally, we provide an upper bound for the last two terms in (30) and (31). Note that

$$\begin{aligned} &\left| \mathbb{E} \left\{ \tilde{\alpha}_n (\theta_n - \theta_\ell)^T (A_n - A) \theta_\ell \middle| \theta_\ell, x_\ell \right\} + \mathbb{E} \left\{ \tilde{\alpha}_n \theta_\ell^T (A_n - A) (\theta_n - \theta_\ell) \middle| \theta_\ell, x_\ell \right\} \right| \\ &\leq \alpha_n \left| \mathbb{E} \left\{ (\theta_n - \theta_\ell)^T (A_n - A) \theta_\ell \middle| \theta_\ell, x_\ell \right\} \right| + \alpha_n \left| \mathbb{E} \left\{ \theta_\ell^T (A_n - A) (\theta_n - \theta_\ell) \middle| \theta_\ell, x_\ell \right\} \right| \\ &\leq 4c_\lambda \alpha_n \mathbb{E} \left\{ \|\theta_n - \theta_\ell\| \middle| \theta_\ell, x_\ell \right\}, \end{aligned} \quad (34)$$

where we use the triangle inequality with $\tilde{\alpha}_n \leq \alpha_n$ for the first inequality and $\|A_n - A\| \leq 2c_\lambda$ in the second inequality. We now apply Lemma B.4 to (34) and get

$$\begin{aligned} &\left| \mathbb{E} \left\{ \tilde{\alpha}_n (\theta_n - \theta_\ell)^T (A_n - A) \theta_\ell \middle| \theta_\ell, x_\ell \right\} + \mathbb{E} \left\{ \tilde{\alpha}_n \theta_\ell^T (A_n - A) (\theta_n - \theta_\ell) \middle| \theta_\ell, x_\ell \right\} \right| \\ &\leq 8c_\lambda^2 \alpha_n \|\theta_\ell\| \alpha_\ell \tilde{\tau}_{\alpha_n} (\|\theta_\ell\| + S_{\max}) \\ &\leq 8c_\lambda^2 \alpha_\ell^2 \tilde{\tau}_{\alpha_n} (\|\theta_\ell\|^2 + \|\theta_\ell\| S_{\max}) \\ &= 8c_\lambda^2 \alpha_\ell^2 \tilde{\tau}_{\alpha_n} \|\theta_\ell\|^2 + 8c_\lambda^2 \alpha_\ell^2 \tilde{\tau}_{\alpha_n} \|\theta_\ell\| S_{\max}, \end{aligned} \quad (35)$$

where we used $\alpha_n \leq \alpha_\ell$ in the second inequality. Combining (32), (33), (35), we get

$$\begin{aligned} &\left| \mathbb{E} \left\{ \tilde{\alpha}_n \theta_n^T (A_n - A) \theta_n \middle| \theta_\ell, x_\ell \right\} \right| \\ &\leq \alpha_n^2 \|\theta_\ell\|^2 + 2c_\lambda \alpha_n \mathbb{E} \left\{ \|\theta_n - \theta_\ell\|^2 \middle| \theta_\ell, x_\ell \right\} + 8c_\lambda^2 \alpha_\ell^2 \tilde{\tau}_{\alpha_n} \|\theta_\ell\|^2 + 8c_\lambda^2 \alpha_\ell^2 \tilde{\tau}_{\alpha_n} \|\theta_\ell\| S_{\max} \\ &= (\alpha_n^2 + 8c_\lambda^2 \alpha_\ell^2 \tilde{\tau}_{\alpha_n}) \|\theta_\ell\|^2 + 8c_\lambda^2 \alpha_\ell^2 \tilde{\tau}_{\alpha_n} \|\theta_\ell\| S_{\max} + 2c_\lambda \alpha_n \mathbb{E} \left\{ \|\theta_n - \theta_\ell\|^2 \middle| \theta_\ell, x_\ell \right\} \\ &\leq (\alpha_\ell^2 + 8c_\lambda^2 \alpha_\ell^2 \tilde{\tau}_{\alpha_n}) \|\theta_\ell\|^2 + 8c_\lambda^2 \alpha_\ell^2 \tilde{\tau}_{\alpha_n} \|\theta_\ell\| S_{\max} + 2c_\lambda \alpha_\ell \mathbb{E} \left\{ \|\theta_n - \theta_\ell\|^2 \middle| \theta_\ell, x_\ell \right\} \\ &\leq 9c_\lambda^2 \alpha_\ell^2 \tilde{\tau}_{\alpha_n} \|\theta_\ell\|^2 + 8c_\lambda^2 \alpha_\ell^2 \tilde{\tau}_{\alpha_n} \|\theta_\ell\| S_{\max} + 2c_\lambda \alpha_\ell \mathbb{E} \left\{ \|\theta_n - \theta_\ell\|^2 \middle| \theta_\ell, x_\ell \right\}, \end{aligned} \quad (36)$$

where in the second and last inequality, $\alpha_n \leq \alpha_\ell$ and $c_\lambda \tilde{\tau}_{\alpha_n} \geq 1$ was respectively used.

Step 3: Combining bounds obtained in previous steps, given in (27) and (36), we get

$$\begin{aligned} &\mathbb{E} \left\{ \theta_n^T (\theta_{n+1} - \theta_n - \tilde{\alpha}_n A \theta_n) \middle| \theta_\ell, x_\ell \right\} \\ &\leq \alpha_\ell^2 c_{w_*} \|\theta_\ell\| + 2c_\lambda \alpha_\ell^2 \tilde{\tau}_{\alpha_n} (\|\theta_\ell\| + S_{\max}) S_{\max} + 8c_\lambda^2 \alpha_\ell^2 \tilde{\tau}_{\alpha_n} \|\theta_\ell\|^2 + 8c_\lambda^2 \alpha_\ell^2 \tilde{\tau}_{\alpha_n} \|\theta_\ell\| S_{\max} \\ &\quad + 2c_\lambda \alpha_\ell \mathbb{E} \left\{ \|\theta_n - \theta_\ell\|^2 \middle| \theta_\ell, x_\ell \right\} \\ &\leq 9c_\lambda^2 \alpha_\ell^2 \tilde{\tau}_{\alpha_n} \|\theta_\ell\|^2 + (10c_\lambda^2 \alpha_\ell^2 \tilde{\tau}_{\alpha_n} S_{\max} + \alpha_\ell^2 c_{w_*}) \|\theta_\ell\| + 2c_\lambda \alpha_\ell^2 \tilde{\tau}_{\alpha_n} S_{\max}^2 + 2c_\lambda \alpha_\ell \mathbb{E} \left\{ \|\theta_n - \theta_\ell\|^2 \middle| \theta_\ell, x_\ell \right\}, \end{aligned}$$

where in the last inequality, we used the fact $c_\lambda \geq 1$. Since $\|\theta_\ell\| \leq \frac{1}{2} + \frac{1}{2}\|\theta_\ell\|^2$, we get

$$\begin{aligned} & \mathbb{E} \left\{ \theta_n^T (\theta_{n+1} - \theta_n - \tilde{\alpha}_n A \theta_n) \middle| \theta_\ell, x_\ell \right\} \\ & \leq 9c_\lambda^2 \alpha_\ell^2 \tilde{\tau}_{\alpha_n} \|\theta_\ell\|^2 + (10c_\lambda^2 \alpha_\ell^2 \tilde{\tau}_{\alpha_n} S_{\max} + \alpha_\ell^2 c_{w_*}) \left(\frac{1}{2} + \frac{1}{2} \|\theta_\ell\|^2 \right) + 2c_\lambda \alpha_\ell^2 \tilde{\tau}_{\alpha_n} S_{\max}^2 + 2c_\lambda \alpha_\ell \mathbb{E} \left\{ \|\theta_n - \theta_\ell\|^2 \middle| \theta_\ell, x_\ell \right\} \\ & \leq (9c_\lambda^2 \alpha_\ell^2 \tilde{\tau}_{\alpha_n} + 5c_\lambda^2 \alpha_\ell^2 \tilde{\tau}_{\alpha_n} S_{\max} + \alpha_\ell^2 c_{w_*}) \|\theta_\ell\|^2 + (5c_\lambda^2 \alpha_\ell^2 \tilde{\tau}_{\alpha_n} S_{\max} + \alpha_\ell^2 c_{w_*} + 2c_\lambda \alpha_\ell^2 \tilde{\tau}_{\alpha_n} S_{\max}^2) \\ & \quad + 2c_\lambda \alpha_\ell \mathbb{E} \left\{ \|\theta_n - \theta_\ell\|^2 \middle| \theta_\ell, x_\ell \right\} \end{aligned} \quad (37)$$

$$\begin{aligned} & \leq (9c_\lambda^2 \alpha_\ell^2 \tilde{\tau}_{\alpha_n} + 5c_\lambda^2 \alpha_\ell^2 \tilde{\tau}_{\alpha_n} + \alpha_\ell^2)(1 + S_{\max}) \|\theta_\ell\|^2 + (5c_\lambda^2 \alpha_\ell^2 \tilde{\tau}_{\alpha_n} S_{\max} + \alpha_\ell^2(1 + S_{\max}) + 2c_\lambda \alpha_\ell^2 \tilde{\tau}_{\alpha_n} S_{\max}^2) \\ & \quad + 2c_\lambda \alpha_\ell \mathbb{E} \left\{ \|\theta_n - \theta_\ell\|^2 \middle| \theta_\ell, x_\ell \right\}, \end{aligned} \quad (38)$$

where in (37), we used $\frac{1}{2} \alpha_\ell^2 c_{w_*} \leq \alpha_\ell^2 c_{w_*}$ and in (38), $1 \leq c_{w_*} \leq S_{\max} + 1$ was used. Since $\tilde{\tau}_{\alpha_n} \geq 1$ and $c_\lambda \geq 1$,

$$\begin{aligned} & \mathbb{E} \left\{ \theta_n^T (\theta_{n+1} - \theta_n - \tilde{\alpha}_n A \theta_n) \middle| \theta_\ell, x_\ell \right\} \\ & \leq 15c_\lambda^2 \alpha_\ell^2 \tilde{\tau}_{\alpha_n} (1 + S_{\max}) \|\theta_\ell\|^2 + 5c_\lambda^2 (\alpha_\ell^2 \tilde{\tau}_{\alpha_n} S_{\max} + \alpha_\ell^2 \tilde{\tau}_{\alpha_n} (1 + S_{\max}) + \alpha_\ell^2 \tilde{\tau}_{\alpha_n} S_{\max}^2) + 2c_\lambda \alpha_\ell \mathbb{E} \left\{ \|\theta_n - \theta_\ell\|^2 \middle| \theta_\ell, x_\ell \right\} \\ & = 15c_\lambda^2 \alpha_\ell^2 \tilde{\tau}_{\alpha_n} (1 + S_{\max}) \|\theta_\ell\|^2 + 5c_\lambda^2 \alpha_\ell^2 \tilde{\tau}_{\alpha_n} (S_{\max}^2 + 2S_{\max} + 1) + 2c_\lambda \alpha_\ell \mathbb{E} \left\{ \|\theta_n - \theta_\ell\|^2 \middle| \theta_\ell, x_\ell \right\} \\ & \leq 30c_\lambda^2 \alpha_\ell^2 \tilde{\tau}_{\alpha_n} (1 + S_{\max}) \mathbb{E} \left\{ \|\theta_n\|^2 \middle| \theta_\ell, x_\ell \right\} + 5c_\lambda^2 \alpha_\ell^2 \tilde{\tau}_{\alpha_n} (S_{\max} + 1)^2 \\ & \quad + (30c_\lambda^2 \alpha_\ell^2 \tilde{\tau}_{\alpha_n} (1 + S_{\max}) + 2c_\lambda \alpha_\ell) \mathbb{E} \left\{ \|\theta_n - \theta_\ell\|^2 \middle| \theta_\ell, x_\ell \right\}, \end{aligned}$$

where in the last inequality, we used the triangle inequality $\|\theta_\ell\|^2 \leq 2\|\theta_n\|^2 + 2\|\theta_n - \theta_\ell\|^2$. Next, we use the identity $\alpha_\ell \tilde{\tau}_{\alpha_n} \leq \frac{1}{4c_\lambda}$. We have

$$\begin{aligned} & \mathbb{E} \left\{ \theta_n^T (\theta_{n+1} - \theta_n - \tilde{\alpha}_n A \theta_n) \middle| \theta_\ell, x_\ell \right\} \\ & \leq 30c_\lambda^2 \alpha_\ell^2 \tilde{\tau}_{\alpha_n} (1 + S_{\max}) \mathbb{E} \left\{ \|\theta_n\|^2 \middle| \theta_\ell, x_\ell \right\} + 5c_\lambda^2 \alpha_\ell^2 \tilde{\tau}_{\alpha_n} (S_{\max} + 1)^2 \\ & \quad + (8c_\lambda \alpha_\ell (1 + S_{\max}) + 2c_\lambda \alpha_\ell) \mathbb{E} \left\{ \|\theta_n - \theta_\ell\|^2 \middle| \theta_\ell, x_\ell \right\} \\ & \leq 30c_\lambda^2 \alpha_\ell^2 \tilde{\tau}_{\alpha_n} (1 + S_{\max}) \mathbb{E} \left\{ \|\theta_n\|^2 \middle| \theta_\ell, x_\ell \right\} + 5c_\lambda^2 \alpha_\ell^2 \tilde{\tau}_{\alpha_n} (S_{\max} + 1)^2 + 10c_\lambda \alpha_\ell (1 + S_{\max}) \mathbb{E} \left\{ \|\theta_n - \theta_\ell\|^2 \middle| \theta_\ell, x_\ell \right\} \\ & \leq 30c_\lambda^2 \alpha_\ell^2 \tilde{\tau}_{\alpha_n} (1 + S_{\max}) \mathbb{E} \left\{ \|\theta_n\|^2 \middle| \theta_\ell, x_\ell \right\} + 5c_\lambda^2 \alpha_\ell^2 \tilde{\tau}_{\alpha_n} (S_{\max} + 1)^2 + 80c_\lambda^2 \alpha_\ell^2 \tilde{\tau}_{\alpha_n} (1 + S_{\max}) \mathbb{E} \left\{ \|\theta_n\|^2 \middle| \theta_\ell, x_\ell \right\} \\ & \quad + 80c_\lambda^2 \alpha_\ell^2 \tilde{\tau}_{\alpha_n} (1 + S_{\max}) S_{\max}^2 \\ & \leq 30c_\lambda^2 \kappa^2 \alpha_n^2 \tilde{\tau}_{\alpha_n} (1 + S_{\max}) \mathbb{E} \left\{ \|\theta_n\|^2 \middle| \theta_\ell, x_\ell \right\} + 5c_\lambda^2 \kappa^2 \alpha_n^2 \tilde{\tau}_{\alpha_n} (S_{\max} + 1)^2 \\ & \quad + 80c_\lambda^2 \kappa^2 \alpha_n^2 \tilde{\tau}_{\alpha_n} (1 + S_{\max}) \mathbb{E} \left\{ \|\theta_n\|^2 \middle| \theta_\ell, x_\ell \right\} + 80c_\lambda^2 \kappa^2 \alpha_n^2 \tilde{\tau}_{\alpha_n} (1 + S_{\max}) S_{\max}^2 \\ & = 110c_\lambda^2 \kappa^2 (1 + S_{\max}) \alpha_n^2 \tilde{\tau}_{\alpha_n} \mathbb{E} \left\{ \|\theta_n\|^2 \middle| \theta_\ell, x_\ell \right\} + (5c_\lambda^2 (S_{\max} + 1)^2 + 80c_\lambda^2 (1 + S_{\max}) S_{\max}^2) \kappa^2 \alpha_n^2 \tilde{\tau}_{\alpha_n}, \end{aligned}$$

where in the second inequality, we used $1 + S_{\max} \geq 1$, in the third inequality, Lemma B.4 was invoked, and the last inequality was due to the condition $\alpha_\ell \leq \kappa \alpha_n$. \square

The last result we need in establishing the asymptotic convergence of TD algorithms is the negative definiteness of the matrix A .

Lemma B.6 (Lemma 6.6 of Bertsekas [3]). *Under Assumptions 4.1, 4.2, 4.3 and 4.4, the matrix*

$$A = \begin{cases} \mathbb{E}_\infty \{ \gamma \phi_n \phi_{n+1}^T - \phi_n \phi_n^T \} & \text{for TD}(0), \\ \mathbb{E}_\infty \{ \gamma e_{-\infty:n} \phi_{n+1}^T - e_{-\infty:n} \phi_n^T \} & \text{for TD}(\lambda), \end{cases}$$

is negative definite, where $e_{-\infty:n} := \sum_{k=0}^{\infty} (\lambda\gamma)^k \phi_{n-k}$ represents the steady-space eligibility trace and \mathbb{E}_∞ represents the expectation with respect to the steady-state distribution of $(x_n)_{n \in \mathbb{N}}$.

We now establish show that $\mathbb{E}\{\|\theta_n\|^2\} = \mathbb{E}\{\|w_n^{\text{im}} - w_*\|^2\}$ converges to zero as n goes to ∞ .

Proof of Theorem 4.7. Note that

$$\begin{aligned} \mathbb{E} \left\{ \theta_{n+1}^\top \theta_{n+1} - \theta_n^\top \theta_n \middle| \theta_\ell, x_\ell \right\} &= \mathbb{E} \left\{ 2\theta_n^\top (\theta_{n+1} - \theta_n) + (\theta_{n+1} - \theta_n)^\top (\theta_{n+1} - \theta_n) \middle| \theta_\ell, x_\ell \right\} \\ &= \mathbb{E} \left\{ 2\theta_n^\top (\theta_{n+1} - \theta_n - \tilde{\alpha}_n A \theta_n) \middle| \theta_\ell, x_\ell \right\} \end{aligned} \quad (39)$$

$$+ \mathbb{E} \left\{ (\theta_{n+1} - \theta_n)^\top (\theta_{n+1} - \theta_n) \middle| \theta_\ell, x_\ell \right\} \quad (40)$$

$$+ \mathbb{E} \left\{ 2\tilde{\alpha}_n \theta_n^\top A \theta_n \middle| \theta_\ell, x_\ell \right\}, \quad (41)$$

where in the second inequality, we add and subtract $\mathbb{E} \{ 2\tilde{\alpha}_n \theta_n^\top A \theta_n | \theta_\ell, x_\ell \}$. Note that from Lemma B.5, we have

$$(39) \leq 2c_1 \alpha_n^2 \tilde{\tau}_{\alpha_n} \mathbb{E} \{ \|\theta_n\|^2 | \theta_\ell, x_\ell \} + 2c_2 \alpha_n^2 \tilde{\tau}_{\alpha_n}.$$

For the term in (40), notice that

$$\begin{aligned} \|\theta_{n+1} - \theta_n\|^2 &= \|\tilde{\alpha}_n (A_n \theta_n + A_n w_* + b_n)\|^2 \leq \alpha_n^2 \|A_n \theta_n + A_n w_* + b_n\|^2 \\ &\leq 2\alpha_n^2 (\|A_n \theta_n\|^2 + \|A_n w_* + b_n\|^2) \\ &\leq 2\alpha_n^2 \{ c_\lambda^2 \|\theta_n\|^2 + S_{\max}^2 \} = 2c_\lambda^2 \alpha_n^2 \|\theta_n\|^2 + 2\alpha_n^2 S_{\max}^2, \end{aligned}$$

where the first inequality is due to Lemma (A.2), the second inequality is from the identity $(a+b)^2 \leq 2a^2 + 2b^2$, and the third inequality is due to Lemma (B.3). For the expression (41), note that

$$\begin{aligned} \mathbb{E} \left\{ \tilde{\alpha}_n \theta_n^\top A \theta_n \middle| \theta_\ell, x_\ell \right\} &\leq \max \left[\alpha_n \mathbb{E} \left\{ \theta_n^\top A \theta_n \middle| \theta_\ell, x_\ell \right\}, \frac{\alpha_n}{1 + \alpha_n} \mathbb{E} \left\{ \theta_n^\top A \theta_n \middle| \theta_\ell, x_\ell \right\} \right], \quad \text{for TD}(0) \\ \mathbb{E} \left\{ \tilde{\alpha}_n \theta_n^\top A \theta_n \middle| \theta_\ell, x_\ell \right\} &\leq \max \left[\alpha_n \mathbb{E} \left\{ \theta_n^\top A \theta_n \middle| \theta_\ell, x_\ell \right\}, \frac{(1 - \lambda\gamma)^2 \alpha_n}{(1 - \lambda\gamma)^2 + \alpha_n} \mathbb{E} \left\{ \theta_n^\top A \theta_n \middle| \theta_\ell, x_\ell \right\} \right], \quad \text{for TD}(\lambda). \end{aligned}$$

Notice that $\frac{\alpha_n}{1 + \alpha_n} \geq \frac{(1 - \lambda\gamma)^2 \alpha_n}{(1 - \lambda\gamma)^2 + \alpha_n} \geq \frac{(1 - \lambda\gamma)^2 \alpha_n}{1 + \alpha_n}$. From Lemma B.6 which states that A is negative definite, for any non-zero θ , we know there exists $\lambda_0 > 0$ such that $\theta^\top A \theta \leq -\lambda_0 \|\theta\|^2 < 0$. Therefore,

we have

$$\mathbb{E} \left\{ \theta_n^\top A \theta_n \middle| \theta_\ell, x_\ell \right\} \leq -\lambda_0 \mathbb{E} \left\{ \|\theta_n\|^2 \middle| \theta_\ell, x_\ell \right\},$$

which gives us (41) $\leq -\frac{2(1-\lambda\gamma)^2\alpha_n\lambda_0}{1+\alpha_n} \mathbb{E} \left\{ \|\theta_n\|^2 \middle| \theta_\ell, x_\ell \right\}$. Combining all three bounds we established, we get

$$\begin{aligned} \mathbb{E} \left\{ \theta_{n+1}^\top \theta_{n+1} - \theta_n^\top \theta_n \middle| \theta_\ell, x_\ell \right\} &\leq \left(2c_1\alpha_n^2\tilde{\tau}_{\alpha_n} + 2c_\lambda^2\alpha_n^2 - \frac{2(1-\lambda\gamma)^2\alpha_n\lambda_0}{1+\alpha_n} \right) \mathbb{E} \left\{ \|\theta_n\|^2 \middle| \theta_\ell, x_\ell \right\} + 2\alpha_n^2 (c_2\tilde{\tau}_{\alpha_n} + S_{\max}^2) \\ &\leq \left(2c_1\alpha_n^2\tilde{\tau}_{\alpha_n} + 2c_\lambda^2\alpha_n^2 - \frac{2(1-\lambda\gamma)^2\alpha_n\lambda_0}{1+\alpha_1} \right) \mathbb{E} \left\{ \|\theta_n\|^2 \middle| \theta_\ell, x_\ell \right\} + 2\alpha_n^2 (c_2\tilde{\tau}_{\alpha_n} + S_{\max}^2) \end{aligned}$$

where the last inequality follows from non-increasingness of $(a_k)_{k \in \mathbb{N}}$. For n large enough, such that

$$2c_1\alpha_n^2\tilde{\tau}_{\alpha_n} + 2c_\lambda^2\alpha_n^2 \leq \frac{(1-\lambda\gamma)^2\alpha_n\lambda_0}{1+\alpha_1},$$

we get

$$\mathbb{E} \left\{ \|\theta_{n+1}\|^2 \middle| \theta_\ell, x_\ell \right\} \leq \left\{ 1 - \frac{(1-\lambda\gamma)^2\alpha_n\lambda_0}{1+\alpha_1} \right\} \mathbb{E} \left\{ \|\theta_n\|^2 \middle| \theta_\ell, x_\ell \right\} + 2\alpha_n^2 (c_2\tilde{\tau}_{\alpha_n} + S_{\max}^2).$$

Taking the expectation with respect to θ_ℓ and x_ℓ , we have

$$\mathbb{E} \left\{ \|\theta_{n+1}\|^2 \right\} \leq \left\{ 1 - \frac{(1-\lambda\gamma)^2\alpha_n\lambda_0}{1+\alpha_1} \right\} \mathbb{E} \left\{ \|\theta_n\|^2 \right\} + 2\alpha_n^2 (c_2\tilde{\tau}_{\alpha_n} + S_{\max}^2).$$

Recursively using this inequality, we get

$$\begin{aligned} \mathbb{E} \left\{ \|\theta_{n+1}\|^2 \right\} &\leq \prod_{k=\ell}^n \left(1 - \frac{(1-\lambda\gamma)^2\alpha_k\lambda_0}{1+\alpha_1} \right) \mathbb{E} \left\{ \|\theta_\ell\|^2 \right\} + \prod_{k=\ell+1}^n \left(1 - \frac{(1-\lambda\gamma)^2\alpha_k\lambda_0}{1+\alpha_1} \right) 2\alpha_\ell^2 (c_2\tilde{\tau}_{\alpha_\ell} + S_{\max}^2) \\ &\quad + \prod_{k=\ell+2}^n \left(1 - \frac{(1-\lambda\gamma)^2\alpha_k\lambda_0}{1+\alpha_1} \right) 2\alpha_{\ell+1}^2 (c_2\tilde{\tau}_{\alpha_{\ell+1}} + S_{\max}^2) + \cdots \\ &\quad + \left(1 - \frac{(1-\lambda\gamma)^2\alpha_n\lambda_0}{1+\alpha_1} \right) 2\alpha_{n-1}^2 (c_2\tilde{\tau}_{\alpha_{n-1}} + S_{\max}^2) + 2\alpha_n^2 (c_2\tilde{\tau}_{\alpha_n} + S_{\max}^2) \\ &= \mathbb{E} \left\{ \|\theta_\ell\|^2 \right\} \prod_{k=\ell}^n \left(1 - \frac{(1-\lambda\gamma)^2\alpha_k\lambda_0}{1+\alpha_1} \right) + \sum_{j=\ell+1}^n \prod_{k=j}^n \left(1 - \frac{(1-\lambda\gamma)^2\alpha_k\lambda_0}{1+\alpha_1} \right) 2\alpha_{j-1}^2 (c_2\tilde{\tau}_{\alpha_{j-1}} + S_{\max}^2) \\ &\quad + 2\alpha_n^2 (c_2\tilde{\tau}_{\alpha_n} + S_{\max}^2). \end{aligned}$$

Using $1 - x \leq \exp(-x)$, we get

$$\begin{aligned}
\mathbb{E} \{ \|\theta_{n+1}\|^2 \} &\leq \mathbb{E} \{ \|\theta_\ell\|^2 \} \prod_{k=\ell}^n \exp \left(-\frac{(1-\lambda\gamma)^2 \alpha_k \lambda_0}{1+\alpha_1} \right) \\
&\quad + \sum_{j=\ell+1}^n \prod_{k=j}^n \exp \left(-\frac{(1-\lambda\gamma)^2 \alpha_k \lambda_0}{1+\alpha_1} \right) 2\alpha_{j-1}^2 (c_2 \tilde{\tau}_{\alpha_{j-1}} + S_{\max}^2) + 2\alpha_n^2 (c_2 \tilde{\tau}_{\alpha_n} + S_{\max}^2) \\
&= \mathbb{E} \{ \|\theta_\ell\|^2 \} \exp \left(-\frac{(1-\lambda\gamma)^2 \lambda_0}{1+\alpha_1} \sum_{k=\ell}^n \alpha_k \right) \\
&\quad + \sum_{j=\ell+1}^n \exp \left(-\frac{(1-\lambda\gamma)^2 \lambda_0}{1+\alpha_1} \sum_{k=\ell}^n \alpha_k \right) 2\alpha_{j-1}^2 (c_2 \tilde{\tau}_{\alpha_{j-1}} + S_{\max}^2) + 2\alpha_n^2 (c_2 \tilde{\tau}_{\alpha_n} + S_{\max}^2).
\end{aligned} \tag{42}$$

For $\alpha_n = \frac{c}{n^s}$, $s \in (0.5, 1]$, we have

$$\lim_{n \rightarrow \infty} \sum_{k=\ell}^n \alpha_k = \infty, \quad \lim_{n \rightarrow \infty} \alpha_n^2 \tilde{\tau}_{\alpha_n} = 0 \quad \text{and} \quad \lim_{n \rightarrow \infty} \alpha_n \rightarrow 0,$$

which implies the convergence of the first and the last term in (42) to zero. Therefore, the rest of the proof is to establish

$$\sum_{j=\ell+1}^n \exp \left(-\frac{(1-\lambda\gamma)^2 \lambda_0}{1+\alpha_1} \sum_{k=\ell}^n \alpha_k \right) 2\alpha_{j-1}^2 (c_2 \tilde{\tau}_{\alpha_{j-1}} + S_{\max}^2) \rightarrow 0, \quad \text{as } n \rightarrow \infty.$$

To this end, note that $\sum_{k=\ell}^n \frac{1}{k} \leq \sum_{k=\ell}^n \frac{1}{k^s}$ for $s \in (0, 1]$, which gives us

$$\exp \left(-\frac{(1-\lambda\gamma)^2 \lambda_0}{1+\alpha_1} \sum_{k=\ell}^n \frac{1}{k^s} \right) \leq \exp \left(-\frac{(1-\lambda\gamma)^2 \lambda_0}{1+\alpha_1} \sum_{k=\ell}^n \frac{1}{k} \right),$$

From the definition of Euler-Mascheroni constant, denoted by $\gamma_* > 0$, we have

$$\log n + \gamma_* + \frac{c'}{n} \leq \sum_{k=1}^n \frac{1}{k} \leq \log n + \gamma_* + \frac{c''}{n},$$

for some constant $c', c'' \in \mathbb{R}$ [16]. Therefore, we get

$$\log n + \gamma_* + \frac{c'}{n} + \tilde{c} \leq \sum_{k=\ell}^n \frac{1}{k} \leq \log n + \gamma_* + \frac{c''}{n} + \tilde{c},$$

where $\tilde{c} = -\sum_{k=1}^{\ell-1} \frac{1}{k}$. This gives us

$$\exp \left(-\frac{(1-\lambda\gamma)^2 \lambda_0}{1+\alpha_1} \sum_{k=\ell}^n \frac{1}{k} \right) \leq \exp \left\{ -\frac{(1-\lambda\gamma)^2 \lambda_0}{1+\alpha_1} \left(\log n + \gamma_* + \frac{c'}{n} + \tilde{c} \right) \right\} = c_n \exp \left(-\frac{(1-\lambda\gamma)^2 \lambda_0}{1+\alpha_1} \log n \right),$$

where $c_n = \exp \left\{ -\frac{(1-\lambda\gamma)^2\lambda_0}{1+\alpha_1} \left(\gamma_* + \frac{c'}{n} + \tilde{c} \right) \right\}$ converges to a finite positive constant as $n \rightarrow \infty$. Therefore, for $s \in (0.5, 1)$, we get

$$\exp \left(-\frac{(1-\lambda\gamma)^2\lambda_0}{1+\alpha_1} \sum_{k=\ell}^n \frac{1}{k^s} \right) \leq \exp \left(-\frac{(1-\lambda\gamma)^2\lambda_0}{1+\alpha_1} \sum_{k=\ell}^n \frac{1}{k} \right) \leq \frac{c_n}{n^{\frac{(1-\lambda\gamma)^2\lambda_0}{1+\alpha_1}}},$$

which converges to zero as $n \rightarrow \infty$. Plugging this upper bound back to (42), we have

$$\begin{aligned} \mathbb{E}\{\|\theta_{n+1}\|^2\} &\leq \mathbb{E}\{\|\theta_\ell\|^2\} \exp \left(-\frac{(1-\lambda\gamma)^2\lambda_0}{1+\alpha_1} \sum_{k=\ell}^n \alpha_k \right) + 2\alpha_n^2 (c_2 \tilde{\tau}_{\alpha_n} + S_{\max}^2) \\ &\quad + \frac{c_n}{n^{\frac{(1-\lambda\gamma)^2\lambda_0}{1+\alpha_1}}} \sum_{j=\ell+1}^n 2\alpha_{j-1}^2 (c_2 \tilde{\tau}_{\alpha_{j-1}} + S_{\max}^2). \end{aligned}$$

Since

$$\sum_{j=1}^n 2\alpha_{j-1}^2 (c_2 \tilde{\tau}_{\alpha_{j-1}} + S_{\max}^2) < \infty,$$

for $\alpha_n = \frac{c}{n^s}$, $s \in (0.5, 1]$, we have

$$\lim_{n \rightarrow \infty} \mathbb{E}\{\|\theta_n\|^2\} = \lim_{n \rightarrow \infty} \mathbb{E}\{\|w_n^{\text{im}} - w_*\|^2\} = 0,$$

which establishes the asymptotic convergence of implicit TD algorithms to w_* . \square

B.2 Finite-time/Convergence analysis for projected implicit TD(0)/TD(λ)

In this section, we establish a finite-time error bound after adding a projection step in the TD algorithm [4]. To this end, we review projections and notations which will be used in this section. Given a radius $R > 0$, at each iteration of the projected TD algorithms proposed in Bhandari et al. [4], we have the following update rule,

$$w_{n+1} = \Pi_R \{w_n + \alpha_n S_n(w_n)\}, \quad (43)$$

where

$$\Pi_R(w) := \underset{w': \|w'\| \leq R}{\operatorname{argmin}} \|w - w'\| = \begin{cases} Rw/\|w\| & \text{if } \|w\| > R \\ w & \text{otherwise.} \end{cases}$$

Therefore, at each n^{th} iteration, projected implicit TD algorithm is defined to be

$$w_{n+1}^{\text{im}} = \Pi_R \{w_n^{\text{im}} + \tilde{\alpha}_n S_n(w_n^{\text{im}})\}.$$

Here is a reminder and introduction to the notations we will use in this section.

- $\xi_n(w) := \{S_n(w) - S(w)\}^\top (w - w_*)$, $\forall w \in \mathbb{R}^d$

- $\Gamma := \sum_{x \in \mathcal{X}} \pi(x) \phi(x) \phi(x)^T = \Phi^T D \Phi$, $D := \text{diag} \{ \pi(x) : x \in \mathcal{X} \}$
- $\min\{\text{eig}(\Gamma)\} = \lambda_{\min}$
- $V_{w_*}(x) := \phi(x)^T w_*$, $\forall x \in \mathcal{X}$
- $\|V_w - V_{w'}\|_D = \|w - w'\|_\Gamma$, where $\|u\|_Q := u^T Q u$

We first establish a result, which relates the value function difference to that of the parameter difference.

Lemma B.7. *For all $w, w' \in \mathbb{R}^d$,*

$$\sqrt{\lambda_{\min}} \|w - w'\| \leq \|V_w - V_{w'}\|_D \leq \|w - w'\|.$$

Proof. Note that

$$\|V_w - V_{w'}\|_D = \sqrt{\sum_{x \in \mathcal{X}} \pi(x) (\phi(x)^T (w - w'))^2} = \left((w - w')^T \Gamma (w - w') \right)^{1/2}.$$

By the definition of Γ ,

$$\lambda_{\max}(\Gamma) = \lambda_{\max} \left(\sum_{x \in \mathcal{X}} \pi(x) \phi(x) \phi(x)^T \right) \leq \sum_{x \in \mathcal{X}} \pi(x) \lambda_{\max} \left(\phi(x) \phi(x)^T \right) \leq \sum_{x \in \mathcal{X}} \pi(x) = 1.$$

Therefore, we have

$$(w - w')^T \Gamma (w - w') \leq (w - w')^T (w - w').$$

The lower bound of $\|V_w - V_{w'}\|$ comes from the fact that $\lambda_{\min} = \min_u \frac{u^T \Gamma u}{\|u\|^2}$. By plugging in $u = w - w'$, we get the lower bound. \square

B.2.1 Finite-time/Convergence analysis for projected implicit TD(0)

In this subsection, we present a finite-time error bound for implicit TD(0) with a projection step. Our approach closely follows that of [4], with a few modifications to account for the data-adaptive step size used in implicit TD algorithms. To ensure clarity and completeness, we also restate some of the proofs from [4]. An upshot of our result is that the projection step in combination with an implicit update will yield a finite-time error bound nearly independent of the step size one chooses. We first list results from [4], which will be used in establishing finite-time error bounds for the projected implicit TD(0) algorithm.

Lemma B.8 (Lemma 3 of Bhandari et al. [4]). *For any $w \in \mathbb{R}^d$,*

$$(w_* - w)^T S(w) \geq (1 - \gamma) \|V_{w_*} - V_w\|_D^2 \geq 0$$

Lemma B.9 (Lemma 6 of Bhandari et al. [4]). *For all $n \in \mathbb{N}$, $w \in \{w' : \|w'\| \leq R\}$,*

$$\|S_n(w)\| \leq G := r_{\max} + (\gamma + 1)R,$$

with probability 1.

Lemma B.10 (Lemma 9 of Bhandari et al. [4]). *Consider two random variables U and \tilde{U} such that*

$$U \rightarrow x_n \rightarrow x_{n+\tau} \rightarrow \tilde{U}$$

for some fixed $n \in \{1, 2, \dots\}$ and $\tau > 0$. Assume the Markov chain mixes as stated in Expression 15. Let U' and \tilde{U}' be independent copies drawn from the marginal distributions of U and \tilde{U} . Then, for any bounded function h ,

$$\left| \mathbb{E} \left\{ h(U, \tilde{U}) \right\} - \mathbb{E} \left\{ h(U', \tilde{U}') \right\} \right| \leq 2\|h\|_{\infty} m \rho^{\tau},$$

for some $m > 0$, $\rho \in (0, 1)$. In particular, with $\tilde{U} = x_{n+\tau}$, the above inequality still holds.

Lemma B.11 (Lemma 10 of Bhandari et al. [4]). *With probability 1, for all $w, v \in \{w' : \|w'\| \leq R\}$,*

$$\begin{aligned} |\xi_n(w)| &\leq 2G^2 \\ |\xi_n(w) - \xi_n(v)| &\leq 6G \|w - v\|, \end{aligned}$$

where $\xi_n(w) = (S_n(w) - S(w))^T(w - w_)$.*

Lemma B.12. *For every $n \geq 1$, with $R \geq \|w_*\|$,*

$$\|w_* - w_{n+1}^{\text{im}}\|^2 \leq \|w_* - w_n^{\text{im}}\|^2 - \frac{2\alpha_n(1-\gamma)}{1+\alpha_n} \|V_{w_*} - V_{w_n^{\text{im}}}\|_D^2 + 2\tilde{\alpha}_n \xi_n(w_n^{\text{im}}) + \alpha_n^2 G^2,$$

holds with probability one.

Proof. With probability one, we have

$$\|w_* - w_{n+1}^{\text{im}}\|^2 = \|\Pi_R(w_*) - \Pi_R\{w_n^{\text{im}} + \tilde{\alpha}_n S_n(w_n^{\text{im}})\}\|^2 \quad (44)$$

$$\leq \|w_* - w_n^{\text{im}} - \tilde{\alpha}_n S_n(w_n^{\text{im}})\|^2 \quad (45)$$

$$\begin{aligned} &= \|w_* - w_n^{\text{im}}\|^2 - 2\tilde{\alpha}_n S_n(w_n^{\text{im}})^{\top} (w_* - w_n^{\text{im}}) + \|\tilde{\alpha}_n S_n(w_n^{\text{im}})\|^2 \\ &\leq \|w_* - w_n^{\text{im}}\|^2 - 2\tilde{\alpha}_n S_n(w_n^{\text{im}})^{\top} (w_* - w_n^{\text{im}}) + \alpha_n^2 G^2 \end{aligned} \quad (46)$$

$$\begin{aligned} &= \|w_* - w_n^{\text{im}}\|^2 - 2\tilde{\alpha}_n S(w_n^{\text{im}})^{\top} (w_* - w_n^{\text{im}}) + 2\tilde{\alpha}_n \xi_n(w_n^{\text{im}}) + \alpha_n^2 G^2 \\ &\leq \|w_* - w_n^{\text{im}}\|^2 - 2\tilde{\alpha}_n(1-\gamma) \left\| V_{w_*} - V_{w_n^{\text{im}}} \right\|_D^2 + 2\tilde{\alpha}_n \xi_n(w_n^{\text{im}}) + \alpha_n^2 G^2 \end{aligned} \quad (47)$$

$$\leq \|w_* - w_n^{\text{im}}\|^2 - \frac{2\alpha_n(1-\gamma)}{1+\alpha_n} \left\| V_{w_*} - V_{w_n^{\text{im}}} \right\|_D^2 + 2\tilde{\alpha}_n \xi_n(w_n^{\text{im}}) + \alpha_n^2 G^2, \quad (48)$$

where (44) is due to the fact that $w_* = \Pi_R(w_*)$, (45) is thanks to non-expansiveness of the projection operator on the convex set, (46) comes from the fact $\tilde{\alpha}_n \leq \alpha_n$ with Lemma B.9 and (47) is by Lemma B.8. Finally, the last inequality is a direct consequence of the Lemma A.2. \square

Lemma B.13. *Given a non-increasing sequence $\alpha_1 \geq \dots \geq \alpha_N$, for any fixed $n < N$, we get*

$$\mathbb{E} \{ \tilde{\alpha}_n \xi_n (w_n^{im}) \} \leq 6\alpha_n G^2 \sum_{i=1}^{n-1} \alpha_i, \quad (49)$$

as well as

$$\mathbb{E} \{ \tilde{\alpha}_n \xi_n (w_n^{im}) \} \leq \alpha_n G^2 (4 + 6\tau_{\alpha_N}) \alpha_{\max\{1, n-\tau_{\alpha_N}\}}. \quad (50)$$

Proof. We first establish a bound on $\mathbb{E}_\infty \{ \xi_n (w_n^{im}) \}$. To this end, recall from Lemma B.11 that

$$\xi_n(w_n^{im}) \leq \xi_n(w_{n-1}^{im}) + 6G \|w_n^{im} - w_{n-1}^{im}\|. \quad (51)$$

For $\tau = 1, \dots, n-1$, from the repeated application of (51), we have

$$\begin{aligned} \xi_n(w_n^{im}) &\leq \xi_n(w_{n-2}^{im}) + 6G \|w_{n-1}^{im} - w_{n-2}^{im}\| + 6G \|w_n^{im} - w_{n-1}^{im}\| \\ &\leq \xi_n(w_{n-\tau}^{im}) + 6G \sum_{i=n-\tau}^{n-1} \|w_{i+1}^{im} - w_i^{im}\|. \end{aligned}$$

Note that

$$\|w_{i+1}^{im} - w_i^{im}\| = \|\Pi_R\{w_i^{im} + \tilde{\alpha}_i S_i(w_i^{im})\} - \Pi_R(w_i^{im})\| \leq \|w_i^{im} + \tilde{\alpha}_i S_i(w_i^{im}) - w_i^{im}\| \leq \alpha_i G,$$

where in the first inequality, we have used the non-expansiveness of the projection operator, and for the second inequality, both Lemma A.2 and B.9 were used. Therefore, for $\tau \in \{1, \dots, n-1\}$, we have

$$\xi_n(w_n^{im}) \leq \xi_n(w_{n-\tau}^{im}) + 6G^2 \sum_{i=n-\tau}^{n-1} \alpha_i \quad (52)$$

$$\leq \xi_n(w_{n-\tau}^{im}) + 6G^2 \tau \alpha_{n-\tau}, \quad (53)$$

where (53) follows from non-increasingness of $(\alpha_n)_{n \in \mathbb{N}}$. We first show (49). From (52) with $\tau = n-1$, we have

$$\xi_n(w_n^{im}) \leq \xi_n(w_1^{im}) + 6G^2 \sum_{i=1}^{n-1} \alpha_i.$$

Taking the expectation with respect to the steady state distribution, we get

$$\mathbb{E} \{ \xi_n (w_n^{im}) \} \leq 6G^2 \sum_{i=1}^{n-1} \alpha_i,$$

since $\mathbb{E}_\infty \{ \xi_n (w) \} = 0$, for any fixed w . From Lemma A.2,

$$\mathbb{E} \{ \tilde{\alpha}_n \xi_n (w_n^{im}) \} \leq \max \left[\alpha_n \mathbb{E} \{ \xi_n (w_n^{im}) \}, \frac{\alpha_n}{1 + \alpha_n} \mathbb{E} \{ \xi_n (w_n^{im}) \} \right], \quad (54)$$

we have

$$\mathbb{E} \{ \tilde{\alpha}_n \xi_n (w_n^{im}) \} \leq 6\alpha_n G^2 \sum_{i=1}^{n-1} \alpha_i,$$

as we desired. We next show (50). We consider two different cases.

Case 1: We first consider when $n \leq \tau_{\alpha_N}$. Setting $\tau = n - 1$ in (53), we get

$$\xi_n (w_n^{im}) \leq \xi_n (w_1^{im}) + 6G^2(n-1)\alpha_1 \leq \xi_n (w_1^{im}) + 6G^2 n \alpha_1.$$

Taking the expectation with respect to the steady-state distribution, we get

$$\mathbb{E} \{ \xi_n (w_n^{im}) \} \leq \mathbb{E} \{ \xi_n (w_1^{im}) \} + 6G^2 n \alpha_1.$$

Since $\mathbb{E} \{ \xi_n (w) \} = 0$, for any fixed w , we get

$$\mathbb{E} \{ \xi_n (w_n^{im}) \} \leq 6G^2 \tau_{\alpha_N} \alpha_1$$

Case 2: We next consider when $n > \tau_{\alpha_N}$. Setting $\tau = \tau_{\alpha_N}$ in (53), we get

$$\xi_n (w_n^{im}) \leq \xi_n (w_{n-\tau_{\alpha_N}}^{im}) + 6G^2 \tau_{\alpha_N} \alpha_{n-\tau_{\alpha_N}}. \quad (55)$$

Recall that $\xi_n(w) = \{S_n(w) - S(w)\}^\top (w - w_*)$, which can be viewed as a function of $u_n = \{x_n, r(x_n), x_{n+1}\}$ and w . Notice that u_n is a Markov process with the same transition probability as x_n . Furthermore, we can view $w_{n-\tau_{\alpha_N}}^{im}$ as a function of $\{u_1, \dots, u_{n-\tau_{\alpha_N}-1}\}$. Now consider $\xi_n (w_{n-\tau_{\alpha_N}}^{im})$, which is a function of both $U = \{u_1, \dots, u_{n-\tau_{\alpha_N}-1}\}$ and $\tilde{U} = u_n$. We set $h(U, \tilde{U}) = \xi_n (w_{n-\tau_{\alpha_N}}^{im})$ to invoke Lemma B.10. The condition for Lemma B.10 is met since $U = \{u_1, \dots, u_{n-\tau_{\alpha_N}-1}\} \rightarrow u_{n-\tau_{\alpha_N}} \rightarrow u_n = \{x_n, r(x_n), x_{n+1}\} = \tilde{U}$ forms a Markov chain. Therefore, we get

$$\mathbb{E} \{ h(U, \tilde{U}) \} - \mathbb{E} \{ h(U', \tilde{U}') \} \leq 2\|h\|_\infty m \rho^{\tau_{\alpha_N}},$$

where $U' = \{u'_1, \dots, u'_{n-\tau_{\alpha_N}-1}\}$ and $\tilde{U}' = \{x'_n, r(x'_n), x'_{n+1}\}$ are independent and have the same marginal distribution as U and \tilde{U} . Let us denote the $(n - \tau_{\alpha_N})^{th}$ implicit TD(0) iterate computed

using U' as $w'_{n-\tau_{\alpha_N}}$. Conditioning on U' , we know $w'_{n-\tau_{\alpha_N}}$ is fixed and hence we get

$$\mathbb{E} \left\{ h(U', \tilde{U}') \right\} = \mathbb{E} \left[\mathbb{E} \left\{ \xi_n \left(w'_{n-\tau_{\alpha_N}} \right) \middle| U' \right\} \right] = 0,$$

since $\mathbb{E} \{ \xi_n(w) \} = 0$, for any fixed w . Combined with Lemma B.11, which states that $\|h\|_\infty \leq 2G^2$ we have

$$\mathbb{E} \left\{ \xi_n \left(w_{n-\tau_{\alpha_N}}^{\text{im}} \right) \right\} \leq 4G^2 m \rho^{\tau_{\alpha_N}}.$$

Taking the expectation of (55) with respect to the stationary distribution, we get

$$\mathbb{E} \{ \xi_n(w_n^{\text{im}}) \} \leq \mathbb{E} \left\{ \xi_n \left(w_{n-\tau_{\alpha_N}}^{\text{im}} \right) \right\} + 6G^2 \tau_{\alpha_N} \alpha_{n-\tau_{\alpha_N}} \leq 4G^2 m \rho^{\tau_{\alpha_N}} + 6G^2 \tau_{\alpha_N} \alpha_{n-\tau_{\alpha_N}}.$$

Therefore, again from (54), we have

$$\begin{aligned} \mathbb{E} \{ \tilde{\alpha}_n \xi_n(w_n^{\text{im}}) \} &\leq \alpha_n \left(4G^2 m \rho^{\tau_{\alpha_N}} + 6G^2 \tau_{\alpha_N} \alpha_{n-\tau_{\alpha_N}} \right) \leq \alpha_n \left(4G^2 \alpha_N + 6G^2 \tau_{\alpha_N} \alpha_{n-\tau_{\alpha_N}} \right) \\ &\leq \alpha_n G^2 (4 + 6\tau_{\alpha_N}) \alpha_{n-\tau_{\alpha_N}}, \end{aligned}$$

where the second inequality follows from the definition of the mixing time and the last inequality is due to non-increasingness of step size, i.e., $\alpha_N \leq \alpha_{n-\tau_{\alpha_N}}$. \square

We now establish a finite-time error bound for TD(0) with a constant step size. **Proof.** of **Theorem 4.10:** Starting from Lemma B.12 with a constant step size, we have

$$\begin{aligned} &\mathbb{E} \left\{ \|w_* - w_{n+1}^{\text{im}}\|^2 \right\} \\ &\leq \mathbb{E} \left\{ \|w_* - w_n^{\text{im}}\|^2 \right\} - \frac{2\alpha(1-\gamma)}{1+\alpha} \mathbb{E} \left\{ \|V_{w_*} - V_{w_n^{\text{im}}}\|_D^2 \right\} + 2\mathbb{E} \{ \tilde{\alpha}_n \xi_n(w_n^{\text{im}}) \} + \alpha^2 G^2 \\ &\leq \mathbb{E} \left\{ \|w_* - w_n^{\text{im}}\|^2 \right\} - \frac{2\alpha(1-\gamma)\lambda_{\min}}{1+\alpha} \mathbb{E} \left\{ \|w_* - w_n^{\text{im}}\|^2 \right\} + 2\mathbb{E} \{ \tilde{\alpha}_n \xi_n(w_n^{\text{im}}) \} + \alpha^2 G^2 \\ &\leq \mathbb{E} \left\{ \|w_* - w_n^{\text{im}}\|^2 \right\} - \frac{2\alpha(1-\gamma)\lambda_{\min}}{1+\alpha} \mathbb{E} \left\{ \|w_* - w_n^{\text{im}}\|^2 \right\} + 2\alpha^2 G^2 (4 + 6\tau_\alpha) + \alpha^2 G^2 \\ &= \left\{ 1 - \frac{2\alpha(1-\gamma)\lambda_{\min}}{1+\alpha} \right\} \mathbb{E} \left\{ \|w_* - w_n^{\text{im}}\|^2 \right\} + \alpha^2 G^2 (9 + 12\tau_\alpha), \end{aligned} \tag{56}$$

where the second inequality is due to Lemma B.7, which gives us $\|V_{w_*} - V_{w_n}\|_D^2 \geq \lambda_{\min} \|w_* - w_n\|^2$ and the third one is thanks to Lemma B.13 with a constant step size. Then, the projected implicit

TD(0) iterates with $R \geq \|w_*\|$ achieves

$$\begin{aligned}
& \mathbb{E} \left\{ \|w_* - w_{N+1}^{\text{im}}\|^2 \right\} \\
& \leq \left\{ 1 - \frac{2\alpha(1-\gamma)\lambda_{\min}}{1+\alpha} \right\} \mathbb{E} \left\{ \|w_* - w_N^{\text{im}}\|^2 \right\} + \alpha^2 G^2 (9 + 12\tau_\alpha) \\
& \leq \left\{ 1 - \frac{2\alpha(1-\gamma)\lambda_{\min}}{1+\alpha} \right\}^N \|w_* - w_1^{\text{im}}\|^2 + (\alpha^2 G^2 (9 + 12\tau_\alpha)) \sum_{t=0}^{\infty} \left(1 - \frac{2\alpha(1-\gamma)\lambda_{\min}}{1+\alpha} \right)^t \\
& \leq e^{-\frac{2\alpha(1-\gamma)\lambda_{\min}}{1+\alpha} N} \|w_* - w_1^{\text{im}}\|^2 + \frac{\alpha(1+\alpha)G^2(9+12\tau_\alpha)}{2(1-\gamma)\lambda_{\min}},
\end{aligned}$$

where in the second inequality, we have recursively used the upper bound in (56) and further bounded the finite sum by an infinite sum. In the last inequality, we used $1 - x \leq \exp(-x)$, and an assumption $\frac{2\alpha(1-\gamma)\lambda_{\min}}{1+\alpha} \in (0, 1)$ to obtain a closed form expression of the infinite sum. \square

We next establish convergence of the projected TD(0) algorithm with a sequence of decreasing step sizes.

Proof of Theorem 4.14. Rearranging terms in Lemma B.12, we have

$$\begin{aligned}
& \frac{\alpha_n(1-\gamma)}{1+\alpha_n} \|V_{w_*} - V_{w_n^{\text{im}}}\|_D^2 \\
& \leq \|w_* - w_n^{\text{im}}\|^2 - \frac{\alpha_n(1-\gamma)}{1+\alpha_n} \|V_{w_*} - V_{w_n^{\text{im}}}\|_D^2 - \|w_* - w_{n+1}^{\text{im}}\|^2 + 2\tilde{\alpha}_n \xi_n(w_n^{\text{im}}) + \alpha_n^2 G^2 \\
& \leq \left(1 - \frac{\alpha_n(1-\gamma)\lambda_{\min}}{1+\alpha_n} \right) \|w_* - w_n^{\text{im}}\|^2 - \|w_* - w_{n+1}^{\text{im}}\|^2 + 2\tilde{\alpha}_n \xi_n(w_n^{\text{im}}) + \alpha_n^2 G^2, \tag{57}
\end{aligned}$$

where in the second inequality, we have used Lemma B.7. Dividing both sides by $\frac{\alpha_n(1-\gamma)}{1+\alpha_n}$ and from the non-negativeness of $\|V_{w_*} - V_{w_n^{\text{im}}}\|_D^2$, we have

$$\begin{aligned}
0 & \leq \frac{1+\alpha_n}{\alpha_n(1-\gamma)} \left\{ \left(1 - \frac{\alpha_n(1-\gamma)\lambda_{\min}}{1+\alpha_n} \right) \|w_* - w_n^{\text{im}}\|^2 - \|w_* - w_{n+1}^{\text{im}}\|^2 + 2\tilde{\alpha}_n \xi_n(w_n^{\text{im}}) + \alpha_n^2 G^2 \right\} \\
& = \left(\frac{1+\alpha_n}{\alpha_n(1-\gamma)} - \lambda_{\min} \right) \|w_* - w_n^{\text{im}}\|^2 - \frac{1+\alpha_n}{\alpha_n(1-\gamma)} \|w_* - w_{n+1}^{\text{im}}\|^2 + \frac{2(1+\alpha_n)}{\alpha_n(1-\gamma)} \tilde{\alpha}_n \xi_n(w_n^{\text{im}}) + \frac{\alpha_n(1+\alpha_n)}{(1-\gamma)} G^2 \\
& \tag{58}
\end{aligned}$$

With the choice of $\alpha_n = \frac{\alpha_1}{\alpha_1 \lambda_{\min}(1-\gamma)(n-1)+1}$, one can show that $\frac{1+\alpha_n}{\alpha_n(1-\gamma)} - \lambda_{\min} = \frac{1+\alpha_{n-1}}{\alpha_{n-1}(1-\gamma)}$. Summing (58) over $n = 1, \dots, N$, we have

$$\begin{aligned}
0 & \leq \left(\frac{1+\alpha_1}{\alpha_1(1-\gamma)} - \lambda_{\min} \right) \|w_* - w_1^{\text{im}}\|^2 - \frac{1+\alpha_N}{\alpha_N(1-\gamma)} \|w_* - w_{N+1}^{\text{im}}\|^2 \\
& \quad + \sum_{n=1}^N \frac{2(1+\alpha_n)}{\alpha_n(1-\gamma)} \tilde{\alpha}_n \xi_n(w_n^{\text{im}}) + \sum_{n=1}^N \frac{\alpha_n(1+\alpha_n)}{(1-\gamma)} G^2.
\end{aligned}$$

Rearranging terms and dividing both sides by $\frac{1+\alpha_N}{\alpha_N(1-\gamma)}$, we have

$$\begin{aligned}\|w_* - w_{N+1}^{\text{im}}\|^2 &\leq \frac{\alpha_N(1-\gamma)}{1+\alpha_N} \left(\frac{1+\alpha_1}{\alpha_1(1-\gamma)} - \lambda_{\min} \right) \|w_* - w_1^{\text{im}}\|^2 \\ &\quad + \frac{\alpha_N(1-\gamma)}{1+\alpha_N} \sum_{n=1}^N \frac{2(1+\alpha_n)}{\alpha_n(1-\gamma)} \tilde{\alpha}_n \xi_n(w_n^{\text{im}}) + \frac{\alpha_N(1-\gamma)}{1+\alpha_N} \sum_{n=1}^N \frac{\alpha_n(1+\alpha_n)}{(1-\gamma)} G^2.\end{aligned}$$

Taking expectations on both sides and canceling out terms, we get

$$\begin{aligned}\mathbb{E} \{ \|w_* - w_{N+1}^{\text{im}}\|^2 \} &\leq \frac{\alpha_N(1-\gamma)}{1+\alpha_N} \left(\frac{1+\alpha_1}{\alpha_1(1-\gamma)} - \lambda_{\min} \right) \|w_* - w_1^{\text{im}}\|^2 \\ &\quad + \frac{2\alpha_N}{1+\alpha_N} \sum_{n=1}^N \left(\frac{1+\alpha_n}{\alpha_n} \right) \mathbb{E} \{ \tilde{\alpha}_n \xi_n(w_n^{\text{im}}) \} + \frac{\alpha_N}{1+\alpha_N} \sum_{n=1}^N \alpha_n(1+\alpha_n) G^2 \quad (59)\end{aligned}$$

We will obtain upper bounds for the second and last terms in (59). We first establish an upper bound for the second term. For N large enough such that $N > \tau_{\alpha_N}$, we have

$$\begin{aligned}\sum_{n=1}^N \left(\frac{1+\alpha_n}{\alpha_n} \right) \mathbb{E} \{ \tilde{\alpha}_n \xi_n(w_n^{\text{im}}) \} &= \sum_{n=1}^{\tau_{\alpha_N}} \left(\frac{1+\alpha_n}{\alpha_n} \right) \mathbb{E} \{ \tilde{\alpha}_n \xi_n(w_n^{\text{im}}) \} + \sum_{n=\tau_{\alpha_N}+1}^N \left(\frac{1+\alpha_n}{\alpha_n} \right) \mathbb{E} \{ \tilde{\alpha}_n \xi_n(w_n^{\text{im}}) \} \\ &\leq \sum_{n=1}^{\tau_{\alpha_N}} \left(\frac{1+\alpha_n}{\alpha_n} \right) 6\alpha_n G^2 \sum_{i=1}^{n-1} \alpha_i + \sum_{n=\tau_{\alpha_N}+1}^N \left(\frac{1+\alpha_n}{\alpha_n} \right) \alpha_n G^2 (4 + 6\tau_{\alpha_N}) \alpha_{n-\tau_{\alpha_N}} \\ &\leq 6(1+\alpha_1) G^2 \sum_{n=1}^{\tau_{\alpha_N}} \sum_{i=1}^{n-1} \alpha_i + (1+\alpha_1) G^2 (4 + 6\tau_{\alpha_N}) \sum_{n=\tau_{\alpha_N}+1}^N \alpha_{n-\tau_{\alpha_N}} \\ &\leq 6(1+\alpha_1) G^2 \tau_{\alpha_N} \sum_{n=1}^N \alpha_i + (1+\alpha_1) G^2 (4 + 6\tau_{\alpha_N}) \sum_{n=1}^N \alpha_i \\ &= (1+\alpha_1) G^2 (4 + 12\tau_{\alpha_N}) \sum_{n=1}^N \alpha_n\end{aligned}$$

where the second inequality is due to Lemma B.13, and in the third inequality, we used $\alpha_n \leq \alpha_1$, and the last inequality is thanks to non-negativity of the sequence $(\alpha_n)_{n \in \mathbb{N}}$. Note that

$$\begin{aligned}\sum_{n=1}^N \alpha_n &= \alpha_1 + \sum_{n=2}^N \frac{\alpha_1}{\alpha_1 \lambda_{\min}(1-\gamma)(n-1) + 1} \leq \alpha_1 + \sum_{n=2}^N \frac{\alpha_1}{\alpha_1 \lambda_{\min}(1-\gamma)(n-1)} \\ &\leq \alpha_1 + \frac{1}{\lambda_{\min}(1-\gamma)} \sum_{n=1}^N \frac{1}{n} \\ &\leq \alpha_1 + \frac{(\log N + 1)}{\lambda_{\min}(1-\gamma)}, \quad (60)\end{aligned}$$

where the first inequality holds due to a smaller positive denominator, the second inequality comes from an additional positive term, and the last inequality is thanks to $\sum_{n=1}^N \frac{1}{n} \leq \log N + 1$. Therefore,

we have

$$\frac{2\alpha_N}{1+\alpha_N} \sum_{n=1}^N \left(\frac{1+\alpha_n}{\alpha_n} \right) \mathbb{E} \{ \tilde{\alpha}_n \xi_n(w_n^{\text{im}}) \} \leq \frac{2\alpha_N(1+\alpha_1)G^2(4+12\tau_{\alpha_N})}{1+\alpha_N} \left\{ \alpha_1 + \frac{(\log N + 1)}{\lambda_{\min}(1-\gamma)} \right\}. \quad (61)$$

For the third term in (59), notice that

$$\begin{aligned} \sum_{n=1}^N \alpha_n^2 &= \alpha_1^2 + \sum_{n=2}^N \left(\frac{\alpha_1}{\alpha_1 \lambda_{\min}(1-\gamma)(n-1) + 1} \right)^2 \leq \alpha_1^2 + \sum_{n=2}^N \left(\frac{\alpha_1}{\alpha_1 \lambda_{\min}(1-\gamma)(n-1)} \right)^2 \\ &\leq \alpha_1^2 + \frac{1}{\lambda_{\min}^2(1-\gamma)^2} \sum_{n=1}^N \frac{1}{n^2} \\ &\leq \alpha_1^2 + \frac{\pi^2}{6\lambda_{\min}^2(1-\gamma)^2}, \end{aligned} \quad (62)$$

where the first inequality again holds due to a smaller positive denominator, the second inequality comes from an additional positive term, and the last inequality is thanks to $\sum_{n=1}^{\infty} \frac{1}{n^2} \leq \sum_{n=1}^{\infty} \frac{1}{n^2} = \frac{\pi^2}{6}$. Utilizing (60) and (62), we observe that

$$G^2 \sum_{n=1}^N \alpha_n + G^2 \sum_{n=1}^N \alpha_n^2 \leq G^2 \left(\alpha_1 + \frac{(\log N + 1)}{\lambda_{\min}(1-\gamma)} \right) + G^2 \left(\alpha_1^2 + \frac{\pi^2}{6\lambda_{\min}^2(1-\gamma)^2} \right)$$

Therefore, the last term in (59) admits the following upper bound,

$$\frac{\alpha_N G^2}{1+\alpha_N} \left(\sum_{n=1}^N \alpha_n + \sum_{n=1}^N \alpha_n^2 \right) \leq \frac{\alpha_N G^2}{1+\alpha_N} \left\{ \alpha_1 + \frac{(\log N + 1)}{\lambda_{\min}(1-\gamma)} + \alpha_1^2 + \frac{\pi^2}{6\lambda_{\min}^2(1-\gamma)^2} \right\} \quad (63)$$

Combining (61) and (63), we get the following upperbound of (59), given by

$$\begin{aligned} \mathbb{E} \{ \|w_* - w_{N+1}^{\text{im}}\|^2 \} &\leq \frac{\alpha_N(1-\gamma)}{1+\alpha_N} \left(\frac{1+\alpha_1}{\alpha_1(1-\gamma)} - \lambda_{\min} \right) \|w_* - w_1\|^2 \\ &\quad + \frac{2\alpha_N(1+\alpha_1)G^2(4+12\tau_{\alpha_N})}{1+\alpha_N} \left\{ \alpha_1 + \frac{(\log N + 1)}{\lambda_{\min}(1-\gamma)} \right\} \\ &\quad + \frac{\alpha_N G^2}{1+\alpha_N} \left\{ \alpha_1 + \frac{(\log N + 1)}{\lambda_{\min}(1-\gamma)} + \alpha_1^2 + \frac{\pi^2}{6\lambda_{\min}^2(1-\gamma)^2} \right\}. \end{aligned}$$

The first term is of $O(\alpha_N)$, the second term is of $O(\alpha_N \log^2 N)$, and the last term is of $O(\alpha_N \log N)$. Combining all and suppressing the logarithmic complexity, the upper bound above is $\tilde{O}(1/N)$. As N goes to ∞ , we observe that $\mathbb{E} \{ \|w_* - w_{N+1}^{\text{im}}\|^2 \}$ tends to zero. \square

B.2.2 Finite-time/Convergence analysis for projected implicit TD(λ)

Recall that, in TD(λ) algorithm, we defined

$$\begin{aligned} S_n(w) &:= r_n e_n + \gamma e_n \phi_{n+1}^T w - e_n \phi_n^T w, \\ S(w) &:= \mathbb{E}_\infty [r_n e_{-\infty:n}] + \mathbb{E}_\infty [\gamma e_{-\infty:n} \phi_{n+1}^T] w - \mathbb{E}_\infty [e_{-\infty:n} \phi_n^T] w, \end{aligned}$$

where $e_{-\infty:n} := \sum_{k=0}^{\infty} (\lambda\gamma)^k \phi_{n-k}$. In addition to these notations, we also define

$$\begin{aligned} S_{\ell:n}(w) &:= r_n e_{\ell:n} + \gamma e_{\ell:n} \phi_{n+1}^T w - e_{\ell:n} \phi_n^T w, \\ \xi_n(w) &:= \{S_n(w) - S(w)\}^\top (w - w_*), \quad \forall w \in \mathbb{R}^d \\ \xi_{\ell:n}(w) &:= \{S_{\ell:n}(w) - S(w)\}^\top (w - w_*), \quad \forall w \in \mathbb{R}^d \end{aligned}$$

where $e_{\ell:n} := \sum_{k=0}^{n-\ell} (\lambda\gamma)^k \phi_{n-k}$. The following results from [4] will be used to both establish the finite-time error bound and asymptotic convergence.

Lemma B.14 (Lemma 16 of Bhandari et al. [4]). *For any $w \in \mathbb{R}^d$,*

$$(w_* - w)^\top S(w) \geq (1 - \kappa) \|V_{w_*} - V_w\|_D^2.$$

Lemma B.15 (Lemma 17 of Bhandari et al. [4]). *With probability 1, for all $w \in \{w' : \|w'\| \leq R\}$, $\|S_n(w)\| \leq B$, $\|S(w)\| \leq B$, where $B := \frac{r_{max} + 2R}{1 - \lambda\gamma}$.*

Lemma B.16. *With probability 1, for every $n \in \mathbb{N}$,*

$$\|w_* - w_{n+1}^{im}\|^2 \leq \|w_* - w_n^{im}\|^2 - \frac{2\alpha_n(1 - \lambda\gamma)^2(1 - \kappa)}{1 + \alpha_n} \|V_{w_*} - V_{w_n^{im}}\|_D^2 + 2\tilde{\alpha}_n \xi_n(w_n) + \alpha_n^2 B^2,$$

where $\kappa = \frac{\gamma(1-\lambda)}{1-\lambda\gamma}$ and $B = \frac{r_{max} + 2R}{1-\lambda\gamma}$.

Proof. *With probability one, the following derivations hold.*

$$\begin{aligned} \|w_* - w_{n+1}^{im}\|^2 &= \|w_* - \Pi_R\{w_n^{im} + \tilde{\alpha}_n S_n(w_n^{im})\}\|^2 \\ &= \|\Pi_R(w_*) - \Pi_R\{w_n^{im} + \tilde{\alpha}_n S_n(w_n^{im})\}\|^2 \end{aligned} \tag{64}$$

$$\leq \|w_* - w_n^{im} - \tilde{\alpha}_n S_n(w_n^{im})\|^2 \tag{65}$$

$$\begin{aligned} &= \|w_* - w_n^{im}\|^2 - 2\tilde{\alpha}_n S_n(w_n^{im})^\top (w_* - w_n^{im}) + \|\tilde{\alpha}_n S_n(w_n^{im})\|^2 \\ &\leq \|w_* - w_n^{im}\|^2 - 2\tilde{\alpha}_n S_n(w_n^{im})^\top (w_* - w_n^{im}) + \alpha_n^2 B^2 \end{aligned} \tag{66}$$

$$\begin{aligned} &= \|w_* - w_n^{im}\|^2 - 2\tilde{\alpha}_n S(w_n^{im})^\top (w_* - w_n^{im}) + 2\tilde{\alpha}_n \xi_n(w_n^{im}) + \alpha_n^2 B^2 \\ &\leq \|w_* - w_n^{im}\|^2 - 2\tilde{\alpha}_n(1 - \kappa) \|V_{w_*} - V_{w_n^{im}}\|_D^2 + 2\tilde{\alpha}_n \xi_n(w_n^{im}) + \alpha_n^2 B^2 \end{aligned} \tag{67}$$

$$\leq \|w_* - w_n^{im}\|^2 - \frac{2\alpha_n(1 - \lambda\gamma)^2(1 - \kappa)}{(1 - \lambda\gamma)^2 + \alpha_n} \|V_{w_*} - V_{w_n^{im}}\|_D^2 + 2\tilde{\alpha}_n \xi_n(w_n^{im}) + \alpha_n^2 B^2, \tag{68}$$

$$\leq \|w_* - w_n^{im}\|^2 - \frac{2\alpha_n(1 - \lambda\gamma)^2(1 - \kappa)}{1 + \alpha_n} \|V_{w_*} - V_{w_n^{im}}\|_D^2 + 2\tilde{\alpha}_n \xi_n(w_n^{im}) + \alpha_n^2 B^2, \tag{69}$$

where (64) is due to the fact that $w_* = \Pi_R(w_*)$, (65) is thanks to non-expansiveness of the projection operator on the convex set, (66) comes from Lemma B.15 with $\tilde{\alpha}_n \leq \alpha_n$, and (67) is obtained through Lemma B.14. Finally, (68) is the direct consequence of Lemma A.2 and (69) is due to $(1 - \lambda\gamma)^2 < 1$. \square

Lemma B.17 (Lemma 19 of Bhandari et al. [4]). *Given any $\ell \leq n$, for any arbitrary $w, v \in \{w' : \|w'\| \leq R\}$, with probability 1,*

1. $|\xi_{\ell:n}(w)| \leq 2B^2$.
2. $|\xi_{\ell:n}(w) - \xi_{\ell:n}(v)| \leq 6B\|w - v\|$.
3. $|\xi_n(w) - \xi_{n-\tau:n}(w)| \leq B^2(\lambda\gamma)^\tau$, for all $\tau \leq n$.
4. $|\xi_n(w) - \xi_{-\infty:n}(w)| \leq B^2(\lambda\gamma)^n$.

Definition B.18. *Given $\epsilon > 0$, we define a modified mixing time τ_{λ, α_N} to be*

$$\tau_\epsilon^\lambda = \min \{n \in \mathbb{N} \mid (\lambda\gamma)^n \leq \epsilon\},$$

$$\tau_{\lambda, \alpha_N} = \max \left\{ \tau_{\alpha_N}, \tau_{\alpha_N}^\lambda \right\}.$$

Lemma B.19. *Given a non-increasing sequence $\alpha_1 \geq \dots \geq \alpha_N$, for any fixed $n < N$, the following hold.*

1. For $2\tau_{\lambda, \alpha_N} < n$,

$$\mathbb{E} \left\{ \tilde{\alpha}_n \xi_n(w_n^{\text{im}}) \right\} \leq \alpha_n B^2 (12\tau_{\lambda, \alpha_N} + 7) \alpha_{n-2\tau_{\lambda, \alpha_N}}.$$

2. For $n \leq 2\tau_{\lambda, \alpha_N}$,

$$\mathbb{E} \left\{ \tilde{\alpha}_n \xi_n(w_n^{\text{im}}) \right\} \leq 6\alpha_n B^2 \sum_{i=1}^{n-1} \alpha_i + \alpha_n B^2 (\lambda\gamma)^n.$$

3. For all $n < N$,

$$\mathbb{E} \left\{ \tilde{\alpha}_n \xi_n(w_n^{\text{im}}) \right\} \leq \alpha_n B^2 (12\tau_{\lambda, \alpha_N} + 7) \alpha_1 + \alpha_n B^2 (\lambda\gamma)^n.$$

Proof. Claim 1: We first consider the case where $n > 2\tau_{\lambda, \alpha_N}$ and obtain a bound for $\mathbb{E} \left\{ \xi_n(w_n^{\text{im}}) \right\}$. Notice that

$$\mathbb{E} \left\{ \xi_n(w_n^{\text{im}}) \right\} \leq \left| \mathbb{E} \left\{ \xi_n(w_n^{\text{im}}) \right\} - \mathbb{E} \left\{ \xi_n(w_{n-2\tau_{\lambda, \alpha_N}}^{\text{im}}) \right\} \right| \quad (70)$$

$$+ \left| \mathbb{E} \left\{ \xi_n(w_{n-2\tau_{\lambda, \alpha_N}}^{\text{im}}) \right\} - \mathbb{E} \left\{ \xi_{n-\tau_{\lambda, \alpha_N}:n}(w_{n-2\tau_{\lambda, \alpha_N}}^{\text{im}}) \right\} \right| \quad (71)$$

$$+ \left| \mathbb{E} \left\{ \xi_{n-\tau_{\lambda, \alpha_N}:n}(w_{n-2\tau_{\lambda, \alpha_N}}^{\text{im}}) \right\} \right|. \quad (72)$$

To get an upper bound of the term in (70), notice that

$$\left| \xi_n(w_n^{\text{im}}) - \xi_n(w_{n-2\tau_{\lambda, \alpha_N}}^{\text{im}}) \right| \leq 6B \left\| w_n^{\text{im}} - w_{n-2\tau_{\lambda, \alpha_N}}^{\text{im}} \right\| \leq 6B \sum_{i=n-2\tau_{\lambda, \alpha_N}}^{n-1} \|w_{i+1}^{\text{im}} - w_i^{\text{im}}\|$$

where the second inequality comes from Lemma B.17 and the third inequality is thanks to the triangle inequality. Note that

$$\|w_{i+1}^{\text{im}} - w_i^{\text{im}}\| = \|\Pi_R(w_i^{\text{im}} + \tilde{\alpha}_i S_i(w_i^{\text{im}})) - \Pi_R(w_i^{\text{im}})\| \leq \|w_i^{\text{im}} + \tilde{\alpha}_i S_i(w_i^{\text{im}}) - w_i^{\text{im}}\| \leq \alpha_i B,$$

where in the first inequality, we have used the non-expansiveness of the projection operator, and for the second inequality, both Lemma A.2 and B.15 were used. Therefore, we have

$$\left| \xi_n(w_n^{\text{im}}) - \xi_n(w_{n-2\tau_{\lambda, \alpha_N}}^{\text{im}}) \right| \leq 6B^2 \sum_{i=n-2\tau_{\lambda, \alpha_N}}^{n-1} \alpha_i, \quad (73)$$

which leads to

$$\left| \mathbb{E} \left\{ \xi_n(w_n^{\text{im}}) \right\} - \mathbb{E} \left\{ \xi_n(w_{n-2\tau_{\lambda, \alpha_N}}^{\text{im}}) \right\} \right| \leq \mathbb{E} \left\{ \left| \xi_n(w_n^{\text{im}}) - \xi_n(w_{n-2\tau_{\lambda, \alpha_N}}^{\text{im}}) \right| \right\} \leq 6B^2 \sum_{i=n-2\tau_{\lambda, \alpha_N}}^{n-1} \alpha_i, \quad (74)$$

where the first inequality is due to the Jensen's inequality [19] and the second inequality is thanks to (73). Next, we obtain an upper bound of (71). From the third claim of Lemma B.17, we have

$$\left| \mathbb{E} \left\{ \xi_n(w_{n-2\tau_{\lambda, \alpha_N}}^{\text{im}}) \right\} - \mathbb{E} \left\{ \xi_{n-\tau_{\lambda, \alpha_N}:n}(w_{n-2\tau_{\lambda, \alpha_N}}^{\text{im}}) \right\} \right| \leq B^2 (\lambda \gamma)^{\tau_{\lambda, \alpha_N}} \leq B^2 \alpha_N, \quad (75)$$

where the last inequality is due to the definition of the modified mixing time τ_{λ, α_N} .

Next, we aim to obtain an upper bound of (72). Notice that for a fixed $w \in \{w' : \|w'\| \leq R\}$, $\xi_{n-\tau_{\lambda, \alpha_N}:n}(w)$ is a function of $u_{n-\tau_{\lambda, \alpha_N}}, \dots, u_{n-1}$, where $u_k = (x_k, r(x_k), x_{k+1})$ for $k = n - \tau_{\lambda, \alpha_N}, \dots, n$. Furthermore, we can view $w_{n-2\tau_{\lambda, \alpha_N}}^{\text{im}}$ as a function of $\{u_1, \dots, u_{n-2\tau_{\lambda, \alpha_N}-1}\}$. Now consider $\xi_{n-\tau_{\lambda, \alpha_N}:n}(w_{n-2\tau_{\lambda, \alpha_N}}^{\text{im}})$, which is a function of both $U = \{u_1, \dots, u_{n-2\tau_{\lambda, \alpha_N}-1}\}$ and $\tilde{U} = \{u_{n-\tau_{\lambda, \alpha_N}}, \dots, u_{n-1}\}$. We set $h(U, \tilde{U}) = \xi_{n-\tau_{\lambda, \alpha_N}:n}(w_{n-2\tau_{\lambda, \alpha_N}}^{\text{im}})$ to invoke Lemma B.10. The condition for Lemma B.10 is met since

$$U = \{u_1, \dots, u_{n-2\tau_{\lambda, \alpha_N}-1}\} \rightarrow \{u_{n-2\tau_{\lambda, \alpha_N}}, \dots, u_{n-\tau_{\lambda, \alpha_N}-1}\} \rightarrow \{u_{n-\tau_{\lambda, \alpha_N}}, \dots, u_{n-1}\} = \tilde{U}$$

forms a Markov chain. Therefore, we get

$$\mathbb{E} \left\{ h(U, \tilde{U}) \right\} - \mathbb{E} \left\{ h(U', \tilde{U}') \right\} \leq 2\|h\|_{\infty} m \rho^{\tau_{\lambda, \alpha_N}}, \quad (76)$$

where $U' = \{u'_1, \dots, u'_{n-2\tau_{\lambda, \alpha_N}-1}\}$ and $\tilde{U}' = \{u'_{n-\tau_{\lambda, \alpha_N}}, \dots, u'_{n-1}\}$ are independent and have the

same marginal distribution as U and \tilde{U} . Let us denote the $(n - 2\tau_{\lambda, \alpha_N})^{\text{th}}$ implicit TD(λ) iterate computed using U' as $w'_{n-2\tau_{\lambda, \alpha_N}}$. From the law of iterated expectations, we have

$$\mathbb{E} \left\{ h(U', \tilde{U}') \right\} = \mathbb{E} \left[\mathbb{E} \left\{ \xi_{n-\tau_{\lambda, \alpha_N}:n} \left(w'_{n-2\tau_{\lambda, \alpha_N}} \right) \middle| U' \right\} \right].$$

Now, for any fixed w , by the definition of $\xi_{n-\tau_{\lambda, \alpha_N}:n}(\cdot)$, we know

$$\begin{aligned} \mathbb{E} \left\{ \xi_{n-\tau_{\lambda, \alpha_N}:n} (w) \right\} &= \left[\mathbb{E} \left\{ S_{n-\tau_{\lambda, \alpha_N}:n}(w) \right\} - S(w) \right]^\top (w - w_*) \\ &= \mathbb{E} \left\{ S_{n-\tau_{\lambda, \alpha_N}:n}(w) - S_{-\infty:n}(w) \right\}^\top (w - w_*). \end{aligned}$$

The second equality follows from

$$\mathbb{E} \left\{ S_{n-\tau_{\lambda, \alpha_N}:n}(w) \right\} - S(w) = \mathbb{E} \left\{ S_{n-\tau_{\lambda, \alpha_N}:n}(w) \right\} - \mathbb{E} \left\{ S_{-\infty:n}(w) \right\} = \mathbb{E} \left\{ S_{n-\tau_{\lambda, \alpha_N}:n}(w) - S_{-\infty:n}(w) \right\}.$$

Notice that

$$\begin{aligned} \left| \left\{ S_{n-\tau_{\lambda, \alpha_N}:n}(w) - S_{-\infty:n}(w) \right\}^\top (w - w_*) \right| &= \left| \xi_{n-\tau_{\lambda, \alpha_N}:n}(w) - \xi_{-\infty:n}(w) \right| \\ &\leq \left| \xi_{n-\tau_{\lambda, \alpha_N}:n}(w) - \xi_n(w) \right| + \left| \xi_n(w) - \xi_{-\infty:n}(w) \right| \\ &\leq 2B^2(\lambda\gamma)^{\tau_{\lambda, \alpha_N}}, \end{aligned}$$

where the first inequality is due to the triangle inequality and the last inequality follows from combining claims 3 and 4 of Lemma B.17 with $\tau_{\lambda, \alpha_N} \leq n$. This yields

$$\mathbb{E} \left\{ h(U', \tilde{U}') \right\} \leq 2B^2(\lambda\gamma)^{\tau_{\lambda, \alpha_N}}. \quad (77)$$

Combining (76) and (77), we arrive at

$$\begin{aligned} \mathbb{E} \left\{ \xi_{n-\tau_{\lambda, \alpha_N}:n} \left(w_{n-\tau_{\lambda, \alpha_N}}^{\text{im}} \right) \right\} &= \mathbb{E} \left\{ h(U, \tilde{U}) \right\} \leq 2\|h\|_\infty m\rho^{\tau_{\lambda, \alpha_N}} + 2B^2(\lambda\gamma)^{\tau_{\lambda, \alpha_N}} \\ &\leq 4B^2m\rho^{\tau_{\lambda, \alpha_N}} + 2B^2(\lambda\gamma)^{\tau_{\lambda, \alpha_N}} \\ &\leq 6B^2\alpha_N \end{aligned} \quad (78)$$

where the second inequality is due to the first claim of Lemma B.17 and the last inequality is due to the definition of modified mixing time τ_{λ, α_N} .

Combining (74), (75) and (78), we get

$$\begin{aligned}
\mathbb{E}\{\xi_n(w_n^{\text{im}})\} &\leq 6B^2 \sum_{i=n-2\tau_{\lambda,\alpha_N}}^{n-1} \alpha_i + 7B^2\alpha_N \\
&\leq 12B^2\tau_{\lambda,\alpha_N}\alpha_{n-2\tau_{\lambda,\alpha_N}} + 7B^2\alpha_N \\
&\leq B^2(12\tau_{\lambda,\alpha_N} + 7)\alpha_{n-2\tau_{\lambda,\alpha_N}},
\end{aligned}$$

where both the second and third inequalities are due to non-increasingness of $(\alpha_n)_{n \in \mathbb{N}}$. Combined with Lemma A.2, we get the first claim.

Claim 2: We next consider the case where $n \leq 2\tau_{\lambda,\alpha_N}$. Using the triangle inequality, we get that

$$\mathbb{E}\{\xi_n(w_n^{\text{im}})\} \leq |\mathbb{E}\{\xi_n(w_n^{\text{im}})\} - \mathbb{E}\{\xi_n(w_1^{\text{im}})\}| \quad (79)$$

$$+ |\mathbb{E}\{\xi_n(w_1^{\text{im}})\} - \mathbb{E}\{\xi_{-\infty:n}(w_1^{\text{im}})\}| \quad (80)$$

$$+ |\mathbb{E}\{\xi_{-\infty:n}(w_1^{\text{im}})\}|. \quad (81)$$

An analogous argument in the proof for the first claim can be applied to obtain a bound for (79). Specifically, we have

$$|\xi_n(w_n^{\text{im}}) - \xi_n(w_1^{\text{im}})| \leq 6B \|w_n^{\text{im}} - w_1^{\text{im}}\| \leq 6B \sum_{i=1}^{n-1} \|w_{i+1}^{\text{im}} - w_i^{\text{im}}\|,$$

where the first inequality comes from Lemma B.17 and the second inequality is thanks to the triangle inequality. Recall that

$$\|w_{i+1}^{\text{im}} - w_i^{\text{im}}\| = \|\Pi_R\{w_i^{\text{im}} + \tilde{\alpha}_i S_i(w_i^{\text{im}})\} - \Pi_R(w_i^{\text{im}})\| \leq \|w_i^{\text{im}} + \tilde{\alpha}_i S_i(w_i^{\text{im}}) - w_i^{\text{im}}\| \leq \alpha_i B,$$

where in the first inequality, we have used the non-expansiveness of the projection operator, and for the second inequality, both Lemma A.2 and B.15 were used. Therefore, we have

$$|\xi_n(w_n^{\text{im}}) - \xi_n(w_1^{\text{im}})| \leq 6B^2 \sum_{i=1}^{n-1} \alpha_i, \quad (82)$$

which leads to

$$|\mathbb{E}\{\xi_n(w_n^{\text{im}})\} - \mathbb{E}\{\xi_n(w_1^{\text{im}})\}| \leq \mathbb{E}\{|\xi_n(w_n^{\text{im}}) - \xi_n(w_1^{\text{im}})|\} \leq 6B^2 \sum_{i=1}^{n-1} \alpha_i, \quad (83)$$

where the first inequality is due to the Jensen's inequality [19] and the second inequality is thanks to (82). Furthermore, from the fourth claim of Lemma B.17, we can obtain an upper bound of (80)

as follows

$$|\mathbb{E} \{ \xi_n (w_1^{\text{im}}) \} - \mathbb{E} \{ \xi_{-\infty:n} (w_1^{\text{im}}) \}| \leq B^2 (\lambda \gamma)^n. \quad (84)$$

Lastly, by definition, since w_1^{im} is fixed, we have $\mathbb{E} \{ \xi_{-\infty:n} (w_1^{\text{im}}) \} = 0$. Combining (83) and (84), we have

$$\mathbb{E} \{ \xi_n (w_n^{\text{im}}) \} \leq 6B^2 \sum_{i=1}^{n-1} \alpha_i + B^2 (\lambda \gamma)^n.$$

Combined with Lemma A.2, we get the second claim.

Claim 3: For $n \leq 2\tau_{\lambda, \alpha_N}$, observe that the bound we obtained in the previous claim admits the following upper bound, given by

$$6B^2 \sum_{i=1}^{n-1} \alpha_i + B^2 (\lambda \gamma)^n \leq 12B^2 \tau_{\lambda, \alpha_N} \alpha_1 + B^2 (\lambda \gamma)^n.$$

Since $\max \left\{ 12B^2 \tau_{\lambda, \alpha_N} \alpha_1 + B^2 (\lambda \gamma)^n, B^2 (12\tau_{\lambda, \alpha_N} + 7) \alpha_{n-2\tau_{\lambda, \alpha_N}} \right\} \leq B^2 \{ (12\tau_{\lambda, \alpha_N} + 7) \alpha_1 + (\lambda \gamma)^n \}$, the third claim directly follows from Lemma A.2. \square

We now establish a finite-time error bound of projected implicit TD(λ).

Proof of Theorem 4.12. Starting from Lemma B.16 with a constant step size, we have

$$\begin{aligned} \mathbb{E} \left\{ \|w_* - w_{n+1}^{\text{im}}\|^2 \right\} &\leq \mathbb{E} \left\{ \|w_* - w_n^{\text{im}}\|^2 \right\} - \frac{2\alpha(1-\lambda\gamma)^2(1-\kappa)}{1+\alpha} \mathbb{E} \left\{ \|V_{w_*} - V_{w_n^{\text{im}}} \|_D^2 \right\} \\ &\quad + 2\mathbb{E} \{ \tilde{\alpha}_n \xi_n (w_n^{\text{im}}) \} + \alpha^2 B^2. \end{aligned}$$

Then, for all $n < N$, we have

$$\begin{aligned} \mathbb{E} \left\{ \|w_* - w_{n+1}^{\text{im}}\|^2 \right\} &\leq \mathbb{E} \left\{ \|w_* - w_n^{\text{im}}\|^2 \right\} - \frac{2\alpha(1-\lambda\gamma)^2(1-\kappa)\lambda_{\min}}{1+\alpha} \mathbb{E} \left\{ \|w_* - w_n^{\text{im}}\|^2 \right\} \\ &\quad + 2\mathbb{E} \{ \tilde{\alpha}_n \xi_n (w_n^{\text{im}}) \} + \alpha^2 B^2 \\ &\leq \mathbb{E} \left\{ \|w_* - w_n^{\text{im}}\|^2 \right\} - \frac{2\alpha(1-\lambda\gamma)^2(1-\kappa)\lambda_{\min}}{1+\alpha} \mathbb{E} \left\{ \|w_* - w_n^{\text{im}}\|^2 \right\} \\ &\quad + \alpha^2 B^2 (24\tau_{\lambda, \alpha} + 14) + 2\alpha B^2 (\lambda \gamma)^n + \alpha^2 B^2 \\ &\leq \left\{ 1 - \frac{2\alpha(1-\lambda\gamma)^2(1-\kappa)\lambda_{\min}}{1+\alpha} \right\} \mathbb{E} \left\{ \|w_* - w_n^{\text{im}}\|^2 \right\} + \alpha^2 B^2 (24\tau_{\lambda, \alpha} + 15) + 2\alpha B^2, \end{aligned}$$

where the first inequality is due to Lemma B.7, which gives us $\|V_{w_*} - V_{w_n}\|_D^2 \geq \lambda_{\min} \|w_* - w_n\|^2$ and the second one is thanks to Lemma B.19 with a constant step size. In the final inequality, we

merged $\alpha_1^2 B^2$ terms and used the fact $\lambda\gamma \leq 1$. Then, we have

$$\begin{aligned}
& \mathbb{E} \left\{ \|w_* - w_{N+1}^{\text{im}}\|^2 \right\} \\
& \leq \left\{ 1 - \frac{2\alpha(1-\kappa)(1-\lambda\gamma)^2\lambda_{\min}}{1+\alpha} \right\} \mathbb{E} \left\{ \|w_* - w_n^{\text{im}}\|^2 \right\} + \alpha^2 B^2 (24\tau_{\lambda,\alpha} + 15) + 2\alpha B^2 \quad (85) \\
& \leq \left\{ 1 - \frac{2\alpha(1-\lambda\gamma)^2(1-\kappa)\lambda_{\min}}{1+\alpha} \right\}^N \|w_* - w_1^{\text{im}}\|^2 \\
& \quad + (\alpha^2 B^2 (24\tau_{\lambda,\alpha} + 15) + 2\alpha B^2) \sum_{t=0}^{\infty} \left\{ 1 - \frac{2\alpha(1-\lambda\gamma)^2(1-\kappa)\lambda_{\min}}{1+\alpha} \right\}^t \\
& \leq e^{-\frac{2\alpha(1-\lambda\gamma)^2(1-\kappa)\lambda_{\min}}{1+\alpha} N} \|w_* - w_1^{\text{im}}\|^2 + \frac{(1+\alpha) \{ \alpha B^2 (24\tau_{\lambda,\alpha} + 15) + 2B^2 \}}{2(1-\kappa)(1-\lambda\gamma)^2\lambda_{\min}},
\end{aligned}$$

where in the second inequality, we have recursively used the upper bound in (85) and further bounded the finite sum through an infinite sum. In the last inequality, we used $1 - x \leq \exp(-x)$, and an assumption $\frac{2\alpha(1-\lambda\gamma)^2(1-\kappa)\lambda_{\min}}{1+\alpha} \in (0, 1)$. \square

We next establish a convergence of the projected implicit TD(λ) with a sequence of decreasing step sizes.

Proof of Theorem 4.15. Rearranging terms in Lemma B.16, we have

$$\begin{aligned}
& \frac{\alpha_n(1-\lambda\gamma)^2(1-\kappa)}{1+\alpha_n} \|V_{w_*} - V_{w_n^{\text{im}}}\|_D^2 \\
& \leq \|w_* - w_n^{\text{im}}\|^2 - \frac{\alpha_n(1-\lambda\gamma)^2(1-\kappa)}{1+\alpha_n} \|V_{w_*} - V_{w_n^{\text{im}}}\|_D^2 - \|w_* - w_{n+1}^{\text{im}}\|^2 + 2\tilde{\alpha}_n \xi_n(w_n^{\text{im}}) + \alpha_n^2 B^2 \\
& \leq \left(1 - \frac{\alpha_n(1-\lambda\gamma)^2(1-\kappa)\lambda_{\min}}{1+\alpha_n} \right) \|w_* - w_n^{\text{im}}\|^2 - \|w_* - w_{n+1}^{\text{im}}\|^2 + 2\tilde{\alpha}_n \xi_n(w_n^{\text{im}}) + \alpha_n^2 B^2, \quad (86)
\end{aligned}$$

where we have used Lemma B.7 in (86). Dividing both sides by $\frac{\alpha_n(1-\lambda\gamma)^2(1-\kappa)}{1+\alpha_n}$ and from non-negativity of $\|V_{w_*} - V_{w_n^{\text{im}}}\|_D^2$, we have

$$\begin{aligned}
& \frac{1+\alpha_n}{\alpha_n(1-\lambda\gamma)^2(1-\kappa)} \left\{ \left(1 - \frac{\alpha_n(1-\lambda\gamma)^2(1-\kappa)\lambda_{\min}}{1+\alpha_n} \right) \|w_* - w_n^{\text{im}}\|^2 - \|w_* - w_{n+1}^{\text{im}}\|^2 + 2\tilde{\alpha}_n \xi_n(w_n^{\text{im}}) + \alpha_n^2 B^2 \right\} \\
& = \left(\frac{1+\alpha_n}{\alpha_n(1-\lambda\gamma)^2(1-\kappa)} - \lambda_{\min} \right) \|w_* - w_n^{\text{im}}\|^2 - \frac{1+\alpha_n}{\alpha_n(1-\lambda\gamma)^2(1-\kappa)} \|w_* - w_{n+1}^{\text{im}}\|^2 \\
& \quad + \frac{2(1+\alpha_n)}{\alpha_n(1-\lambda\gamma)^2(1-\kappa)} \tilde{\alpha}_n \xi_n(w_n^{\text{im}}) + \frac{\alpha_n(1+\alpha_n)}{(1-\lambda\gamma)^2(1-\kappa)} B^2 \geq 0 \quad (87)
\end{aligned}$$

With the choice of $\alpha_n = \frac{\alpha_1}{\alpha_1 \lambda_{\min}(1-\lambda\gamma)^2(1-\kappa)(n-1)+1}$, one can show that $\frac{1+\alpha_n}{\alpha_n(1-\lambda\gamma)^2(1-\kappa)} - \lambda_{\min} =$

$\frac{1+\alpha_{n-1}}{\alpha_{n-1}(1-\lambda\gamma)^2(1-\kappa)}$. Summing (87) over $n = 1, \dots, N$, we have

$$0 \leq \left(\frac{1+\alpha_1}{\alpha_1(1-\lambda\gamma)^2(1-\kappa)} - \lambda_{\min} \right) \|w_* - w_1^{\text{im}}\|^2 - \frac{1+\alpha_N}{\alpha_N(1-\lambda\gamma)^2(1-\kappa)} \|w_* - w_{N+1}^{\text{im}}\|^2 \\ + \sum_{n=1}^N \frac{2(1+\alpha_n)}{\alpha_n(1-\lambda\gamma)^2(1-\kappa)} \tilde{\alpha}_n \xi_n(w_n^{\text{im}}) + \sum_{n=1}^N \frac{\alpha_n(1+\alpha_n)}{(1-\lambda\gamma)^2(1-\kappa)} B^2.$$

Rearranging terms and dividing both sides by $\frac{1+\alpha_N}{\alpha_N(1-\lambda\gamma)^2(1-\kappa)}$, we have

$$\|w_* - w_{N+1}^{\text{im}}\|^2 \leq \frac{\alpha_N(1-\lambda\gamma)^2(1-\kappa)}{1+\alpha_N} \left(\frac{1+\alpha_1}{\alpha_1(1-\lambda\gamma)^2(1-\kappa)} - \lambda_{\min} \right) \|w_* - w_1^{\text{im}}\|^2 \\ + \frac{\alpha_N(1-\lambda\gamma)^2(1-\kappa)}{1+\alpha_N} \sum_{n=1}^N \frac{2(1+\alpha_n)}{\alpha_n(1-\lambda\gamma)^2(1-\kappa)} \tilde{\alpha}_n \xi_n(w_n^{\text{im}}) \\ + \frac{\alpha_N(1-\lambda\gamma)^2(1-\kappa)}{1+\alpha_N} \sum_{n=1}^N \frac{\alpha_n(1+\alpha_n)}{(1-\lambda\gamma)^2(1-\kappa)} B^2.$$

Taking expectations on both sides and canceling out terms, we get

$$\mathbb{E} \{ \|w_* - w_{N+1}^{\text{im}}\|^2 \} \leq \frac{\alpha_N(1-\lambda\gamma)^2(1-\kappa)}{1+\alpha_N} \left(\frac{1+\alpha_1}{\alpha_1(1-\lambda\gamma)^2(1-\kappa)} - \lambda_{\min} \right) \|w_* - w_1^{\text{im}}\|^2 \\ + \frac{2\alpha_N}{1+\alpha_N} \sum_{n=1}^N \left(\frac{1+\alpha_n}{\alpha_n} \right) \mathbb{E} \{ \tilde{\alpha}_n \xi_n(w_n^{\text{im}}) \} + \frac{\alpha_N}{1+\alpha_N} \sum_{n=1}^N \alpha_n(1+\alpha_n) B^2 \quad (88)$$

We will establish upper bounds for both the second and third terms in (88). To this end, first consider the second term in (88). For N large enough such that $N > 2\tau_{\lambda, \alpha_N}$, we have

$$\sum_{n=1}^N \left(\frac{1+\alpha_n}{\alpha_n} \right) \mathbb{E} \{ \tilde{\alpha}_n \xi_n(w_n^{\text{im}}) \} \quad (89) \\ = \sum_{n=1}^{2\tau_{\lambda, \alpha_N}} \left(\frac{1+\alpha_n}{\alpha_n} \right) \mathbb{E} \{ \tilde{\alpha}_n \xi_n(w_n^{\text{im}}) \} + \sum_{n=2\tau_{\lambda, \alpha_N}+1}^N \left(\frac{1+\alpha_n}{\alpha_n} \right) \mathbb{E} \{ \tilde{\alpha}_n \xi_n(w_n^{\text{im}}) \} \\ \leq \sum_{n=1}^{2\tau_{\lambda, \alpha_N}} \left(\frac{1+\alpha_n}{\alpha_n} \right) \alpha_n \left\{ 6B^2 \sum_{i=1}^{n-1} \alpha_i + B^2(\lambda\gamma)^n \right\} + \sum_{n=2\tau_{\lambda, \alpha_N}+1}^N \left(\frac{1+\alpha_n}{\alpha_n} \right) \alpha_n B^2 (12\tau_{\lambda, \alpha_N} + 7) \alpha_{n-2\tau_{\lambda, \alpha_N}} \\ = 6B^2 \sum_{n=1}^{2\tau_{\lambda, \alpha_N}} (1+\alpha_n) \left(\sum_{i=1}^{n-1} \alpha_i \right) + B^2 \sum_{n=1}^{2\tau_{\lambda, \alpha_N}} (1+\alpha_n)(\lambda\gamma)^n + B^2(12\tau_{\lambda, \alpha_N} + 7) \sum_{n=2\tau_{\lambda, \alpha_N}+1}^N (1+\alpha_n) \alpha_{n-2\tau_{\lambda, \alpha_N}} \\ \leq 12(1+\alpha_1)B^2\tau_{\lambda, \alpha_N} \sum_{i=1}^N \alpha_i + \frac{(1+\alpha_1)B^2}{1-\lambda\gamma} + B^2(12\tau_{\lambda, \alpha_N} + 7)(1+\alpha_1) \sum_{i=1}^N \alpha_i \\ = B^2(24\tau_{\lambda, \alpha_N} + 7)(1+\alpha_1) \sum_{i=1}^N \alpha_i + \frac{(1+\alpha_1)B^2}{1-\lambda\gamma} \quad (90)$$

where in the first inequality, we used Lemma B.19 and Lemma A.2, and in the second inequality where we used non-negativity and decreasing property of the sequence $(\alpha_n)_{n \in \mathbb{N}}$ as well as the fact $\sum_{n=1}^{2\tau_{\lambda, \alpha_N}} (\lambda\gamma)^n \leq \sum_{n=0}^{\infty} (\lambda\gamma)^n = \frac{1}{1-\lambda\gamma}$. Since

$$\begin{aligned}
\sum_{n=1}^N \alpha_i &\leq \sum_{n=1}^N \frac{\alpha_1}{\alpha_1 \lambda_{\min}(1-\kappa)(1-\lambda\gamma)^2(n-1)+1} \\
&= \alpha_1 + \sum_{n=2}^N \frac{1}{\lambda_{\min}(1-\kappa)(1-\lambda\gamma)^2(n-1)} \\
&\leq \alpha_1 + \frac{1}{\lambda_{\min}(1-\kappa)(1-\lambda\gamma)^2} \sum_{n=1}^N \frac{1}{n} \\
&\leq \alpha_1 + \frac{(\log N + 1)}{\lambda_{\min}(1-\kappa)(1-\lambda\gamma)^2}
\end{aligned} \tag{91}$$

where the first inequality holds due to a smaller positive denominator, the second inequality comes from an additional positive term, and the last inequality is thanks to $\sum_{n=1}^N \frac{1}{n} \leq \log N + 1$. Therefore, plugging (91) in (90), we get

$$\begin{aligned}
&\frac{2\alpha_N}{1+\alpha_N} \sum_{n=1}^N \left(\frac{1+\alpha_n}{\alpha_n} \right) \mathbb{E} \{ \tilde{\alpha}_n \xi_n(w_n^{\text{im}}) \} \\
&\leq \frac{\alpha_N B^2 (48\tau_{\lambda, \alpha_N} + 14)(1+\alpha_1)}{1+\alpha_N} \left(\alpha_1 + \frac{(\log N + 1)}{\lambda_{\min}(1-\kappa)(1-\lambda\gamma)^2} \right) + \frac{2\alpha_N(1+\alpha_1)B^2}{(1+\alpha_N)(1-\lambda\gamma)}.
\end{aligned} \tag{92}$$

For the third term in (88), notice that

$$\begin{aligned}
\sum_{n=1}^N \alpha_n^2 &= \alpha_1^2 + \sum_{n=2}^N \left(\frac{\alpha_1}{\alpha_1 \lambda_{\min}(1-\kappa)(1-\lambda\gamma)^2(n-1)+1} \right)^2 \\
&\leq \alpha_1^2 + \sum_{n=2}^N \left(\frac{\alpha_1}{\alpha_1 \lambda_{\min}(1-\kappa)(1-\lambda\gamma)^2(n-1)} \right)^2 \\
&\leq \alpha_1^2 + \frac{1}{\lambda_{\min}^2(1-\kappa)^2(1-\lambda\gamma)^4} \sum_{n=1}^N \frac{1}{n^2} \\
&\leq \alpha_1^2 + \frac{\pi^2}{6\lambda_{\min}^2(1-\kappa)^2(1-\lambda\gamma)^4}
\end{aligned} \tag{93}$$

where the first inequality again holds due to a smaller positive denominator, the second inequality comes from an additional positive term, and the last inequality is thanks to $\sum_{n=1}^{\infty} \frac{1}{n^2} \leq \sum_{n=1}^{\infty} \frac{1}{n^2} = \frac{\pi^2}{6}$. Utilizing (91) and (93), we observe that

$$B^2 \sum_{n=1}^N \alpha_n + B^2 \sum_{n=1}^N \alpha_n^2 \leq B^2 \left(\alpha_1 + \frac{(\log N + 1)}{\lambda_{\min}(1-\kappa)(1-\lambda\gamma)^2} \right) + B^2 \left(\alpha_1^2 + \frac{\pi^2}{6\lambda_{\min}^2(1-\kappa)^2(1-\lambda\gamma)^4} \right).$$

Therefore, the last term in (88) admits the following upper bound,

$$\frac{\alpha_N B^2}{1 + \alpha_N} \left(\sum_{n=1}^N \alpha_n + \sum_{n=1}^N \alpha_n^2 \right) \leq \frac{\alpha_N B^2}{1 + \alpha_N} \left\{ \alpha_1 + \frac{(\log N + 1)}{\lambda_{\min}(1 - \kappa)(1 - \lambda\gamma)^2} + \alpha_1^2 + \frac{\pi^2}{6\lambda_{\min}^2(1 - \kappa)^2(1 - \lambda\gamma)^4} \right\}. \quad (94)$$

Combining (92) and (94), we get the following upper bound of (88), given by

$$\begin{aligned} \mathbb{E} \{ \|w_* - w_{N+1}^{\text{im}}\|^2 \} &\leq \frac{\alpha_N(1 - \kappa)(1 - \lambda\gamma)^2}{1 + \alpha_N} \left(\frac{1 + \alpha_1}{\alpha_1(1 - \kappa)(1 - \lambda\gamma)^2} - \lambda_{\min} \right) \|w_* - w_1^{\text{im}}\|^2 \\ &\quad + \frac{\alpha_N B^2(48\tau_{\lambda, \alpha_N} + 14)(1 + \alpha_1)}{1 + \alpha_N} \left(\alpha_1 + \frac{(\log N + 1)}{\lambda_{\min}(1 - \kappa)(1 - \lambda\gamma)^2} \right) + \frac{2\alpha_N(1 + \alpha_1)B^2}{(1 + \alpha_N)(1 - \lambda\gamma)} \\ &\quad + \frac{\alpha_N B^2}{1 + \alpha_N} \left\{ \alpha_1 + \frac{(\log N + 1)}{\lambda_{\min}(1 - \kappa)(1 - \lambda\gamma)^2} + \alpha_1^2 + \frac{\pi^2}{6\lambda_{\min}^2(1 - \kappa)^2(1 - \lambda\gamma)^4} \right\}. \end{aligned}$$

The first term is of $O(\alpha_N)$, the second term is of $O(\alpha_N \log^2 N)$, and the last term is of $O(\alpha_N \log N)$.

Combining all and suppressing the logarithmic complexity, we observe that the upper bound above is $\tilde{O}(1/N)$. As N goes to ∞ , we observe that $\mathbb{E} \{ \|w_* - w_{N+1}^{\text{im}}\|^2 \}$ tends to zero. \square

C Theoretical analysis for implicit TDC

For the ease of presentation, we abbreviate the superscript for the implicit update and consider the following implicit TDC updates given by

$$w_{n+1} = w_n + \alpha'_n \rho_n (r_n + \gamma \phi_{n+1}^T w_n - \phi_n^T w_n) \phi_n - \alpha_n \rho_n \gamma (\phi_n^T u_n) \{ \phi_{n+1} - \alpha'_n \rho_n (\phi_n^T \phi_{n+1}) \phi_n \} \quad (95)$$

$$u_{n+1} = u_n + \beta'_n \rho_n (r_n + \gamma \phi_{n+1}^T w_n - \phi_n^T w_n) \phi_n - \beta'_n \rho_n \phi_n \phi_n^T u_n \quad (96)$$

where $\alpha'_n = \frac{\alpha_n}{1 + \alpha_n \rho_n \|\phi_n\|^2}$ and $\beta'_n = \frac{\beta_n}{1 + \beta_n \rho_n \|\phi_n\|^2}$. We first list notations and establish a linear stochastic approximation form of the implicit TDC update.

- π_b : behavioral policy / π_* : target policy
- μ_{π_b} = stationary distribution of the Markov chain $\{(x_n, a_n, x_{n+1})\}_{n \geq 0}$ under behavioral policy π_b
- μ_{π_*} = stationary distribution of the Markov chain $\{(x_n, a_n, x_{n+1})\}_{n \geq 0}$ under target policy π_*
- $\rho(x, a) = \pi(a|x)/\pi_b(a|x)$, $\rho_{\max} = \max_{x \in \mathcal{X}, a \in \mathcal{A}} \rho(x, a)$
- $O_n = (x_n, a_n, r_n, x_{n+1})$ denotes the observation at time n
- $A = \mathbb{E}_{\mu_{\pi_b}} \left[\rho(x, a) \phi(x) \{ \gamma \phi(x') - \phi(x) \}^\top \right]$, $A_n = \rho_n \phi_n (\gamma \phi_{n+1} - \phi_n)^\top$

- $B = -\gamma \mathbb{E}_{\mu_{\pi_b}} [\rho(x, a) \phi(x') \phi(x)^T]$
- $B_n^s = -\gamma \rho_n \phi_{n+1} \phi_n^T$, $B_n = -\gamma \rho_n \{\phi_{n+1} - \alpha'_n \rho_n (\phi_n^T \phi_{n+1}) \phi_n\} \phi_n^T$
- $C = -\mathbb{E}_{\mu_{\pi_b}} [\rho(x, a) \phi(x) \phi(x)^T]$, $C_n = -\rho_n \phi_n \phi_n^T$
- $b = \mathbb{E}_{\mu_{\pi_b}} [\rho(x, a) r(x) \phi(x)]$, $b_n = \rho_n r_n \phi_n$
- Tracking error vector: $v_n = u_n + C^{-1}(b + A w_n)$
- $f_1(w_n, O_n) = (A_n - B_n^s C^{-1} A) w_n + (b_n - B_n^s C^{-1} b)$,
- $\bar{f}_1(w_n) = (A - B C^{-1} A) w_n + (b - B C^{-1} b)$
- $g_1(v_n, O_n) = B_n^s v_n$, $\bar{g}_1(v_n) = B v_n$
- $f_2(w_n, O_n) = (A_n - C_n C^{-1} A) w_n + (b_n - C_n C^{-1} b)$
- $g_2(v_n, O_n) = C_n v_n$, $\bar{g}_2(v_n) = C v_n$
- $\lambda_c = \text{minimum absolute eigenvalue of the matrix } C$
- $\tau_{\alpha_t} = \min\{i \geq 0 : m \rho^i \leq \alpha_t\}$, $\tau_{\beta_t} = \min\{i \geq 0 : m \rho^i \leq \beta_t\}$

Based on the introduced notations, we can rewrite (95) and (96) as

$$\begin{aligned} w_{n+1} &= w_n + \alpha'_n (b_n + A_n w_n) + \alpha_n B_n u_n \\ u_{n+1} &= u_n + \beta'_n (b_n + A_n w_n + C_n u_n). \end{aligned}$$

Corresponding projected implicit TDC algorithms are provided below

$$w_{n+1} = \prod_{R_w} \{w_n + \alpha'_n (b_n + A_n w_n) + \alpha_n B_n u_n\} \quad (97)$$

$$u_{n+1} = \prod_{R_u} \{u_n + \beta'_n (b_n + A_n w_n + C_n u_n)\}. \quad (98)$$

To facilitate theoretical analysis, we rewrite the above projected linear stochastic approximation form into the following form:

$$\begin{aligned} w_{n+1} &= \Pi_{R_w} \{w_n + \alpha'_n (b_n + A_n w_n) + \alpha_n B_n u_n\} \\ &= \Pi_{R_w} \{w_n + \alpha'_n (b_n + A_n w_n + B_n^s u_n) + (\alpha_n B_n - \alpha'_n B_n^s) u_n\} \\ &= \Pi_{R_w} [w_n + \alpha'_n \{f_1(w_n, O_n) + g_1(v_n, O_n)\} + (\alpha_n B_n - \alpha'_n B_n^s) u_n] \end{aligned}$$

and introduce a tracking error vector $v_n = u_n + C^{-1}(b + Aw_n)$, whose iterative update rule is given below

$$\begin{aligned}
v_{n+1} &= \Pi_{R_u} \{v_n - C^{-1}(b + Aw_n) + \beta'_n(b_n + A_n w_n + C_n u_n)\} + C^{-1}(b + Aw_{n+1}) \\
&= \Pi_{R_u} [v_n - C^{-1}(b + Aw_n) + \beta'_n \{b_n - C_n C^{-1}b + A_n w_n - C_n C^{-1}Aw_n + C_n u_n + C_n C^{-1}(b + Aw_n)\}] \\
&\quad + C^{-1}(b + Aw_{n+1}) \\
&= \Pi_{R_u} [v_n + \beta'_n \{f_2(w_n, O_n) + g_2(v_n, O_n)\} - C^{-1}(b + Aw_n)] + C^{-1}(b + Aw_{n+1}).
\end{aligned}$$

C.1 Technical Lemmas for finite-time analysis for projected implicit TDC

In this section, we establish preliminary lemmas used in the proof of projected implicit TDC's finite-time error bounds.

Lemma C.1. *For all $n \geq 1$,*

$$\begin{aligned}
(a) \quad & \|B_n^s\| \leq \gamma \rho_{\max} \\
(b) \quad & \|B_n\| \leq \gamma \rho_{\max}(1 + c_\alpha \rho_{\max}) \\
(c) \quad & \|B_n - B_n^s\| \leq \gamma \rho_{\max}^2 \alpha_n \\
(d) \quad & \|\alpha_n B_n - \alpha'_n B_n^s\| \leq K_c \alpha_n^2
\end{aligned}$$

where K_c is a positive constant independent of n .

Proof. Recall the definition $B_n^s = -\gamma \rho_n \phi_{n+1} \phi_n^T$, $B_n = -\gamma \rho_n \{\phi_{n+1} - \alpha'_n \rho_n (\phi_n^T \phi_{n+1}) \phi_n\} \phi_n^T$. Part (a) follows from the normalized feature assumption with the Cauchy-Schwarz inequality. For part (b),

$$\|-\gamma \rho_n \{\phi_{n+1} - \alpha'_n \rho_n (\phi_n^T \phi_{n+1}) \phi_n\} \phi_n^T\| \leq \gamma \rho_{\max} (1 + \alpha'_n \rho_{\max}) \leq \gamma \rho_{\max} (1 + c_\alpha \rho_{\max}).$$

For part (c),

$$\|B_n - B_n^s\| = \|\gamma \rho_n^2 \alpha'_n (\phi_n^T \phi_{n+1}) \phi_n \phi_n^T\| \leq \gamma \rho_{\max}^2 \alpha_n.$$

For part (d),

$$\begin{aligned}
\|\alpha_n B_n - \alpha'_n B_n^s\| &= \alpha_n \left\| B_n - \frac{1}{1 + \alpha_n \rho_n \|\phi_n\|^2} B_n^s \right\| \leq \alpha_n \|B_n - B_n^s\| + \alpha_n \left(1 - \frac{1}{1 + \alpha_n \rho_n \|\phi_n\|^2}\right) \|B_n^s\| \\
&\leq \gamma \rho_{\max}^2 \alpha_n^2 + \frac{\alpha_n^2 \rho_n \|\phi_n\|^2}{1 + \alpha_n \rho_n \|\phi_n\|^2} \gamma \rho_{\max} \\
&\leq \gamma \rho_{\max}^2 \alpha_n^2 + \gamma \rho_{\max}^2 \alpha_n^2 =: K_c \alpha_n^2,
\end{aligned}$$

where the second inequality is due to parts (a) and (c). \square

Lemma C.2. For any $w \in \mathbb{R}^d$ such that $\|w\| \leq R_w$,

$$\|f_1(w, O_n)\| \leq K_{f_1}$$

for all $n \geq 1$. Here, K_{f_1} is a positive constant independent of w .

Proof. By the definition of $f_1(w, O_n)$, and $\lambda_c = \min |\lambda(C)|$, we obtain

$$\begin{aligned} \|f_1(w, O_n)\| &= \|(A_n - B_n^s C^{-1} A) w + (b_n - B_n^s C^{-1} b)\| \\ &\leq \|(A_n - B_n^s C^{-1} A) w\| + \|(b_n - B_n^s C^{-1} b)\| \\ &\leq (\|A_n\| + \|B_n^s\| \|C^{-1}\| \|A\|) \|w\| + \|b_n\| + \|B_n^s\| \|C^{-1}\| \|b\| \\ &\leq \left\{ (1 + \gamma) \rho_{\max} + \frac{1}{\lambda_c} \gamma (1 + \gamma) \rho_{\max}^2 \right\} R_w + \rho_{\max} r_{\max} + \frac{1}{\lambda_c} \gamma \rho_{\max}^2 r_{\max} =: K_{f_1}. \end{aligned}$$

□

Lemma C.3. For any $w \in \mathbb{R}^d$ such that $\|w\| \leq R_w$,

$$\|f_2(w, O_n)\| \leq K_{f_2}$$

for all $n \geq 1$. Here, K_{f_2} is a positive constant independent of w .

Proof. By the definition of $f_2(w, O_n)$, and $\lambda_c = \min |\lambda(C)|$, we obtain

$$\begin{aligned} \|f_2(w, O_n)\| &= \|(A_n - C_n C^{-1} A) w + (b_n - C_n C^{-1} b)\| \\ &\leq \|(A_n - C_n C^{-1} A) w\| + \|(b_n - C_n C^{-1} b)\| \\ &\leq (\|A_n\| + \|C_n\| \|C^{-1}\| \|A\|) \|w\| + \|b_n\| + \|C_n\| \|C^{-1}\| \|b\| \\ &\leq \left[(1 + \gamma) \rho_{\max} + \frac{1}{\lambda_c} (1 + \gamma) \rho_{\max}^2 \right] R_w + \rho_{\max} r_{\max} + \frac{1}{\lambda_c} \rho_{\max}^2 r_{\max} =: K_{f_2} \end{aligned}$$

□

Lemma C.4. Let $v := u + C^{-1}(b + Aw)$. Then for any $u, w \in \mathbb{R}^d$ such that $\|u\| \leq R_u$ and $\|w\|_2 \leq R_w$,

- (a) $\|v\| \leq R_v$
- (b) $\|g_1(v, O_n)\| \leq K_{g_1}$
- (c) $\|g_2(v, O_n)\| \leq K_{g_2}$

for all $n \geq 1$. Here, R_v , K_{g_1} and K_{g_2} are some positive constants independent of w and u .

Proof. For (a),

$$\begin{aligned}
\|v\| &= \|u + C^{-1}(b + Aw)\| \\
&\leq \|u\| + \|C^{-1}(b + Aw)\| \\
&\leq R_u + \frac{\rho_{\max} r_{\max} + \rho_{\max}(\gamma + 1)R_w}{\lambda_c} =: R_v
\end{aligned}$$

For (b), by the definition of $g_1(w, O_n)$, we obtain

$$\|g_1(v, O_n)\| = \|B_n^s v\| \leq \|B_n^s\| \|v\| \leq \gamma \rho_{\max} R_v =: K_{g_1}.$$

For (c), by the definition of $g_2(w, O_n)$, we obtain

$$\|g_2(v, O_n)\| = \|C_n v\| \leq \|C_n\| \|v\| \leq \rho_{\max} R_v =: K_{g_2}$$

□

Lemma C.5. Let $\zeta_{f_1}(w, O_n) := \langle f_1(w, O_n) - \bar{f}_1(w), w - w^* \rangle$. For any $w, w' \in \mathbb{R}^d$ such that $\|w\| \leq R_w$ and $\|w'\| \leq R_w$, we have

$$\begin{aligned}
(a) \quad & \|\zeta_{f_1}(w, O_n)\| \leq 4R_w K_{f_1} \\
(b) \quad & |\zeta_{f_1}(w, O_n) - \zeta_{f_1}(w', O_n)| \leq L_{f_1} \|w - w'\|
\end{aligned}$$

for all $n \geq 1$. Here L_{f_1} is a positive constant independent of w and w' .

Proof. For (a), following the same steps in Lemma C.2, we have $\|\bar{f}_1(w)\| \leq K_{f_1}$. Therefore, we get

$$\|\zeta_{f_1}(w, O_n)\| \leq (\|f_1(w, O_n)\| + \|\bar{f}_1(w)\|) (\|w\| + \|w^*\|) \leq 4R_w K_{f_1}.$$

For (b), we derive the bound as follows

$$\begin{aligned}
& |\zeta_{f_1}(w, O_n) - \zeta_{f_1}(w', O_n)| \\
&= |\langle f_1(w, O_n) - \bar{f}_1(w), w - w^* \rangle - \langle f_1(w', O_n) - \bar{f}_1(w'), w' - w^* \rangle| \\
&\leq \|w - w^*\| \|f_1(w, O_n) - \bar{f}_1(w) - f_1(w', O_n) + \bar{f}_1(w')\| + \|f_1(w', O_n) - \bar{f}_1(w')\| \|w - w'\| \\
&\leq \|w - w^*\| (\|f_1(w, O_n) - f_1(w', O_n)\| + \|\bar{f}_1(w') - \bar{f}_1(w)\|) + \|f_1(w', O_n) - \bar{f}_1(w')\| \|w - w'\| \\
&\leq 2R_w (\|(A_n - B_n^s C^{-1}A)(w - w')\| + \|(A - BC^{-1}A)(w' - w)\|) + 2K_{f_1} \|w - w'\| \\
&\leq 4R_w(1 + \gamma)\rho_{\max} \left(1 + \frac{1}{\lambda_{c,1}}\gamma\rho_{\max}\right) \|w - w'\| + 2K_{f_1} \|w - w'\| =: L_{f_1} \|w - w'\|.
\end{aligned}$$

□

Lemma C.6. Let $\zeta_{f_2}(w, v, O_n) := \langle f_2(w, O_n), v \rangle$. For any $w, w', v, v' \in \mathbb{R}^d$ such that $\|w\| \leq R_w$, $\|w'\| \leq R_w$, $\|v\| \leq R_v$ and $\|v'\| \leq R_v$,

$$\begin{aligned} (a) \quad & \|\zeta_{f_2}(w, v, O_n)\| \leq K_{f_2} R_v \\ (b) \quad & |\zeta_{f_2}(w, v, O_n) - \zeta_{f_2}(w', v', O_n)| \leq L_{f_2, w} \|w - w'\| + L_{f_2, v} \|v - v'\| \end{aligned}$$

for all $n \geq 1$. Here $L_{f_2, w}$ and $L_{f_2, v}$ are positive constants independent of w, w', v , and v' .

Proof. For (a), by the definition, we have

$$\|\zeta_{f_2}(w, v, O_n)\| = \|\langle f_2(w, O_n), v \rangle\| \leq \|f_2(w, O_n)\| \|v\| \leq K_{f_2} R_v.$$

For (b), we derive the bound as follows

$$\begin{aligned} |\zeta_{f_2}(w, v, O_n) - \zeta_{f_2}(w', v', O_n)| &= |\langle f_2(w, O_n), v \rangle - \langle f_2(w', O_n), v' \rangle| \\ &\leq \|v\| \|f_2(w, O_n) - f_2(w', O_n)\| + \|f_2(w', O_n)\| \|v - v'\| \\ &\leq R_v \|(A_n - C_n C^{-1} A)(w - w')\| + K_{f_2} \|v - v'\| \\ &\leq R_v \left\{ (1 + \gamma) \rho_{\max} + \frac{1}{\lambda_c} (1 + \gamma) \rho_{\max}^2 \right\} \|w - w'\| + K_{f_2} \|v - v'\| \\ &=: L_{f_2, w} \|w - w'\| + L_{f_2, v} \|v - v'\|. \end{aligned}$$

□

Lemma C.7. Let $\zeta_{g_2}(v, O_n) := \langle g_2(v, O_n) - \bar{g}_2(v), v \rangle$, where $v = u + C^{-1}(b + Aw)$. For all $v, v' \in \mathbb{R}^d$ such that $\|v\| \leq R_v$ and $\|v'\| \leq R_v$, we have

$$\begin{aligned} (a) \quad & \|\zeta_{g_2}(v, O_n)\| \leq 2K_{g_2} R_v \\ (b) \quad & |\zeta_{g_2}(v, O_n) - \zeta_{g_2}(v', O_n)| \leq L_{g_2} \|v - v'\| \end{aligned}$$

for all $n \geq 1$. Here, L_{g_2} is a positive constant independent of v and v' .

Proof. Following the same steps in Lemma C.4, we know $\|\bar{g}_2(v)\| \leq K_{g_2}$. For part (a), by the definition of $\zeta_{g_2}(v, O_n)$, we have

$$\|\zeta_{g_2}(v, O_n)\| = \|\langle g_2(v, O_n) - \bar{g}_2(v), v \rangle\| \leq (\|g_2(v, O_n)\| + \|\bar{g}_2(v)\|) \|v\| \leq 2K_{g_2} R_v.$$

For part (b), we derive the bound as follows.

$$\begin{aligned}
& |\zeta_{g_2}(v, O_n) - \zeta_{g_2}(v', O_n)| \\
&= |\langle g_2(v, O_n) - \bar{g}_2(v), v \rangle - \langle g_2(v', O_n) - \bar{g}_2(v'), v' \rangle| \\
&\leq \|v\| \|g_2(v, O_n) - \bar{g}_2(v) - g_2(v', O_n) + \bar{g}_2(v')\| + \|g_2(v', O_n) - \bar{g}_2(v')\| \|v - v'\| \\
&= \|v\| \|(C_n - C)(v - v')\| + \|g_2(v', O_n) - \bar{g}_2(v')\| \|v - v'\| \\
&\leq R_v (\|C_n\| + \|C\|) \|v - v'\| + 2K_{g_2} \|v - v'\| \\
&\leq 2R_v \rho_{\max} \|v - v'\| + 2K_{g_2} \|v - v'\| =: L_{g_2} \|v - v'\|
\end{aligned}$$

□

Lemma C.8. *For a positive integer t , suppose $i \leq t$ and $(\alpha_n)_{n \in \mathbb{N}}$ is a non-increasing sequence with $\alpha_1 = c_\alpha$. If $i \leq \tau_{\alpha_t}$,*

$$\mathbb{E} \{\zeta_{f_1}(w_i, O_i)\} \leq L_{f_1} \{(K_{f_1} + K_{g_1}) c_\alpha + K_c R_u c_\alpha^2\} \tau_{\alpha_t}.$$

Otherwise,

$$\mathbb{E} \{\zeta_{f_1}(w_i, O_i)\} \leq 8R_w K_{f_1} \alpha_t + L_{f_1} \{(K_{f_1} + K_{g_1}) \alpha_{i-\tau_{\alpha_t}} + K_c R_u \alpha_{i-\tau_{\alpha_t}}^2\} \tau_{\alpha_t}.$$

Proof. Note that for any $i \geq 1$,

$$\begin{aligned}
\|w_{i+1} - w_i\| &= \|\Pi_{R_w} [w_i + \alpha'_i \{f_1(w_i, O_i) + g_1(v_i, O_i)\} + \alpha_i B_i u_i - \alpha'_i B_i^s u_i] - \Pi_{R_w} w_i\| \\
&\leq \|w_i + \alpha'_i \{f_1(w_i, O_i) + g_1(v_i, O_i)\} + \alpha_i B_i u_i - \alpha'_i B_i^s u_i - w_i\| \\
&\leq \alpha_i \|f_1(w_i, O_i) + g_1(v_i, O_i)\| + \|\alpha_i B_i u_i - \alpha'_i B_i^s u_i\| \\
&\leq \alpha_i (K_{f_1} + K_{g_1}) + K_c R_u \alpha_i^2,
\end{aligned} \tag{99}$$

where the last inequality follows from Lemma C.1, Lemma C.2, and Lemma C.4. Applying the Lipschitz continuous property of $\zeta_{f_1}(w, O_i)$, obtained in part (b) of Lemma C.5, for $i > \tau_{\alpha_t}$, it follows that

$$|\zeta_{f_1}(w_i, O_i) - \zeta_{f_1}(w_{i-\tau_{\alpha_t}}, O_i)| \leq L_{f_1} \|w_i - w_{i-\tau_{\alpha_t}}\| \leq L_{f_1} (K_{f_1} + K_{g_1}) \sum_{k=i-\tau_{\alpha_t}}^{i-1} \alpha_k + L_{f_1} K_c R_u \sum_{k=i-\tau_{\alpha_t}}^{i-1} \alpha_k^2$$

We now provide an upper bound for $\mathbb{E} \{\zeta_{f_1}(w_{i-\tau_{\alpha_t}}, O_i)\}$. To this end, we define $w'_{i-\tau_{\alpha_t}}$ and $O'_i = (s'_i, a'_i, r'_i, s'_{i+1})$, which are drawn independently from the marginal distributions of $w_{i-\tau_{\alpha_t}}$ and O_i . From part (a) of Lemma C.5 and Lemma B.10, we have

$$\mathbb{E} \{\zeta_{f_1}(w_{i-\tau_{\alpha_t}}, O_i)\} \leq \left| \mathbb{E} \{\zeta_{f_1}(w_{i-\tau_{\alpha_t}}, O_i)\} - \mathbb{E} \{\zeta_{f_1}(w'_{i-\tau_{\alpha_t}}, O'_i)\} \right| \leq 8R_w K_{f_1} m \rho^{\tau_{\alpha_t}}.$$

It follows that

$$\begin{aligned}
\mathbb{E} \{ \zeta_{f_1} (w_i, O_i) \} &\leq \mathbb{E} \{ \zeta_{f_1} (w_{i-\tau_{\alpha_t}}, O_i) \} + L_{f_1} (K_{f_1} + K_{g_1}) \sum_{k=i-\tau_{\alpha_t}}^{i-1} \alpha_k + L_{f_1} K_c R_u \sum_{k=i-\tau_{\alpha_t}}^{i-1} \alpha_k^2 \\
&\leq 8R_w K_{f_1} m \rho^{\tau_{\alpha_t}} + L_{f_1} \left\{ (K_{f_1} + K_{g_1}) \alpha_{i-\tau_{\alpha_t}} + K_c R_u \alpha_{i-\tau_{\alpha_t}}^2 \right\} \tau_{\alpha_t} \\
&\leq 8R_w K_{f_1} \alpha_t + L_{f_1} \left\{ (K_{f_1} + K_{g_1}) \alpha_{i-\tau_{\alpha_t}} + K_c R_u \alpha_{i-\tau_{\alpha_t}}^2 \right\} \tau_{\alpha_t}.
\end{aligned}$$

On the other hand, if $i \leq \tau_{\alpha_t}$,

$$\begin{aligned}
\mathbb{E} \{ \zeta_{f_1} (w_i, O_i) \} &\leq \mathbb{E} \{ \zeta_{f_1} (w_0, O_i) \} + L_{f_1} (K_{f_1} + K_{g_1}) \sum_{k=0}^{i-1} \alpha_k + L_{f_1} K_c R_u \sum_{k=0}^{i-1} \alpha_k^2 \\
&\leq L_{f_1} (K_{f_1} + K_{g_1}) \tau_{\alpha_t} c_\alpha + L_{f_1} K_c R_u \tau_{\alpha_t} c_\alpha^2 \\
&= L_{f_1} \left\{ (K_{f_1} + K_{g_1}) c_\alpha + K_c R_u c_\alpha^2 \right\} \tau_{\alpha_t}.
\end{aligned}$$

□

Lemma C.9. *Given a positive integer t , suppose $i \leq t$. Furthermore, $(\alpha_n)_{n \in \mathbb{N}}$ and $(\beta_n)_{n \in \mathbb{N}}$ are non-increasing sequences with $\alpha_1 = c_\alpha$, $\beta_1 = c_\beta$. If α_n/β_n is a non-increasing sequence, for $i \leq \tau_{\beta_t}$,*

$$\mathbb{E} \{ \zeta_{f_2} (w_i, v_i, O_i) \} \leq c_\beta K_{r_3} \tau_{\beta_t},$$

and for $i > \tau_{\beta_t}$,

$$\mathbb{E} \{ \zeta_{f_2} (w_i, v_i, O_i) \} \leq 2R_v K_{f_2} \beta_t + K_{r_3} \tau_{\beta_t} \beta_{i-\tau_{\beta_t}},$$

for some constant $K_{r_3} > 0$.

Proof. Notice that

$$\begin{aligned}
\|v_{i+1} - v_i\| &= \|\Pi_{R_u} [v_i + \beta'_i \{f_2(w_i, O_i) + g_2(v_i, O_i)\} - C^{-1}(b + Aw_i)] + C^{-1}(b + Aw_{i+1}) - v_i\| \\
&= \|\Pi_{R_u} [v_i + \beta'_i \{f_2(w_i, O_i) + g_2(v_i, O_i)\} - C^{-1}(b + Aw_i)] + C^{-1}(b + Aw_i) - v_i + C^{-1}A(w_{i+1} - w_i)\| \\
&\leq \|\Pi_{R_u} [v_i + \beta'_i \{f_2(w_i, O_i) + g_2(v_i, O_i)\} - C^{-1}(b + Aw_i)] - \Pi_{R_u} \{v_i - C^{-1}(b + Aw_i)\}\| \\
&\quad + \|C^{-1}A(w_{i+1} - w_i)\| \\
&\leq \|v_i + \beta'_i \{f_2(w_i, O_i) + g_2(v_i, O_i)\} - C^{-1}(b + Aw_i) - \{v_i - C^{-1}(b + Aw_i)\}\| + \|C^{-1}A(w_{i+1} - w_i)\| \\
&\leq \beta_i \|f_2(w_i, O_i) + g_2(v_i, O_i)\| + \|C^{-1}A(w_{i+1} - w_i)\| \\
&\leq \beta_i (K_{f_2} + K_{g_2}) + \alpha_i \|C^{-1}\| \|A\| (K_{f_1} + K_{g_1}) + \|C^{-1}\| \|A\| K_c R_u \alpha_i^2 \\
&\leq \beta_i \left\{ K_{f_2} + K_{g_2} + \frac{\alpha_i (1 + \gamma) \rho_{\max}}{\beta_i \lambda_{c,1}} (K_{f_1} + K_{g_1}) + \frac{\alpha_i^2 (1 + \gamma) \rho_{\max}}{\beta_i \lambda_{c,1}} K_c R_u \right\} \\
&\leq \beta_i \left\{ K_{f_2} + K_{g_2} + \frac{c_\alpha (1 + \gamma) \rho_{\max}}{c_\beta \lambda_{c,1}} (K_{f_1} + K_{g_1}) + \frac{c_\alpha^2 (1 + \gamma) \rho_{\max}}{c_\beta \lambda_{c,1}} K_c R_u \right\} =: \beta_i K_{r_2} \tag{100}
\end{aligned}$$

where the first inequality follows from the fact that $\|u_i\|_2 = \|v_i - C^{-1}(b + Aw_i)\|_2 \leq R_u$, the third inequality is thanks to (99) and the last inequality is due to $\frac{\alpha_i}{\beta_i} \leq \frac{c_\alpha}{c_\beta}$ and $\frac{\alpha_i^2}{\beta_i} \leq \frac{c_\alpha^2}{c_\beta}$. Applying the Lipschitz continuous property of $\zeta_{f_2}(w, v, O_i)$ in part (b) of Lemma C.6, for $i > \tau_{\beta_t}$ it follows that

$$\begin{aligned} & \left| \zeta_{f_2}(w_i, v_i, O_i) - \zeta_{f_2}(w_{i-\tau_{\beta_t}}, v_{i-\tau_{\beta_t}}, O_i) \right| \leq L_{f_2, w} \|w_i - w_{i-\tau_{\beta_t}}\| + L_{f_2, v} \|v_i - v_{i-\tau_{\beta_t}}\| \\ & \leq L_{f_2, w} (K_{f_1} + K_{g_1}) \sum_{k=i-\tau_{\beta_t}}^{i-1} \alpha_k + L_{f_2, w} K_c R_u \sum_{k=i-\tau_{\beta_t}}^{i-1} \alpha_k^2 + L_{f_2, v} K_{r_2} \sum_{k=i-\tau_{\beta_t}}^{i-1} \beta_k, \end{aligned}$$

where the second inequality follows from (99) and (100). The next step is to provide an upper bound for $\mathbb{E} \left\{ \zeta_{f_2}(w_{i-\tau_{\beta_t}}, v_{i-\tau_{\beta_t}}, O_i) \right\}$. To this end, we define $(w'_{i-\tau_{\beta_t}}, v'_{i-\tau_{\beta_t}})$ and $O'_i = (s'_i, a'_i, r'_i, s'_{i+1})$ which are independently drawn from the marginal distributions of $(w_{i-\tau_{\beta_t}}, v_{i-\tau_{\beta_t}})$ and O_i . It can be shown that $\mathbb{E} \left\{ \zeta_{f_2}(w'_{i-\tau_{\beta_t}}, v'_{i-\tau_{\beta_t}}, O'_i) \right\} = 0$. By Lemma C.6 and Lemma B.10, we get

$$\mathbb{E} \left\{ \zeta_{f_2}(w_{i-\tau}, v_{i-\tau}, O_i) \right\} \leq \left| \mathbb{E} \left\{ \zeta_{f_2}(w_{i-\tau}, v_{i-\tau}, O_i) \right\} - \mathbb{E} \left\{ \zeta_{f_2}(w'_{i-\tau}, v'_{i-\tau}, O'_i) \right\} \right| \leq 2R_v K_{f_2} m \rho^{\tau_{\beta_t}}.$$

It follows that

$$\begin{aligned} & \mathbb{E} \left\{ \zeta_{f_2}(w_i, z_i, O_i) \right\} \\ & \leq \mathbb{E} \left\{ \zeta_{f_2}(w_{i-\tau_{\beta_t}}, z_{i-\tau_{\beta_t}}, O_i) \right\} + L_{f_2, w} (K_{f_1} + K_{g_1}) \sum_{k=i-\tau_{\beta_t}}^{i-1} \alpha_k + L_{f_2, w} K_c R_u \sum_{k=i-\tau_{\beta_t}}^{i-1} \alpha_k^2 + L_{f_2, v} K_{r_2} \sum_{k=i-\tau_{\beta_t}}^{i-1} \beta_k \\ & \leq 2R_v K_{f_2} m \rho^{\tau_{\beta_t}} + L_{f_2, w} (K_{f_1} + K_{g_1}) \tau_{\beta_t} \alpha_{i-\tau_{\beta_t}} + L_{f_2, w} K_c R_u \tau_{\beta_t} \alpha_{i-\tau_{\beta_t}}^2 + L_{f_2, v} K_{r_2} \tau_{\beta_t} \beta_{i-\tau_{\beta_t}} \\ & \leq 2R_v K_{f_2} \beta_t + \left[\max \left\{ 1, \frac{c_\alpha}{c_\beta} \right\} L_{f_2, w} (K_{f_1} + K_{g_1}) + \max \left\{ 1, \frac{c_\alpha^2}{c_\beta} \right\} L_{f_2, w} K_c R_u + L_{f_2, v} K_{r_2} \right] \tau_{\beta_t} \beta_{i-\tau_{\beta_t}} \\ & = 2R_v K_{f_2} \beta_t + K_{r_3} \tau_{\beta_t} \beta_{i-\tau_{\beta_t}} \end{aligned}$$

where the last inequality is thanks to $\frac{\alpha_{i-\tau_{\beta_t}}^2}{\beta_{i-\tau_{\beta_t}}} \leq \frac{c_\alpha^2}{c_\beta}$. Similarly, for $i \leq \tau_{\beta_t}$, it follows that

$$\begin{aligned} \mathbb{E} \left\{ \zeta_{f_2}(w_i, v_i, O_i) \right\} & \leq \mathbb{E} \left\{ \zeta_{f_2}(w_0, v_0, O_i) \right\} + L_{f_2, w} (K_{f_1} + K_{g_1}) \sum_{k=0}^{i-1} \alpha_k + L_{f_2, w} K_c R_u \sum_{k=0}^{i-1} \alpha_k^2 + L_{f_2, v} K_{r_2} \sum_{k=0}^{i-1} \beta_k \\ & \leq L_{f_2, w} (K_{f_1} + K_{g_1}) \tau_{\beta_t} c_\alpha + L_{f_2, w} K_c R_u \tau_{\beta_t} c_\alpha^2 + L_{f_2, v} K_{r_2} \tau_{\beta_t} c_\beta \\ & \leq c_\beta \left[\max \left\{ 1, \frac{c_\alpha}{c_\beta} \right\} L_{f_2, w} (K_{f_1} + K_{g_1}) + \max \left\{ 1, \frac{c_\alpha^2}{c_\beta} \right\} L_{f_2, w} K_c R_u + L_{f_2, v} K_{r_2} \right] \tau_{\beta_t} \\ & =: c_\beta K_{r_3} \tau_{\beta_t}. \end{aligned}$$

□

Lemma C.10. For a positive integer t , suppose $i \leq t$ and $(\beta_n)_{n \in \mathbb{N}}$ is a non-increasing sequence

with $\beta_1 = c_\beta$. If $i \leq \tau_{\beta_t}$,

$$\mathbb{E} \{ \zeta_{g_2}(v_i, O_i) \} \leq c_\beta L_{g_2} K_{r_2} \tau_{\beta_t}.$$

Otherwise,

$$\mathbb{E} \{ \zeta_{g_2}(v_i, O_i) \} \leq 4R_w K_{g_2} \beta_t + L_{g_2} K_{r_2} \tau_{\beta_t} \beta_{i-\tau_{\beta_t}}.$$

Proof. Applying the Lipschitz continuous property of $\zeta_{g_2}(v, O_i)$ established in part (b) of Lemma C.7 with (100), for $i > \tau_{\beta_t}$, we have

$$\left| \zeta_{g_2}(v_i, O_i) - \zeta_{g_2}(v_{i-\tau_{\beta_t}}, O_i) \right| \leq L_{g_2} \left\| v_i - v_{i-\tau_{\beta_t}} \right\| \leq L_{g_2} K_{r_2} \sum_{k=i-\tau_{\beta_t}}^{i-1} \beta_k.$$

Like the previous two Lemmas, we provide an upper bound for $\mathbb{E} \left\{ \zeta_{g_2}(v_{i-\tau_{\beta_t}}, O_i) \right\}$. To this end, we define an independent $v'_{i-\tau_{\beta_t}}$ and $O'_i = (s'_i, a'_i, r'_i, s'_{i+1})$ which are independently drawn from marginal distributions of $v_{i-\tau_{\beta_t}}$ and O_i . Using part (a) of Lemma C.7 and Lemma B.10, we obtain

$$\mathbb{E} \left\{ \zeta_{g_2}(v_{i-\tau_{\beta_t}}, O_i) \right\} \leq \left| \mathbb{E} \left\{ \zeta_{g_2}(v_{i-\tau_{\beta_t}}, O_i) \right\} - \mathbb{E} \left\{ \zeta_{f_2}(v'_{i-\tau_{\beta_t}}, O'_i) \right\} \right| \leq 4R_w K_{g_2} m \rho^{\tau_{\beta_t}}.$$

Therefore, it follows that

$$\begin{aligned} \mathbb{E} \{ \zeta_{g_2}(v_i, O_i) \} &\leq \mathbb{E} \left\{ \zeta_{g_2}(v_{i-\tau_{\beta_t}}, O_i) \right\} + L_{g_2} K_{r_2} \sum_{k=i-\tau_{\beta_t}}^{i-1} \beta_k \leq 4R_w K_{g_2} m \rho^{\tau_{\beta_t}} + L_{g_2} K_{r_2} \tau_{\beta_t} \beta_{i-\tau_{\beta_t}} \\ &\leq 4R_w K_{g_2} \beta_t + L_{g_2} K_{r_2} \tau_{\beta_t} \beta_{i-\tau_{\beta_t}}. \end{aligned}$$

For $i \leq \tau_{\beta_t}$, we have that

$$\mathbb{E} \{ \zeta_{g_2}(v_i, O_i) \} \leq \mathbb{E} \{ \zeta_{g_2}(v_0, O_i) \} + L_{g_2} K_{r_2} \sum_{k=0}^{i-1} \beta_k \leq L_{g_2} K_{r_2} i c_\beta \leq c_\beta L_{g_2} K_{r_2} \tau_{\beta_t}.$$

□

Lemma C.11. Suppose $\beta_n = \frac{c_\beta}{n^\nu}$, $n \in \mathbb{N}$ and $\lambda > 0$. Let $T_n = \sum_{k=1}^{n-1} \beta_k$, then for all $n \in \mathbb{N}$,

$$\begin{aligned} (a) \quad & \sum_{i=1}^{\tau_{\beta_t}} e^{-\lambda \sum_{k=i+1}^t \beta_k} \beta_i \leq \frac{e^{\lambda c_\beta}}{\lambda} e^{-\frac{\lambda c_\beta}{1-\nu}} \left[(1+t)^{1-\nu} - (1+\tau_{\beta_t})^{1-\nu} \right] \\ (b) \quad & \beta_t \sum_{i=\tau_{\beta_t}+1}^t e^{-\lambda \sum_{k=i+1}^t \beta_k} \beta_i \leq \frac{e^{\lambda c_\beta} c_\beta}{\lambda t^\nu} \\ (c) \quad & \sum_{i=\tau_{\beta_t}+1}^t e^{-\lambda \sum_{k=i+1}^t \beta_k} \beta_{i-\tau_{\beta_t}} \beta_i \leq \left(e^{\frac{-\lambda c_\beta}{2(1-\nu)}} [(t+1)^{1-\nu} - 1] D_\beta \mathbb{I}_{\tau_{\beta_t}+1 < i_{f_\beta}} + \beta_{t-\tau_{\beta_t}} \right) \frac{2e^{\lambda c_\beta/2}}{\lambda} \end{aligned}$$

where $D_\beta = e^{(\lambda/2) \sum_{k=1}^{i_{f_\beta}} \beta_k} c_\beta$ for some $i_{f_\beta} \in \mathbb{N}$.

Proof. Let us use the convention $\sum_{k=n+1}^n \beta_k = 0$. For part (a), we have

$$\begin{aligned} \sum_{i=1}^{\tau_{\beta_t}} e^{-\lambda \sum_{k=i+1}^t \beta_k} \beta_i &\leq \max_{i \geq 1} \left\{ e^{\lambda \beta_i} \right\} \sum_{i=1}^{\tau_{\beta_t}} e^{-\lambda \sum_{k=i}^t \beta_k} \beta_i = e^{\lambda c_\beta} \sum_{i=1}^{\tau_{\beta_t}} e^{-\lambda (T_{t+1} - T_i)} \beta_i \\ &\leq e^{\lambda c_\beta} \int_0^{T_{\tau_{\beta_t}+1}} e^{-\lambda (T_{t+1} - s)} ds \leq \frac{e^{\lambda c_\beta}}{\lambda} e^{-\lambda (T_{t+1} - T_{\tau_{\beta_t}+1})} \\ &= \frac{e^{\lambda c_\beta}}{\lambda} e^{-\lambda c_\beta \sum_{k=\tau_{\beta_t}+1}^t 1/k^\nu} \leq \frac{e^{\lambda c_\beta}}{\lambda} e^{-\frac{\lambda c_\beta}{1-\nu} [(1+t)^{1-\nu} - (1+\tau_{\beta_t})^{1-\nu}]} \end{aligned}$$

For part (b), we have

$$\begin{aligned} \beta_t \sum_{i=\tau_{\beta_t}+1}^t e^{-\lambda \sum_{k=i+1}^t \beta_k} \beta_i &\leq \max_{i \geq 1} \left\{ e^{\lambda \beta_i} \right\} \beta_t \sum_{i=\tau_{\beta_t}+1}^t e^{-\lambda \sum_{k=i}^t \beta_k} \beta_i = e^{\lambda c_\beta} \beta_t \sum_{i=\tau_{\beta_t}+1}^t e^{-\lambda (T_{t+1} - T_i)} \beta_i \\ &\leq e^{\lambda c_\beta} \beta_t \int_{T_{\tau_{\beta_t}+1}}^{T_{t+1}} e^{-\lambda (T_{t+1} - s)} ds = \frac{e^{\lambda c_\beta}}{\lambda} \beta_t \left(1 - e^{-\lambda (T_{t+1} - T_{\tau_{\beta_t}+1})} \right) \leq \frac{e^{\lambda c_\beta} c_\beta}{\lambda t^\nu} \end{aligned}$$

For part (c), we have

$$\begin{aligned} \sum_{i=\tau_{\beta_t}+1}^t e^{-\lambda \sum_{k=i+1}^t \beta_k} \beta_{i-\tau_{\beta_t}} \beta_i &\leq \max_{i \in [\tau_{\beta_t}+1, t]} \left\{ e^{(-\lambda/2) \sum_{k=i+1}^t \beta_k} \beta_{i-\tau_{\beta_t}} \right\} \sum_{i=\tau_{\beta_t}+1}^t e^{(-\lambda/2) \sum_{k=i+1}^t \beta_k} \beta_i \\ &\leq \max_{i \in [\tau_{\beta_t}+1, t]} \left\{ e^{(-\lambda/2) \sum_{k=i+1}^t \beta_k} \beta_{i-\tau_{\beta_t}} \right\} \frac{2e^{\lambda c_\beta/2}}{\lambda} \end{aligned} \quad (101)$$

where the second inequality follows from the same argument as in part (b). To bound the first term in (101), note that $e^{(-\lambda/2) \sum_{k=i+1}^t \beta_k} \beta_{i-\tau_{\beta_t}}$ is eventually increasing. In other words, there exists $i_{f_\beta} \in \mathbb{N}$ such that,

$$\max_{i \in [\tau_{\beta_t}+1, t]} \left\{ e^{(-\lambda/2) \sum_{k=i+1}^t \beta_k} \beta_{i-\tau_{\beta_t}} \right\} = \beta_{t-\tau_{\beta_t}} \quad \text{if } \tau_{\beta_t} + 1 \geq i_{f_\beta}.$$

If $\tau_{\beta_t} + 1 < i_{f_\beta}$, then

$$\begin{aligned}
& \max_{i \in [\tau_{\beta_t} + 1, t]} \left\{ e^{(-\lambda/2) \sum_{k=i+1}^t \beta_k} \beta_{i-\tau_{\beta_t}} \right\} \\
& \leq \max_{i \in [\tau_{\beta_t} + 1, i_{f_\beta}]} \left\{ e^{(-\lambda/2) \sum_{k=i+1}^t \beta_k} \beta_{i-\tau_{\beta_t}} \right\} + \max_{i \in [i_{f_\beta} + 1, t]} \left\{ e^{(-\lambda/2) \sum_{k=i+1}^t \beta_k} \beta_{i-\tau_{\beta_t}} \right\} \\
& \leq e^{(-\lambda/2) \sum_{k=1}^t \beta_k} \max_{i \in [\tau_{\beta_t} + 1, i_{f_\beta}]} \left\{ e^{(\lambda/2) \sum_{k=1}^i \beta_k} \beta_{i-\tau_{\beta_t}} \right\} + \beta_{t-\tau_{\beta_t}} \\
& \leq e^{(-\lambda/2) \sum_{k=1}^t \beta_k} e^{(\lambda/2) \sum_{k=1}^{i_{f_\beta}} \beta_k} \beta_1 + \beta_{t-\tau_{\beta_t}} \\
& \leq e^{\frac{-\lambda c_\beta}{2(1-\nu)} [(t+1)^{1-\nu} - 1]} D_\beta + \beta_{t-\tau_{\beta_t}}
\end{aligned}$$

where $D_\beta = e^{(\lambda/2) \sum_{k=1}^{i_{f_\beta}} \beta_k} c_\beta$. Combining everything, we get

$$\sum_{i=\tau_{\beta_t}+1}^t e^{-\lambda \sum_{k=i+1}^t \beta_k} \beta_{i-\tau_{\beta_t}} \beta_i \leq \left(e^{\frac{-\lambda c_\beta}{2(1-\nu)} [(t+1)^{1-\nu} - 1]} D_\beta \mathbb{I}_{\tau_{\beta_t} + 1 < i_{f_\beta}} + \beta_{t-\tau_{\beta_t}} \right) \frac{2e^{\lambda c_\beta/2}}{\lambda}.$$

□

Lemma C.12. Suppose $\alpha_n = \frac{c_\alpha}{n^\sigma}$, $n \in \mathbb{N}$ and $\lambda > 0$. Let $T_n = \sum_{i=1}^{n-1} \alpha_i$, then for all $n \in \mathbb{N}$,

$$\sum_{i=1}^n \left(e^{-\lambda \sum_{k=i+1}^n \alpha_k} \right) \alpha_i^2 \leq \left\{ K_b e^{-\frac{\lambda}{2} \sum_{k=1}^n \alpha_k} + \alpha_n \right\} \frac{2e^{\frac{\lambda c_\alpha}{2}}}{\lambda},$$

where $K_b = c_\alpha e^{\frac{\lambda}{2} \sum_{k=1}^{i_0} \alpha_k}$ for some $i_0 \in \mathbb{N}$.

Proof. Let us use the convention $\sum_{k=n+1}^n \alpha_k = 0$ and $\sum_{k=n+1}^n \alpha_k^2 = 0$. Notice that

$$\begin{aligned}
\sum_{i=1}^n \left(e^{-\frac{\lambda}{2} \sum_{k=i+1}^n \alpha_k} \right) \alpha_i & \leq \left(\sup_{i \geq 1} e^{\frac{\lambda}{2} \alpha_i} \right) \left\{ \sum_{i=1}^n \left(e^{-\frac{\lambda}{2} \sum_{k=i}^n \alpha_k} \right) \alpha_i \right\} = \left(\sup_{i \geq 1} e^{\frac{\lambda}{2} \alpha_i} \right) \left\{ \sum_{i=1}^n \left(e^{-\frac{\lambda}{2} (T_{n+1} - T_i)} \right) \alpha_i \right\} \\
& \leq \left(\sup_{i \geq 1} e^{\frac{\lambda}{2} \alpha_i} \right) \int_0^{T_{n+1}} e^{-\frac{\lambda}{2} (T_{n+1} - s)} ds \\
& \leq \left(\sup_{i \geq 1} e^{\frac{\lambda}{2} \alpha_i} \right) \frac{2}{\lambda} \leq \frac{2e^{\frac{\lambda c_\alpha}{2}}}{\lambda}
\end{aligned} \tag{102}$$

Now consider

$$\begin{aligned} \sum_{i=1}^n \left(e^{-\lambda \sum_{k=i+1}^n \alpha_k} \right) \alpha_i^2 &\leq \sup_{1 \leq i \leq n} \left(\alpha_i e^{-\frac{\lambda}{2} \sum_{k=i+1}^n \alpha_k} \right) \left\{ \sum_{i=1}^n \left(e^{-\frac{\lambda}{2} \sum_{k=i+1}^n \alpha_k} \right) \alpha_i \right\} \\ &\leq \sup_{1 \leq i \leq n} \left(\alpha_i e^{-\frac{\lambda}{2} \sum_{k=i+1}^n \alpha_k} \right) \frac{2e^{\frac{\lambda c_\alpha}{2}}}{\lambda} \end{aligned} \quad (103)$$

where the last inequality follows from (102). Note that $\alpha_i e^{-\frac{\lambda}{2} \sum_{k=i+1}^n \alpha_k}$ is monotonically increasing after some time $i_0 \in \mathbb{N}$, i.e., for $n \geq i_0$, we have

$$\sup_{i_0 \leq i \leq n} \left\{ \alpha_i \exp \left(-\frac{\lambda}{2} \sum_{k=i+1}^n \alpha_k \right) \right\} \leq \alpha_n = \frac{c_\alpha}{n^\sigma}.$$

Therefore, we have

$$\begin{aligned} (103) &\leq \left\{ \sup_{1 \leq i \leq i_0} \left(\alpha_i e^{-\frac{\lambda}{2} \sum_{k=i+1}^n \alpha_k} \right) + \alpha_n \right\} \frac{2e^{\frac{\lambda c_\alpha}{2}}}{\lambda} \leq \left\{ e^{-\frac{\lambda}{2} \sum_{k=1}^n \alpha_k} \sup_{1 \leq i \leq i_0} \left(\alpha_i e^{\frac{\lambda}{2} \sum_{k=1}^i \alpha_k} \right) + \alpha_n \right\} \frac{2e^{\frac{\lambda c_\alpha}{2}}}{\lambda} \\ &\leq \left\{ K_b e^{-\frac{\lambda}{2} \sum_{k=1}^n \alpha_k} + \alpha_n \right\} \frac{2e^{\frac{\lambda c_\alpha}{2}}}{\lambda}, \end{aligned}$$

where $K_b = c_\alpha e^{\frac{\lambda}{2} \sum_{k=1}^{i_0} \alpha_k}$. □

Lemma C.13. For $0 < \sigma < 1$, $\lambda_w < 0$, let $\alpha_t = \frac{c_\alpha}{t^\sigma}$, $\alpha'_t = \frac{\alpha_t}{1 + \alpha_t \rho_t \|\phi_t\|^2}$ and $\underline{\alpha}_t = \frac{\alpha_t}{1 + c_\alpha \rho_{\max}}$. Then

$$\sum_{i=1}^t e^{\lambda_w \sum_{k=i+1}^t \alpha_k} \mathbb{E} [\alpha'_i \zeta_{f_1}(w_i, O_i)] = O(\log t / t^\sigma).$$

Proof. Applying Lemma C.8 combined with the fact that $0 < \alpha'_i \leq \alpha_i$ holds almost surely, it

follows that

$$\begin{aligned}
& \sum_{i=1}^t e^{\lambda_w \sum_{k=i+1}^t \underline{\alpha}_k} \mathbb{E} [\alpha'_i \zeta_{f_1}(w_i, O_i)] \\
& \leq L_{f_1} (K_{f_1} + K_{g_1}) c_\alpha \tau_{\alpha_t} \sum_{i=1}^{\tau_{\alpha_t}} e^{\lambda_w \sum_{k=i+1}^t \underline{\alpha}_k} \alpha_i + L_{f_1} K_c R_u c_\alpha^2 \tau_{\alpha_t} \sum_{i=1}^{\tau_{\alpha_t}} e^{\lambda_w \sum_{k=i+1}^t \underline{\alpha}_k} \alpha_i \\
& \quad + 8R_w K_{f_1} \alpha_t \sum_{i=\tau_\alpha+1}^t e^{\lambda_w \sum_{k=i+1}^t \underline{\alpha}_k} \alpha_i + L_{f_1} (K_{f_1} + K_{g_1}) \tau_{\alpha_t} \sum_{i=\tau_\alpha+1}^t e^{\lambda_w \sum_{k=i+1}^t \underline{\alpha}_k} \alpha_{i-\tau_{\alpha_t}} \alpha_i \\
& \quad + L_{f_1} K_c R_u \tau_{\alpha_t} \sum_{i=\tau_\alpha+1}^t e^{\lambda_w \sum_{k=i+1}^t \underline{\alpha}_k} \alpha_{i-\tau_{\alpha_t}}^2 \alpha_i \\
& \leq L_{f_1} \{ (K_{f_1} + K_{g_1}) c_\alpha + K_c R_u c_\alpha^2 \} \tau_{\alpha_t} \sum_{i=1}^{\tau_{\alpha_t}} e^{\lambda_w \sum_{k=i+1}^t \underline{\alpha}_k} \alpha_i + 8R_w K_{f_1} \alpha_t \sum_{i=\tau_\alpha+1}^t e^{\lambda_w \sum_{k=i+1}^t \underline{\alpha}_k} \alpha_i \\
& \quad + L_{f_1} \{ (K_{f_1} + K_{g_1}) + K_c R_u c_\alpha \} \tau_{\alpha_t} \sum_{i=\tau_\alpha+1}^t e^{\lambda_w \sum_{k=i+1}^t \underline{\alpha}_k} \alpha_{i-\tau_{\alpha_t}} \alpha_i
\end{aligned}$$

Applying Lemma C.11, we obtain:

$$\begin{aligned}
& \sum_{i=1}^{\tau_{\alpha_t}} e^{\lambda_w \sum_{k=i+1}^t \underline{\alpha}_k} \alpha_i \leq \frac{e^{-\frac{\lambda_w c_\alpha}{1+c_\alpha \rho_{\max}}} (1 + c_\alpha \rho_{\max})}{-\lambda_w} e^{\frac{\lambda_w c_\alpha}{(1+c_\alpha \rho_{\max})(1-\sigma)}} [(1+t)^{1-\sigma} - (1+\tau_{\alpha_t})^{1-\sigma}] \\
& \alpha_t \sum_{i=\tau_{\alpha_t}+1}^t e^{\lambda_w \sum_{k=i+1}^t \underline{\alpha}_k} \alpha_i \leq \left\{ \frac{e^{-\frac{\lambda_w c_\alpha}{1+c_\alpha \rho_{\max}}} (1 + c_\alpha \rho_{\max})}{-\lambda_w} \right\} \frac{c_\alpha}{t^\sigma} \\
& \quad \sum_{i=\tau_{\alpha_t}+1}^t e^{\lambda_w \sum_{k=i+1}^t \underline{\alpha}_k} \alpha_{i-\tau_{\alpha_t}} \alpha_i \\
& \leq \left\{ \frac{2e^{-\frac{\lambda_w c_\alpha}{2(1+c_\alpha \rho_{\max})}} (1 + c_\alpha \rho_{\max})}{-\lambda_w} \right\} \left\{ e^{\frac{\lambda_w c_\alpha}{2(1+c_\alpha \rho_{\max})(1-\sigma)}} [(t+1)^{1-\sigma} - 1] D_\alpha \mathbb{I}_{\{\tau_{\alpha_t}+1 < i_{f_\alpha}\}} + \alpha_{t-\tau_{\alpha_t}} \right\}
\end{aligned}$$

where $D_\alpha = e^{(\lambda_w/2) \sum_{k=1}^{i_{f_\alpha}} \underline{\alpha}_k} c_\alpha$ for some $i_{f_\alpha} \in \mathbb{N}$. Combined with the fact that $\tau_{\alpha_t} = O(\log t)$, we obtain the desired result. \square

Lemma C.14. For $0 < \nu < 1$, $\lambda_u < 0$, let $\beta_t = \frac{c_\beta}{t^\nu}$, $\beta'_t = \frac{\beta_t}{1+\beta_t \rho_t \|\phi_t\|^2}$ and $\underline{\beta}_t = \frac{\beta_t}{1+c_\beta \rho_{\max}}$. Then

$$\sum_{i=1}^t e^{\lambda_u \sum_{k=i+1}^t \underline{\beta}_k} \mathbb{E} [\beta'_i \zeta_{f_2}(w_i, v_i, O_i)] = O(\log t / t^\nu).$$

Proof. Applying Lemma C.9 combined with the fact that $0 < \beta'_i \leq \beta_i$ holds almost surely, it

follows that

$$\begin{aligned} \sum_{i=1}^t e^{\lambda_u \sum_{k=i+1}^t \underline{\beta}_k} \mathbb{E} [\beta'_i \zeta_{f_2}(w_i, v_i, O_i)] &\leq c_\beta K_{r_3} \tau_{\beta_t} \sum_{i=1}^{\tau_{\beta_t}} e^{\lambda_u \sum_{k=i+1}^t \underline{\beta}_k} \beta_i + 2R_w K_{f_2} \beta_t \sum_{i=\tau_{\beta_t}+1}^t e^{\lambda_u \sum_{k=i+1}^t \underline{\beta}_k} \beta_i \\ &\quad + K_{r_3} \tau_{\beta_t} \sum_{i=\tau_{\beta_t}+1}^t e^{\lambda_u \sum_{k=i+1}^t \underline{\beta}_k} \beta_{i-\tau_{\beta_t}} \beta_i \end{aligned}$$

Applying Lemma C.11, we have

$$\begin{aligned} &\sum_{i=1}^t e^{\lambda_u \sum_{k=i+1}^t \underline{\beta}_k} \mathbb{E} [\beta'_i \zeta_{f_2}(w_i, v_i, O_i)] \\ &\leq c_\beta K_{r_3} \tau_{\beta_t} (1 + c_\beta \rho_{\max}) \left(\frac{e^{-\frac{\lambda_u c_\beta}{1+c_\beta \rho_{\max}}}}{-\lambda_u} \right) e^{\frac{\lambda_u c_\beta}{(1-\nu)(1+c_\beta \rho_{\max})}} [(1+t)^{1-\nu} - (1+\tau_{\beta_t})^{1-\nu}] \\ &\quad + 2R_w K_{f_2} \left\{ \frac{e^{-\frac{\lambda_u c_\beta}{1+c_\beta \rho_{\max}}} (c_\beta + c_\beta^2 \rho_{\max})}{-\lambda_u t^\nu} \right\} \\ &\quad + K_{r_3} \tau_{\beta_t} \left\{ e^{\frac{\lambda_u c_\beta}{2(1-\nu)(1+c_\beta \rho_{\max})}} [(t+1)^{1-\nu} - 1] D_\beta \mathbb{I}_{\tau_{\beta_t}+1 < i_{f_1}} + \beta_{t-\tau_{\beta_t}} \right\} \left\{ \frac{2e^{-\frac{\lambda_u c_\beta}{2(1+c_\beta \rho_{\max})}} (1 + c_\beta \rho_{\max})}{-\lambda_u} \right\}. \end{aligned}$$

Combined with the fact that $\tau_{\beta_t} = O(\log t)$, we obtain the desired result. \square

Lemma C.15. For $0 < \nu < 1, \lambda_u < 0$, let $\beta_t = \frac{c_\beta}{t^\nu}$, $\beta'_t = \frac{\beta_t}{1+\beta_t \rho_t \|\phi_t\|^2}$ and $\underline{\beta}_t = \frac{\beta_t}{1+c_\beta \rho_{\max}}$. Then

$$\sum_{i=1}^t e^{\lambda_u \sum_{k=i+1}^t \underline{\beta}_k} \mathbb{E} [\beta'_i \zeta_{g_2}(v_i, O_i)] = O(\log t / t^\nu).$$

Proof. Applying Lemma C.10 combined with the fact that $0 < \beta'_i \leq \beta_i$ holds almost surely, it follows that

$$\begin{aligned} \sum_{i=1}^t e^{\lambda_u \sum_{k=i+1}^t \underline{\beta}_k} \mathbb{E} [\beta'_i \zeta_{g_2}(v_i, O_i)] &\leq c_\beta L_{g_2} K_{r_2} \tau_{\beta_t} \sum_{i=1}^{\tau_{\beta_t}} e^{\lambda_u \sum_{k=i+1}^t \underline{\beta}_k} \beta_i + 4R_w K_{g_2} \beta_t \sum_{i=\tau_{\beta_t}+1}^t e^{\lambda_u \sum_{k=i+1}^t \underline{\beta}_k} \beta_i \\ &\quad + L_{g_2} K_{r_2} \tau_{\beta_t} \sum_{i=\tau_{\beta_t}+1}^t e^{\lambda_u \sum_{k=i+1}^t \underline{\beta}_k} \beta_{i-\tau_{\beta_t}} \beta_i \end{aligned}$$

Now we invoke Lemma C.11, and get

$$\begin{aligned}
& \sum_{i=1}^t e^{\lambda_u \sum_{k=i+1}^t \underline{\beta}_k} \mathbb{E} [\beta'_i \zeta_{g_2}(v_i, O_i)] \\
& \leq c_\beta L_{g_2} K_{r_2} \tau_{\beta_t} (1 + c_\beta \rho_{\max}) \left(\frac{e^{-\frac{\lambda_u c_\beta}{1+c_\beta \rho_{\max}}}}{-\lambda_u} \right) e^{\frac{\lambda_u c_\beta}{(1-\nu)(1+c_\beta \rho_{\max})}} [(1+t)^{1-\nu} - (1+\tau_{\beta_t})^{1-\nu}] \\
& \quad + 4R_w K_{g_2} (1 + c_\beta \rho_{\max}) \left(\frac{e^{-\frac{\lambda_u c_\beta}{1+c_\beta \rho_{\max}}} c_\beta}{-\lambda_u t^\nu} \right) \\
& \quad + L_{g_2} K_{r_2} \tau_{\beta_t} (1 + c_\beta \rho_{\max}) \left\{ e^{\frac{\lambda_u c_\beta}{2(1-\nu)(1+c_\beta \rho_{\max})}} [(t+1)^{1-\nu} - 1] D_\beta \mathbb{I}_{\tau_{\beta_t} + 1 < i_{f_1}} + \beta_{t-\tau_{\beta_t}} \right\} \left\{ \frac{2e^{-\frac{\lambda_u c_\beta}{2(1+c_\beta \rho_{\max})}}}{-\lambda_u} \right\}.
\end{aligned}$$

Combined with the fact that $\tau_{\alpha_t} = O(\log t)$, we obtain the desired result. \square

Lemma C.16. For given $0 < \nu < \sigma < 1$, $\lambda_u < 0$, let $\beta_t = \frac{c_\beta}{t^\nu}$, $\alpha_t = \frac{c_\alpha}{t^\sigma}$ and $\underline{\beta}_t = \frac{\beta_t}{1+c_\beta \rho_{\max}}$. Then

$$\sum_{i=1}^t e^{\lambda_u \sum_{k=i+1}^t \underline{\beta}_k} \mathbb{E} \langle C^{-1} A(w_{i+1} - w_i), v_i \rangle = O\left(\frac{1}{t^{\sigma-\nu}}\right).$$

Proof. Using the same arguments as in Lemma C.8, it follows that

$$\begin{aligned}
& \sum_{i=1}^t e^{\lambda_u \sum_{k=i+1}^t \underline{\beta}_k} \mathbb{E} \langle C^{-1} A(w_{i+1} - w_i), v_i \rangle \\
& \leq \sum_{i=1}^t e^{\lambda_u \sum_{k=i+1}^t \underline{\beta}_k} \mathbb{E} \{ \|C^{-1}\| \|A\| \|w_{i+1} - w_i\| \|v_i\| \} \\
& \leq \|C^{-1}\| \|A\| (K_{f_1} + K_{g_1} + K_c R_u c_\alpha) R_v \sum_{i=1}^t e^{\lambda_u \sum_{k=i+1}^t \underline{\beta}_k} \alpha_i \\
& \leq \frac{(1+\gamma)\rho_{\max}}{\lambda_{c,1}} (K_{f_1} + K_{g_1} + K_c R_u c_\alpha) R_v \sum_{i=1}^t e^{\lambda_u \sum_{k=i+1}^t \underline{\beta}_k} \beta_i \frac{\alpha_i}{\beta_i} \\
& \leq \frac{c_\alpha(1+\gamma)\rho_{\max}}{c_\beta \lambda_c} (K_{f_1} + K_{g_1} + K_c R_u c_\alpha) R_v \max_{i \in [1, t]} \left\{ e^{(\lambda_u/2) \sum_{k=i+1}^t \underline{\beta}_k} \frac{1}{i^{\sigma-\nu}} \right\} \sum_{i=1}^t e^{(\lambda_u/2) \sum_{k=i+1}^t \underline{\beta}_k} \beta_i
\end{aligned}$$

Following the arguments used in the proof of Lemma C.12, we obtain the following upper bound

$$\sum_{i=1}^t e^{(\lambda_u/2) \sum_{k=i+1}^t \underline{\beta}_k} \beta_i \leq \frac{2(1+c_\beta \rho_{\max}) e^{-\frac{\lambda_u c_\beta}{2(1+c_\beta \rho_{\max})}}}{-\lambda_u}.$$

as well as

$$\max_{i \in [1, t]} \left\{ e^{(\lambda_u/2) \sum_{k=i+1}^t \underline{\beta}_k} \frac{1}{i^{\sigma-\nu}} \right\} \leq e^{\frac{\lambda_u c_\beta}{2(1-\nu)(1+c_\beta \rho_{\max})} [(1+t)^{1-\nu} - 1]} D_{\beta,2} + \frac{1}{t^{\sigma-\nu}}$$

where $D_{\beta,2} = \max_{i \in [1, i_{d_{\beta,2}}]} \left\{ e^{-(\lambda_u/2) \sum_{k=1}^i \underline{\beta}_k} \frac{1}{i^{\sigma-\nu}} \right\}$ for some $i_{d_{\beta,2}} \in \mathbb{N}$. This gives us the desired result. \square

Lemma C.17. Suppose $\mathbb{E} \|v_i\|^2 = O\left(\frac{\log i}{i^\nu}\right) + O\left(\frac{1}{i^{\sigma-\nu}}\right)$. If $\sigma > \frac{3}{2}\nu$, we have

$$\sum_{i=1}^t e^{\lambda_u \sum_{k=i+1}^t \underline{\beta}_k} \mathbb{E} \langle C^{-1} A(w_{i+1} - w_i), v_i \rangle = O\left(\frac{1}{t^\nu}\right),$$

and if $\nu < \sigma \leq \frac{3}{2}\nu$, we have

$$\sum_{i=1}^t e^{\lambda_u \sum_{k=i+1}^t \underline{\beta}_k} \mathbb{E} \langle C^{-1} A(w_{i+1} - w_i), z_i \rangle = O\left(\frac{1}{t^{2(\sigma-\nu)-\epsilon}}\right),$$

where ϵ is any constant in $(0, \sigma - \nu]$.

Proof. If $\sigma \geq 2\nu$, Lemma C.16 implies that

$$\sum_{i=1}^t e^{\lambda_u \sum_{k=i+1}^t \underline{\beta}_k} \mathbb{E} \langle C^{-1} A(w_{i+1} - w_i), v_i \rangle = O\left(\frac{1}{t^\nu}\right)$$

If $\sigma < 2\nu$, it follows that $\mathbb{E} \|v_t\|^2 = O\left(\frac{1}{t^{\sigma-\nu}}\right)$. Hence there exists a constant $0 < C < \infty$ and $T > 0$ such that

$$\mathbb{E} \|v_t\|^2 \leq R_z^2 \quad \text{for all } 0 \leq t \leq T \quad (104)$$

$$\mathbb{E} \|v_t\|^2 \leq \frac{C}{t^{(\sigma-\nu)}} \quad \text{for all } t > T \quad (105)$$

Now consider

$$\begin{aligned} & \sum_{i=1}^t e^{\lambda_u \sum_{k=i+1}^t \underline{\beta}_k} \mathbb{E} \langle C^{-1} A(w_{i+1} - w_i), v_i \rangle \\ & \leq \|C^{-1}\| \|A\| (K_{f_1} + K_{g_1} + K_c R_u c_\alpha) \sum_{i=1}^t e^{\lambda_u \sum_{k=i+1}^t \underline{\beta}_k} \alpha_i \sqrt{\mathbb{E} \|v_i\|^2} \\ & \leq \frac{(1+\gamma)\rho_{\max}}{\lambda_{c,1}} (K_{f_1} + K_{g_1} + K_c R_u c_\alpha) \left(\sum_{i=1}^T e^{\lambda_u \sum_{k=i+1}^t \underline{\beta}_k} \alpha_i \sqrt{\mathbb{E} \|v_i\|^2} + \sum_{i=T+1}^t e^{\lambda_u \sum_{k=i+1}^t \underline{\beta}_k} \alpha_i \sqrt{\mathbb{E} \|v_i\|^2} \right) \\ & \leq \frac{c_\alpha(1+\gamma)\rho_{\max}}{c_\beta \lambda_{c,1}} (K_{f_1} + K_{g_1} + K_c R_u c_\alpha) \left(R_v \sum_{i=1}^T e^{\lambda_u \sum_{k=i+1}^t \underline{\beta}_k} \beta_i \frac{1}{i^{(\sigma-\nu)}} + C \sum_{i=T+1}^t e^{\lambda_u \sum_{k=i+1}^t \underline{\beta}_k} \beta_i \frac{1}{i^{1.5(\sigma-\nu)}} \right) \end{aligned}$$

where the last inequality follows from (104) and (105). Following arguments used in proof for Lemma C.11, we have

$$\sum_{i=1}^T e^{\lambda_u \sum_{k=i+1}^t \beta_k} \beta_i \frac{1}{i^{(\sigma-\nu)}} \leq \sum_{i=1}^T e^{\lambda_u \sum_{k=i+1}^t \beta_k} \beta_i \leq \frac{e^{-\frac{\lambda_u c_\beta}{1+c_\beta \rho_{\max}} (1+c_\beta \rho_{\max})}}{-\lambda_u} e^{\frac{\lambda_u c_\beta}{(1-\nu)(1+c_\beta \rho_{\max})}} [(1+t)^{1-\nu} - (1+T)^{1-\nu}]$$

and

$$\sum_{i=T+1}^t e^{\lambda_u \sum_{k=i+1}^t \beta_k} \beta_i \frac{1}{i^{1.5(\sigma-\nu)}} \leq \frac{2e^{-\frac{\lambda_u c_\beta}{2(1+c_\beta \rho_{\max})}} (1+c_\beta \rho_{\max})}{-\lambda_u} \left\{ e^{\frac{\lambda_u c_\beta}{2(1-\nu)(1+c_\beta \rho_{\max})}} [(1+t)^{1-\nu} - 1] D_{\beta,3} + \frac{1}{t^{1.5(\sigma-\nu)}} \right\}$$

where $D_{\beta,3} = \max_{i \in [1, i_{d_{\beta,3}}]} \left\{ e^{-(\lambda_u/2) \sum_{k=1}^i \beta_k} \frac{1}{i^{1.5(\sigma-\nu)}} \right\}$ for some $i_{d_{\beta,3}} \in \mathbb{N}$. It follows that

$$\sum_{i=1}^t e^{\lambda_u \sum_{k=i+1}^t \beta_k} \mathbb{E} \langle C^{-1} A(w_{i+1} - w_i), v_i \rangle = O\left(\frac{1}{t^{1.5(\sigma-\nu)}}\right).$$

If $\frac{3}{2}\nu < \sigma \leq 2\nu$, we have $\mathbb{E} \|v_t\|^2 = O\left(\frac{1}{t^{1.5(\sigma-\nu)}}\right)$. Following the same steps above, we have

$$\sum_{i=1}^t e^{\lambda_u \sum_{k=i+1}^t \beta_k} \mathbb{E} \langle C^{-1} A(w_{i+1} - w_i), v_i \rangle = O\left(\frac{1}{t^{1.75(\sigma-\nu)}}\right)$$

and $\mathbb{E} \|v_t\|^2 = O\left(\frac{1}{t^{1.75(\sigma-\nu)}}\right)$. Repeating analogous steps for a total of $N = \lceil -\log_2 \left(2 - \frac{\nu}{\sigma-\nu}\right) \rceil$ times, we have

$$\sum_{i=1}^t e^{\lambda_u \sum_{k=i+1}^t \beta_k} \mathbb{E} \langle C^{-1} A(w_{i+1} - w_i), v_i \rangle = O\left(\frac{1}{t^{(2-2^{-N})(\sigma-\nu)}}\right) = O\left(\frac{1}{t^\nu}\right).$$

If $\nu < \sigma \leq \frac{3}{2}\nu$, then we repeat previous steps for a total number $N = \lceil \log_2 \left(\frac{\sigma-\nu}{\epsilon}\right) \rceil$ of times, we have

$$\sum_{i=1}^t e^{\lambda_u \sum_{k=i+1}^t \beta_k} \mathbb{E} \langle C^{-1} A(w_{i+1} - w_i), v_i \rangle = O\left(\frac{1}{t^{2(\sigma-\nu)-\epsilon}}\right).$$

□

Lemma C.18. For $0 < \sigma < 1$, suppose $c_\alpha > 0$, $\alpha_t = \frac{c_\alpha}{(1+t)^\sigma}$, and $0 \leq x \leq 1, 0 \leq y \leq 1$. If $\mathbb{E} \|v_t\|^2 = O\left(\frac{\log t}{t^\nu} + \frac{1}{t^\nu}\right)^x$ and $\mathbb{E} \|w_t - w_*\|^2 = O\left(\frac{\log t}{t^\nu} + \frac{1}{t^\nu}\right)^y$, we have

$$\sum_{i=1}^t e^{\lambda_w \sum_{k=i+1}^t \alpha_k} \mathbb{E} \left\{ \alpha'_i \langle B_i^s v_i, w_i - w_* \rangle \right\} = O\left(\frac{\log t}{t^\nu} + \frac{1}{t^\nu}\right)^{0.5(x+y)}.$$

If $\mathbb{E} \|v_t\|^2 = O\left(\frac{\log t}{t^\nu} + \frac{1}{t^{2(\sigma-\nu)-\epsilon}}\right)^x$ and $\mathbb{E} \|w_t - w_*\|^2 = O\left(\frac{\log t}{t^\nu} + \frac{1}{t^{2(\sigma-\nu)-\epsilon}}\right)^y$, we have

$$\sum_{i=1}^t e^{\lambda_w \sum_{k=i+1}^t \alpha_k} \mathbb{E} \{\alpha'_i \langle B_i^s v_i, w_i - w_* \rangle\} = O\left(\frac{\log t}{t^\nu} + \frac{1}{t^{2(\sigma-\nu)-\epsilon}}\right)^{0.5(x+y)}.$$

Proof. Consider the first case. Without loss of generality, we assume that there exist constants $0 < C_1, C_2 < \infty, T > 0$ such that

$$\begin{aligned} \mathbb{E} \|v_t\|^2 &\leq R_v^2, \quad \mathbb{E} \|w_t - w_*\|^2 \leq R_w^2 \quad \text{for all } 0 \leq t \leq T \\ \mathbb{E} \|v_t\|^2 &\leq C_1^2 \left(\frac{\log t + 1}{t^\nu}\right)^x, \quad \mathbb{E} \|w_t - w_*\|^2 \leq C_2^2 \left(\frac{\log t + 1}{t^\nu}\right)^y \quad \text{for all } t > T \end{aligned}$$

Since $0 \leq \alpha'_i \leq \alpha_i$ holds almost surely,

$$\begin{aligned} \sum_{i=1}^t e^{\lambda_w \sum_{k=i+1}^t \alpha_k} \mathbb{E} \{\alpha'_i \langle B_i^s v_i, w_i - w_* \rangle\} &\leq \sum_{i=1}^t e^{\lambda_w \sum_{k=i+1}^t \alpha_k} \alpha_i \sqrt{\mathbb{E} \|B_i^s v_i\|^2} \sqrt{\mathbb{E} \|w_i - w_*\|^2} \\ &\leq \gamma \rho_{\max} \sum_{i=1}^t e^{\lambda_w \sum_{k=i+1}^t \alpha_k} \alpha_i \sqrt{\mathbb{E} \|v_i\|^2} \sqrt{\mathbb{E} \|w_i - w_*\|^2} \\ &\leq \gamma \rho_{\max} \left(\sum_{i=1}^T e^{\lambda_w \sum_{k=i+1}^t \alpha_k} \alpha_i \sqrt{\mathbb{E} \|v_i\|^2} \sqrt{\mathbb{E} \|w_i - w_*\|^2} + \sum_{i=T+1}^t e^{\lambda_w \sum_{k=i+1}^t \alpha_k} \alpha_i \sqrt{\mathbb{E} \|v_i\|^2} \sqrt{\mathbb{E} \|w_i - w_*\|^2} \right) \\ &\leq \gamma \rho_{\max} \left\{ 2R_v R_w \sum_{i=1}^T e^{\lambda_w \sum_{k=i+1}^t \alpha_k} \alpha_i + C_1 C_2 \sum_{i=T+1}^t e^{\lambda_w \sum_{k=i+1}^t \alpha_k} \alpha_i \left(\frac{\log i + 1}{i^\nu}\right)^{0.5(x+y)} \right\}. \end{aligned}$$

Following arguments used in the proof for part (a) of Lemma C.11, we obtain

$$\sum_{i=1}^T e^{\lambda_w \sum_{k=i+1}^t \alpha_k} \alpha_i \leq \frac{e^{-\frac{\lambda_w c_\alpha}{1+c_\alpha \rho_{\max}}} (1 + c_\alpha \rho_{\max})}{-\lambda_w} e^{\frac{\lambda_w c_\alpha}{(1-\sigma)(1+c_\alpha \rho_{\max})}} [(1+t)^{1-\sigma} - (1+T)^{1-\sigma}]$$

and

$$\begin{aligned} &\sum_{i=T+1}^t e^{\lambda_w \sum_{k=i+1}^t \alpha_k} \alpha_i \left(\frac{\log i + 1}{i^\nu}\right)^{0.5(x+y)} \\ &\leq \frac{2e^{-\frac{\lambda_w c_\alpha}{2(1+c_\alpha \rho_{\max})}} (1 + c_\alpha \rho_{\max})}{-\lambda_w} \left\{ e^{\frac{\lambda_w c_\alpha}{2(1-\sigma)(1+c_\alpha \rho_{\max})}} [(1+t)^{1-\sigma} - 1] D_{\alpha,2} + \left(\frac{\log t + 1}{t^\nu}\right)^{0.5(x+y)} \right\}, \end{aligned}$$

for some constant $D_{\alpha,2}$. The proof for the second case can be done with analogous arguments. \square

Lemma C.19. Suppose $\mathbb{E} \|w_{t+1} - w_*\|^2 = O\left(\frac{\log t}{t^\nu} + \frac{1}{t^\nu}\right)^{0.5}$ and $\epsilon' \in (0, 0.5]$. For $0 < \frac{3}{2}\nu < \sigma < 1$,

if $\mathbb{E} \|v_t\|^2 = O\left(\frac{\log t}{t^\nu}\right) + O\left(\frac{1}{t^\nu}\right)$, we have

$$\sum_{i=1}^t e^{\lambda_w \sum_{k=i+1}^t \alpha_k} \mathbb{E} \left\{ \alpha'_i \langle B_i^s v_i, w_i - w_* \rangle \right\} = O\left(\frac{\log t}{t^\nu} + \frac{1}{t^\nu}\right)^{1-\epsilon'}.$$

For $0 < \nu < \sigma \leq \frac{3}{2}\nu < 1$, if $\mathbb{E} \|v_t\|^2 = O\left(\frac{\log t}{t^\nu}\right) + O\left(\frac{1}{t^{2(\sigma-\nu)-\epsilon}}\right)$, we have

$$\sum_{i=1}^t e^{\lambda_w \sum_{k=i+1}^t \alpha_k} \mathbb{E} \left\{ \alpha'_i \langle B_i^s v_i, w_i - w_* \rangle \right\} = O\left(\frac{\log t}{t^\nu} + \frac{1}{t^{2(\sigma-\nu)-\epsilon}}\right)^{1-\epsilon'}.$$

Proof. Consider the first case. First, $\mathbb{E} \|w_t - w_*\|^2 \leq 4R_w^2 = O(1)$, applying Lemma C.18, we immediately have

$$\sum_{i=1}^t e^{\lambda_w \sum_{k=i+1}^t \alpha_k} \mathbb{E} \left\{ \alpha'_i \langle B_i^s v_i, w_i - w_* \rangle \right\} = O\left(\frac{\log t}{t^\nu} + \frac{1}{t^\nu}\right)^{0.5}.$$

From the assumption $\mathbb{E} \|w_{t+1} - w_*\|^2 = O\left(\frac{\log t}{t^\nu} + \frac{1}{t^\nu}\right)^{0.5}$ with Lemma C.18, we obtain

$$\sum_{i=1}^t e^{\lambda_w \sum_{k=i+1}^t \alpha_k} \mathbb{E} \left\{ \alpha'_i \langle B_i^s v_i, w_i - w_* \rangle \right\} = O\left(\frac{\log t}{t^\nu} + \frac{1}{t^\nu}\right)^{0.75}.$$

Repeating the above steps for a total number of $N = \lceil \log_2 \left(\frac{1}{\epsilon}\right) \rceil$ times, we have

$$\sum_{i=1}^t e^{\lambda_w \sum_{k=i+1}^t \alpha_k} \mathbb{E} \left\{ \alpha'_i \langle B_i^s v_i, w_i - w_* \rangle \right\} = O\left(\frac{\log t}{t^\nu} + \frac{1}{t^\nu}\right)^{1-\frac{1}{2^N}} = O\left(\frac{\log t}{t^\nu} + \frac{1}{t^\nu}\right)^{1-\epsilon'}.$$

The proof for the second case is analogous. □

C.2 Finite-time analysis for projected implicit TDC

We finally provide proofs for the finite-time error bounds of the implicit TDC algorithm under both a sequence of decreasing step sizes as well as the constant step size.

Proof of Theorem 4.18. From the recursion of w_n , for $n \geq 1$, we know

$$\begin{aligned}
\|w_{n+1} - w_*\|^2 &= \|\Pi_{R_w} [w_n + \alpha'_n \{f_1(w_n, O_n) + g_1(v_n, O_n)\} + (\alpha_n B_n - \alpha'_n B_n^s)u_n] - w_*\|^2 \\
&= \|\Pi_{R_w} [w_n + \alpha'_n \{f_1(w_n, O_n) + g_1(v_n, O_n)\} + (\alpha_n B_n - \alpha'_n B_n^s)u_n] - \Pi_{R_w} w_*\|^2 \\
&\leq \|w_n - w_* + \alpha'_n \{f_1(w_n, O_n) + g_1(v_n, O_n)\} + (\alpha_n B_n - \alpha'_n B_n^s)u_n\|^2 \\
&= \|w_n - w_* + \alpha'_n \{f_1(w_n, O_n) + g_1(v_n, O_n)\}\|^2 + \|(\alpha_n B_n - \alpha'_n B_n^s)u_n\|^2 \\
&\quad + 2\langle w_n - w_* + \alpha'_n \{f_1(w_n, O_n) + g_1(v_n, O_n)\}, (\alpha_n B_n - \alpha'_n B_n^s)u_n \rangle \\
&\leq \|w_n - w_*\|^2 + 2\alpha'_n \langle f_1(w_n, O_n), w_n - w_* \rangle + 2\alpha'_n \langle g_1(v_n, O_n), w_n - w_* \rangle \\
&\quad + \alpha_n'^2 \|f_1(w_n, O_n) + g_1(v_n, O_n)\|^2 + \|(\alpha_n B_n - \alpha'_n B_n^s)u_n\|^2 \\
&\quad + 2\|w_n - w_* + \alpha'_n \{f_1(w_n, O_n) + g_1(v_n, O_n)\}\| \|(\alpha_n B_n - \alpha'_n B_n^s)u_n\| \\
&\leq \|w_n - w_*\|^2 + 2\alpha'_n \langle \bar{f}_1(w_n), w_n - w_* \rangle + 2\alpha'_n \langle f_1(w_n, O_n) - \bar{f}_1(w_n), w_n - w_* \rangle \\
&\quad + 2\alpha'_n \langle g_1(v_n, O_n), w_n - w_* \rangle + 2\alpha_n'^2 K_{f_1}^2 + 2\alpha_n'^2 K_{g_1}^2 + \|(\alpha_n B_n - \alpha'_n B_n^s)u_n\|^2 \\
&\quad + 2\{2R_w + \alpha_n(K_{f_1} + K_{g_1})\} \|(\alpha_n B_n - \alpha'_n B_n^s)u_n\|.
\end{aligned}$$

Applying Lemma C.1 we have,

$$\begin{aligned}
\|w_{n+1} - w_*\|^2 &\leq \|w_n - w_*\|^2 + 2\alpha'_n \left\langle \left(A^\top C^{-1} A \right) (w_n - w_*), w_n - w_* \right\rangle \\
&\quad + 2\alpha'_n \langle f_1(w_n, O_n) - \bar{f}_1(w_n), w_n - w_* \rangle + 2\alpha'_n \langle g_1(v_n, O_n), w_n - w_* \rangle \\
&\quad + 2\alpha_n'^2 K_{f_1}^2 + 2\alpha_n'^2 K_{g_1}^2 + K_c^2 \alpha_n^4 + 2\{2R_w + \alpha_n(K_{f_1} + K_{g_1})\} K_c \alpha_n^2 \\
&\leq (1 - \alpha'_n |\lambda_w|) \|w_n - w_*\|^2 + 2\alpha'_n \zeta_{f_1}(w_n, O_n) + 2\alpha'_n \langle B_n^s v_n, w_n - w_* \rangle \\
&\quad + 2\alpha_n'^2 K_{f_1}^2 + 2\alpha_n'^2 K_{g_1}^2 + K_c^2 \alpha_n^4 + 2\{2R_w + c_\alpha(K_{f_1} + K_{g_1})\} K_c \alpha_n^2,
\end{aligned}$$

where in the second inequality we used facts $2\lambda_{\max}(A^\top C^{-1} A) \leq \lambda_w < 0$ and $\zeta_{f_1}(w_n, O_n) = \langle f_1(w_n, O_n) - \bar{f}_1(w_n), w_n - w_* \rangle$. Now note that $\alpha'_n \geq \alpha_n/(1 + c_\alpha \rho_{\max}) =: \underline{\alpha}_n$ holds almost surely. Telescoping the above inequality and taking the expectation on both sides yields that

$$\begin{aligned}
\mathbb{E} \|w_{n+1} - w_*\|^2 &\leq \left\{ \prod_{i=1}^n (1 - \underline{\alpha}_i |\lambda_w|) \right\} \|w_1 - w_*\|^2 + 2 \sum_{i=1}^n \left\{ \prod_{k=i+1}^n (1 - \underline{\alpha}_k |\lambda_w|) \right\} \mathbb{E} \{ \alpha'_i \zeta_{f_1}(\theta_i, O_i) \} \\
&\quad + 2 \sum_{i=1}^n \left\{ \prod_{k=i+1}^n (1 - \underline{\alpha}_k |\lambda_w|) \right\} \mathbb{E} \{ \alpha'_i \langle B_i^s v_i, w_i - w_* \rangle \} \\
&\quad + 2 [K_{f_1}^2 + K_{g_1}^2 + K_c^2 c_\alpha^2 + 2\{2R_w + c_\alpha(K_{f_1} + K_{g_1})\} K_c] \sum_{i=1}^n \left\{ \prod_{k=i+1}^n (1 - \underline{\alpha}_k |\lambda_w|) \right\} \alpha_i^2
\end{aligned} \tag{106}$$

Applying $1 - \underline{\alpha}_i |\lambda_w| \leq e^{-\underline{\alpha}_i |\lambda_w|}$, we obtain

$$\begin{aligned} \mathbb{E} \|w_{n+1} - w_*\|^2 &\leq e^{-|\lambda_w| \sum_{i=1}^n \underline{\alpha}_i} \|w_1 - w_*\|^2 + 2 \sum_{i=1}^n e^{-|\lambda_w| \sum_{k=i+1}^n \underline{\alpha}_k} \mathbb{E} \{ \alpha'_i \zeta_{f_1}(w_i, O_i) \} \\ &\quad + 2 \sum_{i=1}^n e^{-|\lambda_w| \sum_{k=i+1}^t \underline{\alpha}_k} \mathbb{E} \{ \alpha'_i \langle B_i^s z_i, w_i - w_* \rangle \} \\ &\quad + 2 [K_{f_1}^2 + K_{g_1}^2 + K_c^2 c_\alpha^2 + 2 \{2R_w + c_\alpha (K_{f_1} + K_{g_1})\} K_c] \sum_{i=1}^n e^{-|\lambda_w| \sum_{k=i+1}^n \underline{\alpha}_k} \alpha_i^2 \end{aligned}$$

Combining Lemma C.8, Lemma C.12 and Lemma C.19, we obtain

$$\begin{aligned} &\mathbb{E} \|w_{n+1} - w_*\|^2 \\ &\leq e^{\frac{-|\lambda_w| c_\alpha}{(1+c_\alpha \rho_{\max})(1-\sigma)}} [(1+n)^{1-\sigma} - 1] \|w_1 - w_*\|^2 + O\left(\frac{\log n}{n^\sigma}\right) + O\left(\frac{\log n}{n^\nu} + h(\sigma, \nu)\right)^{1-\epsilon'} + O\left(\frac{1}{n^\sigma}\right) \\ &= O\left(e^{\frac{-|\lambda_w| c_\alpha n^{1-\sigma}}{(1+c_\alpha \rho_{\max})(1-\sigma)}}\right) + O\left(\frac{\log n}{n^\sigma}\right) + O\left(\frac{\log n}{n^\nu} + h(\sigma, \nu)\right)^{1-\epsilon'} \end{aligned}$$

where $h(\sigma, \nu) = \begin{cases} \frac{1}{n^\nu}, & \sigma > 1.5\nu \\ \frac{1}{n^{2(\sigma-\nu)-\epsilon}}, & \nu < \sigma \leq 1.5\nu \end{cases}$ for an arbitrarily small constant $\epsilon \in (0, \sigma - \nu]$. Justification of the condition of Lemma C.19 is provided in the analysis for the tracking error vector, which is given below.

We next bound the recursion of the tracking error vector v_n as follows. For any $n \geq 1$, note that

$$\begin{aligned}
\|v_{n+1}\|^2 &= \|\Pi_{R_u} [v_n + \beta'_n \{f_2(w_n, O_n) + g_2(v_n, O_n)\} - C^{-1}(b + Aw_n)] + C^{-1}(b + Aw_{n+1})\|^2 \\
&= \|\Pi_{R_u} [v_n + \beta'_n \{f_2(w_n, O_n) + g_2(v_n, O_n)\} - C^{-1}(b + Aw_n)] + \Pi_{R_u} \{C^{-1}(b + Aw_{n+1})\}\|^2 \\
&\leq \|v_n + \beta'_n \{f_2(w_n, O_n) + g_2(v_n, O_n)\} + C^{-1}A(w_{n+1} - w_n)\|^2 \\
&= \|v_n\|^2 + 2\beta'_n \langle f_2(w_n, O_n), v_n \rangle + 2\beta'_n \langle g_2(v_n, O_n), v_n \rangle + 2\langle C^{-1}A(w_{n+1} - w_n), v_n \rangle \\
&\quad + \|\beta'_n f_2(w_n, O_n) + \beta'_n g_2(v_n, O_n) + C^{-1}A(w_{n+1} - w_n)\|^2 \\
&\leq \|v_n\|^2 + 2\beta'_n \langle \bar{g}_2(v_n), v_n \rangle + 2\beta'_n \langle f_2(w_n, O_n), v_n \rangle + 2\beta'_n \langle g_2(v_n, O_n) - \bar{g}_2(v_n), v_n \rangle \\
&\quad + 2\langle C^{-1}A(w_{n+1} - w_n), v_n \rangle + 3\beta_n'^2 \|f_2(w_n, O_n)\|^2 + 3\beta_n'^2 \|g_2(v_n, O_n)\|^2 \\
&\quad + 3\|C^{-1}A(w_{n+1} - w_n)\|^2 \\
&\leq \|v_n\|^2 + 2\beta'_n \langle Cv_n, v_n \rangle + 2\beta'_n \langle f_2(w_n, O_n), v_n \rangle + 2\beta'_n \langle g_2(v_n, O_n) - \bar{g}_2(v_n), v_n \rangle \\
&\quad + 2\langle C^{-1}A(w_{n+1} - w_n), v_n \rangle + 3\beta_n'^2 \|f_2(w_n, O_n)\|^2 + 3\beta_n'^2 \|g_2(v_n, O_n)\|^2 \\
&\quad + 3\|C^{-1}\|^2 \|A\|^2 \|\alpha'_n \{f_1(w_n, O_n) + g_1(v_n, O_n)\} + (\alpha_n B_n - \alpha'_n B_n^s)u_n\|^2 \\
&\leq \|v_n\|^2 + 2\beta'_n \langle Cv_n, v_n \rangle + 2\beta'_n \langle f_2(w_n, O_n), v_n \rangle + 2\beta'_n \langle g_2(v_n, O_n) - \bar{g}_2(v_n), v_n \rangle \\
&\quad + 2\langle C^{-1}A(w_{n+1} - w_n), v_n \rangle + 3\beta_n'^2 K_{f_2}^2 + 3\beta_n'^2 K_{g_2}^2 \\
&\quad + 12\|C^{-1}\|^2 \|A\|^2 \{\alpha_n'^2 (K_{f_1}^2 + K_{g_1}^2) + K_c^2 \alpha_n^4 R_u^2\} \\
&\leq (1 - \beta'_n |\lambda_u|) \|v_n\|^2 + 2\beta'_n \zeta_{f_2}(w_n, v_n, O_n) + 2\beta'_n \zeta_{g_2}(v_n, O_n) + 2\langle C^{-1}A(w_{n+1} - w_n), v_n \rangle \\
&\quad + 3\beta_n'^2 K_{f_2}^2 + 3\beta_n'^2 K_{g_2}^2 + 12\alpha_n^2 \|C^{-1}\|^2 \|A\|^2 (K_{f_1}^2 + K_{g_1}^2 + K_c^2 c_\alpha^2 R_u^2),
\end{aligned}$$

where $\lambda_{\max}(2C) \leq \lambda_u < 0$, $\zeta_{f_2}(w_n, v_n, O_n) = \langle f_2(w_n, O_n), v_n \rangle$, $\zeta_{g_2}(v_n, O_n) = \langle g_2(v_n, O_n) - \bar{g}_2(v_n), v_n \rangle$. Please refer to Lemma C.2, C.3, C.4 and C.1 for definitions of $K_{f_1}, K_{g_1}, K_{f_2}, K_{g_2}$ and K_c . Furthermore, $\beta'_n \geq \beta_n/(1+\beta_n \rho_{\max}) =: \underline{\beta}_n$ holds almost surely. Let $K_{r_2}^2 = \|C^{-1}\|^2 \|A\|^2 (K_{f_1}^2 + K_{g_1}^2 + K_c^2 c_\alpha^2 R_u^2)$; taking the expectation on both sides, we have

$$\begin{aligned}
\mathbb{E} \|v_{n+1}\|^2 &\leq (1 - \underline{\beta}_n |\lambda_u|) \mathbb{E} \|v_n\|^2 + 2\mathbb{E} \{\beta'_n \zeta_{f_2}(w_n, v_n, O_n)\} + 2\mathbb{E} \{\beta'_n \zeta_{g_2}(v_n, O_n)\} \\
&\quad + 2\mathbb{E} \langle C^{-1}A(w_{n+1} - w_n), v_n \rangle + 3\beta_n'^2 K_{f_2}^2 + 3\beta_n'^2 K_{g_2}^2 + 12\alpha_n^2 K_{r_2}^2
\end{aligned} \tag{107}$$

With $\underline{\beta}_n |\lambda_u| \in (0, 1)$ for all $n \in \mathbb{N}$, telescoping the above inequality yields that

$$\begin{aligned}
\mathbb{E} \|v_{n+1}\|^2 &\leq \left\{ \prod_{i=1}^n (1 - \underline{\beta}_i |\lambda_u|) \right\} \|v_1\|^2 + 2 \sum_{i=1}^n \left\{ \prod_{k=i+1}^n (1 - \underline{\beta}_k |\lambda_u|) \right\} \mathbb{E} \{\beta'_i \zeta_{f_2}(w_i, v_i, O_i)\} \\
&\quad + 2 \sum_{i=1}^n \left\{ \prod_{k=i+1}^n (1 - \underline{\beta}_k |\lambda_u|) \right\} \mathbb{E} \{\beta'_i \zeta_{g_2}(v_i, O_i)\} + 2 \sum_{i=1}^n \left\{ \prod_{k=i+1}^n (1 - \underline{\beta}_k |\lambda_u|) \right\} \mathbb{E} \langle C^{-1}A(w_{i+1} - w_i), v_i \rangle \\
&\quad + 3(K_{f_2}^2 + K_{g_2}^2) \sum_{i=1}^n \left\{ \prod_{k=i+1}^n (1 - \underline{\beta}_k |\lambda_u|) \right\} \beta_i^2 + 12K_{r_2}^2 \sum_{i=1}^n \left\{ \prod_{k=i+1}^n (1 - \underline{\beta}_k |\lambda_u|) \right\} \alpha_i^2
\end{aligned}$$

Since $1 - \underline{\beta}_i |\lambda_u| \leq e^{-\underline{\beta}_i |\lambda_u|}$ and by the fact $(1+i)^{-\nu} \geq (1+i)^{-\sigma}$ for all $i \geq 0$, we have

$$\begin{aligned} \mathbb{E} \|v_{n+1}\|^2 &\leq e^{-|\lambda_u| \sum_{i=1}^n \underline{\beta}_i} \|v_1\|_2^2 + 2 \sum_{i=1}^n e^{-|\lambda_u| \sum_{k=i+1}^n \underline{\beta}_k} \mathbb{E} \{ \beta'_i \zeta_{f_2}(w_i, v_i, O_i) \} \\ &\quad + 2 \sum_{i=1}^n e^{-|\lambda_u| \sum_{k=i+1}^n \underline{\beta}_k} \mathbb{E} \{ \beta'_i \zeta_{g_2}(v_i, O_i) \} + 2 \sum_{i=1}^n e^{-|\lambda_u| \sum_{k=i+1}^n \underline{\beta}_k} \mathbb{E} \langle C^{-1} A(w_{i+1} - w_i), v_i \rangle \\ &\quad + 3 \max \left\{ 1, \frac{c_\alpha^2}{c_\beta^2} \right\} (K_{f_2}^2 + K_{g_2}^2 + 4K_{r_2}^2) \sum_{i=1}^n e^{-|\lambda_u| \sum_{k=i+1}^n \underline{\beta}_k} \beta_i^2 \end{aligned}$$

Combining Lemma C.14, Lemma C.15, Lemma C.16 and applying Lemma C.12, we obtain

$$\begin{aligned} \mathbb{E} \|v_{n+1}\|^2 &= O \left(e^{\frac{-|\lambda_u| c_\beta n^{1-\nu}}{(1-\nu)(1+c_\beta \rho_{\max})}} \right) + O \left(\frac{\log n}{n^\nu} \right) + O \left(\frac{\log n}{n^\nu} \right) + O \left(\frac{1}{n^{\sigma-\nu}} \right) + O \left(\frac{1}{n^\nu} \right) \\ &= O \left(\frac{\log n}{n^\nu} \right) + O \left(\frac{1}{n^{\sigma-\nu}} \right) \end{aligned} \quad (108)$$

Applying Lemma C.17, we can further refine to yield

$$\mathbb{E} \|v_n\|^2 = O \left(\frac{\log n}{n^\nu} \right) + O(h(\sigma, \nu))$$

with

$$h(\sigma, \nu) = \begin{cases} \frac{1}{n^\nu}, & \sigma > 1.5\nu \\ \frac{1}{n^{2(\sigma-\nu)-\epsilon}}, & \nu < \sigma \leq 1.5\nu \end{cases}$$

where $\epsilon \in (0, \sigma - \nu]$ can be an arbitrarily small constant. \square

Lastly, we provide a proof for the finite-time error bounds of implicit TDC with constant step sizes.

Proof of Theorem 4.21. Suppose $\alpha_i = c_\alpha$, $\beta_i = c_\beta$ for all $i \in \mathbb{N}$ and recall $\underline{\alpha} = \frac{c_\alpha}{1+c_\alpha \rho_{\max}}$ and $\underline{\beta} = \frac{c_\beta}{1+c_\beta \rho_{\max}}$. Note that $0 \leq \alpha'_i \leq c_\alpha$ and $0 \leq \beta'_i \leq c_\beta$ hold almost surely. Therefore, from the expression (107), we have

$$\begin{aligned} \mathbb{E} \|v_{n+1}\|^2 &\leq (1 - \underline{\beta} |\lambda_u|)^n \|v_1\|^2 + 2c_\beta \sum_{i=1}^n (1 - \underline{\beta} |\lambda_u|)^{n-i} \mathbb{E} \{ \zeta_{f_2}(w_i, v_i, O_i) \} \\ &\quad + 2c_\beta \sum_{i=1}^n (1 - \underline{\beta} |\lambda_u|)^{n-i} \mathbb{E} \{ \zeta_{g_2}(v_i, O_i) \} \\ &\quad + 2 \sum_{i=1}^n (1 - \underline{\beta} |\lambda_u|)^{n-i} \mathbb{E} \langle C^{-1} A(w_{i+1} - w_i), v_i \rangle \\ &\quad + 3 (K_{f_2}^2 + K_{g_2}^2) c_\beta^2 \sum_{i=1}^n (1 - \underline{\beta} |\lambda_u|)^{n-i} + 12K_{r_2}^2 c_\alpha^2 \sum_{i=1}^n (1 - \underline{\beta} |\lambda_u|)^{n-i} \end{aligned} \quad (109)$$

With a constant step size, Lemma C.9 and Lemma C.10 give us

$$\mathbb{E} \{ \zeta_{f_2} (w_i, v_i, O_i) \} \leq (2R_w K_{f_2} + K_{r_3} \tau_\beta) c_\beta \quad (110)$$

$$\mathbb{E} \{ \zeta_{g_2} (v_i, O_i) \} \leq (4R_w K_{g_2} + L_{g_2} K_{r_2} \tau_\beta) c_\beta \quad (111)$$

and from arguments used in Lemma C.16, we have

$$\begin{aligned} & \sum_{i=1}^n (1 - \underline{\beta} |\lambda_u|)^{n-i} \mathbb{E} \langle C^{-1} A (w_{i+1} - w_i), v_i \rangle \\ & \leq \sum_{i=1}^n (1 - \underline{\beta} |\lambda_u|)^{n-i} \mathbb{E} \{ \|C^{-1}\| \|A\| \|w_{i+1} - w_i\| \|v_i\| \} \\ & \leq \|C^{-1}\| \|A\| (K_{f_1} + K_{g_1} + K_c R_u c_\alpha) R_v c_\alpha \sum_{i=1}^n (1 - \underline{\beta} |\lambda_u|)^{n-i} \\ & \leq \rho_{\max}(\gamma + 1) (K_{f_1} + K_{g_1} + K_c R_u c_\alpha) R_v c_\alpha / (\underline{\beta} |\lambda_u| \lambda_c) \end{aligned} \quad (112)$$

Plugging (110), (111) and (112) into (109), we have

$$\mathbb{E} \|v_{n+1}\|^2 \leq (1 - \underline{\beta} |\lambda_u|)^n \|v_1\|^2 + C_v,$$

where

$$\begin{aligned} C_v &= \frac{8c_\beta^2 \{R_w (K_{f_2} + K_{g_2})\}}{|\lambda_u| \underline{\beta}} + \frac{8c_\beta^2 (K_{r_3} + L_{g_2} K_{r_2}) \tau_\beta}{|\lambda_u| \underline{\beta}} + \frac{3 (K_{f_2}^2 + K_{g_2}^2) c_\beta^2}{|\lambda_u| \underline{\beta}} + \frac{12 K_{r_2}^2 c_\alpha^2}{|\lambda_u| \underline{\beta}} \\ &+ \frac{2 \rho_{\max}(\gamma + 1) (K_{f_1} + K_{g_1} + K_c R_u c_\alpha) R_v c_\alpha}{|\lambda_u| \lambda_c \underline{\beta}} \\ &= O \left(\max \{c_\beta \tau_{c_\beta}, c_\beta^2 \tau_{c_\beta}\} + \max \{c_\alpha, c_\alpha^2\} + \max \{c_\alpha / c_\beta, c_\alpha^2 / c_\beta\} \right). \end{aligned}$$

Therefore, for all $n \geq \tilde{n} := \frac{\log C_v / \|v_1\|^2}{\log(1 - \underline{\beta} |\lambda_u|)}$, we have $\mathbb{E} \|v_n\|^2 \leq 2C_v$. Otherwise, we have $\mathbb{E} \|v_n\|^2 \leq R_v^2$.

For the primary parameter, from the expression (106), we have

$$\begin{aligned} \mathbb{E} \|w_{n+1} - w_*\|^2 &\leq (1 - \underline{\alpha} |\lambda_w|)^n \|w_1 - w_*\|^2 + 2c_\alpha \sum_{i=1}^n (1 - \underline{\alpha} |\lambda_w|)^{n-i} \mathbb{E} \{ \zeta_{f_1} (w_i, O_i) \} \\ &+ 2c_\alpha \sum_{i=1}^n (1 - \underline{\alpha} |\lambda_w|)^{n-i} \mathbb{E} \{ \langle B_i^s v_i, w_i - w_* \rangle \} \\ &+ 2 [K_{f_1}^2 + K_{g_1}^2 + K_c^2 c_\alpha^2 + 2 \{2R_w + c_\alpha (K_{f_1} + K_{g_1})\} K_c] c_\alpha^2 \sum_{i=1}^n (1 - \underline{\alpha} |\lambda_w|)^{n-i} \end{aligned}$$

From Lemma C.8, for a constant step size, we have

$$\mathbb{E} \{ \zeta_{f_1} (w_i, O_i) \} \leq 8R_w K_{f_1} c_\alpha + L_{f_1} \{ (K_{f_1} + K_{g_1}) c_\alpha + K_c R_u c_\alpha^2 \} \tau_\alpha.$$

Furthermore, recall that $\mathbb{E}\|v_n\|^2 \leq 2C_v$, for $n \geq \tilde{n}$ and $\mathbb{E}\|v_n\|^2 \leq R_w^2$, otherwise. With the fact that $\|w_i - w_*\| \leq 2R_w$, we have

$$\begin{aligned}
& \sum_{i=1}^n (1 - \underline{\alpha}|\lambda_w|)^{n-i} \mathbb{E} \{ \langle B_i^s v_i, w_i - w_* \rangle \} \\
& \leq 2\gamma\rho_{\max} R_w \left\{ R_v \sum_{i=1}^{\tilde{n}-1} (1 - \underline{\alpha}|\lambda_w|)^{n-i} + \sqrt{2C_v} \sum_{i=\tilde{n}}^n (1 - \underline{\alpha}|\lambda_w|)^{n-i} \right\} \\
& \leq 2\gamma\rho_{\max} R_w (1 - \underline{\alpha}|\lambda_w|)^n \left\{ R_v \sum_{i=1}^{\tilde{n}} (1 - \underline{\alpha}|\lambda_w|)^{-i} + \sqrt{2C_v} \sum_{i=\tilde{n}}^n (1 - \underline{\alpha}|\lambda_w|)^{-i} \right\} \\
& \leq 2\gamma\rho_{\max} R_w (1 - \underline{\alpha}|\lambda_w|)^n \left\{ \frac{R_v (1 - \underline{\alpha}|\lambda_w|)^{-\tilde{n}}}{\underline{\alpha}|\lambda_w|} + \frac{\sqrt{2C_v} (1 - \underline{\alpha}|\lambda_w|)^{-n}}{\underline{\alpha}|\lambda_w|} \right\} \\
& \leq \frac{2\gamma\rho_{\max} R_w \left(R_v (1 - \underline{\alpha}|\lambda_w|)^{n-\tilde{n}} + \sqrt{2C_v} \right)}{\underline{\alpha}|\lambda_w|}
\end{aligned}$$

Combining everything, we have

$$\mathbb{E} \|w_{n+1} - w_*\|^2 \leq (1 - \underline{\alpha}|\lambda_w|)^n \|w_1 - w_*\|^2 + C_w,$$

where

$$\begin{aligned}
C_w &= \frac{2c_\alpha [8R_w K_{f_1} c_\alpha + L_{f_1} \{ (K_{f_1} + K_{g_1}) c_\alpha + K_c R_u c_\alpha^2 \} \tau_\alpha]}{\underline{\alpha}|\lambda_w|} + \frac{4c_\alpha \gamma \rho_{\max} R_w \left(R_v (1 - \underline{\alpha}|\lambda_w|)^{n-\tilde{n}} + \sqrt{2C_v} \right)}{\underline{\alpha}|\lambda_w|} \\
&+ \frac{2 \left[K_{f_1}^2 + K_{g_1}^2 + K_c^2 c_\alpha^2 + 2 \{ 2R_w + c_\alpha (K_{f_1} + K_{g_1}) \} K_c \right] c_\alpha^2}{\underline{\alpha}|\lambda_w|} \\
&= O(\max\{c_\alpha, c_\alpha^4\}) + O\left(\sqrt{C_v} + c_\alpha \sqrt{C_v}\right) + O(\max\{c_\alpha, c_\alpha^3\} \tau_{c_\alpha}).
\end{aligned}$$

□

D Additional numerical experiments

In this section, we provide numerical experimental results of implicit TD(0.5) along with standard TD(0.5) methods with and without projection on 11-state random walk environments, considered in Subsection 5.1. Similar to what we have observed for TD(0) algorithms, both implicit TD(0.5) and projected implicit TD(0.5) were much more robust to standard TD(0.5) counterparts in terms of the step size choice. In terms of numerical stability, for a moderately large step size, TD(0.5) was more stable than TD(0). However, the quality of the value function approximation was distinctively inferior to that of implicit TD(0.5), which can be observed in Figure 9.

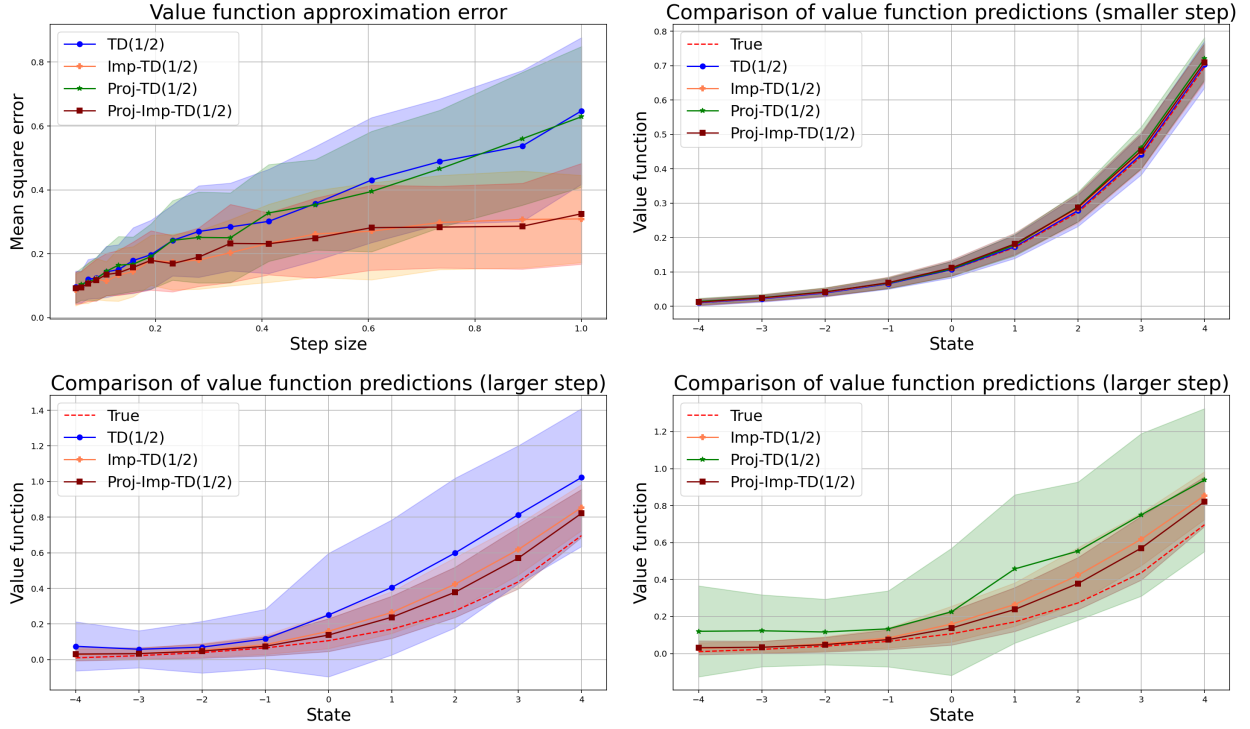


Figure 9: All figures pertain to the random walk environment. **Top left:** Value approximation error versus constant step size over the interval $[0.05, 1]$. Implicit TD(0.5) exhibits a more gradual increase in value approximation error as the step size grows, reflecting its enhanced robustness to large step sizes. **Top right:** Value function approximation with $\alpha_n = 0.05$. Both standard and implicit TD(0.5) algorithms accurately recover the true value function, with tight confidence bands. **Bottom left:** Value function approximation with $\alpha_n = 1.5$. Implicit TD(0.5) achieves closer alignment with the true value function and reduced variance compared to standard TD(0.5). **Bottom right:** Value function approximation with $\alpha_n = 1.5$ using projected TD(0.5). The standard projected TD(0.5) algorithm exhibits larger approximation error and wider confidence band than the implicit TD(0.5) algorithms.

References

- [1] Leemon Baird et al. Residual algorithms: Reinforcement learning with function approximation. In *Proceedings of the 12th International Conference on Machine Learning*, pages 30–37, 1995.
- [2] Albert Benveniste, Michel Métivier, and Pierre Priouret. *Adaptive Algorithms and Stochastic Approximations*, volume 22. Springer, 2012.
- [3] Dimitri P. Bertsekas. Neuro-dynamic programming. *Athena Scientific*, 1996.
- [4] Jalaj Bhandari, Daniel Russo, and Raghav Singal. A finite-time analysis of temporal difference learning with linear function approximation. In *Conference on Learning Theory*, pages 1691–1692. PMLR, 2018.

- [5] Vivek S. Borkar. Stochastic approximation with two time scales. *Systems & Control Letters*, 29(5):291–294, 1997.
- [6] Vivek S. Borkar. *Stochastic Approximation: A Dynamical Systems Viewpoint*, volume 9. Springer, 2008.
- [7] Léon Bottou. Large-scale machine learning with stochastic gradient descent. In *Proceedings of COMPSTAT’2010*, pages 177–186. Springer, 2010.
- [8] Léon Bottou. Stochastic gradient descent tricks. In *Neural Networks: Tricks of the Trade, 2nd ed.*, pages 421–436. Springer, 2012.
- [9] Jerry Chee, Hwanwoo Kim, and Panos Toulis. “Plus/Minus the Learning Rate”: Easy and scalable statistical inference with SGD. In *International Conference on Artificial Intelligence and Statistics*, pages 2285–2309. PMLR, 2023.
- [10] William Dabney and Andrew Barto. Adaptive step-size for online temporal difference learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 26, pages 872–878. AAAI Press, 2012.
- [11] Gal Dalal, Balázs Szörényi, Gagan Thoppe, and Shie Mannor. Finite sample analyses for TD(0) with function approximation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 32. AAAI Press, 2018.
- [12] Gal Dalal, Gagan Thoppe, Balázs Szörényi, and Shie Mannor. Finite sample analysis of two-timescale stochastic approximation with applications to reinforcement learning. In *Conference on Learning Theory*, pages 1199–1233. PMLR, 2018.
- [13] 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.
- [14] Abraham P. George and Warren B. Powell. Adaptive stepsizes for recursive estimation with applications in approximate dynamic programming. *Machine Learning*, 65:167–198, 2006.
- [15] Sina Ghiassian, Andrew Patterson, Shivam Garg, Dhawal Gupta, Adam White, and Martha White. Gradient temporal-difference learning with regularized corrections. In *International Conference on Machine Learning*, pages 3524–3534. PMLR, 2020.
- [16] Xavier Gourdon and Pascal Sebah. The euler constant: γ . *Young*, 1:2n, 2004.
- [17] Harsh Gupta, Rayadurgam Srikant, and Lei Ying. Finite-time performance bounds and adaptive learning-rate selection for two time-scale reinforcement learning. *Advances in Neural Information Processing Systems*, 32, 2019.

- [18] Marcus Hutter and Shane Legg. Temporal difference updating without a learning rate. *Advances in Neural Information Processing Systems*, 20, 2007.
- [19] Olav Kallenberg. *Foundations of Modern Probability*, volume 2. Springer, 1997.
- [20] Prasenjit Karmakar and Shalabh Bhatnagar. Two time-scale stochastic approximation with controlled markov noise and off-policy temporal-difference learning. *Mathematics of Operations Research*, 43(1):130–151, 2018.
- [21] Vijay R. Konda and John N. Tsitsiklis. Convergence rate of linear two-time-scale stochastic approximation. *The Annals of Applied Probability*, 14(2):796 – 819, 2004.
- [22] Chandrashekar Lakshminarayanan and Csaba Szepesvári. Linear stochastic approximation: How far does constant step-size and iterate averaging go? In *International Conference on Artificial Intelligence and Statistics*, pages 1347–1355. PMLR, 2018.
- [23] David A. Levin and Yuval Peres. *Markov Chains and Mixing Times*, volume 107. American Mathematical Society, 2017.
- [24] Lennart Ljung, Georg Pflug, and Harro Walk. *Stochastic Approximation and Optimization of Random Systems*, volume 17. Birkhäuser, 2012.
- [25] Hamid Reza Maei. Gradient temporal-difference learning algorithms. *PhD thesis, University of Alberta*, 2011.
- [26] Ashique R. Mahmood, Richard S. Sutton, Thomas Degris, and Patrick M. Pilarski. Tuning-free step-size adaptation. In *2012 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 2121–2124. IEEE, 2012.
- [27] Aritra Mitra. A simple finite-time analysis of td learning with linear function approximation. *IEEE Transactions on Automatic Control*, 2024.
- [28] Joseph Modayil, Adam White, and Richard S. Sutton. Multi-timescale nexting in a reinforcement learning robot. *Adaptive Behavior*, 22(2):146–160, 2014.
- [29] John P. O’Doherty, Peter Dayan, Karl Friston, Hugo Critchley, and Raymond J. Dolan. Temporal difference models and reward-related learning in the human brain. *Neuron*, 38(2):329–337, 2003.
- [30] Gandharv Patil, L. A. Prashanth, Dheeraj Nagaraj, and Doina Precup. Finite-time analysis of temporal difference learning with linear function approximation: Tail averaging and regularisation. In *International Conference on Artificial Intelligence and Statistics*, pages 5438–5448. PMLR, 2023.
- [31] Herbert Robbins and Sutton Monro. A stochastic approximation method. *The Annals of Mathematical Statistics*, pages 400–407, 1951.

- [32] Benjamin Van Roy. Temporal-difference learning and applications in finance. *Computational Finance 1999*, page 447, 2000.
- [33] Rayadurgam Srikant and Lei Ying. Finite-time error bounds for linear stochastic approximation and TD learning. In *Conference on Learning Theory*, pages 2803–2830. PMLR, 2019.
- [34] Richard S. Sutton. Learning to predict by the methods of temporal differences. *Machine Learning*, 3:9–44, 1988.
- [35] Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*, volume 1. MIT press Cambridge, 1998.
- [36] Richard S. Sutton, Hamid Maei, and Csaba Szepesvári. A convergent $o(n)$ temporal-difference algorithm for off-policy learning with linear function approximation. *Advances in Neural Information Processing Systems*, 21, 2008.
- [37] Richard S. Sutton, Hamid R. 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, 2009.
- [38] Richard S. Sutton, A. Rupam Mahmood, and Martha White. An emphatic approach to the problem of off-policy temporal-difference learning. *Journal of Machine Learning Research*, 17 (73):1–29, 2016.
- [39] Aviv Tamar, Panos Toulis, Shie Mannor, and Edoardo M. Airolidi. Implicit temporal differences. *arXiv preprint arXiv:1412.6734*, 2014.
- [40] Panagiotis Toulis, Edo Airolidi, and Jason Rennie. Statistical analysis of stochastic gradient methods for generalized linear models. In *International Conference on Machine Learning*, pages 667–675. PMLR, 2014.
- [41] Panos Toulis and Edoardo M. Airolidi. Scalable estimation strategies based on stochastic approximations: Classical results and new insights. *Statistics and Computing*, 25:781–795, 2015.
- [42] Panos Toulis and Edoardo M. Airolidi. Asymptotic and finite-sample properties of estimators based on stochastic gradients. *The Annals of Statistics*, 45(4):1694–1727, 2017.
- [43] Mark Towers, Ariel Kwiatkowski, Jordan Terry, John U. Balis, Gianluca De Cola, Tristan Deleu, Manuel Goulao, Andreas Kallinteris, Markus Krimmel, Arjun K. G., et al. Gymnasium: A standard interface for reinforcement learning environments. *arXiv preprint arXiv:2407.17032*, 2024.
- [44] John Tsitsiklis and Benjamin Van Roy. Analysis of temporal-difference learning with function approximation. *Advances in Neural Information Processing Systems*, 9, 1996.

- [45] Yue Wang, Wei Chen, Yuting Liu, Zhi-Ming Ma, and Tie-Yan Liu. Finite sample analysis of the gtd policy evaluation algorithms in the markov setting. *Advances in Neural Information Processing Systems*, 30, 2017.
- [46] Tengyu Xu, Shaofeng Zou, and Yingbin Liang. Two time-scale off-policy TD learning: Non-asymptotic analysis over markovian samples. *Advances in Neural Information Processing Systems*, 32, 2019.
- [47] Sheng Zhang, Zhe Zhang, and Siva Theja Maguluri. Finite sample analysis of average-reward TD learning and Q-learning. *Advances in Neural Information Processing Systems*, 34:1230–1242, 2021.