
Finite-Time Analysis of Whittle Index based Q-Learning for Restless Multi-Armed Bandits with Neural Network Function Approximation

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Whittle index policy is a heuristic to the intractable restless multi-armed bandits
2 (RMAB) problem. Although it is provably asymptotically optimal, finding Whittle
3 indices remains difficult. In this paper, we present `Neural-Q-Whittle`, a Whittle
4 index based Q-learning algorithm for RMAB with neural network function approx-
5 imation, which is an example of nonlinear two-timescale stochastic approximation
6 with Q-function values updated on a faster timescale and Whittle indices on a slower
7 timescale. Despite the empirical success of deep Q-learning, the non-asymptotic
8 convergence rate of `Neural-Q-Whittle`, which couples neural networks with
9 two-timescale Q-learning largely remains unclear. This paper provides a finite-time
10 analysis of `Neural-Q-Whittle`, where data are generated from a Markov chain,
11 and Q-function is approximated by a ReLU neural network. Our analysis leverages
12 a Lyapunov drift approach to capture the evolution of two coupled parameters,
13 and the nonlinearity in value function approximation further requires us to charac-
14 terize the approximation error. Combining these provide `Neural-Q-Whittle` with
15 $\mathcal{O}(1/k^{2/3})$ convergence rate, where k is the number of iterations.

16 1 Introduction

17 We consider the restless multi-armed bandits (RMAB) problem [55], where the decision maker
18 (DM) repeatedly activates K out of N arms at each decision epoch. Each arm is described by a
19 Markov decision process (MDP) [44], and evolves stochastically according to two different transition
20 kernels, depending on whether the arm is activated or not. Rewards are generated with each transition.
21 Although RMAB has been widely used to study constrained sequential decision making problems
22 [5, 37, 12, 60, 28, 35, 34], it is notoriously intractable due to the explosion of state space [43]. A
23 celebrated heuristic is the Whittle index policy [55], which computes the Whittle index for each arm
24 given its current state as the cost to pull the arm. Whittle index policy then activates the K highest
25 indexed arms at each decision epoch, and is provably asymptotically optimal [53].

26 However, the computation of Whittle index requires full knowledge of the underlying MDP associated
27 with each arm, which is often unavailable in practice. To this end, many recent efforts have focused
28 on learning Whittle indices for making decisions in an online manner. First, model-free reinforcement
29 learning (RL) solutions have been proposed [10, 23, 52, 8, 27, 56, 3], among which [3] developed
30 a Whittle index based Q-learning algorithm, which we call `Q-Whittle` for ease of exposition, and
31 provided the first-ever rigorous asymptotic analysis. However, `Q-Whittle` suffers from slow conver-
32 gence since it only updates the Whittle index of a specific state when that state is visited. In addition,
33 `Q-Whittle` needs to store the Q-function values for all state-action pairs, which limits its applicability
34 only to problems with small state space. Second, deep RL methods have been leveraged to predict

Whittle indices via training neural networks [39, 40]. Though these methods are capable of dealing with large state space, there is no asymptotic or finite-time performance guarantee. Furthermore, training neural networks requires tuning hyper-parameters. This introduces an additional layer of complexity to predict Whittle indices. Third, to address aforementioned deficiencies, [57] proposed Q-Whittle-LFA by coupling Q-Whittle with linear function approximation and provided a finite-time convergence analysis. One key limitation of Q-Whittle-LFA is the unrealistic assumption that all data used in Q-Whittle-LFA are sampled i.i.d. from a fixed stationary distribution.

To tackle the aforementioned limitations and inspired by the empirical success of deep Q-learning in numerous applications, we develop Neural-Q-Whittle, a Whittle index based Q-learning algorithm with *neural network function approximation under Markovian observations*. Like [3, 57], the updates of Q-function values and Whittle indices form a two-timescale stochastic approximation (2TSA) with the former operating on a faster timescale and the later on a slower timescale. Unlike [3, 57], our Neural-Q-Whittle uses a deep neural network with the ReLU activation function to approximate the Q-function. However, Q-learning with neural network function approximation can in general diverge [2], and the theoretical convergence of Q-learning with neural network function approximation has been limited to special cases such as fitted Q-iteration with i.i.d. observations [22], which fails to capture the practical setting of Q-learning with neural network function approximation.

In this paper, we study the non-asymptotic convergence of Neural-Q-Whittle with data generated from a Markov decision process. Compared with recent theoretical works for Q-learning with neural network function approximation [13, 22, 58], our Neural-Q-Whittle involves a two-timescale update between two coupled parameters, i.e., Q-function values and Whittle indices. This renders existing finite-time analysis in [13, 22, 58] not applicable to our Neural-Q-Whittle due to the fact that [13, 22, 58] only contains a single-timescale update on Q-function values. Furthermore, [13, 22, 58] required an additional projection step for the update of parameters of neural network function so as to guarantee the boundedness between the unknown parameter at any time step with the initialization. This in some cases is impractical. Hence, a natural question that arises is

Is it possible to provide a non-asymptotic convergence rate analysis of Neural-Q-Whittle with two coupled parameters updated in two timescales under Markovian observations without the extra projection step?

The theoretical convergence guarantee of two-timescale Q-learning with neural network function approximation under Markovian observations remains largely an open problem, and in this paper, we provide an affirmative answer to this question. Our main contributions are summarized as follows:

- We propose Neural-Q-Whittle, a novel Whittle index based Q-learning algorithm with neural network function approximation for RMAB. Inspired by recent work on TD learning [47] and Q-learning [15] with linear function approximation, our Neural-Q-Whittle removes the additional impractical projection step in the neural network function parameter update.
- We establish the first finite-time analysis of Neural-Q-Whittle under Markovian observations. Due to the two-timescale nature for the updates of two coupled parameters (i.e., Q-function values and Whittle indices) in Neural-Q-Whittle, we focus on the convergence rate of these parameters rather than the convergence rate of approximated Q-functions as in [13, 22, 58]. Our key technique is to view Neural-Q-Whittle as a 2TSA for finding the solution of suitable nonlinear equations. Different from recent works on finite-time analysis of a general 2TSA [20] or with linear function approximation [57], the nonlinear parameterization of Q-function in Neural-Q-Whittle under Markovian observations imposes significant difficulty in finding the global optimum of the corresponding nonlinear equations. To mitigate this, we first approximate the original neural network function with a collection of local linearization and focus on finding a surrogate Q-function in the neural network function class that well approximates the optimum. Our finite-time analysis then requires us to consider two Lyapunov functions that carefully characterize the coupling between iterates of Q-function values and Whittle indices, with one Lyapunov function defined with respect to the true neural network function, and the other defined with respect to the locally linearized neural network function. We then characterize the errors between these two Lyapunov functions. Putting them together, we prove that Neural-Q-Whittle achieves a convergence in expectation at a rate $\mathcal{O}(1/k^{2/3})$, where k is the number of iterations.
- Finally, we conduct experiments to validate the convergence performance of Neural-Q-Whittle, and verify the sufficiency of our proposed condition for the stability of Neural-Q-Whittle.

88 2 Preliminaries

89 **RMAB.** We consider an infinite-horizon average-reward RMAB with each arm $n \in \mathcal{N}$ described by
 90 a unichain MDP [44] $\mathcal{M}_n := (\mathcal{S}, \mathcal{A}, P_n, r_n)$, where \mathcal{S} is the state space with cardinality $S < \infty$, \mathcal{A}
 91 is the action space with cardinality A , $P_n(s'|s, a)$ is the transition probability of reaching state s' by
 92 taking action a in state s , and $r_n(s, a)$ is the reward associated with state-action pair (s, a) . At each
 93 time slot t , the DM activates K out of N arms. Arm n is “active” at time t when it is activated, i.e.,
 94 $A_n(t) = 1$; otherwise, arm n is “passive”, i.e., $A_n(t) = 0$. Let Π be the set of all possible policies
 95 for RMAB, and $\pi \in \Pi$ is a feasible policy, satisfying $\pi : \mathcal{F}_t \mapsto \mathcal{A}^N$, where \mathcal{F}_t is the sigma-algebra
 96 generated by random variables $\{S_n(h), A_n(h) : \forall n \in \mathcal{N}, h \leq t\}$. The objective of the DM is to
 97 maximize the expected long-term average reward subject to an instantaneous constraint that only K
 98 arms can be activated at each time slot, i.e.,

$$\text{RMAB: } \max_{\pi \in \Pi} \liminf_{T \rightarrow \infty} \frac{1}{T} \mathbb{E}_\pi \left(\sum_{t=0}^T \sum_{n=1}^N r_n(t) \right), \quad \text{s.t. } \sum_{n=1}^N A_n(t) = K, \quad \forall t. \quad (1)$$

99 **Whittle Index Policy.** It is well known that RMAB (1) suffers from the curse of dimensionality [43].
 100 To address this challenge, Whittle [55] proposed an index policy through decomposition. Specifically,
 101 Whittle relaxed the constraint in (1) to be satisfied on average and obtained a unconstrained prob-
 102 lem: $\max_{\pi \in \Pi} \liminf_{T \rightarrow \infty} \frac{1}{T} \mathbb{E}_\pi \sum_{t=1}^T \sum_{n=1}^N \{r_n(t) + \lambda(1 - A_n(t))\}$, where λ is the Lagrangian
 103 multiplier associated with the constraint. The key observation of Whittle is that this problem can
 104 be decomposed and its solution is obtained by combining solutions of N independent problems via
 105 solving the associated dynamic programming (DP): $V_n(s) = \max_{a \in \{0,1\}} Q_n(s, a)$, $\forall n \in \mathcal{N}$, where

$$Q_n(s, a) + \beta = a \left(r_n(s, a) + \sum_{s'} p_n(s'|s, 1) V_n(s') \right) + (1-a) \left(r_n(s, a) + \lambda + \sum_{s'} p_n(s'|s, 0) V_n(s') \right), \quad (2)$$

106 where β is unique and equals to the maximal long-term average reward of the unichain MDP, and
 107 $V_n(s)$ is unique up to an additive constant, both of which depend on the Lagrangian multiplier λ . The
 108 optimal decision a^* in state s then is the one which maximizes the right hand side of the above DP.
 109 The Whittle index associated with state s is defined as the value $\lambda_n^*(s) \in \mathbb{R}$ such that actions 0 and 1
 110 are equally favorable in state s for arm n [3, 23], satisfying

$$\lambda_n^*(s) := Q_n(s, 0) - Q_n(s, 1) = r_n(s, 1) + \sum_{s'} p_n(s'|s, 1) V_n(s') - r_n(s, 0) - \sum_{s'} p_n(s'|s, 0) V_n(s'). \quad (3)$$

111 Whittle index policy then activates K arms with the largest Whittle indices at each time slot. Addi-
 112 tional discussions are provided in Section B in supplementary materials.

113 **Q-Learning for Whittle Index.** Since the underlying MDPs are often unknown, [3] proposed
 114 Q-Whittle, a tabular Whittle index based Q-learning algorithm, where the updates of Q-function
 115 values and Whittle indices form a 2TSA, with the former operating on a faster timescale for a given
 116 λ_n and the later on a slower timescale. Specifically, the Q-function values for $\forall n \in \mathcal{N}$ are updated as

$$Q_{n,k+1}(s, a) := Q_{n,k}(s, a) + \alpha_{n,k} \mathbb{1}_{\{S_{n,k}=s, A_{n,k}=a\}} \left(r_n(s, a) + (1-a)\lambda_{n,k}(s) \right. \\ \left. + \max_a Q_{n,k}(S_{n,k+1}, a) - I_n(Q_k) - Q_{n,k}(s, a) \right), \quad (4)$$

117 where $I_n(Q_k) = \frac{1}{2S} \sum_{s \in \mathcal{S}} (Q_{n,k}(s, 0) + Q_{n,k}(s, 1))$ is standard in the relative Q-learning for long-
 118 term average MDP setting [1], which differs significantly from the discounted reward setting [44, 1].
 119 $\{\alpha_{n,k}\}$ is a step-size sequence satisfying $\sum_k \alpha_{n,k} = \infty$ and $\sum_k \alpha_{n,k}^2 < \infty$.

120 Accordingly, the Whittle index is updated as

$$\lambda_{n,k+1}(s) = \lambda_{n,k}(s) + \eta_{n,k} (Q_{n,k}(s, 1) - Q_{n,k}(s, 0)), \quad (5)$$

121 with the step-size sequence $\{\eta_{n,k}\}$ satisfying $\sum_k \eta_{n,k} = \infty$, $\sum_k \eta_{n,k}^2 < \infty$ and $\eta_{n,k} = o(\alpha_{n,k})$.
 122 The coupled iterates (4) and (5) form a 2TSA, and [3] provided an asymptotic convergence analysis.

123 3 Neural Q-Learning for Whittle Index

124 A closer look at (5) reveals that Q-Whittle only updates the Whittle index of a specific state when
 125 that state is visited. This makes Q-Whittle suffers from slow convergence. In addition, Q-Whittle

Algorithm 1 Neural-Q-Whittle: Neural Q-Learning for Whittle Index

```

1: Input:  $\phi(s, a)$  for  $\forall s \in \mathcal{S}, a \in \mathcal{A}$ , and learning rates  $\{\alpha_k\}_{k=1, \dots, T}, \{\eta_k\}_{k=1, \dots, T}$ 
2: Initialization:  $b_r \sim \text{Unif}(\{-1, 1\}), \mathbf{w}_{r,0} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d/d), \forall r \in [1, m]$  and  $\lambda(s) = 0, \forall s \in \mathcal{S}$ 
3: for  $s \in \mathcal{S}$  do
4:   for  $k = 1, \dots, T$  do
5:     Sample  $(S_k, A_k, S_{k+1})$  according to the  $\epsilon$ -greedy policy;
6:      $\Delta_k = r(S_k, A_k) + (1 - A_k)\lambda_k(s) + \max_a f(\boldsymbol{\theta}_k; \boldsymbol{\phi}(S_{k+1}, a)) - I(\boldsymbol{\theta}_k) - f(\boldsymbol{\theta}_k; \boldsymbol{\phi}(S_k, A_k));$ 
7:      $\boldsymbol{\theta}_{k+1} = \boldsymbol{\theta}_k + \alpha_k \Delta_k \nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}_k; \boldsymbol{\phi}(S_k, A_k));$ 
8:      $\lambda_{k+1}(s) = \lambda_k(s) + \eta_k (f(\boldsymbol{\theta}_k; \boldsymbol{\phi}(s, 1)) - f(\boldsymbol{\theta}_k; \boldsymbol{\phi}(s, 0)));$ 
9:   end for
10: end for
11: Return:  $\lambda(s), \forall s \in \mathcal{S}$ .

```

needs to store the Q-function values for all state-action pairs, which limits its applicability only to problems with small state space. To address this challenge and inspired by the empirical success of deep Q-learning, we develop Neural-Q-Whittle through coupling Q-Whittle with neural network function approximation by using low-dimensional feature mapping and leveraging the strong representation power of neural networks. For ease of presentation, we drop the subscript n in (4) and (5), and discussions in the rest of the paper apply to any arm $n \in \mathcal{N}$.

Specifically, given a set of basis functions $\phi_\ell : \mathcal{S} \times \mathcal{A} \mapsto \mathbb{R}, \forall \ell = 1, \dots, d$ with $d \ll SA$, the approximation of Q-function $Q_\theta(s, a)$ parameterized by a unknown weight vector $\boldsymbol{\theta} \in \mathbb{R}^{md}$, is given by $Q_\theta(s, a) = f(\boldsymbol{\theta}; \boldsymbol{\phi}(s, a)), \forall s \in \mathcal{S}, a \in \mathcal{A}$, where f is a nonlinear neural network function parameterized by $\boldsymbol{\theta}$ and $\boldsymbol{\phi}(s, a)$, with $\boldsymbol{\phi}(s, a) = (\phi_1(s, a), \dots, \phi_d(s, a))^\top$. The feature vectors are assumed to be linearly independent and are normalized so that $\|\boldsymbol{\phi}(s, a)\| \leq 1, \forall s \in \mathcal{S}, a \in \mathcal{A}$. In particular, we parameterize the Q-function by using a two-layer neural network [13, 58]

$$f(\boldsymbol{\theta}; \boldsymbol{\phi}(s, a)) := \frac{1}{\sqrt{m}} \sum_{r=1}^m b_r \sigma(\mathbf{w}_r^\top \boldsymbol{\phi}(s, a)), \quad (6)$$

where $\boldsymbol{\theta} = (b_1, \dots, b_m, \mathbf{w}_1^\top, \dots, \mathbf{w}_m^\top)^\top$ with $b_r \in \mathbb{R}$ and $\mathbf{w}_r \in \mathbb{R}^{d \times 1}, \forall r \in [1, m]$. $b_r, \forall r$ are uniformly initialized in $\{-1, 1\}$ and $\mathbf{w}_r, \forall r$ are initialized as a zero mean Gaussian distribution according to $\mathcal{N}(\mathbf{0}, \mathbf{I}_d/d)$. During training process, only $\mathbf{w}_r, \forall r$ are updated while $b_r, \forall r$ are fixed as the random initialization. Hence, we use $\boldsymbol{\theta}$ and $\mathbf{w}_r, \forall r$ interchangeably throughout this paper. $\sigma(x) = \max(0, x)$ is the rectified linear unit (ReLU) activation function.

Given (6), we can rewrite the Q-function value updates in (4) as

$$\boldsymbol{\theta}_{k+1} = \boldsymbol{\theta}_k + \alpha_k \Delta_k \nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}_k; \boldsymbol{\phi}(S_k, A_k)), \quad (7)$$

with Δ_k being the temporal difference (TD) error defined as $\Delta_k := r(S_k, A_k) + (1 - A_k)\lambda_k(s) - I(\boldsymbol{\theta}_k) + \max_a f(\boldsymbol{\theta}_k; \boldsymbol{\phi}(S_{k+1}, a)) - f(\boldsymbol{\theta}_k; \boldsymbol{\phi}(S_k, A_k))$, where $I(\boldsymbol{\theta}_k) = \frac{1}{2S} \sum_{s \in \mathcal{S}} [f(\boldsymbol{\theta}_k; \boldsymbol{\phi}(s, 0)) + f(\boldsymbol{\theta}_k; \boldsymbol{\phi}(s, 1))]$. Similarly, the Whittle index update (5) can be rewritten as

$$\lambda_{k+1}(s) = \lambda_k(s) + \eta_k (f(\boldsymbol{\theta}_k; \boldsymbol{\phi}(s, 1)) - f(\boldsymbol{\theta}_k; \boldsymbol{\phi}(s, 0))). \quad (8)$$

The coupled iterates in (7) and (8) form Neural-Q-Whittle as summarized in Algorithm 1, which aims to learn the coupled parameters $(\boldsymbol{\theta}^*, \lambda^*(s))$ such that $f(\boldsymbol{\theta}^*; \boldsymbol{\phi}(s, 1)) = f(\boldsymbol{\theta}^*; \boldsymbol{\phi}(s, 0)), \forall s \in \mathcal{S}$.

Remark 1. Unlike recent works for Q-learning with linear [6, 36, 63] or neural network function approximations [13, 22, 58], we do not assume an additional projection step of the updates of unknown parameters $\boldsymbol{\theta}_k$ in (7) to confine $\boldsymbol{\theta}_k, \forall k$ into a bounded set. This projection step is often used to stabilize the iterates related to the unknown stationary distribution of the underlying Markov chain, which in some cases is impractical. More recently, [47] removed the extra projection step and established the finite-time convergence of TD learning, which is treated as a linear stochastic approximation algorithm. [15] extended it to the Q-learning with linear function approximation. However, these state-of-the-art works only contained a single-timescale update on Q-function values, i.e., with the only unknown parameter $\boldsymbol{\theta}$, while our Neural-Q-Whittle involves a two-timescale update between two coupled unknown parameters $\boldsymbol{\theta}$ and λ as in (7) and (8). Our goal in this paper is to expand the frontier by providing a finite-time bound for Neural-Q-Whittle under Markovian noise without requiring an additional projection step.

4 Finite-Time Analysis of Neural-Q-Whittle

In this section, we present the finite-time analysis of Neural-Q-Whittle for learning Whittle index $\lambda(s)$ of any state $s \in \mathcal{S}$ when data are generated from a MDP. To simplify notation, we abbreviate $\lambda(s)$ as λ in the rest of the paper. We start by first rewriting the updates of Neural-Q-Whittle in (7) and (8) as a nonlinear two-timescale stochastic approximation (2TSA) in Section 4.1.

4.1 A Nonlinear 2TSA Formulation with Neural Network Function

We first show that Neural-Q-Whittle can be rewritten as a variant of the nonlinear 2TSA. For any fixed policy π , since the state of each arm $\{S_k\}$ evolves according to a Markov chain, we can construct a new variable $X_k = (S_k, A_k, S_{k+1})$, which also forms a Markov chain with state space $\mathcal{X} := \{(s, a, s') | s \in \mathcal{S}, \pi(a|s) \geq 0, p(s'|s, a) > 0\}$. Therefore, the coupled updates (7) and (8) of Neural-Q-Whittle can be rewritten in the form of a nonlinear 2TSA [20]:

$$\theta_{k+1} = \theta_k + \alpha_k h(X_k, \theta_k, \lambda_k), \quad \lambda_{k+1} = \lambda_k + \eta_k g(X_k, \theta_k, \lambda_k), \quad (9)$$

where θ_0 and λ_0 being arbitrarily initialized in \mathbb{R}^{md} and \mathbb{R} , respectively; and $h(\cdot)$ and $g(\cdot)$ satisfy

$$h(X_k, \theta_k, \lambda_k) := \nabla_{\theta} f(\theta_k; \phi(S_k, A_k)) \Delta_k, \quad \theta_k \in \mathbb{R}^{md}, \lambda_k \in \mathbb{R}, \quad (10)$$

$$g(X_k, \theta_k, \lambda_k) := f(\theta_k; \phi(s, 1)) - f(\theta_k; \phi(s, 0)), \quad \theta_k \in \mathbb{R}^{md}. \quad (11)$$

Since $\eta_k \ll \alpha_k$, the dynamics of θ evolves much faster than those of λ . We aim to establish the finite-time performance of the nonlinear 2TSA in (9), where $f(\cdot)$ is the neural network function defined in (6). This is equivalent to find the root¹ (θ^*, λ^*) of a system with *two coupled* nonlinear equations $h : \mathcal{X} \times \mathbb{R}^{md} \times \mathbb{R} \rightarrow \mathbb{R}^{md}$ and $g : \mathcal{X} \times \mathbb{R}^{md} \times \mathbb{R} \rightarrow \mathbb{R}$ such that

$$H(\theta, \lambda) := \mathbb{E}_{\mu}[h(X, \theta, \lambda)] = 0, \quad G(\theta, \lambda) := \mathbb{E}_{\mu}[g(X, \theta, \lambda)] = 0, \quad (12)$$

where X is a random variable in finite state space \mathcal{X} with unknown distribution μ . For a fixed θ , to study the stability of λ , we assume the condition on the existence of a mapping such that $\lambda = y(\theta)$ is the unique solution of $G(\theta, \lambda) = 0$. In particular, $y(\theta)$ is given as

$$y(\theta) = r(s, 1) + \sum_{s'} p(s'|s, 1) \max_a f(\theta; \phi(s', a)) - r(s, 0) - \sum_{s'} p(s'|s, 0) \max_a f(\theta; \phi(s', a)). \quad (13)$$

4.2 Main Results

As inspired by [20], the finite-time analysis of such a nonlinear 2TSA boils down to the choice of two step sizes $\{\alpha_k, \eta_k, \forall k\}$ and a Lyapunov function that couples the two iterates in (9). To this end, we first define the following two error terms:

$$\tilde{\theta}_k = \theta_k - \theta^*, \quad \tilde{\lambda}_k = \lambda_k - y(\theta_k), \quad (14)$$

which characterize the coupling between θ_k and λ_k . If $\tilde{\theta}_k$ and $\tilde{\lambda}_k$ go to zero simultaneously, the convergence of (θ_k, λ_k) to (θ^*, λ^*) can be established. Thus, to prove the convergence of (θ_k, λ_k) of the nonlinear 2TSA in (9) to its true value (θ^*, λ^*) , we can equivalently study the convergence of $(\tilde{\theta}_k, \tilde{\lambda}_k)$ by providing the finite-time analysis for the mean squared error generated by (9). To couple the fast and slow iterates, we define the following weighted Lyapunov function

$$M(\theta_k, \lambda_k) := \frac{\eta_k}{\alpha_k} \|\tilde{\theta}_k\|^2 + \|\tilde{\lambda}_k\|^2 = \frac{\eta_k}{\alpha_k} \|\theta_k - \theta^*\|^2 + \|\lambda_k - y(\theta_k)\|^2, \quad (15)$$

where $\|\cdot\|$ stands for the the Euclidean norm for vectors throughout the paper. It is clear that the Lyapunov function $M(\theta_k, \lambda_k)$ combines the updates of θ and λ with respect to the true neural network function $f(\theta; \phi(s, a))$ in (6).

To this end, our goal turns to characterize finite-time convergence of $\mathbb{E}[M(\theta_k, \lambda_k)]$. However, it is challenging to directly finding the global optimum of the corresponding nonlinear equations due

¹The root (θ^*, λ^*) of the nonlinear 2TSA (9) can be established by using the ODE method following the solution of suitably defined differential equations [9, 48, 3, 21, 19, 20], i.e., $\dot{\theta} = H(\theta, \lambda)$, $\dot{\lambda} = \frac{\eta}{\alpha} G(\theta, \lambda)$, where a fixed stepsize is assumed for ease of expression at this moment.

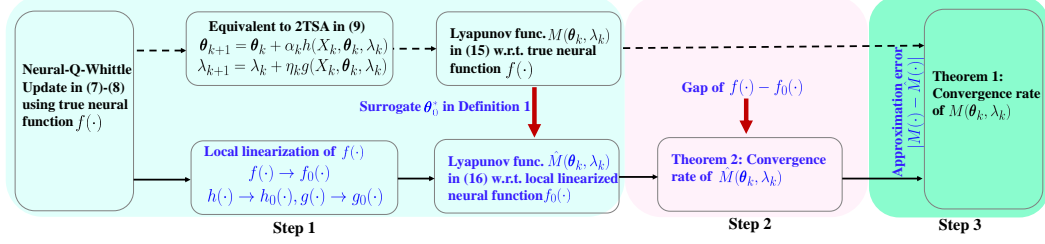


Figure 1: Neural-Q-Whittle operates w.r.t. true neural function $f(\cdot)$ with its finite-time performance given in Theorem 1 (indicated in dashed lines). Our proofs operate in three steps: (i) Step 1: Obtain local linearization $f_0(\cdot)$ and define Lyapunov function $\hat{M}(\cdot)$ w.r.t. $f_0(\cdot)$. (ii) Step 2: Characterize the finite-time performance w.r.t. $\hat{M}(\cdot)$ using Lyapunov drift method. Since Neural-Q-Whittle is updated w.r.t. $f(\cdot)$, we need to characterize the gap between $f(\cdot)$ and $f_0(\cdot)$. (iii) Step 3: Similarly, we characterize the approximation errors between $M(\cdot)$ and $\hat{M}(\cdot)$.

to the nonlinear parameterization of Q-function in Neural-Q-Whittle. In addition, the operators $h(\cdot)$, $g(\cdot)$ and $y(\cdot)$ in (10), (11) and (13) directly relate with the convoluted neural network function $f(\theta; \phi(s, a))$ in (6), which hinders us to characterize the smoothness properties of these operators. Such properties are often required for the analysis of stochastic approximation [15, 19, 21].

To mitigate this, **(Step 1)** we instead approximate the true neural network function $f(\theta, \phi(s, a))$ with a collection of local linearization $f_0(\theta; \phi(s, a))$ at the initial point θ_0 . Based on the surrogate stationary point θ_0^* of $f_0(\theta; \phi(s, a))$, we correspondingly define a modified Lyapunov function $\hat{M}(\theta_k, \lambda_k)$ combining updates of θ and λ with respect to such local linearization. Specifically, we have

$$\hat{M}(\theta_k, \lambda_k) := \frac{\eta_k}{\alpha_k} \|\theta_k - \theta_0^*\|^2 + \|\lambda_k - y_0(\theta_k)\|^2, \quad (16)$$

where $y_0(\cdot)$ is in the same expression as $y(\cdot)$ in (13) by replacing $f(\cdot)$ with $f_0(\cdot)$, and we will describe this in details below. **(Step 2)** We then study the convergence rate of the nonlinear 2TSA using this modified Lyapunov function under general conditions. **(Step 3)** Finally, since the two coupled parameters θ and λ in (9) are updated with respect to the true neural network function $f(\theta; \phi(s, a))$ in (6) in Neural-Q-Whittle, while we characterize their convergence using the approximated neural network function in Step 2. Hence, this further requires us to characterize the approximation errors. We visualize the above three steps in Figure 1 and provide a proof sketch in Section 4.3. Combining them together gives rise to our main theoretical results on the finite-time performance of Neural-Q-Whittle, which is formally stated in the following theorem.

Theorem 1. Consider iterates $\{\theta_k\}$ and $\{\lambda_k\}$ generated by Neural-Q-Whittle in (7) and (8). Given $\alpha_k = \frac{\alpha_0}{(k+1)}$, $\eta_k = \frac{\eta_0}{(k+1)^{4/3}}$, we have for $\forall k \geq \tau$

$$\begin{aligned} \mathbb{E}[M(\theta_{k+1}, \lambda_{k+1}) | \mathcal{F}_{k-\tau}] &\leq \frac{2\tau^2 \mathbb{E}[\hat{M}(\theta_\tau, \lambda_\tau)]}{(k+1)^2} + \frac{1200\alpha_0^3 (C_1 + \|\hat{\theta}_0\|)^2 + (2C_1 + \|\hat{\lambda}_0\|)^2}{\eta_0 (k+1)^{2/3}} \\ &+ \frac{2\eta_0 c_0^2}{\alpha_0(1-\kappa)^2} \|\text{span}(\Pi_{\mathcal{F}} f(\theta^*) - f(\theta^*))\|^2 + \left(\frac{2}{(k+1)^{2/3}} + 2 \right) \mathcal{O}\left(\frac{c_1^3 (\|\theta_0\| + \|\lambda_0\| + 1)^3}{m^{1/2}} \right), \end{aligned} \quad (17)$$

where $C_1 := c_1(\|\theta_0\| + \|\lambda_0\| + 1)$ with c_1 being a proper chosen constant, c_0 is a constant defined in Assumption 3, τ is the mixing time defined in (22), and $\Pi_{\mathcal{F}}$ represents the projection to the set of \mathcal{F} containing all possible $f_0(\theta; \phi(s, a))$ in (18).

The first term on the right hand side (17) corresponds to the bias compared to the Lyapunov function at the mixing time τ , which goes to zero at a rate of $\mathcal{O}(1/k^2)$. The second term corresponds to the accumulated estimation error of the nonlinear 2TSA due to Markovian noise, which vanishes at the rate $\mathcal{O}(1/k^{2/3})$. Hence it dominates the overall convergence rate in (17). The third term captures the distance between the optimal solution (θ^*, λ^*) to the true neural network function $f(\theta_k; \phi(s, a))$ in (6) and the optimal one $(\theta_0^*, y_0(\theta_0^*))$ with local linearization $f_0(\theta_k; \phi(s, a))$ in (18), which quantifies the error when $f(\theta^*)$ does not fall into the function class \mathcal{F} . The last term characterizes the distance between $f(\theta_k; \phi(s, a))$ and $f_0(\theta_k; \phi(s, a))$ with any θ_k . Both terms diminish as $m \rightarrow \infty$. Theorem 1 implies the convergence to the optimal value (θ^*, λ^*) is bounded by the approximation error, which

will diminish to zero as representation power of $f_0(\boldsymbol{\theta}_k; \boldsymbol{\phi}(s, a))$ increases when $m \rightarrow \infty$. Need to mention that a constant step size will result in a non-vanishing accumulated error as in [15].

Remark 2. A finite-time analysis of nonlinear 2TSA was presented in [38]. However, [38] required a stability condition that $\lim_{k \rightarrow \infty} (\boldsymbol{\theta}_k, \lambda_k) = (\boldsymbol{\theta}^*, \lambda^*)$, and both h and g are locally approximated as linear functions. [19, 57] relaxed these conditions and provided a finite-time analysis under i.i.d. noise. These results were later extended to Markovian noise [20] under the assumption that H function is strongly monotone in $\boldsymbol{\theta}$ and G function is strongly monotone in λ . Since [20] leveraged the techniques in [19], it needed to explicitly characterize the covariance between the error caused by Markovian noise and the parameters' residual error in (14), leading to the convergence analysis much more intrinsic. [15] exploited the mixing time to avoid the covariance between the error caused by Markovian noise and the parameters' residual error; however, it only considered the single timescale Q -learning with linear function approximation. Though our *Neural-Q-Whittle* can be rewritten as a nonlinear 2TSA, the nonlinear parameterization of Q -function caused by the neural network function approximation makes the aforementioned analysis not directly applicable to ours and requires additional characterization as highlighted in Figure 1. The explicit characterization of approximation errors further distinguish our work.

4.3 Proof Sketch

In this section, we sketch the proofs of the three steps shown in Figure 1 as required for Theorem 1.

4.3.1 Step 1: Approximated Solution of Neural-Q-Whittle

We first approximate the optimal solution by projecting the Q -function in (6) to some function classes parameterized by $\boldsymbol{\theta}$. The common choice of the projected function classes is the local linearization of $f(\boldsymbol{\theta}; \boldsymbol{\phi}(s, a))$ at the initial point $\boldsymbol{\theta}_0$ [13, 58], i.e., $\mathcal{F} := \{f_0(\boldsymbol{\theta}; \boldsymbol{\phi}(s, a)), \forall \boldsymbol{\theta} \in \Theta\}$, where

$$f_0(\boldsymbol{\theta}; \boldsymbol{\phi}(s, a)) = \frac{1}{\sqrt{m}} \sum_{r=1}^m b_r \mathbb{1}\{\mathbf{w}_{r,0}^\top \boldsymbol{\phi}(s, a) > 0\} \mathbf{w}_r^\top \boldsymbol{\phi}(s, a). \quad (18)$$

Then, we define the approximate stationary point $\boldsymbol{\theta}_0^*$ with respect to $f_0(\boldsymbol{\theta}; \boldsymbol{\phi}(s, a))$ as follows.

Definition 1. [[13, 58]] A point $\boldsymbol{\theta}_0^* \in \Theta$ is said to be the approximate stationary point of Algorithm 1 if for all feasible $\boldsymbol{\theta} \in \Theta$ it holds that $\mathbb{E}_{\mu, \pi, \mathcal{P}}[(\Delta_0 \cdot \nabla_{\boldsymbol{\theta}} f_0(\boldsymbol{\theta}; \boldsymbol{\phi}(s, a)))^\top (\boldsymbol{\theta} - \boldsymbol{\theta}_0^*)] \geq 0, \forall \boldsymbol{\theta} \in \Theta$, with $\Delta_0 := [r(s, a) + (1 - a)\lambda^* - I_0(\boldsymbol{\theta}) + \max_{a'} f_0(\boldsymbol{\theta}; \boldsymbol{\phi}(s', a)) - f_0(\boldsymbol{\theta}; \boldsymbol{\phi}(s, a))]$, where $I_0(\boldsymbol{\theta}) = \frac{1}{2S} \sum_{s \in \mathcal{S}} [f_0(\boldsymbol{\theta}; \boldsymbol{\phi}(s, 0)) + f_0(\boldsymbol{\theta}; \boldsymbol{\phi}(s, 1))]$

Though there is a gap between the true neural function (6) and the approximated local linearized function (18), the gap diminishes as the width of neural network i.e., m , becomes large [13, 58].

With the approximated stationary point $\boldsymbol{\theta}_0^*$, we can redefine the two error terms in (14) as

$$\hat{\boldsymbol{\theta}}_k = \boldsymbol{\theta}_k - \boldsymbol{\theta}_0^*, \quad \hat{\lambda}_k = \lambda_k - y_0(\boldsymbol{\theta}_k), \quad (19)$$

using which we correspondingly define a modified Lyapunov function $\hat{M}(\boldsymbol{\theta}_k, \lambda_k)$ in (16), where

$$y_0(\boldsymbol{\theta}) = r(s, 1) + \sum_{s'} p(s'|s, 1) \max_a f_0(\boldsymbol{\theta}; \boldsymbol{\phi}(s', a)) - r(s, 0) - \sum_{s'} p(s'|s, 0) \max_a f_0(\boldsymbol{\theta}; \boldsymbol{\phi}(s', a)). \quad (20)$$

4.3.2 Step 2: Convergence Rate of $\hat{M}(\boldsymbol{\theta}_k, \lambda_k)$ in (16)

Since we approximate the true neural network function $f(\boldsymbol{\theta}; \boldsymbol{\phi}(s, a))$ in (6) with the local linearized function $f_0(\boldsymbol{\theta}; \boldsymbol{\phi}(s, a))$ in (18), the operators $h(\cdot)$ and $g(\cdot)$ in (10)-(11) turn correspondingly to be

$$h_0(X_k, \boldsymbol{\theta}_k, \lambda_k) = \nabla_{\boldsymbol{\theta}} f_0(\boldsymbol{\theta}_k; \boldsymbol{\phi}(S_k, A_k)) \Delta_{k,0}, \quad g_0(\boldsymbol{\theta}_k) := f_0(\boldsymbol{\theta}_k; \boldsymbol{\phi}(s, 1)) - f_0(\boldsymbol{\theta}_k; \boldsymbol{\phi}(s, 0)), \quad (21)$$

with $\Delta_{k,0} := r(S_k, A_k) + (1 - A_k)\lambda_k - I_0(\boldsymbol{\theta}_k) + \max_a f_0(\boldsymbol{\theta}_k; \boldsymbol{\phi}(S_{k+1}, a)) - f_0(\boldsymbol{\theta}_k; \boldsymbol{\phi}(S_k, A_k))$.

Before we present the finite-time error bound of the nonlinear 2TSA (9) under Markovian noise, we first discuss the mixing time of the Markov chain $\{X_k\}$ and our assumptions.

Definition 2 (Mixing time [15]). For any $\delta > 0$, define τ_δ as

$$\tau_\delta = \min\{k \geq 1 : \|\mathbb{E}[h_0(X_k, \boldsymbol{\theta}, \lambda) | X_0 = x] - H_0(\boldsymbol{\theta}, \lambda)\| \leq \delta(\|\boldsymbol{\theta} - \boldsymbol{\theta}_0^*\| + \|\lambda - y_0(\boldsymbol{\theta}_0^*)\|)\}. \quad (22)$$

Assumption 1. The Markov chain $\{X_k\}$ is irreducible and aperiodic. Hence, there exists a unique stationary distribution μ [31], and constants $C > 0$ and $\rho \in (0, 1)$ such that $d_{TV}(P(X_k|X_0 = x), \mu) \leq C\rho^k, \forall k \geq 0, x \in \mathcal{X}$, where $d_{TV}(\cdot, \cdot)$ is the total-variation (TV) distance [31].

Remark 3. Assumption 1 is often assumed to study the asymptotic convergence of stochastic approximation under Markovian noise [4, 9, 15].

Lemma 1. The function $h_0(X, \theta, \lambda)$ defined in (21) is globally Lipschitz continuous w.r.t θ and λ uniformly in X , i.e., $\|h_0(X, \theta_1, \lambda_1) - h_0(X, \theta_2, \lambda_2)\| \leq L_{h,1}\|\theta_1 - \theta_2\| + L_{h,2}\|\lambda_1 - \lambda_2\|, \forall X \in \mathcal{X}$, and $L_{h,1} = 3, L_{h,2} = 1$ are valid Lipschitz constants.

Lemma 2. The function $g_0(\theta)$ defined in (21) is linear and thus Lipschitz continuous in θ , i.e., $\|g_0(\theta_1) - g_0(\theta_2)\| \leq L_g\|\theta_1 - \theta_2\|$, and $L_g = 2$ is a valid Lipschitz constant.

Lemma 3. The function $y_0(\theta)$ defined in (20) is linear and thus Lipschitz continuous in θ , i.e., $\|y_0(\theta_1) - y_0(\theta_2)\| \leq L_y\|\theta_1 - \theta_2\|$, and $L_y = 2$ is a valid Lipschitz constant.

Remark 4. The Lipschitz continuity of h guarantees the existence of a solution θ to the ODE $\dot{\theta}$ for a fixed λ , while the Lipschitz continuity of g and y ensures the existence of a solution λ to the ODE $\dot{\lambda}$ when θ is fixed. These lemmas often serve as assumptions when proving the convergence rate for both linear and nonlinear 2TSA [30, 38, 18, 24, 19, 17, 26].

Lemma 4. For a fixed λ , there exists a constant $\mu_1 > 0$ such that $h_0(X, \theta, \lambda)$ defined in (10) satisfies

$$\mathbb{E}[\hat{\theta}^\top h_0(X, \theta, \lambda)] \leq -\mu_1 \|\hat{\theta}\|^2.$$

For fixed θ , there exists a constant $\mu_2 > 0$ such that $g_0(X, \theta, \lambda)$ defined in (11) satisfies

$$\mathbb{E}[\hat{\lambda} g_0(X, \theta, \lambda)] \leq -\mu_2 \|\hat{\lambda}\|^2.$$

Remark 5. Lemma 4 guarantees the stability and uniqueness of the solution θ to the ODE $\dot{\theta}$ for a fixed λ , and the uniqueness of the solution λ to the ODE $\dot{\lambda}$ for a fixed θ . This assumption can be viewed as a relaxation of the stronger monotone property of nonlinear mappings [19, 15], since it is automatically satisfied if h and g are strong monotone as assumed in [19].

Lemma 5. Under Assumption 1 and Lemma 1, there exist constants $C > 0, \rho \in (0, 1)$ and $L = \max(3, \max_X h_0(X, \theta_0^*), y_0(\theta_0^*))$ such that

$$\tau_\delta \leq \frac{\log(1/\delta) + \log(2LCmd)}{\log(1/\rho)}.$$

Remark 6. τ_δ is equivalent to the mixing time of the underlying Markov chain satisfying $\lim_{\delta \rightarrow 0} \delta \tau_\delta = 0$ [15]. For simplicity, we remove the subscript and denote it as τ .

We now present the finite-time error bound for the Lyapunov function $\hat{M}(\theta_k, \lambda_k)$ in (16).

Theorem 2. Consider iterates $\{\theta_k\}$ and $\{\lambda_k\}$ generated by Neural-Q-Whittle in (7) and (8). Given Lemma 1-4, $\alpha_k = \frac{\alpha_0}{(k+1)}, \eta_k = \frac{\eta_0}{(k+1)^{4/3}}, C_1 := c_1(\|\theta_0\| + \|\lambda_0\| + 1)$ with a constant c_1 ,

$$\begin{aligned} \mathbb{E}[\hat{M}(\theta_{k+1}, \lambda_{k+1}) | \mathcal{F}_{k-\tau}] &\leq \frac{\tau^2 \mathbb{E}[\hat{M}(\theta_\tau, \lambda_\tau)]}{(k+1)^2} + \frac{600\alpha_0^3 (C_1 + \|\hat{\theta}_0\|)^2 + (2C_1 + \|\hat{\lambda}_0\|)^2}{\eta_0 (k+1)^{2/3}} \\ &\quad + \frac{\mathcal{O}\left(c_1^3(\|\theta_0\| + \|\lambda_0\| + 1)^3 m^{-1/2}\right)}{(k+1)^{2/3}}, \quad \forall k \geq \tau. \end{aligned} \quad (23)$$

4.3.3 Step 3: Approximation Error between $M(\theta_k, \lambda_k)$ and $\hat{M}(\theta_k, \lambda_k)$

Finally, we characterize the approximation error between Lyapunov functions $M(\theta_k, \lambda_k)$ and $\hat{M}(\theta_k, \lambda_k)$. Since we are dealing with long-term average MDP, we assume that the total variation of the MDP is bounded [46].

Assumption 2. There exists $0 < \kappa < 1$ such that $\sup_{(s,a),(s',a')} \|p(\cdot|s,a) - p(\cdot|s',a')\|_{TV} = 2\kappa$.

Hence, the Bellman operator is a span-contraction operator [46], i.e.,

$$\text{span}(\mathcal{T}f_0(\theta_0^*) - \mathcal{T}f(\theta^*)) \leq \kappa \text{span}(f_0(\theta_0^*) - f(\theta^*)). \quad (24)$$

Assumption 3. $\|\theta_0^* - \theta^*\| \leq c_0 \|\text{span}(f_0(\theta_0^*) - f(\theta^*))\|$, with c_0 being a positive constant.

Lemma 6. For $M(\theta_k, \lambda_k)$ in (15) and $\hat{M}(\theta_k, \lambda_k)$ in (16), with constants c_1 and c_0 (Assumption 3),

$$M(\theta_k, \lambda_k) \leq 2\hat{M}(\theta_k, \lambda_k) + \frac{2\eta_k c_0^2}{\alpha_k(1-\kappa)} \|\text{span}(\Pi_{\mathcal{F}} f(\theta^*) - f(\theta^*))\| + 2\mathcal{O}\left(\frac{c_1^3(\|\theta_0\| + \|\lambda_0\| + 1)^3}{m^{1/2}}\right).$$

5 Numerical Experiments

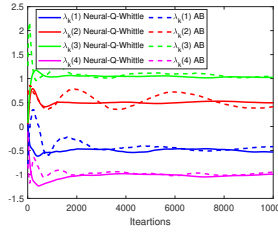
We numerically evaluate the performance of Neural-Q-Whittle using an example of circulant dynamics [23, 3, 8]. The state space is $\mathcal{S} = \{1, 2, 3, 4\}$. Rewards are $r(1, a) = -1, r(2, a) = r(3, a) = 0$, and $r(4, a) = 1$ for $a \in \{0, 1\}$. The dynamics of states are circulant and defined as

$$P^1 = \begin{bmatrix} 0.5 & 0.5 & 0 & 0 \\ 0 & 0.5 & 0.5 & 0 \\ 0 & 0 & 0.5 & 0.5 \\ 0.5 & 0 & 0 & 0.5 \end{bmatrix} \text{ and } P^0 = \begin{bmatrix} 0.5 & 0 & 0 & 0.5 \\ 0.5 & 0.5 & 0 & 0 \\ 0 & 0.5 & 0.5 & 0 \\ 0 & 0 & 0.5 & 0.5 \end{bmatrix}.$$

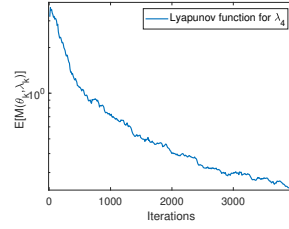
This indicates that the process either remains in its current state or increments if it is active (i.e., $a = 1$), or it either remains the current state or decrements if it is passive (i.e., $a = 0$). The exact value of Whittle indices [23] are $\lambda(1) = -0.5, \lambda(2) = 0.5, \lambda(3) = 1$, and $\lambda(4) = -1$. In our experiments, we set the learning rates as $\alpha_k = 0.5/(k + 1)$ and $\eta_k = 0.1/(k + 1)^{4/3}$. We use ϵ -greedy for the exploration and exploitation tradeoff with $\epsilon = 0.5$. We consider a two-layer neural network with the number of neurons in the hidden layer as $m = 200$. As described in Algorithm 1, $b_r, \forall r$ are uniformly initialized in $\{-1, 1\}$ and $w_r, \forall r$ are initialized as a zero mean Gaussian distribution according to $\mathcal{N}(\mathbf{0}, \mathbf{I}_d/d)$. These results are carried out by Monte Carlo simulations with 100 independent trials.

Convergence to true Whittle index.

We verify that Neural-Q-Whittle converges to true Whittle indices. As illustrated in Figure 2a, Neural-Q-Whittle guarantees the convergence to true Whittle indices and outperforms Q-Whittle [3] in the convergence speed. This is due to the fact that Neural-Q-Whittle updates the Whittle index of a specific state even when the current visited state is not that state. Note that many existing Whittle index based Q-learning algorithms do not guarantee convergence to the true Whittle indices, e.g., WIQL [8] and QWIC [23]. See additional discussions in supplementary materials.



(a) Neural-Q-Whittle vs. Q-Whittle [3].



(b) Convergence of Lyapunov function in (15).

Figure 2: Convergence of Neural-Q-Whittle.

Convergence of the Lyapunov function defined in (15). We also evaluate the convergence of the proposed Lyapunov function defined in (15), which is presented in Figure 2b. It depicts $\mathbb{E}[M(\theta_k, \lambda_k)]$ vs. the number of iterations in logarithmic scale. For ease of presentation, we only take state $s = 4$ as an illustrative example. It is clear that $M(\theta_k, \lambda_k)$ converges to zero as the number of iterations increases, which is in alignment with our theoretical results in Theorem 1.

Verification of Assumption 3. We now verify Assumption 3 that the gap between θ_0^* and θ^* can be bounded by the span of $f_0(\theta_0^*)$ and $f(\theta^*)$ with a constant c_0 . In Figure 3, we show c_0 as a function of the number of neurons in the hidden layer m . It clearly indicates that constant c_0 exists and decreases as the number of neurons grows larger.

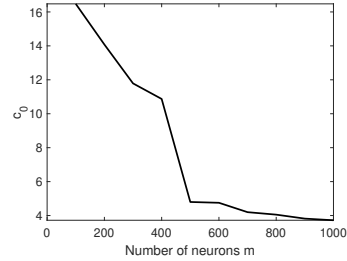


Figure 3: Verification of Assumption 3 w.r.t the constant c_0 .

6 Conclusion

We presented Neural-Q-Whittle, a Whittle index based Q-learning algorithm for the notoriously intractable RMAB with neural network function approximation. We proved that Neural-Q-Whittle achieves an $\mathcal{O}(1/k^{2/3})$ convergence rate, where k is the number of iterations when data are generated from a Markov chain and Q-function is approximated by a ReLU neural network. By viewing Neural-Q-Whittle as 2TSA and leveraging the Lyapunov drift method, we removed the projection step on parameter update of Q-learning with neural network function approximation. Extending the current framework to two-timescale Q-learning (i.e., the coupled iterates between Q-function values and Whittle indices) with general deep neural network approximation is our future work.

References

- [1] Jinane Abounadi, Dimitris Bertsekas, and Vivek S Borkar. Learning algorithms for markov decision processes with average cost. *SIAM Journal on Control and Optimization*, 40(3):681–698, 2001.
- [2] Joshua Achiam, Ethan Knight, and Pieter Abbeel. Towards characterizing divergence in deep q-learning. *arXiv preprint arXiv:1903.08894*, 2019.
- [3] Konstantin E Avrachenkov and Vivek S Borkar. Whittle index based q-learning for restless bandits with average reward. *Automatica*, 139:110186, 2022.
- [4] Dimitri Bertsekas and John N Tsitsiklis. *Neuro-dynamic programming*. Athena Scientific, 1996.
- [5] Dimitris Bertsimas and José Niño-Mora. Restless Bandits, Linear Programming Relaxations, and A Primal-Dual Index Heuristic. *Operations Research*, 48(1):80–90, 2000.
- [6] 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.
- [7] Shalabh Bhatnagar, Richard S Sutton, Mohammad Ghavamzadeh, and Mark Lee. Natural Actor–Critic Algorithms. *Automatica*, 45(11):2471–2482, 2009.
- [8] Arpita Biswas, Gaurav Aggarwal, Pradeep Varakantham, and Milind Tambe. Learn to intervene: An adaptive learning policy for restless bandits in application to preventive healthcare. In *Proc. of IJCAI*, 2021.
- [9] Vivek S Borkar. *Stochastic Approximation: A Dynamical Systems Viewpoint*, volume 48. Springer, 2009.
- [10] Vivek S Borkar and Karan Chadha. A reinforcement learning algorithm for restless bandits. In *2018 Indian Control Conference (ICC)*, pages 89–94. IEEE, 2018.
- [11] Vivek S Borkar and Vijaymohan R Konda. The Actor-Critic Algorithm as Multi-Time-Scale Stochastic Approximation. *Sadhana*, 22(4):525–543, 1997.
- [12] Vivek S Borkar, K Ravikumar, and Krishnakant Saboo. An index policy for dynamic pricing in cloud computing under price commitments. *Applicationes Mathematicae*, 44:215–245, 2017.
- [13] Qi Cai, Zhuoran Yang, Jason D Lee, and Zhaoran Wang. Neural temporal difference and q learning provably converge to global optima. *Mathematics of Operations Research*, 2023.
- [14] Zaiwei Chen, Siva Theja Maguluri, Sanjay Shakkottai, and Karthikeyan Shanmugam. Finite-Sample Analysis of Stochastic Approximation Using Smooth Convex Envelopes. *arXiv e-prints*, pages arXiv–2002, 2020.
- [15] Zaiwei Chen, Sheng Zhang, Thinh T Doan, Siva Theja Maguluri, and John-Paul Clarke. Performance of Q-learning with Linear Function Approximation: Stability and Finite-Time Analysis. *arXiv preprint arXiv:1905.11425*, 2019.
- [16] Wenhan Dai, Yi Gai, Bhaskar Krishnamachari, and Qing Zhao. The Non-Bayesian Restless Multi-Armed Bandit: A Case of Near-Logarithmic Regret. In *Proc. of IEEE ICASSP*, 2011.
- [17] Gal Dalal, Balazs Szorenyi, and Gagan Thoppe. A tale of two-timescale reinforcement learning with the tightest finite-time bound. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 3701–3708, 2020.
- [18] 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.
- [19] Thinh T Doan. Nonlinear two-time-scale stochastic approximation: Convergence and finite-time performance. *arXiv preprint arXiv:2011.01868*, 2020.

- [20] Thinh T Doan. Finite-time convergence rates of nonlinear two-time-scale stochastic approximation under markovian noise. *arXiv preprint arXiv:2104.01627*, 2021.
- [21] Thinh T Doan and Justin Romberg. Linear Two-Time-Scale Stochastic Approximation A Finite-Time Analysis. In *Proc. of Allerton*, 2019.
- [22] Jianqing Fan, Zhaoran Wang, Yuchen Xie, and Zhuoran Yang. A theoretical analysis of deep q-learning. In *Learning for Dynamics and Control*, pages 486–489. PMLR, 2020.
- [23] Jing Fu, Yoni Nazarathy, Sarat Moka, and Peter G Taylor. Towards q-learning the whittle index for restless bandits. In *2019 Australian & New Zealand Control Conference (ANZCC)*, pages 249–254. IEEE, 2019.
- [24] 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.
- [25] Young Hun Jung and Ambuj Tewari. Regret Bounds for Thompson Sampling in Episodic Restless Bandit Problems. *Proc. of NeurIPS*, 2019.
- [26] Maxim Kaledin, Eric Moulines, Alexey Naumov, Vladislav Tadic, and Hoi-To Wai. Finite time analysis of linear two-timescale stochastic approximation with markovian noise. In *Conference on Learning Theory*, pages 2144–2203. PMLR, 2020.
- [27] Jackson A Killian, Arpita Biswas, Sanket Shah, and Milind Tambe. Q-learning lagrange policies for multi-action restless bandits. In *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining*, pages 871–881, 2021.
- [28] Jackson A Killian, Andrew Perrault, and Milind Tambe. Beyond" To Act or Not to Act": Fast Lagrangian Approaches to General Multi-Action Restless Bandits. In *Proc.of AAMAS*, 2021.
- [29] Vijay R Konda and John N Tsitsiklis. Actor-Critic Algorithms. In *Proc. of NIPS*, 2000.
- [30] 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.
- [31] David A Levin and Yuval Peres. *Markov chains and mixing times*, volume 107. American Mathematical Soc., 2017.
- [32] Haoyang Liu, Keqin Liu, and Qing Zhao. Logarithmic Weak Regret of Non-Bayesian Restless Multi-Armed Bandit. In *Proc of IEEE ICASSP*, 2011.
- [33] Haoyang Liu, Keqin Liu, and Qing Zhao. Learning in A Changing World: Restless Multi-Armed Bandit with Unknown Dynamics. *IEEE Transactions on Information Theory*, 59(3):1902–1916, 2012.
- [34] Aditya Mate, Jackson Killian, Haifeng Xu, Andrew Perrault, and Milind Tambe. Collapsing bandits and their application to public health intervention. *Advances in Neural Information Processing Systems*, 33:15639–15650, 2020.
- [35] Aditya Mate, Andrew Perrault, and Milind Tambe. Risk-Aware Interventions in Public Health: Planning with Restless Multi-Armed Bandits. In *Proc.of AAMAS*, 2021.
- [36] Francisco S Melo, Sean P Meyn, and M Isabel Ribeiro. An analysis of reinforcement learning with function approximation. In *Proceedings of the 25th international conference on Machine learning*, pages 664–671, 2008.
- [37] Rahul Meshram, Aditya Gopalan, and D Manjunath. Optimal recommendation to users that react: Online learning for a class of pomdps. In *2016 IEEE 55th Conference on Decision and Control (CDC)*, pages 7210–7215. IEEE, 2016.
- [38] Abdelkader Makkadem and Mariane Pelletier. Convergence rate and averaging of nonlinear two-time-scale stochastic approximation algorithms. *The Annals of Applied Probability*, 16(3):1671–1702, 2006.

- [39] Khaled Nakhleh, Santosh Ganji, Ping-Chun Hsieh, I Hou, Srinivas Shakkottai, et al. Neurwin: Neural whittle index network for restless bandits via deep rl. *Advances in Neural Information Processing Systems*, 34, 2021.
- [40] Khaled Nakhleh, I Hou, et al. Deeptop: Deep threshold-optimal policy for mdps and rmabs. *arXiv preprint arXiv:2209.08646*, 2022.
- [41] Ronald Ortner, Daniil Ryabko, Peter Auer, and Rémi Munos. Regret Bounds for Restless Markov Bandits. In *Proc. of Algorithmic Learning Theory*, 2012.
- [42] Tejas Pagare, Vivek Borkar, and Konstantin Avrachenkov. Full gradient deep reinforcement learning for average-reward criterion. *arXiv preprint arXiv:2304.03729*, 2023.
- [43] Christos H Papadimitriou and John N Tsitsiklis. The Complexity of Optimal Queueing Network Control. In *Proc. of IEEE Conference on Structure in Complexity Theory*, 1994.
- [44] Martin L Puterman. *Markov Decision Processes: Discrete Stochastic Dynamic Programming*. John Wiley & Sons, 1994.
- [45] Guannan Qu and Adam Wierman. Finite-Time Analysis of Asynchronous Stochastic Approximation and Q -Learning. In *Proc. of COLT*, 2020.
- [46] Hiteshi Sharma, Mehdi Jafarnia-Jahromi, and Rahul Jain. Approximate relative value learning for average-reward continuous state mdps. In *Uncertainty in Artificial Intelligence*, pages 956–964. PMLR, 2020.
- [47] 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.
- [48] Wesley Suttle, Kaiqing Zhang, Zhuoran Yang, Ji Liu, and David Kraemer. Reinforcement Learning for Cost-Aware Markov Decision Processes. In *Proc. of ICML*, 2021.
- [49] Cem Tekin and Mingyan Liu. Online Learning of Rested and Restless Bandits. *IEEE Transactions on Information Theory*, 58(8):5588–5611, 2012.
- [50] John N Tsitsiklis and Benjamin Van Roy. Average Cost Temporal-Difference Learning. *Automatica*, 35(11):1799–1808, 1999.
- [51] Yi Wan, Abhishek Naik, and Richard S Sutton. Learning and Planning in Average-Reward Markov Decision Processes. In *Proc. of ICML*, 2021.
- [52] Siwei Wang, Longbo Huang, and John Lui. Restless-UCB, an Efficient and Low-complexity Algorithm for Online Restless Bandits. In *Proc. of NeurIPS*, 2020.
- [53] Richard R Weber and Gideon Weiss. On An Index Policy for Restless Bandits. *Journal of applied probability*, pages 637–648, 1990.
- [54] Chen-Yu Wei, Mehdi Jafarnia Jahromi, Haipeng Luo, Hiteshi Sharma, and Rahul Jain. Model-free Reinforcement Learning in Infinite-Horizon Average-Reward Markov Decision Processes. In *Proc. of ICML*, 2020.
- [55] Peter Whittle. Restless Bandits: Activity Allocation in A Changing World. *Journal of applied probability*, pages 287–298, 1988.
- [56] Guojun Xiong, Jian Li, and Rahul Singh. Reinforcement Learning Augmented Asymptotically Optimal Index Policies for Finite-Horizon Restless Bandits. In *Proc. of AAAI*, 2022.
- [57] Guojun Xiong, Shufan Wang, Jian Li, and Rahul Singh. Whittle index based q-learning for wireless edge caching with linear function approximation. *arXiv preprint arXiv:2202.13187*, 2022.
- [58] Pan Xu and Quanquan Gu. A finite-time analysis of q-learning with neural network function approximation. In *International Conference on Machine Learning*, pages 10555–10565. PMLR, 2020.

- 485 [59] Zhuoran Yang, Yongxin Chen, Mingyi Hong, and Zhaoran Wang. Provably Global Convergence
486 of Actor-Critic: A Case for Linear Quadratic Regulator with Ergodic Cost. In *Proc. of NeurIPS*,
487 2019.
- 488 [60] Zhe Yu, Yunjian Xu, and Lang Tong. Deadline Scheduling as Restless Bandits. *IEEE Transac-*
489 *tions on Automatic Control*, 63(8):2343–2358, 2018.
- 490 [61] Shangdong Zhang, Yi Wan, Richard S Sutton, and Shimon Whiteson. Average-Reward Off-
491 Policy Policy Evaluation with Function Approximation. *arXiv preprint arXiv:2101.02808*,
492 2021.
- 493 [62] Sheng Zhang, Zhe Zhang, and Siva Theja Maguluri. Finite Sample Analysis of Average-Reward
494 TD Learning and Q -Learning. *Proc. of NeurIPS*, 2021.
- 495 [63] Shaofeng Zou, Tengyu Xu, and Yingbin Liang. Finite-sample analysis for sarsa with linear
496 function approximation. *Advances in neural information processing systems*, 32, 2019.

497 A Related Work

498 **Online Restless Bandits.** The online RMAB setting, where the underlying MDPs are unknown,
 499 has been gaining attention, e.g., [16, 32, 33, 49, 41, 25]. However, these methods do not exploit the
 500 special structure available in the problem and contend directly with an extremely high dimensional
 501 state-action space yielding the algorithms to be too slow to be useful. Recently, RL based algorithms
 502 have been developed [23, 3, 8, 27, 52, 56], to explore the problem structure through index policies.
 503 For instance, [23] proposed a Q-learning algorithm for Whittle index under the discounted setting,
 504 which lacks of convergence guarantees. [8] approximated Whittle index using the difference of
 505 $Q(s, 1) - Q(s, 0)$ for any state s , which is not guaranteed to converge to the true Whittle index
 506 in general scenarios. To our best knowledge, the Q-Whittle in (4)-(5) proposed by [3] is the
 507 first algorithm with a rigorous asymptotic analysis. Therefore, [23, 3, 8, 27] lacked finite-time
 508 performance analysis and multi-timescale stochastic approximation algorithms usually suffer from
 509 slow convergence.

510 [52, 56] designed model-based low-complexity policy but is constrained to either a specific Markovian
 511 model or depends on a simulator for a finite-horizon setting which cannot be directly applied here.
 512 Latter on, [57] showed the finite-time convergence performance under the Q-Whittle setting of
 513 [3] with linear function approximation. However, the underlying assumption in [3, 57] is that data
 514 samples are drawn i.i.d per iteration. This is often not the case in practice since data samples of
 515 Q-learning are drawn according to the underlying Markov decision process. Till now, the finite-time
 516 convergence rate of Q-Whittle under the more challenging Markovian setting remains to be an
 517 open problem. Though [42] proposed a novel DQN method and applied it to Whittle index learning,
 518 it lacks of theoretical convergence analysis. To our best knowledge, our work is the first to study
 519 low-complexity model-free Q-learning for RMAB with neural network function approximation and
 520 provide a finite-time performance guarantee.

521 **Two-Timescale Stochastic Approximation.** The theoretical understanding of average-reward rein-
 522 forcement learning (RL) methods is limited. Most existing results focus on asymptotic convergence
 523 [50, 1, 51, 61], or finite-time performance guarantee for discounted Q-learning [15, 45, 14]. How-
 524 ever, the analysis of average-reward RL algorithms is known to be more challenging than their
 525 discounted-reward counterparts [62, 54]. In particular, our Neural-Q-Whittle follows the 2TSA
 526 scheme [11, 29, 7]. The standard technique for analyzing 2TSA is via the ODE method to prove
 527 asymptotic convergence [9]. Building off the importance of asymptotic results, recent years have
 528 witnessed a focus shifted to non-asymptotic, finite-time analysis of 2TSA [24, 21, 19, 59]. The
 529 closest work is [19], which characterized the convergence rate for a general non-linear 2TSA with
 530 i.i.d. noise. We generalize this result to provide a finite-time analysis of our Neural-Q-Whittle
 531 with Markovian noise. In addition, existing finite-time analysis, e.g., sample complexity [62] and
 532 regret [54] of Q-learning with average reward focus on a single-timescale SA, and hence cannot
 533 be directly applied to our Neural-Q-Whittle. Finally, existing Q-learning with linear function
 534 approximation [36, 6, 63] and neural network function approximation [13, 58] requires an additional
 535 projection step onto a bounded set related to the unknown stationary distribution of the underlying
 536 MDPs, or focuses on a single-timescale SA [15].

537 B Review on Whittle Index Policy

538 Whittle index policy addresses the intractable issue of RMAB through decomposition. In each
 539 round t , it first calculates the Whittle index for each arm n independently only based on its current
 540 state $s_n(t)$, and then the Whittle index policy simply selects the K arms with the highest indices to
 541 activate. Following Whittle’s approach[55], we can consider a system with only one arm due to the
 542 decomposition, and the Lagrangian is expressed as

$$L(\pi, \lambda) = \liminf_{T \rightarrow \infty} \frac{1}{T} \mathbb{E}_\pi \sum_{t=1}^T \left\{ r(t) + \lambda \left(1 - a(t) \right) \right\}, \quad (25)$$

543 where λ is the Lagrangian multiplier (or the subsidy for selecting passive action). For a particular λ ,
 544 the optimal activation policy can be expressed by a set of states in which it would activate this arm,
 545 which is denoted $D(\lambda)$.

Definition 3 (Indexability). We denote $D(\lambda)$ as the set of states S for which the optimal action for the arm is to choose a passive action, i.e., $A = 0$. Then the arm is said to be indexable if $D(\lambda)$ increases with λ , i.e., if $\lambda > \lambda'$, then $D(\lambda) \supseteq D(\lambda')$.

Following the indexability property, the Whittle index in a particular state S is defined as follows.

Definition 4 (Whittle Index). The Whittle index in state S for the indexable arm is the smallest value of the Lagrangian multiplier λ such that the optimal policy at state S is indifferent towards actions $A = 0$ and $A = 1$. We denote such a Whittle index as $\lambda(S)$ satisfying $\lambda(S) := \inf_{\lambda \geq 0} \{S \in D(\lambda)\}$.

Definition 5 (Whittle index policy). Whittle index policy is a controlled policy which activates the K arms with the highest whittle index $\lambda_i(S_i(t))$ at each time slot t .

C Proof of Lemmas for “Step 2: Convergence Rate of $\hat{M}(\theta_k, \lambda_k)$ in (16)”

C.1 Proof of Lemma 1

Proof. Recall that

$$f_0(\theta; \phi(s, a)) = \frac{1}{\sqrt{m}} \sum_{r=1}^m b_r \mathbb{1}\{\mathbf{w}_{r,0}^\top \phi(s, a) > 0\} \mathbf{w}_r^\top \phi(s, a).$$

Thus we denote $\nabla_{\theta} f_0(\theta; \phi(s, a))$ as

$$\begin{aligned} \nabla_{\theta} f_0(\theta; \phi(s, a)) &:= \left[\frac{1}{\sqrt{m}} b_1 \mathbb{1}\{\mathbf{w}_{1,0}^\top \phi(s, a) > 0\} \phi(s, a)^\top, \dots, \right. \\ &\quad \left. \frac{1}{\sqrt{m}} b_m \mathbb{1}\{\mathbf{w}_{m,0}^\top \phi(s, a) > 0\} \phi(s, a)^\top \right]^\top. \end{aligned} \quad (26)$$

Since $\|\phi(s, a)\| \leq 1, \forall s \in \mathcal{S}, a \in \mathcal{A}$ and the fact that $b_r, \forall r \in [m]$ is uniformly initialized as 1 and -1 , we have $\|\nabla_{\theta} f_0(\theta; \phi(s, a))\| \leq 1$.

Therefore, we have the following inequality for any parameter pairs (θ_1, λ_1) and (θ_2, λ_2) with $X = (s, a, s') \in \mathcal{X}$,

$$\begin{aligned} &\|h_0(X, \theta_1, \lambda_1) - h_0(X, \theta_2, \lambda_2)\| \\ &= \left\| \nabla_{\theta} f_0(\theta_1; \phi(s, a)) \left[r(s, a) + (1-a)\lambda_1 - I_0(\theta_1) + \max_{a_1} f_0(\theta_1; \phi(s', a_1)) - f_0(\theta_1; \phi(s, a)) \right] \right. \\ &\quad \left. - \nabla_{\theta} f_0(\theta_2; \phi(s, a)) \left[r(s, a) + (1-a)\lambda_2 - I_0(\theta_2) + \max_{a_2} f_0(\theta_2; \phi(s', a_2)) - f_0(\theta_2; \phi(s, a)) \right] \right\| \\ &\stackrel{(a_1)}{=} \left\| \nabla_{\theta} f_0(\theta_1; \phi(s, a)) \left[(1-a)(\lambda_1 - \lambda_2) + I_0(\theta_2) - I_0(\theta_1) + f_0(\theta_2; \phi(s, a)) - f_0(\theta_1; \phi(s, a)) \right] \right. \\ &\quad \left. + \max_{a_1} \left(f_0(\theta_1; \phi(s', a_1)) - \max_{a_2} f_0(\theta_2; \phi(s', a_2)) \right) \right\| \\ &\stackrel{(a_2)}{\leq} \|(1-a)(\lambda_1 - \lambda_2)\| + \|f_0(\theta_2; \phi(s, a)) - f_0(\theta_1; \phi(s, a))\| \\ &\quad + \left\| \frac{1}{2S} \sum_{\tilde{s} \in \mathcal{S}} f_0(\theta_2; \phi(\tilde{s}, 0)) - f_0(\theta_1; \phi(\tilde{s}, 0)) + f_0(\theta_2; \phi(\tilde{s}, 1)) - f_0(\theta_1; \phi(\tilde{s}, 1)) \right\| \\ &\quad + \left\| \max_{a_1} \left(f_0(\theta_1; \phi(s', a_1)) - \max_{a_2} f_0(\theta_2; \phi(s', a_2)) \right) \right\| \\ &\stackrel{(a_3)}{\leq} \|(1-a)(\lambda_1 - \lambda_2)\| + \|\nabla_{\theta} f_0(\theta_1; \phi(s, a))(\theta_2 - \theta_1)\| \\ &\quad + \left\| \frac{1}{2S} \sum_{\tilde{s} \in \mathcal{S}} \nabla_{\theta} f_0(\theta_1; \phi(\tilde{s}, 0))(\theta_2 - \theta_1) + \nabla_{\theta} f_0(\theta_1; \phi(\tilde{s}, 1))(\theta_2 - \theta_1) \right\| \\ &\quad + \left\| \max_{a_1} \left(f_0(\theta_1; \phi(s', a_1)) - \max_{a_2} f_0(\theta_2; \phi(s', a_2)) \right) \right\| \\ &\stackrel{(a_4)}{\leq} \|(\lambda_1 - \lambda_2)\| + 2\|\theta_1 - \theta_2\| + \left\| \max_{a_1} \left(f_0(\theta_1; \phi(s', a_1)) - \max_{a_2} f_0(\theta_2; \phi(s', a_2)) \right) \right\| \end{aligned}$$

$$\begin{aligned}
&\stackrel{(a_5)}{\leq} \|(\lambda_1 - \lambda_2)\| + 2\|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2\| + \left\| \max_{a'} f_0(\boldsymbol{\theta}_1; \boldsymbol{\phi}(s', a')) - f_0(\boldsymbol{\theta}_2; \boldsymbol{\phi}(s', a')) \right\| \\
&\stackrel{(a_6)}{\leq} \|(\lambda_1 - \lambda_2)\| + 3\|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2\|.
\end{aligned}$$

Specifically, (a_1) holds due to the fact that $\nabla_{\boldsymbol{\theta}} f_0(\boldsymbol{\theta}_1; \boldsymbol{\phi}(s, a)) = \nabla_{\boldsymbol{\theta}} f_0(\boldsymbol{\theta}_2; \boldsymbol{\phi}(s, a))$ as in (26). Since

$$I_0(\boldsymbol{\theta}_k) = \frac{1}{2S} \sum_{\tilde{s} \in \mathcal{S}} \left[f_0(\boldsymbol{\theta}_k; \boldsymbol{\phi}(\tilde{s}, 0)) + f_0(\boldsymbol{\theta}_k; \boldsymbol{\phi}(\tilde{s}, 1)) \right],$$

562 (a_2) is due to the fact that $\|\mathbf{x} + \mathbf{y}\| \leq \|\mathbf{x}\| + \|\mathbf{y}\|, \forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^{md}$ and $\|\mathbf{x} \cdot \mathbf{y}\| \leq \|\mathbf{x}\| \cdot \|\mathbf{y}\|, \forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^{md}$
563 and $\|\boldsymbol{\phi}(s, a)\| \leq 1, \forall s, a$. (a_3) holds since

$$f_0(\boldsymbol{\theta}_2; \boldsymbol{\phi}(s, a)) - f_0(\boldsymbol{\theta}_1; \boldsymbol{\phi}(s, a)) = \nabla_{\boldsymbol{\theta}} f_0(\boldsymbol{\theta}_1; \boldsymbol{\phi}(s, a))(\boldsymbol{\theta}_2 - \boldsymbol{\theta}_1), \forall s \in \mathcal{S}, a \in \mathcal{A}. \quad (27)$$

564 (a_4) holds for the same reason as (a_2) . (a_5) is due to the fact that

$$\begin{aligned}
\left\| \max_{a'} f_0(\boldsymbol{\theta}_1; \boldsymbol{\phi}(s', a')) - f_0(\boldsymbol{\theta}_2; \boldsymbol{\phi}(s', a')) \right\| &\leq \max \left(\left\| \max_{a'} f_0(\boldsymbol{\theta}_1; \boldsymbol{\phi}(s', a')) - f_0(\boldsymbol{\theta}_2; \boldsymbol{\phi}(s', a')) \right\|, \right. \\
&\quad \left. \left\| \min_{a'} f_0(\boldsymbol{\theta}_1; \boldsymbol{\phi}(s', a')) - f_0(\boldsymbol{\theta}_2; \boldsymbol{\phi}(s', a')) \right\| \right). \quad (28)
\end{aligned}$$

565 (a_6) holds for the same reason as (a_3) and (a_4) . \square

566 C.2 Proof of Lemma 2

567 *Proof.* Since $g_0(\cdot)$ is irrelevant with X and λ , in the following, we write $g_0(X, \boldsymbol{\theta}, \lambda)$ with $g_0(\lambda)$
568 interchangeably. For any $\boldsymbol{\theta}_1 \in \mathbb{R}^{md}$ and $\boldsymbol{\theta}_2 \in \mathbb{R}^{md}$, we have

$$\begin{aligned}
&\|g_0(\boldsymbol{\theta}_1) - g_0(\boldsymbol{\theta}_2)\| \\
&= \|f_0(\boldsymbol{\theta}_1; \boldsymbol{\phi}(s, 1)) - f_0(\boldsymbol{\theta}_1; \boldsymbol{\phi}(s, 0)) - f_0(\boldsymbol{\theta}_2; \boldsymbol{\phi}(s, 1)) + f_0(\boldsymbol{\theta}_2; \boldsymbol{\phi}(s, 0))\| \\
&\leq \|f_0(\boldsymbol{\theta}_1; \boldsymbol{\phi}(s, 1)) - f_0(\boldsymbol{\theta}_2; \boldsymbol{\phi}(s, 1))\| + \|f_0(\boldsymbol{\theta}_1; \boldsymbol{\phi}(s, 0)) - f_0(\boldsymbol{\theta}_2; \boldsymbol{\phi}(s, 0))\| \\
&= \|\nabla_{\boldsymbol{\theta}} f_0(\boldsymbol{\theta}_1; \boldsymbol{\phi}(s, 1))(\boldsymbol{\theta}_2 - \boldsymbol{\theta}_1)\| + \|\nabla_{\boldsymbol{\theta}} f_0(\boldsymbol{\theta}_1; \boldsymbol{\phi}(s, 0))(\boldsymbol{\theta}_2 - \boldsymbol{\theta}_1)\| \\
&\leq \|\nabla_{\boldsymbol{\theta}} f_0(\boldsymbol{\theta}_1; \boldsymbol{\phi}(s, 1))\| \cdot \|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2\| + \|\nabla_{\boldsymbol{\theta}} f_0(\boldsymbol{\theta}_1; \boldsymbol{\phi}(s, 0))\| \cdot \|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2\| \\
&\leq 2\|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2\|,
\end{aligned}$$

569 where the first inequality is due to the fact that $\|\mathbf{x} + \mathbf{y}\| \leq \|\mathbf{x}\| + \|\mathbf{y}\|, \forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^{md}$, the second
570 inequality holds due to $\|\mathbf{x} \cdot \mathbf{y}\| \leq \|\mathbf{x}\| \cdot \|\mathbf{y}\|, \forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^{md}$, and the last inequality holds since
571 $\|\nabla_{\boldsymbol{\theta}} f_0(\boldsymbol{\theta}_1; \boldsymbol{\phi}(s, a))\| \leq 1, \forall s \in \mathcal{S}, a \in \mathcal{A}$. \square

572 C.3 Proof of Lemma 3

573 *Proof.* For any $\boldsymbol{\theta}_1 \in \mathbb{R}^{md}$ and $\boldsymbol{\theta}_2 \in \mathbb{R}^{md}$, we have

$$\begin{aligned}
&\|y_0(\boldsymbol{\theta}_1) - y_0(\boldsymbol{\theta}_2)\| \\
&= \left\| r(s, 1) - r(s, 0) + \sum_{s'} P(s'|s, 1) \max_a f_0(\boldsymbol{\theta}_1; \boldsymbol{\phi}(s', a)) - \sum_{s'} P(s'|s, 0) \max_a f_0(\boldsymbol{\theta}_1; \boldsymbol{\phi}(s', a)) \right. \\
&\quad \left. - r(s, 1) - r(s, 0) + \sum_{s'} P(s'|s, 1) \max_a f_0(\boldsymbol{\theta}_2; \boldsymbol{\phi}(s', a)) - \sum_{s'} P(s'|s, 0) \max_a f_0(\boldsymbol{\theta}_2; \boldsymbol{\phi}(s', a)) \right\| \\
&= \left\| \sum_{s'} P(s'|s, 1) \max_a f_0(\boldsymbol{\theta}_1; \boldsymbol{\phi}(s', a)) - \sum_{s'} P(s'|s, 1) \max_a f_0(\boldsymbol{\theta}_2; \boldsymbol{\phi}(s', a)) \right. \\
&\quad \left. - \sum_{s'} P(s'|s, 0) \max_a f_0(\boldsymbol{\theta}_1; \boldsymbol{\phi}(s', a)) + \sum_{s'} P(s'|s, 0) \max_a f_0(\boldsymbol{\theta}_2; \boldsymbol{\phi}(s', a)) \right\| \\
&\leq \left\| \sum_{s'} P(s'|s, 1) \max_a f_0(\boldsymbol{\theta}_1; \boldsymbol{\phi}(s', a)) - \sum_{s'} P(s'|s, 1) \max_a f_0(\boldsymbol{\theta}_2; \boldsymbol{\phi}(s', a)) \right\| \\
&\quad + \left\| \sum_{s'} P(s'|s, 0) \max_a f_0(\boldsymbol{\theta}_1; \boldsymbol{\phi}(s', a)) - \sum_{s'} P(s'|s, 0) \max_a f_0(\boldsymbol{\theta}_2; \boldsymbol{\phi}(s', a)) \right\| \\
&\leq 2\|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2\|,
\end{aligned}$$

574 with the last inequality holds due to (27) and (28). \square

575 **C.4 Proof of Lemma 4**

576 *Proof.* 1) We first show that there exists a constant $\mu_1 > 0$ such that $\mathbb{E}[\hat{\boldsymbol{\theta}}^\top h_0(X, \boldsymbol{\theta}, \lambda)] \leq -\mu_1 \|\hat{\boldsymbol{\theta}}\|^2$.
 577 According to the definition of $\boldsymbol{\theta}_0^*$ given in Definition 1, $\mathbb{E}[h_0(X, \boldsymbol{\theta}_0^*, y_0(\boldsymbol{\theta}_0^*))] = 0$. Hence, we have

$$\begin{aligned}
 & \mathbb{E} \left[\hat{\boldsymbol{\theta}}^\top (h_0(X, \boldsymbol{\theta}, \lambda) - h_0(X, \boldsymbol{\theta}_0^*, y_0(\boldsymbol{\theta}_0^*))) \right] \\
 &= \hat{\boldsymbol{\theta}}^\top \mathbb{E} [h_0(X, \boldsymbol{\theta}, \lambda) - h_0(X, \boldsymbol{\theta}_0^*, y_0(\boldsymbol{\theta}_0^*))] \\
 &= \hat{\boldsymbol{\theta}}^\top \mathbb{E} \left[\nabla_{\boldsymbol{\theta}} f_0(\boldsymbol{\theta}; \boldsymbol{\phi}(s, a)) \left[r(s, a) + (1-a)\lambda - I_0(\boldsymbol{\theta}) + \max_{a_1} f_0(\boldsymbol{\theta}; \boldsymbol{\phi}(s', a_1)) - f_0(\boldsymbol{\theta}; \boldsymbol{\phi}(s, a)) \right] \right. \\
 &\quad \left. - \nabla_{\boldsymbol{\theta}} f_0(\boldsymbol{\theta}_0^*; \boldsymbol{\phi}(s, a)) \left[r(s, a) + (1-a)y_0(\boldsymbol{\theta}_0^*) - I_0(\boldsymbol{\theta}_0^*) + \max_{a_2} f_0(\boldsymbol{\theta}_0^*; \boldsymbol{\phi}(s', a_2)) - f_0(\boldsymbol{\theta}_0^*; \boldsymbol{\phi}(s, a)) \right] \right] \\
 &\stackrel{(b_1)}{=} \hat{\boldsymbol{\theta}}^\top \mathbb{E} \left[\nabla_{\boldsymbol{\theta}} f_0(\boldsymbol{\theta}; \boldsymbol{\phi}(s, a)) \left[(1-a)(\lambda - y_0(\boldsymbol{\theta}_0^*)) + I_0(\boldsymbol{\theta}_0^*) - I_0(\boldsymbol{\theta}) + f_0(\boldsymbol{\theta}_0^*; \boldsymbol{\phi}(s, a)) - f_0(\boldsymbol{\theta}; \boldsymbol{\phi}(s, a)) \right. \right. \\
 &\quad \left. \left. + \max_{a_1} f_0(\boldsymbol{\theta}; \boldsymbol{\phi}(s', a_1)) - \max_{a_2} f_0(\boldsymbol{\theta}_0^*; \boldsymbol{\phi}(s', a_2)) \right] \right] \\
 &= \hat{\boldsymbol{\theta}}^\top \mathbb{E} \left[\nabla_{\boldsymbol{\theta}} f_0(\boldsymbol{\theta}; \boldsymbol{\phi}(s, a)) \left[\max_{a_1} f_0(\boldsymbol{\theta}; \boldsymbol{\phi}(s', a_1)) - \max_{a_2} f_0(\boldsymbol{\theta}_0^*; \boldsymbol{\phi}(s', a_2)) \right] \right] \\
 &\quad - \hat{\boldsymbol{\theta}}^\top \mathbb{E} \left[\nabla_{\boldsymbol{\theta}} f_0(\boldsymbol{\theta}; \boldsymbol{\phi}(s, a)) [I_0(\boldsymbol{\theta}) - I_0(\boldsymbol{\theta}_0^*)] \right] - \hat{\boldsymbol{\theta}}^\top \mathbb{E} \left[\nabla_{\boldsymbol{\theta}} f_0(\boldsymbol{\theta}; \boldsymbol{\phi}(s, a)) [f_0(\boldsymbol{\theta}; \boldsymbol{\phi}(s, a)) - f_0(\boldsymbol{\theta}_0^*; \boldsymbol{\phi}(s, a))] \right] \\
 &\quad + \hat{\boldsymbol{\theta}}^\top \mathbb{E} \left[\nabla_{\boldsymbol{\theta}} f_0(\boldsymbol{\theta}; \boldsymbol{\phi}(s, a)) [(1-a)(\lambda - y_0(\boldsymbol{\theta}_0^*))] \right] \\
 &\stackrel{(b_2)}{\leq} \hat{\boldsymbol{\theta}}^\top \mathbb{E} \left[\nabla_{\boldsymbol{\theta}} f_0(\boldsymbol{\theta}; \boldsymbol{\phi}(s, a)) \max_{a'} \left[f_0(\boldsymbol{\theta}; \boldsymbol{\phi}(s', a')) - f_0(\boldsymbol{\theta}_0^*; \boldsymbol{\phi}(s', a')) \right] \right] \\
 &\quad - \hat{\boldsymbol{\theta}}^\top \mathbb{E} \left[\nabla_{\boldsymbol{\theta}} f_0(\boldsymbol{\theta}; \boldsymbol{\phi}(s, a)) [I_0(\boldsymbol{\theta}) - I_0(\boldsymbol{\theta}_0^*)] \right] - \hat{\boldsymbol{\theta}}^\top \mathbb{E} \left[\nabla_{\boldsymbol{\theta}} f_0(\boldsymbol{\theta}; \boldsymbol{\phi}(s, a)) [f_0(\boldsymbol{\theta}; \boldsymbol{\phi}(s, a)) - f_0(\boldsymbol{\theta}_0^*; \boldsymbol{\phi}(s, a))] \right] \\
 &\quad + \hat{\boldsymbol{\theta}}^\top \mathbb{E} \left[\nabla_{\boldsymbol{\theta}} f_0(\boldsymbol{\theta}; \boldsymbol{\phi}(s, a)) [(1-a)(\lambda - y_0(\boldsymbol{\theta}_0^*))] \right] \\
 &\stackrel{(b_3)}{\leq} \hat{\boldsymbol{\theta}}^\top \mathbb{E} \left[\nabla_{\boldsymbol{\theta}} f_0(\boldsymbol{\theta}; \boldsymbol{\phi}(s, a)) \max_{a'} \left[f_0(\boldsymbol{\theta}; \boldsymbol{\phi}(s', a')) - f_0(\boldsymbol{\theta}_0^*; \boldsymbol{\phi}(s', a')) \right] \right] \\
 &\quad - \hat{\boldsymbol{\theta}}^\top \mathbb{E} \left[\nabla_{\boldsymbol{\theta}} f_0(\boldsymbol{\theta}; \boldsymbol{\phi}(s, a)) [I_0(\boldsymbol{\theta}) - I_0(\boldsymbol{\theta}_0^*)] \right] - \hat{\boldsymbol{\theta}}^\top \mathbb{E} \left[\nabla_{\boldsymbol{\theta}} f_0(\boldsymbol{\theta}; \boldsymbol{\phi}(s, a)) [f_0(\boldsymbol{\theta}; \boldsymbol{\phi}(s, a)) - f_0(\boldsymbol{\theta}_0^*; \boldsymbol{\phi}(s, a))] \right] \\
 &\stackrel{(b_4)}{=} \|\hat{\boldsymbol{\theta}}\|^2 \mathbb{E} \left[\nabla_{\boldsymbol{\theta}} f_0(\boldsymbol{\theta}; \boldsymbol{\phi}(s, a))^\top \nabla_{\boldsymbol{\theta}} f_0(\boldsymbol{\theta}; \boldsymbol{\phi}(s', \tilde{a})) \right] \\
 &\quad - \|\hat{\boldsymbol{\theta}}\|^2 \mathbb{E} \left[\nabla_{\boldsymbol{\theta}} f_0(\boldsymbol{\theta}; \boldsymbol{\phi}(s, a))^\top \left[\frac{1}{2S} \sum_{\tilde{s} \in S} \nabla_{\boldsymbol{\theta}} f_0(\boldsymbol{\theta}; \boldsymbol{\phi}(\tilde{s}, 0)) + f_0(\boldsymbol{\theta}; \boldsymbol{\phi}(\tilde{s}, 1)) \right] \right] \\
 &\quad - \|\hat{\boldsymbol{\theta}}\|^2 \mathbb{E} \left[\nabla_{\boldsymbol{\theta}} f_0(\boldsymbol{\theta}; \boldsymbol{\phi}(s, a))^\top \nabla_{\boldsymbol{\theta}} f_0(\boldsymbol{\theta}; \boldsymbol{\phi}(s, a)) \right]
 \end{aligned}$$

578 where (b_1) holds since $\nabla_{\boldsymbol{\theta}} f_0(\boldsymbol{\theta}; \boldsymbol{\phi}(s, a)) = \nabla_{\boldsymbol{\theta}} f_0(\boldsymbol{\theta}_0^*; \boldsymbol{\phi}(s, a))$ as in (26), (b_2) is due to the fact that
 579 $\max_{a_1} f_0(\boldsymbol{\theta}; \boldsymbol{\phi}(s', a_1)) - \max_{a_2} f_0(\boldsymbol{\theta}_0^*; \boldsymbol{\phi}(s', a_2)) \leq \max_{a'} [f_0(\boldsymbol{\theta}; \boldsymbol{\phi}(s', a')) - f_0(\boldsymbol{\theta}_0^*; \boldsymbol{\phi}(s', a'))]$,
 580 and (b_3) holds due to the fact that $\hat{\boldsymbol{\theta}}^\top \mathbb{E} [\nabla_{\boldsymbol{\theta}} f_0(\boldsymbol{\theta}; \boldsymbol{\phi}(s, a)) [(1-a)(\lambda - y_0(\boldsymbol{\theta}_0^*))]] \leq 0$ since a larger
 581 Whittle index λ will choose the action $a = 1$. Notice that the \tilde{a} in (b_4) represents the action a' which
 582 maximizes $f_0(\boldsymbol{\theta}; \boldsymbol{\phi}(s', a')) - f_0(\boldsymbol{\theta}_0^*; \boldsymbol{\phi}(s', a'))$. Due to the definition of $\nabla_{\boldsymbol{\theta}} f_0(\boldsymbol{\theta}; \boldsymbol{\phi}(s, a))$ in (26), we
 583 show that $\mathbb{E}[\hat{\boldsymbol{\theta}}^\top h_0(X, \boldsymbol{\theta}, \lambda)] \leq 0$.

584 2) Next, we show that there exists a constant $\mu_2 > 0$ such that $\mathbb{E}[\hat{\lambda} g_0(X, \boldsymbol{\theta}, \lambda)] \leq -\mu_2 \|\hat{\lambda}\|^2$.
 585 According to the definition of $g_0(\boldsymbol{\theta})$, i.e., $g_0(\boldsymbol{\theta}) := f_0(\boldsymbol{\theta}; \boldsymbol{\phi}(s, 1)) - f_0(\boldsymbol{\theta}; \boldsymbol{\phi}(s, 0))$. Since $y_0(\boldsymbol{\theta})$
 586 is the solution of λ such that $f_0(\boldsymbol{\theta}; \boldsymbol{\phi}(s, 1)) = f_0(\boldsymbol{\theta}; \boldsymbol{\phi}(s, 0))$, the signs of $\hat{\lambda} := \lambda - y_0(\boldsymbol{\theta})$ and
 587 $f_0(\boldsymbol{\theta}; \boldsymbol{\phi}(s, 1)) - f_0(\boldsymbol{\theta}; \boldsymbol{\phi}(s, 0))$ are always opposite. Hence, we have $\mathbb{E}[\hat{\lambda} g_0(X, \boldsymbol{\theta}, \lambda)] \leq 0$, which
 588 completes the proof.

589 □

590 C.5 Proof of Lemma 5

591 *Proof.* Under Lemma 1, we have

$$\|h_0(X, \boldsymbol{\theta}, \lambda) - h_0(X, \boldsymbol{\theta}^*, \lambda^*)\| \leq 3\|\boldsymbol{\theta} - \boldsymbol{\theta}^*\| + \|\lambda - \lambda^*\|. \quad (29)$$

592 Let $L = \max(3, \max_X h_0(X, \boldsymbol{\theta}^*, \lambda^*))$, then according to (29), we have

$$\|h_0(X, \boldsymbol{\theta}, \lambda)\| \leq L(\|\boldsymbol{\theta} - \boldsymbol{\theta}^*\| + \|\lambda - \lambda^*\| + 1).$$

593 Denote $h_0^i(X, \boldsymbol{\theta}, \lambda)$ as the i -th element of $h_0(X, \boldsymbol{\theta}, \lambda)$. Following [15], we can show that $\boldsymbol{\theta} \in \mathbb{R}^{md}$,
 594 $\lambda \in \mathbb{R}^1$, and $x \in \mathcal{X}$,

$$\begin{aligned} & \|\mathbb{E}[h_0(X_k, \boldsymbol{\theta}, \lambda) | X_0 = x] - \mathbb{E}_\mu[h_0(X, \boldsymbol{\theta}, \lambda)]\| \\ & \leq \sum_{i=1}^{md} |\mathbb{E}[h_i(X_k, \boldsymbol{\theta}, \lambda) | X_0 = x] - \mathbb{E}_\mu[h_0^i(X, \boldsymbol{\theta}, \lambda)]| \\ & \leq 2L(\|\boldsymbol{\theta} - \boldsymbol{\theta}^*\| + \|\lambda - \lambda^*\| + 1) \sum_{i=1}^{md} \left| \mathbb{E} \left[\frac{h_0^i(X_k, \boldsymbol{\theta}, \lambda)}{2L(\|\boldsymbol{\theta} - \boldsymbol{\theta}^*\| + \|\lambda - \lambda^*\| + 1)} \middle| X_0 = x \right] \right. \\ & \quad \left. - \mathbb{E}_\mu \left[\frac{h_0^i(X, \boldsymbol{\theta}, \lambda)}{2L(\|\boldsymbol{\theta} - \boldsymbol{\theta}^*\| + \|\lambda - \lambda^*\| + 1)} \right] \right| \\ & \leq 2L(\|\boldsymbol{\theta} - \boldsymbol{\theta}^*\| + \|\lambda - \lambda^*\| + 1) md C \rho^k, \end{aligned}$$

595 where the last inequality holds due to Assumption 1. To guarantee $2L(\|\boldsymbol{\theta} - \boldsymbol{\theta}^*\| + \|\lambda - \lambda^*\| +$
 596 $1) md C \rho^k \leq \delta(\|\boldsymbol{\theta} - \boldsymbol{\theta}^*\| + \|\lambda - \lambda^*\| + 1)$, we have

$$\tau_\delta \leq \frac{\log(1/\delta) + \log(2LCmd)}{\log(1/\rho)},$$

597 which completes the proof.

598 □

599 C.6 Proof of Lemma 6

600 *Proof.* Based on the definition of $M(\boldsymbol{\theta}_k, \lambda_k)$ in (15), we have

$$\begin{aligned} M(\boldsymbol{\theta}_k, \lambda_k) &:= \frac{\eta_k}{\alpha_k} \|\boldsymbol{\theta}_k - \boldsymbol{\theta}^*\|^2 + \|\lambda_k - y(\boldsymbol{\theta}_k)\|^2 \\ &= \frac{\eta_k}{\alpha_k} \|\boldsymbol{\theta}_k - \boldsymbol{\theta}_0^* + \boldsymbol{\theta}_0^* - \boldsymbol{\theta}^*\|^2 + \|\lambda_k - y_0(\boldsymbol{\theta}_k) + y_0(\boldsymbol{\theta}_k) - y(\boldsymbol{\theta}_k)\|^2 \\ &\leq \frac{2\eta_k}{\alpha_k} (\|\boldsymbol{\theta}_k - \boldsymbol{\theta}_0^*\|^2 + \|\boldsymbol{\theta}_0^* - \boldsymbol{\theta}^*\|^2) + 2(\|\lambda_k - y_0(\boldsymbol{\theta}_k)\|^2 + \|y_0(\boldsymbol{\theta}_k) - y(\boldsymbol{\theta}_k)\|^2) \\ &= 2\hat{M}(\boldsymbol{\theta}_k, \lambda_k) + \frac{2\eta_k}{\alpha_k} \|\boldsymbol{\theta}_0^* - \boldsymbol{\theta}^*\|^2 + 2\|y_0(\boldsymbol{\theta}_k) - y(\boldsymbol{\theta}_k)\|^2 \\ &\leq 2\hat{M}(\boldsymbol{\theta}_k, \lambda_k) + \frac{2\eta_k c_0^2}{\alpha_k} \|\text{span}(f_0(\boldsymbol{\theta}_0^*) - f(\boldsymbol{\theta}^*))\|^2 + 2\|y_0(\boldsymbol{\theta}_k) - y(\boldsymbol{\theta}_k)\|^2, \quad (30) \end{aligned}$$

601 where the first inequality holds based on $\|\mathbf{x} + \mathbf{y}\|^2 \leq 2\|\mathbf{x}\|^2 + 2\|\mathbf{y}\|^2$, and the second inequality
 602 holds based on Assumption 3. Next, we bound $\|\text{span}(f_0(\boldsymbol{\theta}_0^*) - f(\boldsymbol{\theta}^*))\|$ as follows

$$\begin{aligned} \|\text{span}(f_0(\boldsymbol{\theta}_0^*) - f(\boldsymbol{\theta}^*))\| &= \|\text{span}(f_0(\boldsymbol{\theta}_0^*) - \Pi_{\mathcal{F}} f(\boldsymbol{\theta}^*) + \Pi_{\mathcal{F}} f(\boldsymbol{\theta}^*) - f(\boldsymbol{\theta}^*))\| \\ &\leq \|\text{span}(f_0(\boldsymbol{\theta}_0^*) - \Pi_{\mathcal{F}} f(\boldsymbol{\theta}^*))\| + \|\text{span}(\Pi_{\mathcal{F}} f(\boldsymbol{\theta}^*) - f(\boldsymbol{\theta}^*))\| \\ &= \|\text{span}(\Pi_{\mathcal{F}} \mathcal{T} f_0(\boldsymbol{\theta}_0^*) - \Pi_{\mathcal{F}} \mathcal{T} f(\boldsymbol{\theta}^*))\| + \|\text{span}(\Pi_{\mathcal{F}} f(\boldsymbol{\theta}^*) - f(\boldsymbol{\theta}^*))\| \\ &\leq \kappa \|\text{span}(f_0(\boldsymbol{\theta}_0^*) - f(\boldsymbol{\theta}^*))\| + \|\text{span}(\Pi_{\mathcal{F}} f(\boldsymbol{\theta}^*) - f(\boldsymbol{\theta}^*))\|, \quad (31) \end{aligned}$$

603 where the last inequality follows (24). This indicates that

$$\|\text{span}(f_0(\boldsymbol{\theta}_0^*) - f(\boldsymbol{\theta}^*))\|^2 \leq \frac{1}{(1 - \kappa)^2} \|\text{span}(\Pi_{\mathcal{F}} f(\boldsymbol{\theta}^*) - f(\boldsymbol{\theta}^*))\|^2. \quad (32)$$

604 We further bound $\|y_0(\boldsymbol{\theta}_k) - y(\boldsymbol{\theta}_k)\|^2$ as follows

$$\begin{aligned}
\|y_0(\boldsymbol{\theta}_k) - y(\boldsymbol{\theta}_k)\|^2 &= \left\| \sum_{s'} p(s'|s,1) \max_{a_1} f_0(\boldsymbol{\theta}_k; \boldsymbol{\phi}(s', a_1)) - \sum_{s'} p(s'|s,0) \max_{a_2} f_0(\boldsymbol{\theta}_k; \boldsymbol{\phi}(s', a_2)) \right. \\
&\quad \left. - \sum_{s'} p(s'|s,1) \max_{a_3} f(\boldsymbol{\theta}_k; \boldsymbol{\phi}(s', a_3)) + \sum_{s'} p(s'|s,0) \max_{a_4} f(\boldsymbol{\theta}_k; \boldsymbol{\phi}(s', a_4)) \right\|^2 \\
&= \left\| \sum_{s'} p(s'|s,1) (\max_{a_1} f_0(\boldsymbol{\theta}_k; \boldsymbol{\phi}(s', a_1)) - \max_{a_3} f(\boldsymbol{\theta}_k; \boldsymbol{\phi}(s', a_3))) \right. \\
&\quad \left. - \sum_{s'} p(s'|s,0) (\max_{a_2} f_0(\boldsymbol{\theta}_k; \boldsymbol{\phi}(s', a_2)) - \max_{a_4} f(\boldsymbol{\theta}_k; \boldsymbol{\phi}(s', a_4))) \right\|^2 \\
&\leq 2 \left\| \max_{(s,a)} f_0(\boldsymbol{\theta}_k; \boldsymbol{\phi}(s, a)) - f(\boldsymbol{\theta}_k; \boldsymbol{\phi}(s, a)) \right\|^2 \\
&\leq 2\mathcal{O}\left(\frac{c_1^3(\|\boldsymbol{\theta}_0\| + |\lambda_0| + 1)^3}{m^{1/2}}\right), \tag{33}
\end{aligned}$$

605 where the last inequality is due to Lemma 10. Substituting (32) and (33) back to (30) yields the final
606 results.

607 □

608 D Proof of the Theorem 2

609 To prove Theorem 2, we need the following three key lemmas about the error terms defined in (19).

610 **Lemma 7.** *Let $\{\boldsymbol{\theta}_k, \lambda_k\}$ be generated by (9). Then under Lemmas 1-4, for any $k \geq \tau$, we have*

$$\begin{aligned}
\mathbb{E} \left[\left\| \tilde{\boldsymbol{\theta}}_{k+1} \right\|^2 | \mathcal{F}_{k-\tau} \right] &\leq (1 + 150\alpha_k^2 + \eta_k/\alpha_k - 2\alpha_k\mu_1) \mathbb{E} \left[\left\| \hat{\boldsymbol{\theta}}_k \right\|^2 | \mathcal{F}_{k-\tau} \right] + 6\alpha_k^2 \mathbb{E} \left[\left\| \hat{\lambda}_k \right\|^2 | \mathcal{F}_{k-\tau} \right] \\
&\quad + \frac{\alpha_k^3}{\eta_k} \mathcal{O}\left(c_1^3(\|\boldsymbol{\theta}_0\| + |\lambda_0| + 1)^3 \cdot m^{-1/2}\right). \tag{34}
\end{aligned}$$

611 *Proof.* According to (19), we have $\hat{\boldsymbol{\theta}}_{k+1} := \boldsymbol{\theta}_{k+1} - \boldsymbol{\theta}_0^* = \hat{\boldsymbol{\theta}}_k + \alpha_k h(X_k, \boldsymbol{\theta}_k, \lambda_k)$, which leads to

$$\begin{aligned}
\left\| \hat{\boldsymbol{\theta}}_{k+1} \right\|^2 &= \left\| \hat{\boldsymbol{\theta}}_k \right\|^2 + 2\alpha_k \hat{\boldsymbol{\theta}}_k^\top h(X_k, \boldsymbol{\theta}_k, \lambda_k) + \left\| \alpha_k h(X_k, \boldsymbol{\theta}_k, \lambda_k) \right\|^2 \\
&= \left\| \hat{\boldsymbol{\theta}}_k \right\|^2 + 2\alpha_k \hat{\boldsymbol{\theta}}_k^\top (h(X_k, \boldsymbol{\theta}_k, \lambda_k) - h_0(X_k, \boldsymbol{\theta}_k, \lambda_k)) + 2\alpha_k \hat{\boldsymbol{\theta}}_k^\top h_0(X_k, \boldsymbol{\theta}_k, \lambda_k) \\
&\quad + \alpha_k^2 \|h(X_k, \boldsymbol{\theta}_k, \lambda_k) - h_0(X_k, \boldsymbol{\theta}_k, \lambda_k) + h_0(X_k, \boldsymbol{\theta}_k, \lambda_k)\|^2 \\
&\leq \left\| \hat{\boldsymbol{\theta}}_k \right\|^2 + 2\alpha_k \hat{\boldsymbol{\theta}}_k^\top (h(X_k, \boldsymbol{\theta}_k, \lambda_k) - h_0(X_k, \boldsymbol{\theta}_k, \lambda_k)) + 2\alpha_k \hat{\boldsymbol{\theta}}_k^\top h_0(X_k, \boldsymbol{\theta}_k, \lambda_k) \\
&\quad + 2\alpha_k^2 \|h(X_k, \boldsymbol{\theta}_k, \lambda_k) - h_0(X_k, \boldsymbol{\theta}_k, \lambda_k)\|^2 + 2\alpha_k^2 \|h_0(X_k, \boldsymbol{\theta}_k, \lambda_k)\|^2. \tag{35}
\end{aligned}$$

612 The above inequality holds due to the fact that $\|\mathbf{x} + \mathbf{y}\|^2 \leq 2\|\mathbf{x}\|^2 + 2\|\mathbf{y}\|^2$. Taking expectations of
613 $\|\hat{\boldsymbol{\theta}}_{k+1}\|^2$ w.r.t $\mathcal{F}_{k-\tau}$ yields

$$\begin{aligned}
\mathbb{E} \left[\left\| \hat{\boldsymbol{\theta}}_{k+1} \right\|^2 | \mathcal{F}_{k-\tau} \right] &\leq \mathbb{E} \left[\left\| \hat{\boldsymbol{\theta}}_k \right\|^2 | \mathcal{F}_{k-\tau} \right] + 2\alpha_k \mathbb{E} \left[\hat{\boldsymbol{\theta}}_k^\top h_0(X_k, \boldsymbol{\theta}_k, \lambda_k) | \mathcal{F}_{k-\tau} \right] \\
&\quad + \underbrace{2\alpha_k^2 \mathbb{E} \left[\left\| h_0(X_k, \boldsymbol{\theta}_k, \lambda_k) \right\|^2 | \mathcal{F}_{k-\tau} \right]}_{\text{Term}_1} \\
&\quad + \underbrace{2\alpha_k \mathbb{E} \left[\hat{\boldsymbol{\theta}}_k^\top (h(X_k, \boldsymbol{\theta}_k, \lambda_k) - h_0(X_k, \boldsymbol{\theta}_k, \lambda_k)) | \mathcal{F}_{k-\tau} \right]}_{\text{Term}_2} \\
&\quad + \underbrace{2\alpha_k^2 \mathbb{E} \left[\left\| h(X_k, \boldsymbol{\theta}_k, \lambda_k) - h_0(X_k, \boldsymbol{\theta}_k, \lambda_k) \right\|^2 | \mathcal{F}_{k-\tau} \right]}_{\text{Term}_3}
\end{aligned}$$

$$\leq \mathbb{E} \left[\left\| \hat{\boldsymbol{\theta}}_k \right\|^2 | \mathcal{F}_{k-\tau} \right] - 2\alpha_k \mu_1 \mathbb{E} \left[\left\| \tilde{\boldsymbol{\theta}}_k \right\|^2 | \mathcal{F}_{k-\tau} \right] + \text{Term}_1 + \text{Term}_2 + \text{Term}_3, \quad (36)$$

614 where the last inequality is due to Lemma 4. Next, we bound each individual term. Term_1 is bounded
615 as

$$\begin{aligned} \text{Term}_1 &= 2\alpha_k^2 \mathbb{E} \left[\left\| h_0(X_k, \boldsymbol{\theta}_k, \lambda_k) \right\|^2 | \mathcal{F}_{k-\tau} \right] \\ &\stackrel{(c_1)}{=} 2\alpha_k^2 \mathbb{E} \left[\left\| h_0(X_k, \boldsymbol{\theta}_k, \lambda_k) - h_0(X_k, \boldsymbol{\theta}_k, y_0(\boldsymbol{\theta}_k)) + h_0(X_k, \boldsymbol{\theta}_k, y_0(\boldsymbol{\theta}_k)) \right. \right. \\ &\quad \left. \left. - h_0(X_k, \boldsymbol{\theta}_0^*, y_0(\boldsymbol{\theta}_0^*)) + h_0(X_k, \boldsymbol{\theta}_0^*, y_0(\boldsymbol{\theta}_0^*)) - H_0(\boldsymbol{\theta}_0^*, y_0(\boldsymbol{\theta}_0^*)) \right\|^2 | \mathcal{F}_{k-\tau} \right] \\ &\stackrel{(c_2)}{\leq} 6\alpha_k^2 \mathbb{E} \left[\left\| h_0(X_k, \boldsymbol{\theta}_k, \lambda_k) - h_0(X_k, \boldsymbol{\theta}_k, y_0(\boldsymbol{\theta}_k)) \right\|^2 | \mathcal{F}_{k-\tau} \right] \\ &\quad + 6\alpha_k^2 \mathbb{E} \left[\left\| h_0(X_k, \boldsymbol{\theta}_k, y_0(\boldsymbol{\theta}_k)) - h_0(X_k, \boldsymbol{\theta}_0^*, y_0(\boldsymbol{\theta}_0^*)) \right\|^2 | \mathcal{F}_{k-\tau} \right] \\ &\quad + 6\alpha_k^2 \mathbb{E} \left[\left\| h_0(X_k, \boldsymbol{\theta}_0^*, y_0(\boldsymbol{\theta}_0^*)) - H_0(\boldsymbol{\theta}_0^*, y_0(\boldsymbol{\theta}_0^*)) \right\|^2 | \mathcal{F}_{k-\tau} \right] \\ &\stackrel{(c_3)}{\leq} 6\alpha_k^2 \mathbb{E} \left[\left\| \hat{\lambda}_k \right\|^2 | \mathcal{F}_{k-\tau} \right] + 150\alpha_k^2 \mathbb{E} \left[\left\| \hat{\boldsymbol{\theta}}_k \right\|^2 | \mathcal{F}_{k-\tau} \right], \end{aligned} \quad (37)$$

616 where (c_1) holds due to $H_0(\boldsymbol{\theta}_0, y_0(\boldsymbol{\theta}_0^*)) = 0$, (c_2) follows from the triangular inequality, and (c_3)
617 follows from the Lipschitz continuity of $h_0(X, \boldsymbol{\theta}, \lambda)$ in Lemma 1.

618 Term_2 is bounded as

$$\begin{aligned} \text{Term}_2 &= 2\alpha_k \mathbb{E} \left[\hat{\boldsymbol{\theta}}_k^\top (h(X_k, \boldsymbol{\theta}_k, \lambda_k) - h_0(X_k, \boldsymbol{\theta}_k, \lambda_k)) | \mathcal{F}_{k-\tau} \right] \\ &\stackrel{(c_4)}{\leq} \frac{\eta_k}{\alpha_k} \mathbb{E} \left[\left\| \hat{\boldsymbol{\theta}}_k \right\|^2 | \mathcal{F}_{k-\tau} \right] + \frac{\alpha_k^3}{\eta_k} \mathbb{E} \left[\left\| h(X_k, \boldsymbol{\theta}_k, \lambda_k) - h_0(X_k, \boldsymbol{\theta}_k, \lambda_k) \right\|^2 | \mathcal{F}_{k-\tau} \right] \\ &\stackrel{(c_5)}{\leq} \frac{\eta_k}{\alpha_k} \mathbb{E} \left[\left\| \hat{\boldsymbol{\theta}}_k \right\|^2 | \mathcal{F}_{k-\tau} \right] + \frac{\alpha_k^2}{\eta_k} \mathcal{O} \left(c_1^3 (\|\boldsymbol{\theta}_0\| + |\lambda_0| + 1)^3 \cdot m^{-1/2} \right), \end{aligned} \quad (38)$$

619 where (c_4) holds due to the fact that $2\mathbf{x}^\top \mathbf{y} \leq \|\mathbf{x}\|^2 + \|\mathbf{y}\|^2$ and (c_5) is due to Lemma 10.

620 Term_3 is bounded as

$$\begin{aligned} \text{Term}_3 &= 2\alpha_k^2 \mathbb{E} \left[\left\| h(X_k, \boldsymbol{\theta}_k, \lambda_k) - h_0(X_k, \boldsymbol{\theta}_k, \lambda_k) \right\|^2 | \mathcal{F}_{k-\tau} \right] \\ &\stackrel{(c_6)}{\leq} 2\alpha_k^2 \mathcal{O} \left(c_1^3 (\|\boldsymbol{\theta}_0\| + |\lambda_0| + 1)^3 \cdot m^{-1/2} \right), \end{aligned} \quad (39)$$

621 where (c_6) comes from Lemma 10. Substituting Term_1 , Term_2 , and Term_3 back into (36) leads to
622 the desired result in (34), which is

$$\begin{aligned} \mathbb{E} \left[\left\| \hat{\boldsymbol{\theta}}_{k+1} \right\|^2 | \mathcal{F}_{k-\tau} \right] &\leq \mathbb{E} \left[\left\| \hat{\boldsymbol{\theta}}_k \right\|^2 | \mathcal{F}_{k-\tau} \right] - 2\alpha_k \mu_1 \mathbb{E} \left[\left\| \tilde{\boldsymbol{\theta}}_k \right\|^2 | \mathcal{F}_{k-\tau} \right] \\ &\quad + 6\alpha_k^2 \mathbb{E} \left[\left\| \hat{\lambda}_k \right\|^2 | \mathcal{F}_{k-\tau} \right] + 150\alpha_k^2 \mathbb{E} \left[\left\| \hat{\boldsymbol{\theta}}_k \right\|^2 | \mathcal{F}_{k-\tau} \right] \\ &\quad + \frac{\eta_k}{\alpha_k} \mathbb{E} \left[\left\| \hat{\boldsymbol{\theta}}_k \right\|^2 | \mathcal{F}_{k-\tau} \right] + \frac{\alpha_k^3}{\eta_k} \mathcal{O} \left(c_1^3 (\|\boldsymbol{\theta}_0\| + |\lambda_0| + 1)^3 \cdot m^{-1/2} \right) \\ &\quad + 2\alpha_k^2 \mathcal{O} \left(c_1^3 (\|\boldsymbol{\theta}_0\| + |\lambda_0| + 1)^3 \cdot m^{-1/2} \right) \\ &= (1 + 150\alpha_k^2 + \eta_k/\alpha_k - 2\alpha_k \mu_1) \mathbb{E} \left[\left\| \hat{\boldsymbol{\theta}}_k \right\|^2 | \mathcal{F}_{k-\tau} \right] + 6\alpha_k^2 \mathbb{E} \left[\left\| \hat{\lambda}_k \right\|^2 | \mathcal{F}_{k-\tau} \right] \\ &\quad + (\alpha_k^3/\eta_k + 2\alpha_k^2) \mathcal{O} \left(c_1^3 (\|\boldsymbol{\theta}_0\| + |\lambda_0| + 1)^3 \cdot m^{-1/2} \right). \end{aligned}$$

623 By neglecting higher order infinitesimal, we have the inequality in (34). This completes the proof. \square

624 **Lemma 8.** Let $\{\boldsymbol{\theta}_k, \lambda_k\}$ be generated by (9). Then under Lemmas 1-4, for any $k \geq \tau$, we have

$$\begin{aligned} \mathbb{E} \left[\left\| \hat{\lambda}_{k+1} \right\|^2 \middle| \mathcal{F}_{k-\tau} \right] &\leq (1 - 2\eta_k \mu_2 + \alpha_k \eta_k + 24\alpha_k^2 + \frac{\eta_k}{\alpha_k} - \frac{2\eta_k^2 \mu_k}{\alpha_k} + \eta_k^2 + \frac{24\alpha_k^3}{\eta_k}) \mathbb{E} \left[\left\| \hat{\lambda}_k \right\|^2 \middle| \mathcal{F}_{k-\tau} \right] \\ &\quad + (600\alpha_k^2 + 8\eta_k^2 + \frac{8\eta_k^3}{\alpha_k} + \frac{600\alpha_k^3}{\eta_k}) \mathbb{E} \left[\left\| \hat{\boldsymbol{\theta}}_k \right\|^2 \middle| \mathcal{F}_{k-\tau} \right] \\ &\quad + \frac{\eta_k}{\alpha_k} \mathcal{O} \left(c_1^3 (\|\boldsymbol{\theta}_0\| + |\lambda_0| + 1)^3 \cdot m^{-1/2} \right). \end{aligned} \quad (40)$$

625 *Proof.* According to the definition in (14), we have

$$\begin{aligned} \hat{\lambda}_{k+1} &= \lambda_{k+1} - y_0(\boldsymbol{\theta}_{k+1}) \\ &= \hat{\lambda}_k + \eta_k g(\boldsymbol{\theta}_k) + y_0(\boldsymbol{\theta}_k) - y_0(\boldsymbol{\theta}_{k+1}), \end{aligned}$$

626 which leads to

$$\begin{aligned} \left\| \hat{\lambda}_{k+1} \right\|^2 &= \left\| \hat{\lambda}_k + \eta_k g(\boldsymbol{\theta}_k) + y_0(\boldsymbol{\theta}_k) - y_0(\boldsymbol{\theta}_{k+1}) \right\|^2 \\ &= \underbrace{\left\| \hat{\lambda}_k + \eta_k g(\boldsymbol{\theta}_k) \right\|^2}_{\text{Term}_1} + \underbrace{\left\| y_0(\boldsymbol{\theta}_k) - y_0(\boldsymbol{\theta}_{k+1}) \right\|^2}_{\text{Term}_2} \\ &\quad + 2 \underbrace{\left(\hat{\lambda}_k + \eta_k g(\boldsymbol{\theta}_k) \right) \left(y_0(\boldsymbol{\theta}_k) - y_0(\boldsymbol{\theta}_{k+1}) \right)}_{\text{Term}_3}. \end{aligned} \quad (41)$$

627 The second equality is due to the fact that $\|\mathbf{x} + \mathbf{y}\|^2 = \|\mathbf{x}\|^2 + \|\mathbf{y}\|^2 + 2\mathbf{x}^\top \mathbf{y}$. We next analyze the
628 conditional expectation of each term in $\left\| \hat{\lambda}_{k+1} \right\|^2$ on $\mathcal{F}_{k-\tau}$. We first focus on Term₁.

$$\begin{aligned} \mathbb{E} [\text{Term}_1 | \mathcal{F}_{k-\tau}] &= \mathbb{E} \left[\left\| \hat{\lambda}_k \right\|^2 + 2\eta_k \hat{\lambda}_k g(\boldsymbol{\theta}_k) + \left\| \eta_k g(\boldsymbol{\theta}_k) \right\|^2 \middle| \mathcal{F}_{k-\tau} \right] \\ &= \mathbb{E} \left[\left\| \hat{\lambda}_k \right\|^2 + 2\eta_k \hat{\lambda}_k g_0(\boldsymbol{\theta}_k) + 2\eta_k \hat{\lambda}_k (g(\boldsymbol{\theta}_k) - g_0(\boldsymbol{\theta}_k)) + \eta_k^2 \left\| g(\boldsymbol{\theta}_k) - g_0(\boldsymbol{\theta}_k) + g_0(\boldsymbol{\theta}_k) \right\|^2 \middle| \mathcal{F}_{k-\tau} \right] \\ &\leq \mathbb{E} \left[\left\| \hat{\lambda}_k \right\|^2 \middle| \mathcal{F}_{k-\tau} \right] + 2\eta_k \mathbb{E} \left[\hat{\lambda}_k g_0(\boldsymbol{\theta}_k) \middle| \mathcal{F}_{k-\tau} \right] + 2\eta_k \mathbb{E} \left[\hat{\lambda}_k (g(\boldsymbol{\theta}_k) - g_0(\boldsymbol{\theta}_k)) \middle| \mathcal{F}_{k-\tau} \right] \\ &\quad + 2\eta_k^2 \mathbb{E} \left[\left\| g(\boldsymbol{\theta}_k) - g_0(\boldsymbol{\theta}_k) \right\|^2 \middle| \mathcal{F}_{k-\tau} \right] + 2\eta_k^2 \mathbb{E} \left[\left\| g_0(\boldsymbol{\theta}_k) \right\|^2 \middle| \mathcal{F}_{k-\tau} \right] \\ &\stackrel{(d_1)}{=} \mathbb{E} \left[\left\| \hat{\lambda}_k \right\|^2 \middle| \mathcal{F}_{k-\tau} \right] + 2\eta_k \mathbb{E} \left[\hat{\lambda}_k g_0(\boldsymbol{\theta}_k) \middle| \mathcal{F}_{k-\tau} \right] + 2\eta_k \mathbb{E} \left[\hat{\lambda}_k (g(\boldsymbol{\theta}_k) - g_0(\boldsymbol{\theta}_k)) \middle| \mathcal{F}_{k-\tau} \right] \\ &\quad + 2\eta_k^2 \mathbb{E} \left[\left\| g(\boldsymbol{\theta}_k) - g_0(\boldsymbol{\theta}_k) \right\|^2 \middle| \mathcal{F}_{k-\tau} \right] + 2\eta_k^2 \mathbb{E} \left[\left\| g_0(\boldsymbol{\theta}_k) - g_0(\boldsymbol{\theta}_0^*) \right\|^2 \middle| \mathcal{F}_{k-\tau} \right] \\ &\stackrel{(d_2)}{=} \mathbb{E} \left[\left\| \hat{\lambda}_k \right\|^2 \middle| \mathcal{F}_{k-\tau} \right] - 2\eta_k \mu_2 \mathbb{E} \left[\left\| \hat{\lambda}_k \right\|^2 \middle| \mathcal{F}_{k-\tau} \right] + 8\eta_k^2 \mathbb{E} \left[\left\| \hat{\boldsymbol{\theta}}_k \right\|^2 \middle| \mathcal{F}_{k-\tau} \right] \\ &\quad + 2\eta_k \mathbb{E} \left[\hat{\lambda}_k (g(\boldsymbol{\theta}_k) - g_0(\boldsymbol{\theta}_k)) \middle| \mathcal{F}_{k-\tau} \right] + 2\eta_k^2 \mathbb{E} \left[\left\| g(\boldsymbol{\theta}_k) - g_0(\boldsymbol{\theta}_k) \right\|^2 \middle| \mathcal{F}_{k-\tau} \right] \\ &\stackrel{(d_3)}{\leq} \mathbb{E} \left[\left\| \hat{\lambda}_k \right\|^2 \middle| \mathcal{F}_{k-\tau} \right] - 2\eta_k \mu_2 \mathbb{E} \left[\left\| \hat{\lambda}_k \right\|^2 \middle| \mathcal{F}_{k-\tau} \right] + 8\eta_k^2 \mathbb{E} \left[\left\| \hat{\boldsymbol{\theta}}_k \right\|^2 \middle| \mathcal{F}_{k-\tau} \right] \\ &\quad + \alpha_k \eta_k \mathbb{E} \left[\left\| \hat{\lambda}_k \right\|^2 \middle| \mathcal{F}_{k-\tau} \right] + (4\eta_k / \alpha_k + 8\eta_k^2) \mathcal{O} \left(c_1^3 (\|\boldsymbol{\theta}_0\| + |\lambda_0| + 1)^3 \cdot m^{-1/2} \right), \end{aligned}$$

629 where (d_1) follows from $g_0(\boldsymbol{\theta}_0^*) = 0$, (d_2) holds due to Lemma 4 and the Lipschitz continuity of y_0
630 in Lemma 3, and (d_3) comes from Lemma 10. For Term₂, we have

$$\mathbb{E} [\text{Term}_2 | \mathcal{F}_{k-\tau}] = \mathbb{E} \left[\left\| y_0(\boldsymbol{\theta}_k) - y_0(\boldsymbol{\theta}_{k+1}) \right\|^2 \middle| \mathcal{F}_{k-\tau} \right]$$

$$\begin{aligned}
&= 4\mathbb{E} \left[\left\| \boldsymbol{\theta}_k - \boldsymbol{\theta}_{k+1} \right\|^2 | \mathcal{F}_{k-\tau} \right] \\
&= 4\alpha_k^2 \mathbb{E} \left[\left\| h(X_k, \boldsymbol{\theta}_k, \lambda_k) \right\|^2 | \mathcal{F}_{k-\tau} \right] \\
&= 4\alpha_k^2 \mathbb{E} \left[\left\| h(X_k, \boldsymbol{\theta}_k, \lambda_k) - h_0(X_k, \boldsymbol{\theta}_k, \lambda_k) + h_0(X_k, \boldsymbol{\theta}_k, \lambda_k) \right\|^2 | \mathcal{F}_{k-\tau} \right] \\
&= 8\alpha_k^2 \mathbb{E} \left[\left\| h_0(X_k, \boldsymbol{\theta}_k, \lambda_k) \right\|^2 | \mathcal{F}_{k-\tau} \right] + 8\alpha_k^2 \mathbb{E} \left[\left\| h(X_k, \boldsymbol{\theta}_k, \lambda_k) - h_0(X_k, \boldsymbol{\theta}_k, \lambda_k) \right\|^2 | \mathcal{F}_{k-\tau} \right] \\
&\stackrel{(d_4)}{=} 8\alpha_k^2 \mathbb{E} \left[\left\| h_0(X_k, \boldsymbol{\theta}_k, \lambda_k) - h_0(X_k, \boldsymbol{\theta}_k, y_0(\boldsymbol{\theta}_k)) + h_0(X_k, \boldsymbol{\theta}_k, y_0(\boldsymbol{\theta}_k)) \right. \right. \\
&\quad \left. \left. - h_0(X_k, \boldsymbol{\theta}_0^*, y_0(\boldsymbol{\theta}_0^*)) + h_0(X_k, \boldsymbol{\theta}_0^*, y_0(\boldsymbol{\theta}_0^*)) - H_0(\boldsymbol{\theta}_0^*, y_0(\boldsymbol{\theta}_0^*)) \right\|^2 | \mathcal{F}_{k-\tau} \right] \\
&\quad + 8\alpha_k^2 \mathbb{E} \left[\left\| h(X_k, \boldsymbol{\theta}_k, \lambda_k) - h_0(X_k, \boldsymbol{\theta}_k, \lambda_k) \right\|^2 | \mathcal{F}_{k-\tau} \right] \\
&\stackrel{(d_5)}{\leq} 24\alpha_k^2 \mathbb{E} \left[\left\| h_0(X_k, \boldsymbol{\theta}_k, \lambda_k) - h_0(X_k, \boldsymbol{\theta}_k, y_0(\boldsymbol{\theta}_k)) \right\|^2 | \mathcal{F}_{k-\tau} \right] \\
&\quad + 24\alpha_k^2 \mathbb{E} \left[\left\| h_0(X_k, \boldsymbol{\theta}_k, y_0(\boldsymbol{\theta}_k)) - h_0(X_k, \boldsymbol{\theta}_0^*, y_0(\boldsymbol{\theta}_0^*)) \right\|^2 | \mathcal{F}_{k-\tau} \right] \\
&\quad + 24\alpha_k^2 \mathbb{E} \left[\left\| h_0(X_k, \boldsymbol{\theta}_0^*, y_0(\boldsymbol{\theta}_0^*)) - H_0(\boldsymbol{\theta}_0^*, y_0(\boldsymbol{\theta}_0^*)) \right\|^2 | \mathcal{F}_{k-\tau} \right] \\
&\quad + 8\alpha_k^2 \mathbb{E} \left[\left\| h(X_k, \boldsymbol{\theta}_k, \lambda_k) - h_0(X_k, \boldsymbol{\theta}_k, \lambda_k) \right\|^2 | \mathcal{F}_{k-\tau} \right] \\
&\stackrel{(d_6)}{\leq} 24\alpha_k^2 \mathbb{E} \left[\left\| \hat{\lambda}_k \right\|^2 | \mathcal{F}_{k-\tau} \right] + 600\alpha_k^2 \mathbb{E} \left[\left\| \hat{\boldsymbol{\theta}}_k \right\|^2 | \mathcal{F}_{k-\tau} \right] \\
&\quad + 8\alpha_k^2 \mathbb{E} \left[\left\| h(X_k, \boldsymbol{\theta}_k, \lambda_k) - h_0(X_k, \boldsymbol{\theta}_k, \lambda_k) \right\|^2 | \mathcal{F}_{k-\tau} \right] \\
&\stackrel{(d_7)}{\leq} 24\alpha_k^2 \mathbb{E} \left[\left\| \hat{\lambda}_k \right\|^2 | \mathcal{F}_{k-\tau} \right] + 600\alpha_k^2 \mathbb{E} \left[\left\| \hat{\boldsymbol{\theta}}_k \right\|^2 | \mathcal{F}_{k-\tau} \right] \\
&\quad + 8\alpha_k^2 \mathcal{O} \left(c_1^3 (\|\boldsymbol{\theta}_0\| + |\lambda_0| + 1)^3 \cdot m^{-1/2} \right) \tag{42}
\end{aligned}$$

631 where (d_4) is due to the fact that $H_0(\boldsymbol{\theta}_0^*, y_0(\boldsymbol{\theta}_0^*)) = 0$, (d_5) holds according to $\|\mathbf{x} + \mathbf{y} + \mathbf{z}\|^2 \leq$
 632 $3\|\mathbf{x}\|^2 + 3\|\mathbf{y}\|^2 + 3\|\mathbf{z}\|^2$ since $g(X_k, f(\boldsymbol{\lambda}^*), \boldsymbol{\lambda}^*) = \mathbf{0}$, (d_6) holds because of the Lipschitz continuity
 633 of h_0 and y_0 in Lemma 1 and Lemma 3, and (d_7) comes from Lemma 10. Next, we have the
 634 conditional expectation of Term₃ as

$$\begin{aligned}
\mathbb{E} [\text{Term}_3 | \mathcal{F}_{k-\tau}] &= 2\mathbb{E} \left[\left\| \hat{\lambda}_k + \eta_k g(\boldsymbol{\theta}_k) \right\| \cdot \left\| y_0(\boldsymbol{\theta}_k) - y_0(\boldsymbol{\theta}_{k+1}) \right\| | \mathcal{F}_{k-\tau} \right] \\
&\stackrel{(d_8)}{\leq} \frac{\eta_k}{\alpha_k} \text{Term}_1 + \frac{\alpha_k}{\eta_k} \text{Term}_2 \\
&= \frac{\eta_k}{\alpha_k} \mathbb{E} \left[\left\| \hat{\lambda}_k \right\|^2 | \mathcal{F}_{k-\tau} \right] - \frac{2\eta_k^2 \mu_2}{\alpha_k} \mathbb{E} \left[\left\| \hat{\lambda}_k \right\|^2 | \mathcal{F}_{k-\tau} \right] + \frac{8\eta_k^3}{\alpha_k} \mathbb{E} \left[\left\| \hat{\boldsymbol{\theta}}_k \right\|^2 | \mathcal{F}_{k-\tau} \right] \\
&\quad + \eta_k^2 \mathbb{E} \left[\left\| \hat{\lambda}_k \right\|^2 | \mathcal{F}_{k-\tau} \right] + \frac{\eta_k}{\alpha_k} (4\eta_k / \alpha_k + 8\eta_k^2) \mathcal{O} \left(c_1^3 (\|\boldsymbol{\theta}_0\| + |\lambda_0| + 1)^3 \cdot m^{-1/2} \right) \\
&\quad + \frac{24\alpha_k^3}{\eta_k} \mathbb{E} \left[\left\| \hat{\lambda}_k \right\|^2 | \mathcal{F}_{k-\tau} \right] + \frac{600\alpha_k^3}{\eta_k} \mathbb{E} \left[\left\| \hat{\boldsymbol{\theta}}_k \right\|^2 | \mathcal{F}_{k-\tau} \right] \\
&\quad + \frac{8\alpha_k^3}{\eta_k} \mathcal{O} \left(c_1^3 (\|\boldsymbol{\theta}_0\| + |\lambda_0| + 1)^3 \cdot m^{-1/2} \right),
\end{aligned}$$

635 where (d_8) holds because $2\mathbf{x}^T \mathbf{y} \leq 1/\beta \|\mathbf{x}\|^2 + \beta \|\mathbf{y}\|^2, \forall \beta > 0$. Summing Term₁, Term₂, and Term₃
 636 and neglecting higher order infinitesimal yield the desired result. \square

Now we are ready to prove the results in Theorem 2. Providing Lemma 8 and Lemma 7, if $\frac{\eta_k}{\alpha_k}$ is non-increasing, we have the following inequality

$$\begin{aligned}
\mathbb{E} \left[\hat{M}(\boldsymbol{\theta}_{k+1}, \lambda_{k+1}) \middle| \mathcal{F}_{k-\tau} \right] &= \mathbb{E} \left[\frac{\eta_k}{\alpha_k} \left\| \hat{\boldsymbol{\theta}}_{k+1} \right\|^2 + \left\| \hat{\lambda}_{k+1} \right\|^2 \middle| \mathcal{F}_{k-\tau} \right] \\
&\leq \frac{\eta_k}{\alpha_k} (1 - 2\alpha_k \mu_1) \mathbb{E} \left[\left\| \hat{\boldsymbol{\theta}}_k \right\|^2 \middle| \mathcal{F}_{k-\tau} \right] + \frac{600\alpha_k^3}{\eta_k} \mathbb{E} \left[\left\| \hat{\boldsymbol{\theta}}_k \right\|^2 \middle| \mathcal{F}_{k-\tau} \right] \\
&\quad + \frac{8\alpha_k^3}{\eta_k} \mathcal{O} \left(c_1^3 (\|\boldsymbol{\theta}_0\| + |\lambda_0| + 1)^3 \cdot m^{-1/2} \right) \\
&\quad + (1 - 2\eta_k \mu_2) \mathbb{E} \left[\left\| \hat{\lambda}_k \right\|^2 \middle| \mathcal{F}_{k-\tau} \right] \\
&\quad + \frac{600\alpha_k^3}{\eta_k} \mathbb{E} \left[\left\| \hat{\lambda}_k \right\|^2 \middle| \mathcal{F}_{k-\tau} \right]. \tag{43}
\end{aligned}$$

Since $(k+1)^2 \cdot \frac{\alpha_k^3}{\eta} = \frac{\alpha_0^3}{\eta_0} (k+1)^{1/3}$, multiplying both sides of (43) with $(k+1)^2$, we have

$$\begin{aligned}
(k+1)^2 \mathbb{E} \left[\hat{M}(\boldsymbol{\theta}_{k+1}, \lambda_{k+1}) \middle| \mathcal{F}_{k-\tau} \right] &\leq k^2 \mathbb{E} \left[\hat{M}(\boldsymbol{\theta}_k, \lambda_k) \middle| \mathcal{F}_{k-\tau} \right] + \frac{600\alpha_0^3}{\eta_0} (k+1)^{1/3} \left(\left\| \hat{\boldsymbol{\theta}}_k \right\|^2 + \left\| \hat{\lambda}_k \right\|^2 \right) \\
&\quad + \frac{8\alpha_0^3}{\eta_0} (k+1)^{1/3} \mathcal{O} \left(c_1^3 (\|\boldsymbol{\theta}_0\| + |\lambda_0| + 1)^3 \cdot m^{-1/2} \right). \tag{44}
\end{aligned}$$

Summing (44) from time step τ to time step k , we have

$$\begin{aligned}
(k+1)^2 \mathbb{E} \left[\hat{M}(\boldsymbol{\theta}_{k+1}, \lambda_{k+1}) \middle| \mathcal{F}_k \right] &\leq \tau^2 \mathbb{E} \left[\hat{M}(\boldsymbol{\theta}_\tau, \lambda_\tau) \right] + \frac{600\alpha_0^3}{\eta_0} (k+1)^{4/3} \left(\left\| \hat{\boldsymbol{\theta}}_\tau \right\|^2 + \left\| \hat{\lambda}_\tau \right\|^2 \right) \\
&\quad + \frac{8\alpha_0^3}{\eta_0} (k+1)^{4/3} \mathcal{O} \left(c_1^3 (\|\boldsymbol{\theta}_0\| + |\lambda_0| + 1)^3 \cdot m^{-1/2} \right) \\
&\leq \tau^2 \mathbb{E} \left[\hat{M}(\boldsymbol{\theta}_\tau, \lambda_\tau) \right] + \frac{600\alpha_0^3}{\eta_0} \frac{(C_1 + \|\hat{\boldsymbol{\theta}}_0\|)^2 + (2C_1 + \|\hat{\lambda}_0\|)^2}{(k+1)^{-4/3}} \\
&\quad + \frac{8\alpha_0^3}{\eta_0} \frac{\mathcal{O} \left(c_1^3 (\|\boldsymbol{\theta}_0\| + |\lambda_0| + 1)^3 \cdot m^{-1/2} \right)}{(k+1)^{-4/3}}, \tag{45}
\end{aligned}$$

where the second inequality holds due to Lemma 9. Finally, dividing both sides by $(k+1)^2$ and moving the constant term into $\mathcal{O}(\cdot)$ yields the results in Theorem 2.

E Auxiliary Lemmas

In this part, we present several key lemmas which are needed for the major proofs. We first show the parameters update in (9) is bounded in the following lemma.

Lemma 9. *The update of $\boldsymbol{\theta}_k$ and λ_k in (9) is bounded with respect to the initial $\boldsymbol{\theta}_0$ and λ_0 , i.e.,*

$$\|\boldsymbol{\theta}_k - \boldsymbol{\theta}_0\| + |\lambda_k - \lambda_0| \leq c_1 (\|\boldsymbol{\theta}_0\| + |\lambda_0| + 1),$$

with c_1 be the constant, i.e., $c_1 := \frac{1}{2} + \frac{3}{2}(L'_h \alpha_\tau + L'_g \eta_\tau)(L'_h \alpha_\tau + L'_g \eta_\tau + 1)$.

Proof. Without loss of generality, we assume that

$$L'_h \geq \max(3, \max_{X \in \mathcal{X}} \|h_0(X, 0, 0)\|), \quad L'_g \geq \max(2, \max_{X \in \mathcal{X}} \|g_0(X, 0, 0)\|).$$

Then based on triangular inequality and Lemmas 1-2, we have

$$\|h_0(X, \boldsymbol{\theta}, \lambda)\| \leq L'_h (\|\boldsymbol{\theta}\| + |\lambda| + 1), \quad \|g_0(X, \boldsymbol{\theta}, \lambda)\| \leq L'_g (\|\boldsymbol{\theta}\| + |\lambda| + 1), \quad \forall \boldsymbol{\theta}, \lambda, X \in \mathcal{X}. \tag{46}$$

648 Since we have $\boldsymbol{\theta}_{k+1} = \boldsymbol{\theta}_k + \alpha_k h(X_k, \boldsymbol{\theta}_k, \lambda_k)$, we have the following inequality due to Lipschitz
 649 continuity of h in (46)

$$\|\boldsymbol{\theta}_{k+1} - \boldsymbol{\theta}_k\| = \alpha_k \|h(X_k, \boldsymbol{\theta}_k, \lambda_k)\| \leq \alpha_k L'_h (\|\boldsymbol{\theta}_k\| + |\lambda_k| + 1). \quad (47)$$

650 Similarly, we have

$$|\lambda_{k+1} - \lambda_k| = \eta_k |g(X_k, \boldsymbol{\theta}_k, \lambda_k)| \leq \eta_k L'_g (\|\boldsymbol{\theta}_k\| + |\lambda_k| + 1). \quad (48)$$

651 Due to triangular inequality, adding (47) and (48) leads to

$$\begin{aligned} \|\boldsymbol{\theta}_{k+1}\| + |\lambda_{k+1}| + 1 &\leq (L'_h \alpha_k + L'_g \eta_k + 1)(\|\boldsymbol{\theta}_k\| + |\lambda_k| + 1) \\ &\leq (L'_h \alpha_0 + L'_g \eta_0 + 1)(\|\boldsymbol{\theta}_k\| + |\lambda_k| + 1), \end{aligned} \quad (49)$$

652 where the second inequality holds due to the non-increasing learning rates $\{\alpha_k, \eta_k\}$. Rewriting the
 653 above inequality in (49) in a recursive manner yields

$$\|\boldsymbol{\theta}_k\| + \lambda_k + 1 \leq (L'_h \alpha_0 + L'_g \eta_0 + 1)^{k-\tau} (\|\boldsymbol{\theta}_\tau\| + |\lambda_\tau| + 1). \quad (50)$$

654 Hence, we have

$$\begin{aligned} \|\boldsymbol{\theta}_k - \boldsymbol{\theta}_{k-\tau}\| + |\lambda_k - \lambda_{k-\tau}| &\leq \sum_{t=k-\tau}^{k-1} \|\boldsymbol{\theta}_{t+1} - \boldsymbol{\theta}_t\| + |\lambda_{t+1} - \lambda_t| \\ &\leq (L'_h \alpha_0 + L'_g \eta_0) \sum_{t=k-\tau}^{k-1} (\|\boldsymbol{\theta}_t\| + |\lambda_t| + 1) \\ &\leq (L'_h \alpha_0 + L'_g \eta_0) (\|\boldsymbol{\theta}_{k-\tau}\| + |\lambda_{k-\tau}| + 1) \sum_{t=k-\tau}^{k-1} (L'_h \alpha_0 + L'_g \eta_0 + 1)^{t-\tau} \\ &= [(L'_h \alpha_0 + L'_g \eta_0 + 1)^\tau - 1] (\|\boldsymbol{\theta}_{k-\tau}\| + |\lambda_{k-\tau}| + 1) \\ &\leq (e^{(L'_h \alpha_0 + L'_g \eta_0)\tau} - 1) (\|\boldsymbol{\theta}_{k-\tau}\| + |\lambda_{k-\tau}| + 1) \\ &\leq 2(L'_h \alpha_0 + L'_g \eta_0) \tau (\|\boldsymbol{\theta}_{k-\tau}\| + |\lambda_{k-\tau}| + 1), \end{aligned}$$

655 where the last inequality holds when $(L'_h \alpha_0 + L'_g \eta_0) \tau \leq 1/4$. This implies when $k = \tau$, we have

$$\|\boldsymbol{\theta}_\tau - \boldsymbol{\theta}_0\| + |\lambda_\tau - \lambda_0| \leq 2(L'_h \alpha_0 + L'_g \eta_0) \tau (\|\boldsymbol{\theta}_0\| + |\lambda_0| + 1). \quad (51)$$

656 Similarly, we also have

$$\begin{aligned} \|\boldsymbol{\theta}_k - \boldsymbol{\theta}_\tau\| + |\lambda_k - \lambda_\tau| &\leq \sum_{t=\tau}^{k-1} \|\boldsymbol{\theta}_{t+1} - \boldsymbol{\theta}_t\| + |\lambda_{t+1} - \lambda_t| \\ &\leq \sum_{t=\tau}^{k-1} (L'_h \alpha_t + L'_g \eta_t) (\|\boldsymbol{\theta}_t\| + |\lambda_t| + 1) \\ &\leq (\|\boldsymbol{\theta}_\tau\| + |\lambda_\tau| + 1) \sum_{t=\tau}^{k-1} (L'_h \alpha_t + L'_g \eta_t) \prod_{i=0}^{t-\tau} (L'_h \alpha_{\tau+i} + L'_g \eta_{\tau+i} + 1). \end{aligned} \quad (52)$$

657 Therefore, the following inequality holds

$$\begin{aligned} &\|\boldsymbol{\theta}_k - \boldsymbol{\theta}_0\| + |\lambda_k - \lambda_0| \\ &\leq \|\boldsymbol{\theta}_k - \boldsymbol{\theta}_\tau\| + |\lambda_k - \lambda_\tau| + \|\boldsymbol{\theta}_\tau - \boldsymbol{\theta}_0\| + |\lambda_\tau - \lambda_0| \\ &\leq 2(L'_h \alpha_0 + L'_g \eta_0) \tau (\|\boldsymbol{\theta}_0\| + |\lambda_0| + 1) \\ &\quad + (\|\boldsymbol{\theta}_\tau\| + |\lambda_\tau| + 1) \sum_{t=\tau}^{k-1} (L'_h \alpha_t + L'_g \eta_t) \prod_{i=0}^{t-\tau} (L'_h \alpha_{\tau+i} + L'_g \eta_{\tau+i} + 1) \\ &\leq 2(L'_h \alpha_0 + L'_g \eta_0) \tau (\|\boldsymbol{\theta}_0\| + |\lambda_0| + 1) \end{aligned}$$

$$\begin{aligned}
& + (2(L'_h\alpha_0 + L'_g\eta_0)\tau + 1) \sum_{t=\tau}^{k-1} (L'_h\alpha_t + L'_g\eta_t) \prod_{i=0}^{t-\tau} (L'_h\alpha_{\tau+i} + L'_g\eta_{\tau+i} + 1) (\|\boldsymbol{\theta}_0\| + |\lambda_0| + 1) \\
& \leq \left(\frac{1}{2} + \frac{3}{2} \sum_{t=\tau}^{k-1} (L'_h\alpha_t + L'_g\eta_t) \prod_{i=0}^{t-\tau} (L'_h\alpha_{\tau+i} + L'_g\eta_{\tau+i} + 1) \right) (\|\boldsymbol{\theta}_0\| + |\lambda_0| + 1),
\end{aligned}$$

with the last equality holds when $(L'_h\alpha_0 + L'_g\eta_0)\tau \leq 1/4$. When $\sum_{t=\tau}^{k-1} (L'_h\alpha_t + L'_g\eta_t) \prod_{i=0}^{t-\tau} (L'_h\alpha_{\tau+i} + L'_g\eta_{\tau+i} + 1)$ is non-increasing with k , then we can set c_1 as $c_1 := \frac{1}{2} + \frac{3}{2} \cdot (L'_h\alpha_\tau + L'_g\eta_\tau)(L'_h\alpha_\tau + L'_g\eta_\tau + 1)$. This completes the proof. \square

Provided Lemma 9, we have the following lemma related with local linearization of Q functions and the original Q functions.

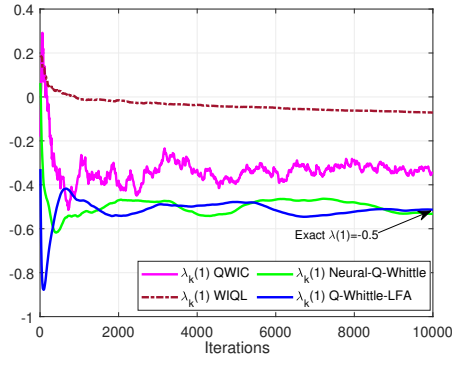
Lemma 10 (Lemma 5.2 in [13]). *There exists a constant c_1 such that*

$$\mathbb{E} \left[\|h(X_k, \boldsymbol{\theta}_k, \lambda_k) - h_0(X_k, \boldsymbol{\theta}_k, \lambda_k)\|^2 | \mathcal{F}_{k-\tau} \right] \leq \mathcal{O} \left(c_1^3 (\|\boldsymbol{\theta}_0\| + |\lambda_0| + 1)^3 \cdot m^{-1/2} \right).$$

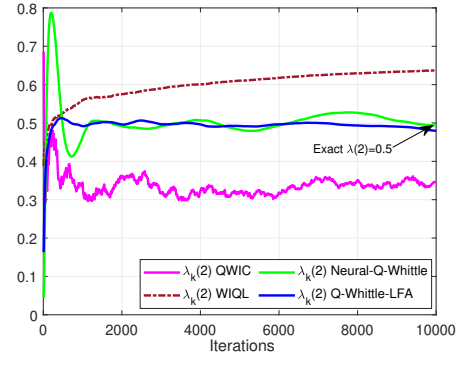
Lemma 10 indicates that if the updated parameter is always bounded in a ball with the initialized one as the center and a fixed radius, the local linearized function $f_0(\cdot)$ in (18) and the original neural network approximated function $f(\cdot)$ in (6) have bounded gap, which tends to be zero as the width of hidden layer m grow large. For interested readers, please refer to [13] for detailed proofs of this lemma.

F Additional Experimental Results

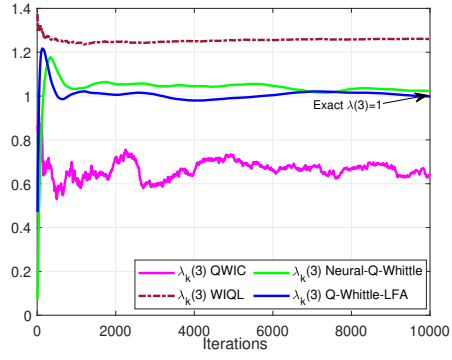
In this section, we provide some additional experimental results related with the convergence comparison of the proposed Neural-Q-Whittle to other Whittle index learning algorithms, i.e., Q-Whittle-LFA [57], WIQL [8] and QWIC [23]. The underlying configuration is the same as that in Section 5. In Figure 4, we present the convergence results of Neural-Q-Whittle and the aforementioned benchmarks. As we observe from Figure 4, only Neural-Q-Whittle and Q-Whittle-LFA in [57] can converge to the true Whittle indices for each state, while the other two benchmarks algorithms do not guarantee the convergence of true Whittle indices. Interestingly, the learning Whittle indices converge and maintain a correct relative order of magnitude, which is still be able to be used in real world problems [57]. Moreover, we observe that Neural-Q-Whittle achieves similar convergence performance as Q-Whittle-LFA in the considered example, whereas the latter has been shown to achieve good performance in real world applications in [57]. Though this work focuses on the theoretical convergence analysis of Q-learning based whittle index under the neural network function approximation, it might be promising to implement it in real-world applications to fully leverage the strong representation ability of neural network functions, which serves as future investigation of this work.



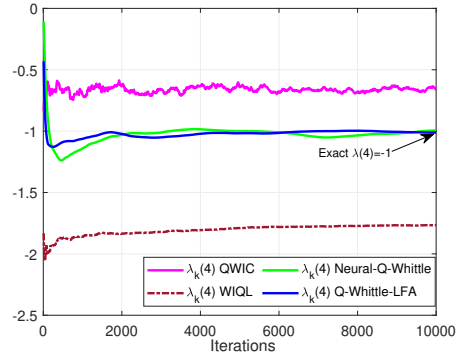
(a) Whittle index $\lambda(1)$ for state $s = 1$.



(b) Whittle index $\lambda(2)$ for state $s = 2$.



(c) Whittle index $\lambda(3)$ for state $s = 3$.



(d) Whittle index $\lambda(4)$ for state $s = 4$.

Figure 4: Convergence comparison between Neural-Q-Whittle and other benchmark algorithms.