

Discrete-Time High Order Tuner With A Time-Varying Learning Rate*

Yingnan Cui[†] and Anuradha M. Annaswamy

Abstract

We propose a new discrete-time online parameter estimation algorithm that combines two different aspects, one that adds momentum, and another that includes a time-varying learning rate. It is well known that recursive least squares based approaches that include a time-varying gain can lead to exponential convergence of parameter errors under persistent excitation, while momentum-based approaches have demonstrated a fast convergence of tracking error towards zero with constant regressors. The question is when combined, will the filter from the momentum method come in the way of exponential convergence. This paper proves that exponential convergence of parameter is still possible with persistent excitation. Simulation results demonstrated competitive properties of the proposed algorithm compared to the recursive least squares algorithm with forgetting.

1 Introduction

An essential part of any adaptive control algorithm is reliable, fast online parameter estimation [1, 2]. Beyond the basic gradient descent method, a large amount of works have focused on proposing provably stable, more efficient algorithms for online parameter estimation in adaptive control [3–7].

It is well known that the introduction of a time-varying learning rate leads to exponential learning of the parameters in the presence of persistent excitation. Both recursive least squares (RLS) and RLS with forgetting have been frequently adopted for parameter estimation [8]. This idea of adopting time-varying learning rate has also led to some major breakthroughs in the optimization community. AdaGrad, for example, adapts the learning rate to the adjustment of parameters, applying larger updates for infrequently adjusted parameters and smaller updates for frequently adjusted parameters [9]. AdaDelta adopts an exponential decaying average of the past gradients to address AdaGrad's aggressive, monotonically decaying learning rate [10].

Yet another recent set of results that leads to accelerated performance, such as fast reduction of a loss function, is through the addition of momentum. It is a well observed fact that gradient descent method often performs badly around saddle points and local optima [11], and provides a convergence rate in $\mathcal{O}(1/k)$, where k is the iteration number. In contrast, Nesterov's acceleration, which adopts the idea of momentum, is a method that helps accelerate gradient descent and can lead to a convergence rate of $\mathcal{O}(1/k^2)$ when the loss function is smooth [12]. In problems of parameter estimation, it has been shown more recently that momentum-based methods, also known as high-order tuners (HT), can lead to acceleration even with time-varying regressors if the loss is strongly convex [4].

In real-time systems, it is of paramount importance to have both acceleration in performance, i.e. in a fast decrease of the loss function, and in learning, i.e. fast convergence of the parameter estimates to their true values. The question therefore is if HT can be combined with time-varying learning rates and lead to both accelerated performance and accelerated learning. Since HT includes an additional

*This work is supported by the Boeing Strategic University Initiative.

[†]Y. Cui and A.M. Annaswamy are with the Department of Mechanical Engineering, Massachusetts Institute of Technology, Cambridge, MA, 02139.

filter, it needs to be ensured that the filtering action does not compromise the property of fast learning in the presence of time-varying gains. In this paper, we show that is not the case and that persistent excitation guarantees exponential convergence of the parameter estimates to the true value.

The specific HT that we consider is that based on Heavy Ball method (HB) that is proposed by Polyak [13]. We add a time-varying gain matrix in addition to the momentum term that is present in the HB method. Through careful adjustment of the time-varying gain matrix, we show that the gain matrix remains bounded, does not go to zero with persistent excitation, and that the parameter estimates converge to their true values exponentially. This is the central contribution of this paper. All results are in the context of a nonlinear ARMA model with unknown parameters that are to be estimated.

Section 2 states the problem we want to solve. Section 3 presents the algorithm. We discuss stability properties of the algorithm in section 4 and show numerical simulations in section 5. Section 6 summarizes the paper and discusses future works.

2 Problem Statement

We consider a class of discrete-time nonlinear plant models of the form

$$y_k = - \sum_{i=1}^n a_i^* y_{k-i} + \sum_{j=1}^m b_j^* u_{k-j-d} + \sum_{\ell=1}^p c_\ell^* f_\ell(y_{k-1}, \dots, y_{k-n}, u_{k-1-d}, \dots, u_{k-m-d}), \quad (1)$$

where a_i^* , b_j^* and c_ℓ^* are unknown parameters that are constant and need to be identified, and d is a known time-delay. The function f_ℓ is an analytic function and is assumed to be such that the system in (1) is bounded-input-bounded-output (BIBO) stable. Denote $z_{k-1} = [y_{k-1}, \dots, y_{k-n}]^\top$ and $v_{k-d-1} = [u_{k-1-d}, \dots, u_{k-m-d}]^\top$. We rewrite (1) in the form of a linear regression

$$y_k = \phi_k^\top \theta^*, \quad (2)$$

where $\phi_k = [z_{k-1}^\top, v_{k-d-1}^\top, f_1(z_{k-1}^\top, v_{k-d-1}^\top), \dots, f_p(z_{k-1}^\top, v_{k-d-1}^\top)]^\top$ is a regressor determined by exogenous signals and $\theta^* = [a_1^*, \dots, a_n^*, b_1^*, \dots, b_m^*, c_1^*, \dots, c_p^*]^\top$ is the underlying unknown parameter vector. We propose to identify the parameter θ^* as θ_k using an estimator

$$\hat{y}_k = \phi_k^\top \theta_k, \quad (3)$$

which leads to a prediction error

$$e_{y,k} = \phi_k^\top \tilde{\theta}_k, \quad (4)$$

where $e_{y,k} = \hat{y}_k - y_k$ is the output prediction error and $\tilde{\theta}_k = \theta_k - \theta^*$ is the parameter error. The goal of parameter identification is to design an iterative procedure such that the parameter error $\|\tilde{\theta}_k\|$ converges to zero exponentially fast.

The iterative procedure for estimating the parameters is based on a squared loss function,

$$L_k(\theta_k) = \frac{1}{2} e_{y,k}^2 = \frac{1}{2} \tilde{\theta}_k^\top \phi_k \phi_k^\top \tilde{\theta}_k, \quad (5)$$

where the subscript k in L_k denotes k th iteration. In the literature, a normalized gradient descent algorithm has been shown to be stable although having a slow convergence rate [2]

$$\theta_{k+1} = \theta_k - \alpha \frac{\nabla L_k(\theta_k)}{\mathcal{N}_k}, \quad 0 < \alpha < 2, \quad (6)$$

where \mathcal{N}_k is a normalizing signal and is defined as $\mathcal{N}_k = 1 + \|\phi_k\|^2$.

The following definitions will be utilized for proving the main results.

Definition 2.1. The regressor ϕ_k is said to satisfy the persistent excitation (PE) condition over an interval ΔT , if for all $k \geq 0$,

$$\epsilon_1 I \leq \sum_{i=k-\Delta T}^{k-1} \phi_i \phi_i^\top \leq \epsilon_2 I. \quad (7)$$

Definition 2.2 (From [14]). For any fixed $p \in [1, \infty)$, a sequence of scalars $\xi = \{\xi_0, \xi_1, \dots\}$ is defined to belong to ℓ_p if

$$\|\xi\|_\infty \equiv \left(\lim_{k \rightarrow \infty} \sum_{i=0}^k \|\xi_i\|^p \right)^{1/p} < \infty. \quad (8)$$

When $p = \infty$, $\xi \in \ell_\infty$ if

$$\|\xi\|_{\ell_\infty} \equiv \sup_{i \geq 0} \|\xi_i\| < \infty \quad (9)$$

Let $k \geq 0$ and consider the following time-varying dynamic system

$$x_{k+1} = f(k, x_k), \quad (10)$$

where $x_k \in \mathcal{D}$, $k \geq 0$, \mathcal{D} is an open set such that $0 \in \mathcal{D}$, $f : \mathbb{N} \times \mathcal{D} \rightarrow \mathbb{R}^n$ is continuous and for all $k \in \mathbb{N}$, $f(k, 0) = 0$. The following definition and theorem of uniform global exponential stability is modified from [15, Page 783-785].

Definition 2.3 (Uniform global exponential stability). The origin in (10) is uniformly globally exponentially stable if there exist scalars $c_1 > 0$ and $c_2 > 1$ such that $\|x_k\| \leq c_1 \|x_0\| \exp(-c_2^{-1}k)$, for all $x_0 \in \mathbb{R}^n$.

Theorem 2.1. *If there exist a continuous function $V : \mathbb{N}^+ \times \mathcal{D} \rightarrow \mathbb{R}$ and positive constants $\bar{\alpha}, \bar{\beta}, \bar{\gamma}$ such that*

$$\bar{\alpha} \|x\|^2 \leq V(k, x) \leq \bar{\beta} \|x\|^2, \quad (k, x) \in \mathbb{N}^+ \times \mathcal{D}, \quad (11)$$

$$\Delta V \leq -\bar{\gamma} \|x\|^2, \quad (k, x) \in \mathbb{N}^+ \times \mathcal{D}, \quad (12)$$

then the origin in (10) is uniformly globally exponentially stable.

3 The Algorithm

The Heavy Ball method, initially proposed in [13], achieves acceleration by adding a momentum term in addition to normalized gradient descent method

$$\theta_{k+1} = \theta_k - \bar{\gamma} \frac{\nabla L_k(\theta_k)}{\mathcal{N}_k} + \bar{\beta}(\theta_k - \theta_{k-1}), \quad (13)$$

where $\bar{\gamma}$ is the learning rate constant and $\bar{\beta}$ is a constant that controls the momentum. In this work, we consider a time-varying matrix F_k as an alternative to the constant $\bar{\gamma}$ in an effort to not only achieve fast convergence of the output error to zero but also have parameter error $\tilde{\theta}_k$ to zero. We propose the resulting algorithm as

$$\vartheta_{k+1} = \vartheta_k - F_k \frac{\nabla L_k(\theta_{k+1})}{\mathcal{N}_k}, \quad (14)$$

$$\theta_{k+1} = \theta_k - \beta(\theta_k - \vartheta_k), \quad (15)$$

where $\mathcal{N}_k = 1 + \eta \phi_k^\top F_{k-1} \phi_k$ and F_k is updated as

$$F_k = \lambda \left(F_{k-1} - \kappa \frac{F_{k-1} \phi_k \phi_k^\top F_{k-1}}{\mathcal{N}_k} \right). \quad (16)$$

In (14), (15) and (16), λ , κ , β and $\eta \geq \kappa$ are positive hyperparameters whose bounds will be defined later. The update of F_k is similar to the covariance matrix update in recursive least squares (RLS) algorithm with forgetting [2] but differs in the choice of the normalization and in the update of F_k . The main contribution of this paper is to show that the algorithm in (14)-(16) results in exponential convergence under PE.

4 Stability Analysis

In this section, we show that the algorithm in (14), (15) and (16) guarantees exponential convergence for suitable choices of the hyperparameters λ , κ , β and η . Let

$$\mu = \min\{c_1, c_2\}, \quad (17)$$

where

$$c_1 = \left(1 - \frac{1}{\lambda}\right) F_{\max}^{-1} \geq 0, \quad (18)$$

$$c_2 = \left\{1 - (1 - \beta)^2 \left[\frac{1}{\lambda} + \frac{\kappa}{\lambda(\eta - \kappa)} + \frac{4\lambda}{\eta^2}\right]\right\} F_{\max}^{-1} \geq 0, \quad (19)$$

and F_{\max} is the upper bound of F_k under the persistent excitation in Definition 2.1. When $\lambda = 1$, there is no forgetting in F_k and from (16), $F_k \leq F_{k-1}$. Therefore $F_{\max} = \sigma_{\max}\{F_0\}$. When $\lambda > 1$, the following lemma gives the upper bound for F_k .

Lemma 4.1. *When the regressor ϕ_k satisfies PE condition in Definition 2.1, the hyperparameters in (13), (14) and (15) satisfy $\frac{\kappa\epsilon_1(\lambda-1)}{\lambda(\lambda^{\Delta T}-1)} > (\eta - \kappa) \max_i \|\phi_i\|^2$ and $F_0 \leq \frac{F_{\max}}{\lambda^{\Delta T}-1} I$, there exists $F_{\max}^{-1} = \frac{\kappa\epsilon_1(\lambda-1)}{\lambda(\lambda^{\Delta T}-1)} - (\eta - \kappa) \max_i \|\phi_i\|^2 \in \mathbb{R}^+$ such that $F_k \leq F_{\max} I$ for all $k \geq 0$.*

Proof. Denote $\Delta_{\max} = 1 + (\eta - \kappa) F_{\max} \max_i \|\phi_i\|^2$. From (16), for all $k \geq 0$, $F_k^{-1} \leq \lambda F_{k+1}^{-1}$ and $\phi_k \phi_k^\top / \Delta_{\max} \leq \lambda F_k^{-1} / \kappa$. Since $F_0 \leq \frac{F_{\max}}{\lambda^{\Delta T}-1} I$, $F_k \leq F_{\max} I$ for all $0 \leq k \leq \Delta T - 1$. For all $k \geq \Delta T$, apply the PE definition, we obtain

$$\begin{aligned} \frac{\epsilon_1 I}{\Delta_{\max}} &\leq \sum_{i=k-\Delta T}^{k-1} \frac{\phi_i \phi_i^\top}{\Delta_{\max}} \\ &\leq \frac{\lambda}{\kappa} (1 + \lambda + \dots + \lambda^{\Delta T-1}) F_{k-1}^{-1} \\ &\leq \frac{\lambda}{\kappa} \frac{\lambda^{\Delta T} - 1}{\lambda - 1} F_{k-1}^{-1} \end{aligned}$$

Therefore $F_{k-1}^{-1} \geq F_{\max}^{-1} I$ for all $k \geq \Delta T$. □

Remark 1. Due to the differences in the update of the learning rates between our algorithm and RLS, certain constraints on the hyperparameters have to be assumed for proof of the upper bound of F_k . This is mainly due to the choice of the denominator \mathcal{N}_k in (16).

When $\lambda = 1$, it can be shown that under PE $F_k \rightarrow 0$ as $k \rightarrow \infty$. The following lemma gives a lower bound on F_k under PE when $\lambda > 1$.

Lemma 4.2. *When the regressor ϕ_k satisfies PE condition in Definition 2.1, there exists $F_{\min} \in \mathbb{R}^+$, where $F_{\min}^{-1} I = F_{\Delta T-1}^{-1} + \frac{\kappa\epsilon_2 I}{\lambda(1-1/\lambda^{\Delta T})}$, such that $F_k \geq F_{\min}$ for all $k \geq 0$.*

Proof. From (16), $F_k^{-1} \leq \lambda F_{k+1}^{-1}$, therefore for all $k \geq \Delta T$,

$$\begin{aligned}
F_k^{-1} &\leq \frac{1-1/\lambda}{1-1/\lambda^{\Delta T}} \sum_{i=k-1}^{k+\Delta T-2} F_{i+1}^{-1} \\
&\leq \frac{1-1/\lambda}{1-1/\lambda^{\Delta T}} \left(\frac{1}{\lambda} \sum_{i=k-1}^{k+\Delta T-2} F_i^{-1} + \frac{\kappa}{\lambda} \epsilon_2 I \right) \\
&\leq \frac{1-1/\lambda}{1-1/\lambda^{\Delta T}} \left(\frac{1}{\lambda^k} \sum_{i=0}^{\Delta T-1} F_i^{-1} + \frac{\kappa}{\lambda} \frac{1-1/\lambda^k}{1-1/\lambda} \epsilon_2 I \right) \\
&\leq \lambda^{\Delta T-k} F_{\Delta T-1}^{-1} + \frac{\kappa}{\lambda} \frac{1-1/\lambda^k}{1-1/\lambda^{\Delta T}} \epsilon_2 I \\
&\leq F_{\Delta T-1}^{-1} + \frac{\kappa \epsilon_2 I}{\lambda(1-1/\lambda^{\Delta T})} \\
&= F_{\min}^{-1} I
\end{aligned}$$

□

Remark 2. When $\lambda > 1$, from the expressions of F_{\max} and F_{\min} , we can observe that F_{\max} and ϵ_1 are inversely correlated, F_{\min} and ϵ_2 are inversely correlated. In the presence of weak excitation signals, F_{\max} increases and can potentially become infinite, which is similar to the covariance matrix update in RLS with forgetting.

The following theorem states accelerated learning properties of the proposed algorithm, and corresponds to the main result of this paper.

Theorem 4.3. *With $\lambda \geq 1$, $\kappa < 2\lambda$, $0 < \beta < 2$ and $\eta \geq \max \left\{ \frac{\lambda(\kappa+2\lambda)+\lambda\sqrt{5\kappa^2-4\lambda\kappa+4\lambda^2}}{2\lambda-\kappa}, \frac{4\lambda(1-\beta)^2}{\lambda-(1-\beta)^2} \right\}$, the update law in (14), (15) and (16) will result in (i) $\vartheta_k - \theta^* \in \ell_\infty$, $\theta_k - \vartheta_k \in \ell_\infty$, and (ii) $\|\vartheta_k - \theta^*\|^2 + \|\theta_k - \vartheta_k\|^2 \leq \exp(-\mu k) V_0$, where μ is defined in (17).*

Proof. Applying matrix inversion lemma to (16), we obtain

$$F_k^{-1} = \frac{1}{\lambda} F_{k-1}^{-1} + \frac{\kappa}{\lambda} \frac{\phi_k \phi_k^\top}{\mathcal{N}_k - \kappa \phi_k^\top F_{k-1} \phi_k} \quad (20)$$

Consider the candidate Lyapunov function

$$V_k = (\vartheta_k - \theta^*)^\top F_{k-1}^{-1} (\vartheta_k - \theta^*) + (\theta_k - \vartheta_k)^\top F_{k-1}^{-1} (\theta_k - \vartheta_k) \quad (21)$$

The increment $\Delta V_k := V_{k+1} - V_k$ may be expanded as

$$\begin{aligned}
\Delta V_k &= (\vartheta_{k+1} - \theta^*)^\top F_k^{-1} (\vartheta_{k+1} - \theta^*) + (\theta_{k+1} - \vartheta_{k+1})^\top F_k^{-1} (\theta_{k+1} - \vartheta_{k+1}) \\
&\quad - (\vartheta_k - \theta^*)^\top F_{k-1}^{-1} (\vartheta_k - \theta^*) - (\theta_k - \vartheta_k)^\top F_{k-1}^{-1} (\theta_k - \vartheta_k) \\
&= \left[\vartheta_k - \theta^* - F_k \frac{\nabla L_k(\theta_{k+1})}{\mathcal{N}_k} \right]^\top F_k^{-1} \left[\vartheta_k - \theta^* - F_k \frac{\nabla L_k(\theta_{k+1})}{\mathcal{N}_k} \right] \\
&\quad + \left[\theta_k - \beta(\theta_k - \vartheta_k) - \vartheta_k + F_k \frac{\nabla L_k(\theta_{k+1})}{\mathcal{N}_k} \right]^\top F_k^{-1} \left[\theta_k - \beta(\theta_k - \vartheta_k) - \vartheta_k + F_k \frac{\nabla L_k(\theta_{k+1})}{\mathcal{N}_k} \right] \\
&\quad - (\vartheta_k - \theta^*)^\top F_{k-1}^{-1} (\vartheta_k - \theta^*) - (\theta_k - \vartheta_k)^\top F_{k-1}^{-1} (\theta_k - \vartheta_k) \\
&= (\vartheta_k - \theta^*)^\top F_k^{-1} (\vartheta_k - \theta^*) + (1-\beta)^2 (\theta_k - \vartheta_k)^\top F_k^{-1} (\theta_k - \vartheta_k)
\end{aligned}$$

$$\begin{aligned}
& -\frac{2}{\mathcal{N}_k}(\vartheta_k - \theta^*)^\top \nabla L_k(\theta_{k+1}) + \frac{2(1-\beta)}{\mathcal{N}_k}(\theta_k - \vartheta_k)^\top \nabla L_k(\theta_{k+1}) \\
& + \frac{2}{\mathcal{N}_k^2} [\nabla L_k(\theta_{k+1})]^\top F_k \nabla L_k(\theta_{k+1}) \\
& - (\vartheta_k - \theta^*)^\top F_{k-1}^{-1}(\vartheta_k - \theta^*) - (\theta_k - \vartheta_k)^\top F_{k-1}^{-1}(\theta_k - \vartheta_k) \\
= & \frac{1}{\lambda}(\vartheta_k - \theta^*)^\top F_{k-1}^{-1}(\vartheta_k - \theta^*) + \frac{\kappa}{\lambda[1 + (\eta - \kappa)\phi_k^\top F_{k-1}\phi_k]}(\vartheta_k - \theta^*)^\top \phi_k \phi_k^\top (\vartheta_k - \theta^*) \\
& + \frac{(1-\beta)^2}{\lambda}(\theta_k - \vartheta_k)^\top F_{k-1}^{-1}(\theta_k - \vartheta_k) + \frac{\kappa(1-\beta)^2}{\lambda[1 + (\eta - \kappa)\phi_k^\top F_{k-1}\phi_k]}(\theta_k - \vartheta_k)^\top \phi_k \phi_k^\top (\theta_k - \vartheta_k) \\
& - \frac{2}{\mathcal{N}_k}(\vartheta_k - \theta^*)^\top \nabla L_k(\theta_{k+1}) + \frac{2(1-\beta)}{\mathcal{N}_k}(\theta_k - \vartheta_k)^\top \nabla L_k(\theta_{k+1}) \\
& + \frac{2}{\mathcal{N}_k^2} [\nabla L_k(\theta_{k+1})]^\top F_k \nabla L_k(\theta_{k+1}) \\
& - (\vartheta_k - \theta^*)^\top F_{k-1}^{-1}(\vartheta_k - \theta^*) - (\theta_k - \vartheta_k)^\top F_{k-1}^{-1}(\theta_k - \vartheta_k)
\end{aligned}$$

Now substitute $\nabla L_k(\theta_{k+1}) = \phi_k \phi_k^\top \tilde{\theta}_{k+1}$ into the above, we get

$$\begin{aligned}
\Delta V_k = & \frac{1}{\lambda}(\vartheta_k - \theta^*)^\top F_{k-1}^{-1}(\vartheta_k - \theta^*) + \frac{\kappa}{\lambda[1 + (\eta - \kappa)\phi_k^\top F_{k-1}\phi_k]} \|(\vartheta_k - \theta^*)^\top \phi_k\|^2 \\
& + \frac{(1-\beta)^2}{\lambda}(\theta_k - \vartheta_k)^\top F_{k-1}^{-1}(\theta_k - \vartheta_k) + \frac{\kappa(1-\beta)^2}{\lambda[1 + (\eta - \kappa)\phi_k^\top F_{k-1}\phi_k]} \|(\theta_k - \vartheta_k)^\top \phi_k\|^2 \\
& - \frac{2}{\mathcal{N}_k}(\vartheta_k - \theta^*)^\top \phi_k \phi_k^\top \tilde{\theta}_{k+1} + \frac{2(1-\beta)}{\mathcal{N}_k}(\theta_k - \vartheta_k)^\top \phi_k \phi_k^\top \tilde{\theta}_{k+1} \\
& + \frac{2}{\mathcal{N}_k^2} [\nabla L_k(\theta_{k+1})]^\top F_k \nabla L_k(\theta_{k+1}) \\
& - (\vartheta_k - \theta^*)^\top F_{k-1}^{-1}(\vartheta_k - \theta^*) - (\theta_k - \vartheta_k)^\top F_{k-1}^{-1}(\theta_k - \vartheta_k)
\end{aligned}$$

Since $\tilde{\theta}_{k+1} = \theta_{k+1} - \vartheta_k + \vartheta_k - \theta^* = (1-\beta)(\theta_k - \vartheta_k) + (\vartheta_k - \theta^*)$,

$$\begin{aligned}
\Delta V_k = & \frac{1}{\lambda}(\vartheta_k - \theta^*)^\top F_{k-1}^{-1}(\vartheta_k - \theta^*) + \frac{\kappa}{\lambda[1 + (\eta - \kappa)\phi_k^\top F_{k-1}\phi_k]} \|(\vartheta_k - \theta^*)^\top \phi_k\|^2 \\
& + \frac{(1-\beta)^2}{\lambda}(\theta_k - \vartheta_k)^\top F_{k-1}^{-1}(\theta_k - \vartheta_k) + \frac{\kappa(1-\beta)^2}{\lambda[1 + (\eta - \kappa)\phi_k^\top F_{k-1}\phi_k]} \|(\theta_k - \vartheta_k)^\top \phi_k\|^2 \\
& - \frac{2}{\mathcal{N}_k} \|(\vartheta_k - \theta^*)^\top \phi_k\|^2 - \frac{2(1-\beta)}{\mathcal{N}_k}(\vartheta_k - \theta^*)^\top \phi_k \phi_k^\top (\theta_k - \vartheta_k) \\
& + \frac{2(1-\beta)^2}{\mathcal{N}_k} \|(\theta_k - \vartheta_k)^\top \phi_k\|^2 + \frac{2(1-\beta)}{\mathcal{N}_k}(\theta_k - \vartheta_k)^\top \phi_k \phi_k^\top (\vartheta_k - \theta^*) \\
& + \frac{2}{\mathcal{N}_k^2} [\nabla L_k(\theta_{k+1})]^\top F_k \nabla L_k(\theta_{k+1}) \\
& - (\vartheta_k - \theta^*)^\top F_{k-1}^{-1}(\vartheta_k - \theta^*) - (\theta_k - \vartheta_k)^\top F_{k-1}^{-1}(\theta_k - \vartheta_k) \\
= & \frac{1}{\lambda}(\vartheta_k - \theta^*)^\top F_{k-1}^{-1}(\vartheta_k - \theta^*) + \frac{\kappa}{\lambda[1 + (\eta - \kappa)\phi_k^\top F_{k-1}\phi_k]} \|(\vartheta_k - \theta^*)^\top \phi_k\|^2 \\
& + \frac{(1-\beta)^2}{\lambda}(\theta_k - \vartheta_k)^\top F_{k-1}^{-1}(\theta_k - \vartheta_k) + \frac{\kappa(1-\beta)^2}{\lambda[1 + (\eta - \kappa)\phi_k^\top F_{k-1}\phi_k]} \|(\theta_k - \vartheta_k)^\top \phi_k\|^2 \\
& - \frac{2}{\mathcal{N}_k} \|(\vartheta_k - \theta^*)^\top \phi_k\|^2 + \frac{2(1-\beta)^2}{\mathcal{N}_k} \|(\theta_k - \vartheta_k)^\top \phi_k\|^2 + \frac{2}{\mathcal{N}_k^2} \tilde{\theta}_{k+1}^\top \phi_k \phi_k^\top F_k \phi_k \phi_k^\top \tilde{\theta}_{k+1} \\
& - (\vartheta_k - \theta^*)^\top F_{k-1}^{-1}(\vartheta_k - \theta^*) - (\theta_k - \vartheta_k)^\top F_{k-1}^{-1}(\theta_k - \vartheta_k)
\end{aligned}$$

Let $A_k = \frac{2\lambda}{\eta^2 \mathcal{N}_k^3} [(\eta - \kappa)\mathcal{N}_k + \kappa](\mathcal{N}_k - 1)$, the above becomes

$$\begin{aligned} \Delta V_k &= \frac{1}{\lambda} (\vartheta_k - \theta^*)^\top F_{k-1}^{-1} (\vartheta_k - \theta^*) + \frac{\kappa}{\lambda[1 + (\eta - \kappa)\phi_k^\top F_{k-1}\phi_k]} \|(\vartheta_k - \theta^*)^\top \phi_k\|^2 \\ &\quad + \frac{(1 - \beta)^2}{\lambda} (\theta_k - \vartheta_k)^\top F_{k-1}^{-1} (\theta_k - \vartheta_k) + \frac{\kappa(1 - \beta)^2}{\lambda[1 + (\eta - \kappa)\phi_k^\top F_{k-1}\phi_k]} \|(\theta_k - \vartheta_k)^\top \phi_k\|^2 \\ &\quad - \frac{2}{\mathcal{N}_k} \|(\vartheta_k - \theta^*)^\top \phi_k\|^2 + \frac{2(1 - \beta)^2}{\mathcal{N}_k} \|(\theta_k - \vartheta_k)^\top \phi_k\|^2 \\ &\quad + A_k [\|(\theta_k - \vartheta_k)^\top \phi_k\|^2 + \|(\vartheta_k - \theta^*)^\top \phi_k\|^2 + 2(\theta_k - \vartheta_k)^\top \phi_k \phi_k^\top (\vartheta_k - \theta^*)] \\ &\quad - (\vartheta_k - \theta^*)^\top F_{k-1}^{-1} (\vartheta_k - \theta^*) - (\theta_k - \vartheta_k)^\top F_{k-1}^{-1} (\theta_k - \vartheta_k) \end{aligned}$$

Combining similar terms,

$$\begin{aligned} \Delta V_k &= \left(\frac{1}{\lambda} - 1\right) (\vartheta_k - \theta^*)^\top F_{k-1}^{-1} (\vartheta_k - \theta^*) + \left[\frac{(1 - \beta)^2}{\lambda} - 1\right] (\theta_k - \vartheta_k)^\top F_{k-1}^{-1} (\theta_k - \vartheta_k) \\ &\quad + \left\{ \frac{\kappa\eta}{\lambda[(\eta - \kappa)\mathcal{N}_k + \kappa]} - \frac{2}{\mathcal{N}_k} + A_k \right\} \|(\vartheta_k - \theta^*)^\top \phi_k\|^2 \\ &\quad + \left\{ \frac{\kappa\eta}{\lambda[(\eta - \kappa)\mathcal{N}_k + \kappa]} + \frac{2}{\mathcal{N}_k} + A_k \right\} (1 - \beta)^2 \|(\theta_k - \vartheta_k)^\top \phi_k\|^2 \\ &\quad + 2A_k(1 - \beta)(\theta_k - \vartheta_k)^\top \phi_k \phi_k^\top (\vartheta_k - \theta^*) \\ &= \left(\frac{1}{\lambda} - 1\right) (\vartheta_k - \theta^*)^\top F_{k-1}^{-1} (\vartheta_k - \theta^*) + \left[\frac{(1 - \beta)^2}{\lambda} - 1\right] (\theta_k - \vartheta_k)^\top F_{k-1}^{-1} (\theta_k - \vartheta_k) \\ &\quad + \left\{ \frac{\kappa\eta}{\lambda[(\eta - \kappa)\mathcal{N}_k + \kappa]} - \frac{2}{\mathcal{N}_k} + 2A_k \right\} \|(\vartheta_k - \theta^*)^\top \phi_k\|^2 \\ &\quad + \left\{ \frac{\kappa\eta}{\lambda[(\eta - \kappa)\mathcal{N}_k + \kappa]} + \frac{2}{\mathcal{N}_k} + 2A_k \right\} (1 - \beta)^2 \|(\theta_k - \vartheta_k)^\top \phi_k\|^2 \\ &\quad - A_k [(\vartheta_k - \theta^*)^\top \phi_k - (1 - \beta)(\theta_k - \vartheta_k)^\top \phi_k]^2 \end{aligned}$$

From Cauchy-Schwarz inequality,

$$\frac{1}{\mathcal{N}_k} \|(\theta_k - \vartheta_k)^\top \phi_k\|^2 \leq \frac{1}{\eta} (\theta_k - \vartheta_k)^\top F_{k-1}^{-1} (\theta_k - \vartheta_k)$$

Therefore

$$\begin{aligned} \Delta V_k &\leq \left(\frac{1}{\lambda} - 1\right) (\vartheta_k - \theta^*)^\top F_{k-1}^{-1} (\vartheta_k - \theta^*) + \left[\frac{(1 - \beta)^2}{\lambda} - 1\right] (\theta_k - \vartheta_k)^\top F_{k-1}^{-1} (\theta_k - \vartheta_k) \\ &\quad + \left\{ \frac{\kappa\eta}{\lambda[(\eta - \kappa)\mathcal{N}_k + \kappa]} - \frac{2}{\mathcal{N}_k} + 2A_k \right\} \|(\vartheta_k - \theta^*)^\top \phi_k\|^2 \\ &\quad + \left\{ \frac{\kappa\eta\mathcal{N}_k}{\lambda[(\eta - \kappa)\mathcal{N}_k + \kappa]} + 2 + 2A_k\mathcal{N}_k \right\} \frac{(1 - \beta)^2}{\eta} (\theta_k - \vartheta_k)^\top F_{k-1}^{-1} (\theta_k - \vartheta_k) \\ &\quad - A_k [(\vartheta_k - \theta^*)^\top \phi_k - (1 - \beta)(\theta_k - \vartheta_k)^\top \phi_k]^2 \end{aligned}$$

Since

$$\frac{\kappa\eta\mathcal{N}_k}{\lambda[(\eta - \kappa)\mathcal{N}_k + \kappa]} \leq \frac{\kappa\eta}{\lambda(\eta - \kappa)},$$

and

$$A_k\mathcal{N}_k \leq \frac{2\lambda}{\eta}$$

and also since $\lambda \geq 1$, $\kappa < 2\lambda$, $0 < \beta < 2$, $\eta \geq \max \left\{ \frac{\lambda(\kappa+2\lambda)+\lambda\sqrt{5\kappa^2-4\lambda\kappa+4\lambda^2}}{2\lambda-\kappa}, \frac{4\lambda(1-\beta)^2}{\lambda-(1-\beta)^2} \right\}$, under persistent excitation and from Lemma 4.1, we get

$$\begin{aligned} \Delta V_k &\leq -c_1 \|\vartheta_k - \theta^*\|^2 - c_2 \|\theta_k - \vartheta_k\|^2 \\ &\quad - A_k [(\vartheta_k - \theta^*)^\top \phi_k - (1-\beta)(\theta_k - \vartheta_k)^\top \phi_k]^2 \\ &\leq 0 \end{aligned}$$

where c_1 and c_2 are defined in (18)-(19). Thus $\vartheta_k - \theta^* \in \ell_\infty$ and $\theta_k - \vartheta_k \in \ell_\infty$. Furthermore,

$$\Delta V_k \leq -\mu (\|\vartheta_k - \theta^*\|^2 + \|\theta_k - \vartheta_k\|^2),$$

where μ is defined in (17).

Since $F_{\max}^{-1}(\|\vartheta_k - \theta^*\|^2 + \|\theta_k - \vartheta_k\|^2) \leq V_k \leq F_{\min}^{-1}(\|\vartheta_k - \theta^*\|^2 + \|\theta_k - \vartheta_k\|^2)$, and $\Delta V_k \leq -\mu (\|\vartheta_k - \theta^*\|^2 + \|\theta_k - \vartheta_k\|^2)$, according to Theorem 2.1, $\|\vartheta_k - \theta^*\| \rightarrow 0$ and $\|\theta_k - \vartheta_k\| \rightarrow 0$ globally uniformly exponentially fast. \square

Remark 3. Theorem 4.3 states that under PE, both $\|\vartheta_k - \theta^*\|$ and $\|\theta_k - \vartheta_k\|$ go to zero exponentially fast. Together, $\|\vartheta_k - \theta^*\|^2 + \|\theta_k - \vartheta_k\|^2$ goes to zero exponentially fast.

Remark 4. Note that directly applying the covariance matrix update in RLS with forgetting to (14) leads to Lyapunov stability in only very limited cases such as when F_0 is small. It is the introduction of new hyperparameters in the update of F_k that gives more flexibility in the choice of hyperparameters and made global uniform exponential convergence possible.

Remark 5. From (16) and (20), under weak or no excitation, the eigenvalues of F_k keep increasing. To avoid the values of F_k getting too large when excitation is weak, a barrier function such as the one in [16] can be applied to the update of F_k in (16), or a variable forgetting factor such as the one introduced in [6] can be considered. A complete proof of this case is beyond the scope of this paper and will be addressed in future work.

Remark 6. Exponential decrease of V_k still happens when there is no persistent excitation in the regressors. However, that does not mean θ_k keeps converging to its true value. As an extreme example, when $\phi_k = 0$, V_k converges to zero exponentially, due to the exponential increase of F_k , but both θ_k and ϑ_k are not changing.

Remark 7. As we will show in Section 5, the additional benefits of the proposed algorithm compared to RLS with forgetting become apparent under weak excitation signals. The presence of momentum helps boost both output error and parameter error convergence towards zero.

5 Numerical Simulations

A linear discrete system is given by

$$G(\mathbf{q}) = \frac{-0.6213\mathbf{q} + 0.5839}{\mathbf{q}^2 - 1.8403\mathbf{q} + 0.8591}. \quad (22)$$

The problem is that the coefficients of (22) are unknown and we use the proposed algorithm to identify them.

5.1 Exponential Parameter Convergence Under Persistent Excitation

For identification, we apply the following signal as an input

$$u(k) = 1 + \sin\left(\frac{3\pi k}{4}\right) + \sin\left(\frac{2\pi k}{5}\right) + \sin\left(\frac{\pi k}{5}\right). \quad (23)$$

The proposed algorithm in (13)-(15) is tested. The forgetting factor in our algorithm is set to be $\lambda = 1.01$ and the initial values of the learning rate matrices are set to be $100I$. The regressors and estimated parameters are all set to be zero at the initial step. For the proposed algorithm, we set $\beta = 0.6$, $\kappa = 0.7$ and $\eta = 3.6$, satisfying the assumptions in Theorem 4.3.

Figure 1 and Figure 2 show output errors and parameter errors in semi-log scale, respectively. Figure 3 shows histories of the learning rate matrix in the algorithm. In this case, both the output error and parameter error converge exponentially towards zero.

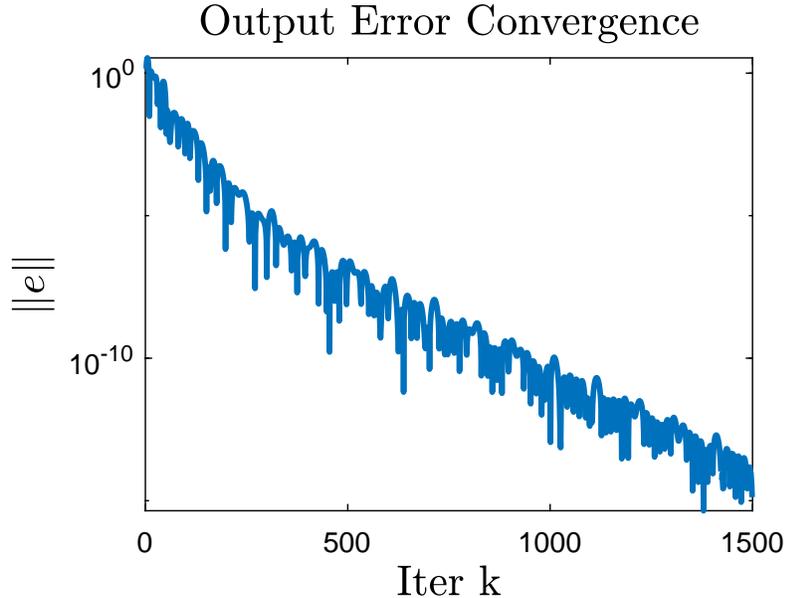


Figure 1: Output error $\|e\|$ of the proposed algorithm under PE.

5.2 Comparison To RLS With Forgetting Under Increasingly Weaker Excitation

RLS algorithm with a forgetting factor (RLS-FF) has been a widely used algorithm for online parameter estimation, and is quite similar to the algorithm proposed in this paper. In RLS-FF, the parameters are updated as follows [2]:

$$P_k = \frac{1}{\bar{\lambda}} P_{k-1} - \frac{P_{k-1} \phi_k \phi_k^\top P_{k-1}}{\bar{\lambda} + \phi_k^\top P_{k-1} \phi_k}, \quad (24)$$

$$\theta_{k+1} = \theta_k + \frac{P_{k-1} \phi_k (y_{k+1} - \phi_k^\top \theta_k)}{\bar{\lambda} + \phi_k^\top P_{k-1} \phi_k}, \quad (25)$$

where P_k is the covariance matrix and $\bar{\lambda}$ is the forgetting factor. The difference between P_k in (24) and F_k in (16) can be seen to be slight, but still makes a distinction as shown below. The other major difference is an HT aspect (see (14)-(15)) in our algorithm, while a gradient descent idea is employed in (25).

To demonstrate the added benefits of our algorithm compared to RLS with forgetting, we apply the following signal as an input to the system in (22):

$$u(k) =$$

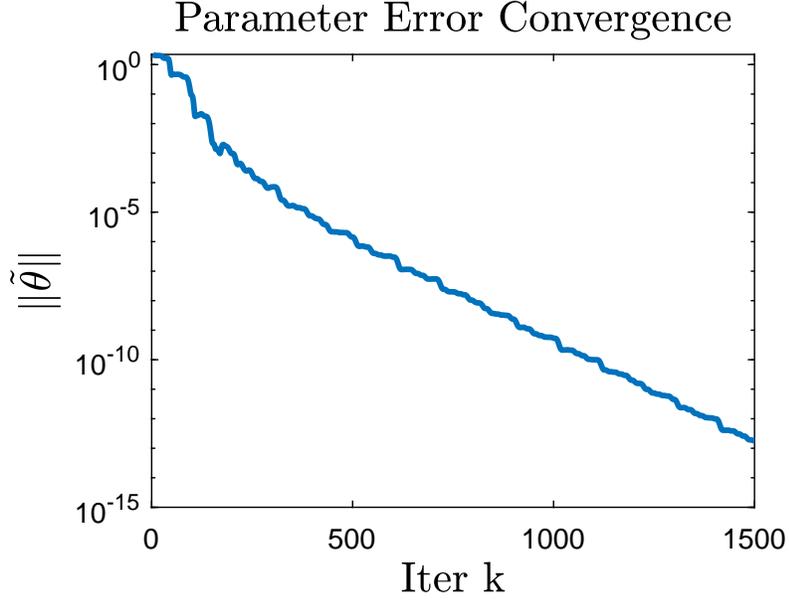


Figure 2: Parameter error $\|\tilde{\theta}\|$ of the proposed algorithm under PE.

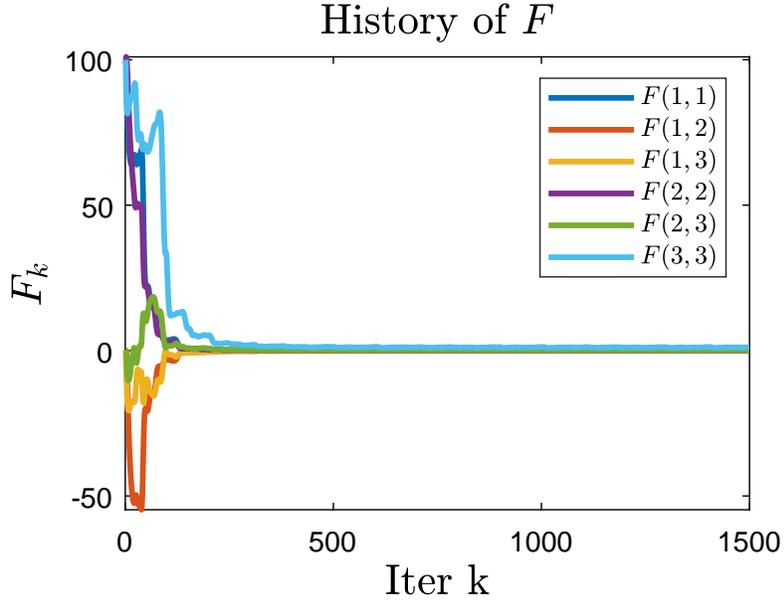


Figure 3: Learning rate matrix histories of the proposed algorithm.

$$1 + \exp(-0.03k) \left[\sin\left(\frac{3\pi k}{4}\right) + \sin\left(\frac{2\pi k}{5}\right) + \sin\left(\frac{\pi k}{5}\right) \right]$$

which is an increasingly weaker excitation signal. To ensure a fair comparison, we choose hyperparameters and initial values to be $\bar{\lambda} = 1/\lambda = 0.99$, $\kappa = 1.06$, $\eta = 3$, $\beta = 0.5$, $F_0 = P_0 = 100I$ and $\theta_0 = 0$ such that the time-varying matrices in the two algorithms are roughly the same as k increases, see Figure 4. Figure 5 shows the output error comparison and Figure 6 shows the parameter error

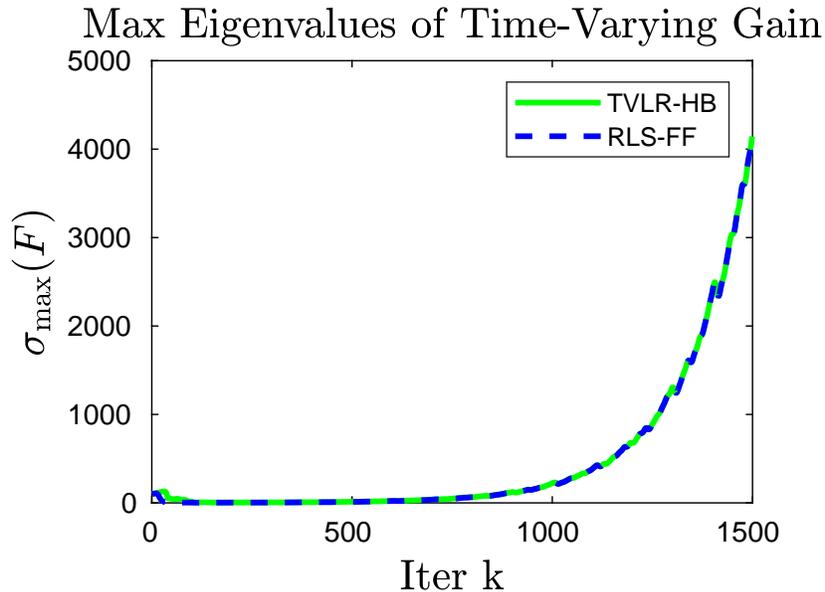


Figure 4: Comparison of the max eigenvalues of time-varying gain matrices between the proposed algorithm and RLS with forgetting.

comparison. Our algorithm demonstrates faster convergence results in this case. We speculate that the main reason for this faster convergence is the momentum term in the HB method, which in turn allows fast decrease in output error e . The time-varying F_k exploits persistent excitation and ensures that this fast decrease in performance error translates into fast decrease in learning error.

6 Conclusion

We introduced an online parameter estimation algorithm that adopts the ideas of momentum and time-varying learning rate. Under persistent excitation, the algorithm results in exponential convergence of the parameter error towards zero. Compared to recursive least squares with a forgetting factor, the presence of momentum in the update provides more flexibility. As shown in the simulation results, this flexibility translates into better performance and learning when the excitation is weak. Similar to recursive least squares with forgetting, one disadvantage of the algorithm is the unboundedness of the learning rate matrix when persistent excitation is not assured. In that case, projection operators need to be included to regulate the behavior of the learning rate matrix, which will be considered in future works.

References

- [1] K. S. Narendra and A. M. Annaswamy, *Stable Adaptive Systems*. NJ: Dover Publications, 2005, (original publication by Prentice-Hall Inc., 1989).
- [2] G. C. Goodwin and K. S. Sin, *Adaptive Filtering Prediction and Control*. Prentice Hall, 1984.
- [3] J. E. Gaudio, A. M. Annaswamy, E. Lavretsky, and M. Bolender, “Parameter estimation in adaptive control of time-varying systems under a range of excitation conditions,” *IEEE Transactions on Automatic Control*, 2021.

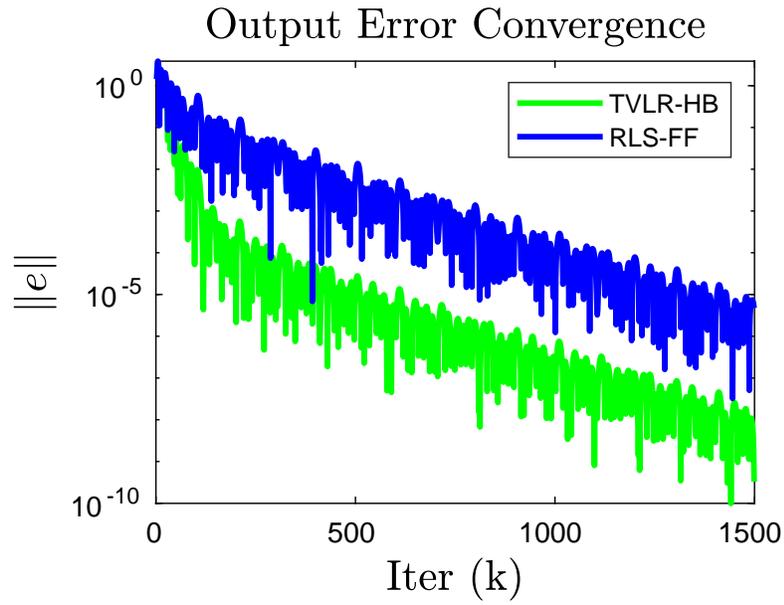


Figure 5: Comparison of output error convergence between the proposed algorithm and RLS with forgetting.

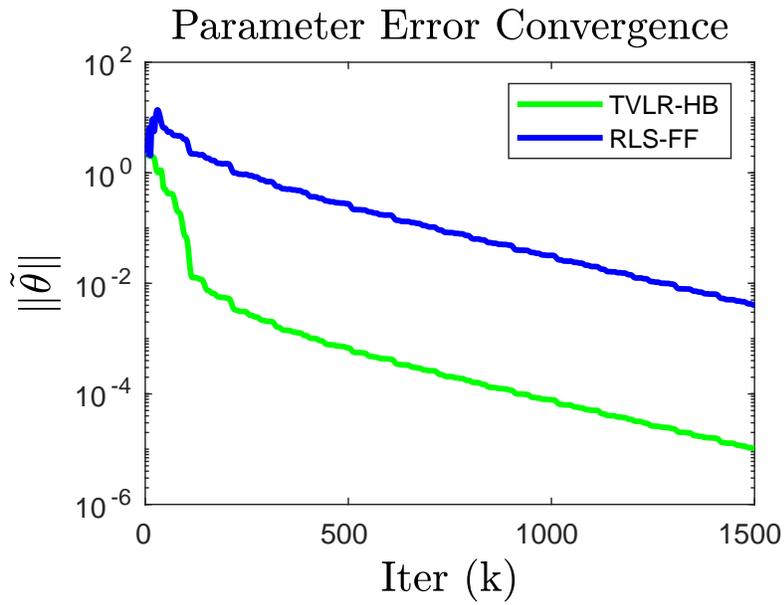


Figure 6: Comparison of parameter error convergence between the proposed algorithm and RLS with forgetting.

- [4] J. E. Gaudio, A. M. Annaswamy, J. M. Moreu, M. A. Bolender, and T. E. Gibson, “Accelerated learning with robustness to adversarial regressors,” *Proceedings of the 3rd Conference on Learning for Dynamics and Control, PMLR 144:636-650*, 2020.
- [5] A. L. Bruce, A. Goel, and D. S. Bernstein, “Recursive least squares with matrix forgetting,” in *2020 American Control Conference (ACC)*. IEEE, 2020, pp. 1406–1410.
- [6] —, “Convergence and consistency of recursive least squares with variable-rate forgetting,” *Automatica*, vol. 119, p. 109052, 2020.
- [7] A. Goel, A. L. Bruce, and D. S. Bernstein, “Recursive least squares with variable-direction forgetting: Compensating for the loss of persistency [lecture notes],” *IEEE Control Systems Magazine*, vol. 40, no. 4, pp. 80–102, 2020.
- [8] R. M. Johnstone, C. R. Johnson Jr, R. R. Bitmead, and B. D. Anderson, “Exponential convergence of recursive least squares with exponential forgetting factor,” *Systems & Control Letters*, vol. 2, no. 2, pp. 77–82, 1982.
- [9] J. Duchi, E. Hazan, and Y. Singer, “Adaptive subgradient methods for online learning and stochastic optimization,” *Journal of Machine Learning Research*, vol. 12, pp. 2121–2159, 2011.
- [10] M. D. Zeiler, “Adadelta: an adaptive learning rate method,” *arXiv preprint arXiv:1212.5701*, 2012.
- [11] S. S. Du, C. Jin, J. D. Lee, M. I. Jordan, A. Singh, and B. Póczos, “Gradient descent can take exponential time to escape saddle points,” *Advances in neural information processing systems*, vol. 30, 2017.
- [12] Y. Nesterov, *Lectures on Convex Optimization*. Springer International Publishing, 2018.
- [13] B. T. Polyak, “Some methods of speeding up the convergence of iteration methods,” *Ussr computational mathematics and mathematical physics*, vol. 4, no. 5, pp. 1–17, 1964.
- [14] D. G. Luenberger, *Optimization by vector space methods*. John Wiley & Sons, 1997.
- [15] V. Chellaboina and W. M. Haddad, *Nonlinear dynamical systems and control: A Lyapunov-based approach*. Princeton University Press, 2008.
- [16] Y. Cui, J. E. Gaudio, and A. M. Annaswamy, “New algorithms for discrete-time parameter estimation,” in *2022 American Control Conference (ACC)*. IEEE, 2022, pp. 3382–3387.