# Efficient model-based reinforcement learning for approximate online optimal control

Rushikesh Kamalapurkar, Joel A. Rosenfeld, and Warren E. Dixon

***Abstract*—In this paper the infinite horizon optimal regulation problem is solved online for a deterministic control-affine nonlinear dynamical system using the state following (StaF) kernel method to approximate the value function. Unlike traditional methods that aim to approximate a function over a large compact set, the StaF kernel method aims to approximate a function in a small neighborhood of a state that travels within a compact set. Simulation results demonstrate that stability and approximate optimality of the control system can be achieved with significantly fewer basis functions than may be required for global approximation methods.**

## I. Introduction

Reinforcement learning (RL) has become a popular tool for determining online solutions of optimal control problems for systems with finite state and action spaces [1]–[3]. Due to technical challenges, implementation of RL in systems with continuous state and action spaces has remained an open problem. In recent years, adaptive dynamic programming (ADP) has been successfully used to implement RL in deterministic autonomous control-affine systems to solve optimal control problems via value function approximation [3]–[13]. ADP techniques employ parametric function approximation (typically by employing neural networks (NNs)) to approximate the value function. Implementation of function approximation in ADP is challenging because the controller is void of pre-designed stabilizing feedback and is completely defined by the estimated parameters. Hence, the error between the optimal and the estimated value function is required to decay to a sufficiently small bound sufficiently fast to establish closed-loop stability. The size of the error bound is determined by the selected basis functions, and the convergence rate is determined by richness of the data used for learning.

Sufficiently accurate approximation of the value function over a sufficiently large neighborhood often requires a large number of basis functions, and hence, introduces a large number of unknown parameters. One way to achieve accurate function approximation with fewer unknown parameters is to use some knowledge about the system to determine the basis functions. However, for general nonlinear systems, prior

Rushikesh Kamalapurkar, Joel A. Rosenfeld, and Warren E. Dixon are with the Department of Mechanical and Aerospace Engineering, University of Florida, Gainesville, FL, USA. Email: {rkamalapurkar, joelar, wdixon}@ufl.edu.

knowledge of the features of the optimal value function is generally not available; hence, a large number of generic basis functions is often the only feasible option.

Sufficiently fast approximation of the value function over a sufficiently large neighborhood requires sufficiently rich data to be available for learning. In traditional ADP methods such as [9], [11], [14], richness of data manifests itself as the amount of excitation in the system. In experience replay-based techniques such as [15]–[18], richness of data is quantified by eigenvalues of the recorded history stack. In model-based RL techniques such as [19]–[21], richness of data corresponds to the eigenvalues of a learning matrix. As the dimension of the system and the number of basis functions increases, the required richness of data increases. In traditional ADP methods, the demand for richer data causes the designer to design increasingly aggressive excitation signals, thereby causing undesirable oscillations. Hence, implementation of traditional ADP techniques such as [3]–[14] in high dimensional systems are seldom found in the literature. In experience replay-based ADP methods and in model-based RL, the demand for richer data causes the required amount of data stored in the history stack, and the number of points selected to construct the learning matrix, respectively, to grow exponentially with the dimension of the system. Hence, implementation of data-driven ADP techniques such as [18]–[23] are scarcely found in the literature.

The contribution of this paper is the development of a novel model-based RL technique to achieve sufficient excitation without causing undesirable oscillations and expenditure of control effort like traditional ADP techniques and at a lower computational cost than state-of-the-art data-driven ADP techniques. Motivated by the fact that the computational effort required to implement ADP and the data-richness required to achieve convergence decrease with decreasing number of basis functions, this paper focuses on reduction of the number of basis functions used for value function approximation. A key contribution of this paper is the observation that online implementation of an ADP-based approximate optimal controller does not require an estimate of the optimal value function over the entire domain of operation of the system. Instead, only an estimate of the slope of the value function evaluated at the current state is required for feedback. Hence, estimation of the value function over a small neighborhood of the current state should be sufficient to implement an ADP-based approximate optimal controller. Furthermore, it is reasonable to postulate that approximation of the value function over a smaller local domain would require fewer basis functions as opposed to approximation over the entire domain of operation.

In this paper, reduction in the number of basis functions required for value function approximation is achieved via selection of basis functions that travel with the system state (referred to as state-following (StaF) kernels) to achieve accurate approximation of the value function over a small neighborhood of the state. The use of StaF kernel introduces a technical challenge owing to the fact that the ideal values of the unknown parameters corresponding to the StaF kernels are functions of the system state. The Lyapunov-based stability analysis presented in Section IV explicitly incorporates this functional relationship using the result that the ideal weights are continuously differentiable functions of the system state.

Sufficient exploration without the addition of an aggressive excitation signal is achieved via model-based RL based on BE extrapolation [19], [20]. The computational load associated with BE extrapolation is reduced via the selection of a single time-varying extrapolation function instead of a large number of autonomous extrapolation functions used in [19], [20]. Stability and convergence to optimality are obtained under a PE condition on the extrapolated regressor. Intuitively, selection of a single time-varying BE extrapolation function results in virtual excitation. That is, instead of using input-output data from a persistently excited system, the dynamic model is used to simulate persistent excitation to facilitate parameter convergence. Simulation results are included to demonstrate the effectiveness of the developed technique.

## II. StaF Kernel Functions

The objective in StaF-based function approximation is to maintain good approximation of the target function in a small region of interest in the neighborhood of a point of interest $x \in \mathbb{R}^n$. In state-of-the-art online approximate control, the optimal value function is approximated using a linear-in-the-parameters approximation scheme, and the approximate control law drives the system along the steepest negative gradient of the approximated value function. To compute the controller at the current state, only the gradient of the value function evaluated at the current state is required. Hence, in this application, the target function is the optimal value function, and the point of interest is the system state.

Since the system state evolves through the state-space with time, the region of interest for function approximation also evolves through the state-space. The StaF technique aims to maintain a uniform approximation of the value function over a small region around the current system state so that the gradient of the value function at the current state, and hence, the optimal controller at the current state, can be approximated.

To facilitate the theoretical development, this section summarizes key results from [24], where the theory of reproducing kernel Hilbert spaces (RKHSs) is used to establish continuous differentiability of the ideal weights with respect to the system state, and the postulate that approximation of the value function over a small neighborhood of the current state would require fewer basis functions is stated and proved.

To facilitate the discussion, let $H$ be a universal RKHS over a compact set $\chi \subset \mathbb{R}^n$ with a continuously differentiable positive definite kernel $k : \chi \times \chi \to \mathbb{R}$. Let $\overline{V}^* : \chi \to \mathbb{R}$ be a function such that $\overline{V}^* \in H$ . Let $c \triangleq [c_1, c_2, \cdots c_L]^T \in \chi^L$ be a set of distinct centers, and let $\sigma : \chi \times \chi^L \to \mathbb{R}^L$ be defined as $\sigma(x, c) = [k(x, c_1), \cdots, k(x, c_L)]^T$. Then, there exists a unique set of weights $W_H$ such that

$$W_H(c) = \arg \min_{a \in \mathbb{R}^L} \left\| a^T \sigma(\cdot, c) - \overline{V}^* \right\|_H,$$

where $\|\cdot\|_H$ denotes the Hilbert space norm.

In the StaF approach, the centers are selected to follow the current state $x$, i.e., $c(x) \triangleq [c_1(x), c_2(x), \cdots c_L(x)]^T : \chi \to \chi^L$. Since the system state evolves in time, the ideal weights are not constant. To approximate the ideal weights using gradient-based algorithms, it is essential that the weights change smoothly with respect to the system state.

Let $B_r(x) \subset \chi$ denote a closed ball of radius $r$ centered at the current state $x$. Let $H_{x,r}$ denote the restriction of the Hilbert space $H$ to $B_r(x)$. Then, $H_{x,r}$ is a Hilbert space with the restricted kernel $k_{x,r} : B_r(x) \times B_r(x) \to \mathbb{R}$ defined as $k_{x,r}(y, z) = k(y, z), \forall (y, z) \in B_r(x) \times B_r(x)$. The following result, first stated and proved in [24] is stated here to motivate the use of StaF kernels.

**Theorem 1.** *[24] Let $K(x, y) = e^{x^T y}$ be the exponential kernel function, which corresponds to an universal RKHS, and let $\epsilon, r > 0$. Then, for each $y \in \chi$, there exists a finite number of centers, $c_1, c_2, ..., c_{M_{y,\epsilon}} \in B_r(y)$ and weights $w_1, w_2, ..., w_{M_{y,\epsilon}}$ such that*

$$\left\| \overline{V}^*(x) - \sum_{i=1}^{M_{y,\epsilon}} w_i e^{x^T c_i} \right\|_{B_r(y), \infty} < \epsilon.$$

*If $p$ is an approximating polynomial that achieves the same approximation over $B_r(y)$ with degree $N_{y,\epsilon}$, then an asymptotically similar bound can be found with $M_{y,\epsilon}$ kernel functions, where $M_{y,\epsilon} < \binom{n + N_{y,\epsilon} + S_{y,\epsilon}}{N_{y,\epsilon} + S_{y,\epsilon}}$ for some constant $S_{y,\epsilon}$. Moreover, $N_{y,\epsilon}$ and $S_{y,\epsilon}$ can be bounded uniformly over $\chi$.*

The Weierstrass theorem indicates that as $r$ decreases, the degree $N_{y,\epsilon}$ of the polynomial needed to achieve the same error $\epsilon$ over $B_r(y)$ decreases [25]. Hence, by Theorem 1, approximation of a function over a smaller domain requires a smaller number of exponential kernels. Furthermore, provided the region of interest is small enough, the number of kernels required to approximate continuous functions with arbitrary accuracy can be reduced to $n + 2$ where $n$ is the state dimension.

The following result, first stated and proved in [24] is stated here to facilitate Lyapunov-based stability analysis of the closed-loop system.

**Theorem 2.** *[24] Let the kernel function $k$ be such that the functions $k(\cdot, c)$ are $l-$times continuously differentiable for all $c \in \chi$. Let $C$ be an ordered collection of $M$ distinct centers, $C = (c_1, c_2, ..., c_M) \in \chi^M$, with associated ideal weights*

$$W_H(C) = \arg \min_{a \in R^M} \left\| \sum_{i=1}^M a_i k(\cdot, c_i) - V(\cdot) \right\|_H.$$

*The function $W(C)$ is $l-$times continuously differentiable with respect to each component of $C$.*

Thus, if the kernels are selected as functions $c_i : \chi \to \chi$ of the state that are $l$−times continuously differentiable, then the ideal weight functions $W : \chi \to \mathbb{R}^L$ defined as $W(x) \triangleq W_{H_{x,r}}(c(x))$ are also $l$−times continuously differentiable.

Theorem 1 motivates the use of StaF kernels for model-based RL, and Theorem 2 facilitates implementation of gradient-based update laws to learn the time-varying ideal weights in real-time. In the following, the StaF-based function approximation approach is used to approximately solve an optimal regulation problem online using exact model knowledge via value function approximation. Selection of an optimal regulation problem and the assumption that the system dynamics are known are motivated by ease of exposition. Using a concurrent learning-based adaptive system identifier and the state augmentation technique developed in [20], the technique developed in this paper can be extended to a class of trajectory tracking problems in the presence of uncertainties in the system drift dynamics. Simulation results in Section V-B demonstrate the performance of such an extension.

## III. STaF Kernel Functions for Online Approximate Optimal Control

### A. Problem Formulation

Consider a control affine nonlinear dynamical system of the form

$$\dot{x}(t) = f(x(t)) + g(x(t)) u(t), \quad (1)$$

$t \in \mathbb{R}_{\geq t_0}$, where $t_0$ denotes the initial time, $x : \mathbb{R}_{\geq t_0} \to \mathbb{R}^n$ denotes the system state $f : \mathbb{R}^n \to \mathbb{R}^n$ and $g : \mathbb{R}^n \to \mathbb{R}^{n \times m}$ denote the drift dynamics and the control effectiveness, respectively, and $u : \mathbb{R}_{\geq 0} \to \mathbb{R}^m$ denotes the control input. The functions $f$ and $g$ are assumed to be locally Lipschitz continuous. Furthermore, $f(0) = 0$ and $\nabla f : \mathbb{R}^n \to \mathbb{R}^{n \times n}$ is continuous. In the following, the notation $\phi^u(t; t_0, x_0)$ denotes the trajectory of the system in (1) under the control signal $u$ with the initial condition $x_0 \in \mathbb{R}^n$ and initial time $t_0 \in \mathbb{R}_{\geq 0}$.

The control objective is to solve the infinite-horizon optimal regulation problem online, i.e., to design a control signal $u$ online to minimize the cost functional

$$J(x, u) \triangleq \int_{t_0}^{\infty} r(x(\tau), u(\tau)) \, d\tau, \quad (2)$$

under the dynamic constraint in (1) while regulating the system state to the origin. In (2), $r : \mathbb{R}^n \times \mathbb{R}^m \to \mathbb{R}_{\geq 0}$ denotes the instantaneous cost defined as

$$r(x^o, u^o) \triangleq Q(x^o) + u^{oT} R u^o, \quad (3)$$

for all $x^o \in \mathbb{R}^n$ and $u^o \in \mathbb{R}^m$, where $Q : \mathbb{R}^n \to \mathbb{R}_{\geq 0}$ is a positive definite function and $R \in \mathbb{R}^{m \times m}$ is a constant positive definite symmetric matrix. In (3) and in the reminder of this paper, the notation $(\cdot)^o$ is used to denote a dummy variable.

### B. Exact Solution

It is well known that since the functions $f$, $g$, and $Q$ are stationary (time-invariant) and the time-horizon is infinite, the optimal control input is a stationary state-feedback policy $u(t) = \xi(x(t))$ for some function $\xi : \mathbb{R}^n \to \mathbb{R}^m$. Furthermore, the function that maps each state to the total accumulated cost starting from that state and following a stationary state-feedback policy, i.e., the value function, is also a stationary function. Hence, the optimal value function $V^* : \mathbb{R}^n \to \mathbb{R}_{\geq 0}$ can be expressed as

$$V^*(x^o) \triangleq \inf_{u(\tau)|\tau \in \mathbb{R}_{\geq t}} \int_{t}^{\infty} r(\phi^u(\tau; t, x^o), u(\tau)) \, d\tau, \quad (4)$$

for all $x^o \in \mathbb{R}^n$, where $U \subset \mathbb{R}^m$ is a compact set. Assuming an optimal controller exists, the optimal value function can be expressed as

$$V^*(x^o) \triangleq \min_{u(\tau)|\tau \in \mathbb{R}_{\geq t}} \int_{t}^{\infty} r(\phi^u(\tau; t, x^o), u(\tau)) \, d\tau. \quad (5)$$

The optimal value function is characterized by the corresponding HJB equation [26]

$$0 = \min_{u^o \in U} \left( \nabla V(x^o) (f(x^o) + g(x^o) u^o) + r(x^o, u^o) \right), \quad (6)$$

for all $x^o \in \mathbb{R}^n$, with the boundary condition $V(0) = 0$. Provided the HJB in (6) admits a continuously differentiable solution, it constitutes a necessary and sufficient condition for optimality, i.e., if the optimal value function in (5) is continuously differentiable, then it is the unique solution to the HJB in (6) [27]. In (6) and in the following development, the notation $\nabla f(x, y, \cdots)$ denotes the partial derivative of $f$ with respect to the first argument. The optimal control policy $u^* : \mathbb{R}^n \to \mathbb{R}^m$ can be determined from (6) as [26]

$$u^*(x^o) \triangleq -\frac{1}{2} R^{-1} g^T(x^o) (\nabla V^*(x^o))^T. \quad (7)$$

The HJB in (6) can be expressed in the open-loop form

$$\nabla V^*(x^o) (f(x^o) + g(x^o) u^*(x^o)) + r(x^o, u^*(x^o)) = 0, \quad (8)$$

and using (7), the HJB in (8) can be expressed in the closed-loop form

$$-\frac{1}{4} \nabla V^*(x^o) g(x^o) R^{-1} g^T(x^o) (\nabla V^*(x^o))^T + \nabla V^*(x^o) f(x^o) + Q(x^o) = 0. \quad (9)$$

The optimal policy can now be obtained using (7) if the HJB in (9) can be solved for the optimal value function $V^*$.

### C. Value Function Approximation

An analytical solution of the HJB equation is generally infeasible; hence, an approximate solution is sought. In an approximate actor-critic-based solution, the optimal value function $V^*(x^o)$ is replaced by a parametric estimate $\hat{V}(x^o, W)$, where $W \in \mathbb{R}^L$ denotes the vector of ideal parameters. Replacing $V^*(x^o)$ by $\hat{V}(x^o, W)$ in (7), an approximation to the optimal policy $u^*(x^o)$ is obtained as $\hat{u}^o(x^o, W)$. The objective of the critic is to learn the parameters $W$, and the objective of the actor is to implement a stabilizing controller based on the parameters learned by the critic. Motivated by the stability analysis, the actor and the critic maintain separate

estimates $\hat{W}_a$ and $\hat{W}_c$, respectively, of the ideal parameters $W$. Substituting the estimates $\hat{V}$ and $\hat{u}$ for $V^*$ and $u^*$ in (8), respectively, a residual error $\delta : \mathbb{R}^n \times \mathbb{R}^L \times \mathbb{R}^L \to \mathbb{R}$, called the Bellman error (BE), is computed as

$$\delta \left( x^o, \hat{W}_c, \hat{W}_a \right) \triangleq r \left( x^o, \hat{u} \left( x^o, \hat{W}_a \right) \right)$$
$$+ \nabla \hat{V} \left( x^o, \hat{W}_c \right) \left( f(x^o) + g(x^o) \hat{u} \left( x^o, \hat{W}_a \right) \right).$$

To solve the optimal control problem, the critic aims to find a set of parameters $\hat{W}_c$ and the actor aims to find a set of parameters $\hat{W}_a$ such that $\delta \left( x^o, \hat{W}_c, \hat{W}_a \right) = 0, \ \forall x^o \in \mathbb{R}^n$. Since an exact basis for value function approximation is generally not available, an approximate set of parameters that minimizes the BE is sought.

The expression for the optimal policy in (7) indicates that to compute the optimal action when the system is at any given state $x^o \in \mathbb{R}^n$, one only needs to evaluate the gradient $\nabla V^*$ at $x^o$. Hence, to compute the optimal policy at any given state $x^o$, one only needs to approximate the value function over a small neighborhood around $x^o$. As established in Theorem 1, the number of basis functions required to approximate the value function is smaller if the approximation space is smaller in the sense of set containment. Hence, in this result, instead of aiming to obtain a uniform approximation of the value function over the entire operating domain, which might require a computationally intractable number of basis functions, the aim is to obtain a uniform approximation of the value function over a small neighborhood around the current system state.

StaF kernels are employed to achieve the aforementioned objective. To facilitate the development, let $\chi \subset \mathbb{R}^n$ be compact. Then, for all $\epsilon > 0$, there exists a function $\overline{V}^* = W^T(x^o) \sigma(x^o, c(x^o)) \in H$ such that $\sup_{x^o \in \chi} \left\| V^*(x^o) - \overline{V}^*(x^o) \right\| < \epsilon$, where $H$ is a universal RKHS, introduced in Section II and $W : \mathbb{R}^n \to \mathbb{R}^L$ denotes the ideal weight function. In the developed StaF-based method, a small compact set $B_r(x^o)$ around the current state $x^o$ is selected for value function approximation by selecting the centers $c^o$ such that $c^o = c(x^o) \in B_r(x^o)$ for some function $c : \chi \to \mathbb{R}^{nL}$. The approximate value function $\hat{V} : \chi \times \mathbb{R}^L \to \mathbb{R}$ and the approximate policy $\hat{u} : \chi \times \mathbb{R}^L \to \mathbb{R}$ can then be expressed as

$$\hat{V} \left( x^o, \hat{W}_c \right) \triangleq \hat{W}_c^T \sigma(x^o, c(x^o)),$$
$$\hat{u} \left( x^o, \hat{W}_a \right) \triangleq -\frac{1}{2} R^{-1} g^T(x^o) \nabla \sigma(x^o, c(x^o))^T \hat{W}_a, \quad (10)$$

where $\sigma : \chi \times \chi^L \to \mathbb{R}^L$ denotes the vector of basis functions introduced in Section II.

It should be noted that since the centers of the kernel functions change as the system state changes, the ideal weights also change as the system state changes. The state-dependent nature of the ideal weights differentiates this approach from state-of-the-art ADP methods in the sense that the stability analysis needs to account for changing ideal weights. Based on Theorem 2, it can be established that the ideal weight function $W$ defined as

$$W(x) \triangleq \arg \min_{a \in \mathbb{R}^L} \left\| a^T \sigma(\cdot, c(x)) - \overline{V}^*(\cdot) \right\|_{H_{x,r}},$$

is continuously differentiable with respect to the system state provided the functions $\sigma$ and $c$ are continuously differentiable.

### D. Online Learning Based on Simulation of Experience

To learn the ideal parameters online, the critic evaluates a form $\delta_t : \mathbb{R}_{\geq t_0} \to \mathbb{R}$ of the BE at each time instance $t$ as

$$\delta_t(t) \triangleq \delta \left( x(t), \hat{W}_c(t), \hat{W}_a(t) \right), \quad (11)$$

where $\hat{W}_a(t)$ and $\hat{W}_c(t)$ denote the estimates of the actor and the critic weights, respectively, at time $t$, and the notation $x(t)$ is used to denote the state the system in (1) at time $t$ when starting from initial time $t_0$, initial state $x_0$, and under the feedback controller

$$u(t) = \hat{u} \left( x(t), \hat{W}_a(t) \right). \quad (12)$$

Since (8) constitutes a necessary and sufficient condition for optimality, the BE serves as an indirect measure of how close the critic parameter estimates $\hat{W}_c$ are to their ideal values; hence, in RL literature, each evaluation of the BE is interpreted as gained experience. Since the BE in (11) is evaluated along the system trajectory, the experience gained is along the system trajectory.

Learning based on simulation of experience is achieved by extrapolating the BE to unexplored areas of the state space. The critic selects a set of functions $\{x_i : \mathbb{R}^n \times \mathbb{R}_{\geq 0} \to \mathbb{R}^n\}_{i=1}^N$ such that each $x_i$ maps the current state $x(t)$ to a point $x_i(x(t), t) \in B_r(x(t))$.

The critic then evaluates a form $\delta_{ti} : \mathbb{R}_{\geq t_0} \to \mathbb{R}$ of the BE for each $x_i$ as

$$\delta_{ti}(t) = \hat{W}_c^T(t) \omega_i(t) + r(x_i(x(t), t), \hat{u}_i(t)), \quad (13)$$

where

$$\hat{u}_i(t) \triangleq -\frac{1}{2} R^{-1} g^T(x_i(x(t), t))$$
$$\cdot \nabla \sigma(x_i(x(t), t), c(x(t)))^T \hat{W}_a(t),$$

and

$$\omega_i(t) \triangleq \nabla \sigma(x_i(x(t), t), c(x(t))) f(x_i(x(t), t))$$
$$- \frac{1}{2} \nabla \sigma(x_i(x(t), t), c(x(t))) g(x_i(x(t), t)) R^{-1} \cdot$$
$$g^T(x_i(x(t), t)) \nabla \sigma^T(x_i(x(t), t), c(x(t))) \hat{W}_a(t).$$

The critic then uses the BEs from (11) and (13) to improve the estimate $\hat{W}_c(t)$ using the recursive least-squares-based update law

$$\dot{\hat{W}}_c = -\eta_{c1} \Gamma(t) \frac{\omega(t)}{\rho(t)} \delta_t(t) - \frac{\eta_{c2}}{N} \Gamma(t) \sum_{i=1}^N \frac{\omega_i(t)}{\rho_i(t)} \delta_{ti}(t), \quad (14)$$

where

$$\omega(t) \triangleq \nabla \sigma(x(t), c(x(t))) f(x(t))$$
$$- \frac{1}{2} \nabla \sigma(x(t), c(x(t))) g(x(t)) R^{-1} g^T(x(t))$$
$$\cdot \nabla \sigma^T(x(t), c(x(t))) \hat{W}_a(t),$$

$\rho_i(t) \triangleq \sqrt{1 + \nu\omega_i^T(t)\omega_i(t)}$, $\rho(t) \triangleq \sqrt{1 + \nu\omega^T(t)\omega(t)}$, $\eta_{c1}, \eta_{c2}, \nu \in \mathbb{R}_{>0}$ are constant learning gains, and $\Gamma(t)$ denotes the least-square learning gain matrix updated according to

$$\dot{\Gamma}(t) = \beta\Gamma(t) - \eta_{c1}\Gamma(t)\frac{\omega(t)\omega^T(t)}{\rho^2(t)}\Gamma(t)$$
$$- \frac{\eta_{c2}}{N}\Gamma(t)\sum_{i=1}^{N}\frac{\omega_i(t)\omega_i^T(t)}{\rho_i^2(t)}\Gamma(t), \ \Gamma(0) = \Gamma_0. \quad (15)$$

In (15), $\beta \in \mathbb{R}_{>0}$ is a constant forgetting factor.

Motivated by a Lyapunov-based stability analysis, the actor improves the estimate $\hat{W}_a(t)$ using the update law

$$\dot{\hat{W}}_a(t) = -\eta_{a1}\left(\hat{W}_a(t) - \hat{W}_c(t)\right) - \eta_{a2}\hat{W}_a(t) +$$
$$\frac{\eta_{c1}G_\sigma^T(t)\hat{W}_a(t)\omega(t)^T}{4\rho(t)}\hat{W}_c(t)$$
$$+ \sum_{i=1}^{N}\frac{\eta_{c2}G_{\sigma i}^T(t)\hat{W}_a(t)\omega_i^T(t)}{4N\rho_i(t)}\hat{W}_c(t), \quad (16)$$

where $\eta_{a1}, \eta_{a2} \in \mathbb{R}_{>0}$ are learning gains,

$$G_\sigma(t) \triangleq \nabla\sigma(x(t), c(x(t)))g(x(t))R^{-1}g^T(x(t))$$
$$\cdot \nabla\sigma^T(x(t), c(x(t))),$$

and

$$G_{\sigma i}(t) \triangleq \nabla\sigma(x_i(x(t), t), c(x(t)))g(x_i(x(t), t))R^{-1}$$
$$\cdot g^T(x_i(x(t), t))\nabla\sigma^T(x_i(x(t), t), c(x(t))).$$

## IV. STABILITY ANALYSIS

For notational brevity, time-dependence of all the signals is suppressed hereafter. Let $B_\zeta \subset \mathbb{R}^{n+2L}$ denote a closed ball with radius $\zeta$ centered at the origin. Let $B_\chi \triangleq B_\zeta \cap \mathbb{R}^n$. Let the notation $\overline{\|(\cdot)\|}$ be defined as $\overline{\|h\|} \triangleq \sup_{\xi \in B_\chi}\|h(\xi)\|$, for some continuous function $h : \mathbb{R}^n \to \mathbb{R}^k$. To facilitate the subsequent stability analysis, the BEs in (11) and (13) are expressed in terms of the weight estimation errors $\tilde{W}_c \triangleq W - \hat{W}_c$ and $\tilde{W}_a = W - \hat{W}_a$ as

$$\delta_t = -\omega^T\tilde{W}_c + \frac{1}{4}\tilde{W}_a G_\sigma \tilde{W}_a + \Delta(x),$$
$$\delta_{ti} = -\omega_i^T\tilde{W}_c + \frac{1}{4}\tilde{W}_a^T G_{\sigma i}\tilde{W}_a + \Delta_i(x). \quad (17)$$

where the functions $\Delta, \Delta_i : \mathbb{R}^n \to \mathbb{R}$ are uniformly bounded over $B_\chi$ such that the bounds $\overline{\|\Delta\|}$ and $\overline{\|\Delta_i\|}$ decreases with decreasing $\overline{\|\nabla\epsilon\|}$. Let a candidate Lyapunov function $V_L : \mathbb{R}^{n+2L} \times \mathbb{R}_{\geq 0} \to \mathbb{R}$ be defined as

$$V_L(Z, t) \triangleq V^*(x) + \frac{1}{2}\tilde{W}_c^T\Gamma^{-1}(t)\tilde{W}_c + \frac{1}{2}\tilde{W}_a^T\tilde{W}_a,$$

where $V^*$ is the optimal value function, and

$$Z = \left[x^T, \tilde{W}_c^T, \tilde{W}_a^T\right]^T.$$

To facilitate learning, the system states $x$ or the selected functions $x_i$ are assumed to satisfy the following.

**Assumption 1.** There exists a positive constant $T \in \mathbb{R}_{>0}$ and nonnegative constants $\underline{c}_1, \underline{c}_2$, and $\underline{c}_3 \in \mathbb{R}_{\geq 0}$ such that

$$\underline{c}_1 I_L \leq \int_t^{t+T}\left(\frac{\omega(\tau)\omega^T(\tau)}{\rho^2(\tau)}\right)d\tau, \ \forall t \in \mathbb{R}_{\geq 0},$$

$$\underline{c}_2 I_L \leq \inf_{t \in \mathbb{R}_{\geq 0}}\left(\frac{1}{N}\sum_{i=1}^{N}\frac{\omega_i(t)\omega_i^T(t)}{\rho_i^2(t)}\right),$$

$$\underline{c}_3 I_L \leq \frac{1}{N}\int_t^{t+T}\left(\sum_{i=1}^{N}\frac{\omega_i(\tau)\omega_i^T(\tau)}{\rho_i^2(\tau)}\right)d\tau, \ \forall t \in \mathbb{R}_{\geq 0}.$$

Furthermore, at least one of $\underline{c}_1, \underline{c}_2$, and $\underline{c}_3$ is strictly positive.

*Remark* 1. Assumption 1 requires either the regressor $\omega$ or the regressor $\omega_i$ to be persistently exciting. The regressor $\omega$ is completely determined by the system state $x$, and the weights $\hat{W}_a$. Hence, excitation in $\omega$ vanishes as the system states and the weights converge. Hence, in general, it is unlikely that $\underline{c}_1 > 0$. However, the regressor $\omega_i$ depends on the functions $x_i$, which can be designed independent of the system state $x$. Hence, heuristically, $\underline{c}_3$ can be made strictly positive if the signal $x_i$ contains enough frequencies, and $\underline{c}_2$ can be made strictly positive by selecting a large number of extrapolation functions.

In previous model-based RL results such as [19], stability and convergence of the developed method relied on $\underline{c}_2$ being strictly positive. In the simulation example in Section V-A the extrapolation algorithm from [19] is used in the sense that large number of extrapolation functions is selected to make $\underline{c}_2$ strictly positive. In this example, the extrapolation algorithm from [19] is rendered computationally feasible by the fact that the value function is a function of only two variables. However, the number of extrapolation functions required to make $\underline{c}_2$ strictly positive increases exponentially with increasing state dimension. Hence, implementation of techniques such as [19] is rendered computationally infeasible in higher dimensions. In this paper, the computational efficiency of model-based RL is improved by allowing time-varying extrapolation functions that ensure that $\underline{c}_3$ is strictly positive, which can be achieved using a single extrapolation trajectory that contains enough frequencies. The performance of the developed extrapolation method is demonstrated in the simulation example in Section (V-B), where the value function is a function of four variables, and a single time-varying extrapolation point is used to improve computational efficiency instead of a large number of fixed extrapolation functions.

The following Lemma facilitates the stability analysis by establishing upper and lower bound on the eigenvalues of the least-squares learning gain matrix $\Gamma$.

**Lemma 1.** *Provided Assumption 1 holds and $\lambda_{\min}\{\Gamma_0^{-1}\} > 0$, the update law in (15) ensures that the least squares gain matrix satisfies*

$$\underline{\Gamma}I_L \leq \Gamma(t) \leq \overline{\Gamma}I_L, \quad (18)$$

*where* $\overline{\Gamma} = \frac{1}{\min\{\eta_{c1}\underline{c}_1 + \eta_{c2}\max\{\underline{c}_2 T, \underline{c}_3\}, \lambda_{\min}\{\Gamma_0^{-1}\}\}e^{-\beta T}}$ *and* $\underline{\Gamma} = \frac{1}{\lambda_{\max}\{\Gamma_0^{-1}\} + \frac{(\eta_{c1}+\eta_{c2})}{\beta\nu}}$. *Furthermore, $\overline{\Gamma} > 0$.*

*Proof:* The proof closely follows the proof of [28, Corollary 4.3.2]. The update law in (15) implies that $\frac{d}{dt}\Gamma^{-1}(t) = -\beta\Gamma^{-1}(t) + \eta_{c1}\frac{\omega(t)\omega^T(t)}{\rho^2(t)} + \frac{\eta_{c2}}{N}\sum_{i=1}^N \frac{\omega_i(t)\omega_i^T(t)}{\rho_i^2(t)}$. Hence,

$$\Gamma^{-1}(t) = e^{-\beta t}\Gamma_0^{-1} + \eta_{c1}\int_0^t e^{-\beta(t-\tau)}\frac{\omega(\tau)\omega^T(\tau)}{\rho^2(\tau)}d\tau$$

$$+ \frac{\eta_{c2}}{N}\int_0^t e^{-\beta(t-\tau)}\sum_{i=1}^N \frac{\omega_i(\tau)\omega_i^T(\tau)}{\rho_i^2(\tau)}d\tau$$

To facilitate the proof, let $t < T$. Then,

$$\Gamma^{-1}(t) \geq e^{-\beta t}\Gamma_0^{-1} \geq e^{-\beta T}\Gamma_0^{-1} \geq \lambda_{\min}\{\Gamma_0^{-1}\}e^{-\beta T}I_L.$$

If $t \geq T$, then since the integrands are positive, $\Gamma^{-1}$ can be bounded as

$$\Gamma^{-1}(t) \geq \eta_{c1}\int_{t-T}^t e^{-\beta(t-\tau)}\frac{\omega(\tau)\omega^T(\tau)}{\rho^2(\tau)}d\tau$$

$$+ \frac{\eta_{c2}}{N}\int_{t-T}^t e^{-\beta(t-\tau)}\sum_{i=1}^N \frac{\omega_i(\tau)\omega_i^T(\tau)}{\rho_i^2(\tau)}d\tau.$$

Hence,

$$\Gamma^{-1}(t) \geq \eta_{c1}e^{-\beta T}\int_{t-T}^t \frac{\omega(\tau)\omega^T(\tau)}{\rho^2(\tau)}d\tau$$

$$+ \frac{\eta_{c2}}{N}e^{-\beta T}\int_{t-T}^t \sum_{i=1}^N \frac{\omega_i(\tau)\omega_i^T(\tau)}{\rho_i^2(\tau)}d\tau.$$

Using Assumption 1,

$$\frac{1}{N}\int_{t-T}^t \sum_{i=1}^N \frac{\omega_i(\tau)\omega_i^T(\tau)}{\rho_i^2(\tau)}d\tau \geq \max\{\underline{c}_2 T, \underline{c}_3\}I_L,$$

$$\int_{t-T}^t \frac{\omega(\tau)\omega^T(\tau)}{\rho^2(\tau)}d\tau \geq \underline{c}_1 I_L.$$

Hence a lower bound for $\Gamma^{-1}$ is obtained as,

$$\Gamma^{-1}(t) \geq \min\Big\{\eta_{c1}\underline{c}_1 + \eta_{c2}\max\{\underline{c}_2 T, \underline{c}_3\},$$

$$\lambda_{\min}\{\Gamma_0^{-1}\}\Big\}e^{-\beta T}I_L. \quad (19)$$

Provided Assumption 1 holds, the lower bound in (19) is strictly positive. Furthermore, using the facts that $\frac{\omega(t)\omega^T(t)}{\rho^2(t)} \leq \frac{1}{\nu}$ and $\frac{\omega_i(t)\omega_i^T(t)}{\rho_i^2(t)} \leq \frac{1}{\nu}$ for all $t \in \mathbb{R}_{\geq 0}$,

$$\Gamma^{-1}(t) \leq e^{-\beta t}\Gamma_0^{-1} + \int_0^t e^{-\beta(t-\tau)}\left(\eta_{c1}\frac{1}{\nu} + \frac{\eta_{c2}}{N}\sum_{i=1}^N \frac{1}{\nu}\right)I_L d\tau,$$

$$\leq \left(\lambda_{\max}\{\Gamma_0^{-1}\} + \frac{(\eta_{c1}+\eta_{c2})}{\beta\nu}\right)I_L.$$

Since inverse of the lower and upper bounds on $\Gamma^{-1}$ are the upper and lower bounds on $\Gamma$, respectively, the proof is complete. ∎

Since the optimal value function is positive definite, (18) and [29, Lemma 4.3] can be used to show that the candidate Lyapunov function satisfies the following bounds

$$\underline{v_l}(\|Z^o\|) \leq V_L(Z^o, t) \leq \overline{v_l}(\|Z^o\|), \quad (20)$$

for all $t \in \mathbb{R}_{\geq t_0}$ and for all $Z^o \in \mathbb{R}^{2+2L}$. In (20), $\underline{v_l}, \overline{v_l} : \mathbb{R}_{\geq 0} \to \mathbb{R}_{\geq 0}$ are class $\mathcal{K}$ functions. To facilitate the analysis, let $\underline{c} \in \mathbb{R}_{>0}$ be a constant defined as

$$\underline{c} \triangleq \frac{\beta}{2\overline{\Gamma}\eta_{c2}} + \frac{c_2}{2}, \quad (21)$$

and let $\iota \in \mathbb{R}_{>0}$ be a constant defined as

$$\iota \triangleq \frac{3\left(\frac{(\eta_{c1}+\eta_{c2})\overline{\|\Delta\|}}{\sqrt{v}} + \frac{\overline{\|\nabla Wf\|}}{\underline{\Gamma}} + \frac{\overline{\|\Gamma^{-1}G_{W\sigma}W\|}}{2}\right)^2}{4\eta_{c2}\underline{c}}$$

$$+ \frac{1}{(\eta_{a1}+\eta_{a2})}\left(\frac{\overline{\|G_{W\sigma}W\|} + \overline{\|G_{V\sigma}\|}}{2} + \eta_{a2}\overline{\|W\|}\right.$$

$$\left. + \overline{\|\nabla Wf\|} + \frac{(\eta_{c1}+\eta_{c2})\overline{\|G_\sigma\|\|W\|}^2}{4\sqrt{v}}\right)^2$$

$$+ \frac{1}{2}\overline{\|G_{V\epsilon}\|}.$$

Let $v_l : \mathbb{R}_{\geq 0} \to \mathbb{R}_{\geq 0}$ be a class $\mathcal{K}$ function such that

$$v_l(\|Z\|) \leq \frac{Q(x)}{2} + \frac{\eta_{c2}\underline{c}}{6}\left\|\tilde{W}_c\right\|^2 + \frac{(\eta_{a1}+\eta_{a2})}{8}\left\|\tilde{W}_a\right\|^2.$$

The sufficient conditions for the subsequent Lyapunov-based stability analysis are given by

$$\frac{\eta_{c2}\underline{c}}{3} \geq \frac{\left(\frac{\overline{\|G_{W\sigma}\|}}{2\underline{\Gamma}} + \frac{(\eta_{c1}+\eta_{c2})\overline{\|W^T G_\sigma\|}}{4\sqrt{v}} + \eta_{a1}\right)^2}{(\eta_{a1}+\eta_{a2})},$$

$$\frac{(\eta_{a1}+\eta_{a2})}{4} \geq \left(\frac{\overline{\|G_{W\sigma}\|}}{2} + \frac{(\eta_{c1}+\eta_{c2})\overline{\|W\|\|G_\sigma\|}}{4\sqrt{v}}\right),$$

$$v_l^{-1}(\iota) < \overline{v_l}^{-1}(\underline{v_l}(\zeta)). \quad (22)$$

Note that the sufficient conditions can be satisfied provided the points for BE extrapolation are selected such that the minimum eigenvalue $\underline{c}$, introduced in (21) is large enough and that the StaF kernels for value function approximation are selected such that $\overline{\|\epsilon\|}$ and $\overline{\|\nabla\epsilon\|}$ are small enough. To improve computational efficiency, the size of the domain around the current state where the StaF kernels provide good approximation of the value function is desired to be small. Smaller approximation domain results in almost identical extrapolated points, which in turn, results in smaller $\underline{c}$. Hence, the approximation domain cannot be selected to be arbitrarily small and needs to be large enough to meet the sufficient conditions in (22).

**Theorem 3.** *Provided Assumption 1 holds and the sufficient gain conditions in (22) are satisfied, the controller in (12) and the update laws in (14) - (16) ensure that the state $x$ and the weight estimation errors $\tilde{W}_c$ and $\tilde{W}_a$ are ultimately bounded.*

*Proof:* The time-derivative of the Lyapunov function is given by

$$\dot{V}_L = \dot{V}^* + \tilde{W}_c^T \Gamma^{-1} \left( \dot{W} - \dot{\hat{W}}_c \right) + \frac{1}{2} \tilde{W}_c^T \dot{\Gamma}^{-1} \tilde{W}_c \\ + \tilde{W}_a^T \left( \dot{W} - \dot{\hat{W}}_a \right).$$

Using Theorem 2, the time derivative of the ideal weights can be expressed as

$$\dot{W} = \nabla W (x) (f (x) + g (x) u). \qquad (23)$$

Using (14) - (17) and (23), the time derivative of the Lyapunov function is expressed as

$$\dot{V}_L = \nabla V^* (x) (f (x) + g (x) u) \\ + \tilde{W}_c^T \Gamma^{-1} \nabla W (x) (f (x) + g (x) u) \\ - \tilde{W}_c^T \Gamma^{-1} \left( -\eta_{c1} \Gamma \frac{\omega}{\rho} \left( -\omega^T \tilde{W}_c + \frac{1}{4} \tilde{W}_a G_\sigma \tilde{W}_a + \Delta (x) \right) \right) \\ - \tilde{W}_c^T \Gamma^{-1} \left( -\frac{\eta_{c2}}{N} \Gamma \sum_{i=1}^N \frac{\omega_i}{\rho_i} \frac{1}{4} \tilde{W}_a^T G_{\sigma i} \tilde{W}_a \right) \\ - \tilde{W}_c^T \Gamma^{-1} \left( -\frac{\eta_{c2}}{N} \Gamma \sum_{i=1}^N \frac{\omega_i}{\rho_i} \left( -\omega_i^T \tilde{W}_c + \Delta_i (x) \right) \right) \\ - \frac{1}{2} \tilde{W}_c^T \Gamma^{-1} \left( \beta \Gamma - \eta_{c1} \Gamma \frac{\omega \omega^T}{\rho} \Gamma \right) \Gamma^{-1} \tilde{W}_c \\ - \frac{1}{2} \tilde{W}_c^T \Gamma^{-1} \left( -\frac{\eta_{c2}}{N} \Gamma \sum_{i=1}^N \frac{\omega_i \omega_i^T}{\rho_i} \Gamma \right) \Gamma^{-1} \tilde{W}_c \\ + \tilde{W}_a^T \left( \nabla W (x) (f (x) + g (x) u) - \dot{\hat{W}}_a \right).$$

Provided the sufficient conditions in (22) hold, the time derivative of the candidate Lyapunov function can be bounded as

$$\dot{V}_L \le -v_l (\|Z\|), \quad \forall \zeta > \|Z\| > v_l^{-1} (\iota). \qquad (24)$$

Using (20), (22), and (24), [29, Theorem 4.18] can be invoked to conclude that $Z$ is ultimately bounded, in the sense that $\limsup_{t \to \infty} \|Z (t)\| \le \underline{v_l}^{-1} (\overline{v_l} (\iota))$. ∎

## V. Simulation

### A. Optimal regulation problem with exact model knowledge

*1) Simulation parameters:* To demonstrate the effectiveness of the StaF kernels, simulations are performed on a two-dimensional nonlinear dynamical system. The system dynamics are given by (1), where $x^o = [x_1^o, x_2^o]^T$,

$$f (x^o) = \begin{bmatrix} -x_1^o + x_2^o \\ -\frac{1}{2} x_1^o - \frac{1}{2} x_2^o \left( \cos (2 x_1^o) + 2 \right)^2 \end{bmatrix},$$
$$g (x^o) = \begin{bmatrix} 0 \\ \cos (2 x_1^o) + 2 \end{bmatrix}. \qquad (25)$$

The control objective is to minimize the cost

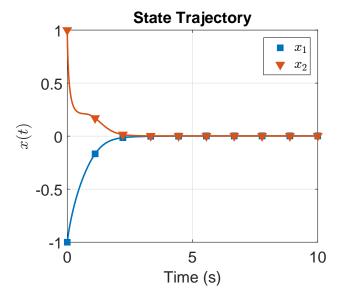$$\int_0^\infty \left( x^T (\tau) x (\tau) + u^2 (\tau) \right) d\tau. \qquad (26)$$



Figure 1. State trajectories generated using StaF kernel-based ADP.

The system in (25) and the cost in (26) are selected because the corresponding optimal control problem has a known analytical solution. The optimal value function is $V^* (x^o) = \frac{1}{2} x_1^{o2} + x_2^{o2}$, and the optimal control policy is $u^* (x^o) = -(\cos(2 x_1^o) + 2) x_2^o$ (cf. [9]).

To apply the developed technique to this problem, the value function is approximated using three exponential StaF kernels, i.e, $\sigma (x^o, c^o) = [\sigma_1 (x^o, c_1^o), \sigma_2 (x^o, c_2^o), \sigma_3 (x^o, c_3^o)]^T$. The kernels are selected to be $\sigma_i (x^o, c_i^o) = e^{x^{oT} c_i^o} - 1$, $i = 1, \cdots, 3$. The centers $c_i^o$ are selected to be on the vertices of a shrinking equilateral triangle around the current state, i.e., $c_i^o = x^o + d_i (x^o)$, $i = 1, \cdots, 3$, where $d_1 (x^o) = 0.7 \nu^o (x^o) \cdot [0, 1]^T$, $d_2 (x^o) = 0.7 \nu^o (x^o) \cdot [0.87, -0.5]^T$, and $d_3 (x^o) = 0.7 \nu^o (x^o) \cdot [-0.87, -0.5]^T$, and $\nu^o (x^o) \triangleq \left( \frac{x^{oT} x^o + 0.01}{1 + \nu_2 x^{oT} x^o} \right)$ denotes the shrinking function. The point for BE extrapolation is selected at random from a uniform distribution over a $2.1 \nu^o (x (t)) \times 2.1 \nu^o (x (t))$ square centered at the current state $x (t)$ so that the function $x_i$ is of the form $x_i (x^o, t) = x^o + a_i (t)$ for some $a_i (t) \in \mathbb{R}^2$.

The system is initialized at the initial conditions

$$x (0) = [-1, 1]^T, \ \hat{W}_c (0) = 0.4 \times \mathbf{1}_{3 \times 1},$$
$$\Gamma (0) = 500 I_3, \ \hat{W}_a (0) = 0.7 \hat{W}_c (0),$$

where $I_3$ denotes a $3 \times 3$ identity matrix and $\mathbf{1}_{3 \times 1}$ denotes a $3 \times 1$ matrix of ones. and the learning gains are selected as

$$\eta_{c1} = 0.001, \ \eta_{c2} = 0.25, \ \eta_{a1} = 1.2, \ \eta_{a2} = 0.01,$$
$$\beta = 0.003, \ v = 0.05, \ \nu_2 = 1.$$

*2) Results:* Figure 1 shows that the developed StaF-based controller drives the system states to the origin while maintaining system stability. Figure 2 shows the implemented control signal compared with the optimal control signal. It is clear that the implemented control converges to the optimal controller. Figure 3 shows that the weight estimates for the StaF-based value function and policy approximation remain bounded and converge as the state converges to the origin. Since the ideal
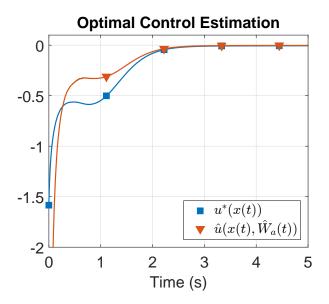
## Optimal Control Estimation



Figure 2. Control trajectory generated using StaF kernel-based ADP compared with the optimal control trajectory.
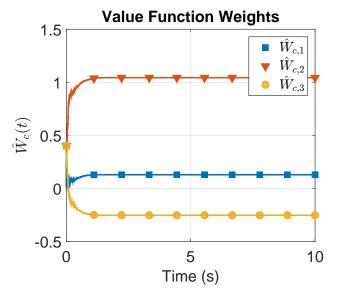
## Value Function Weights



Figure 3. Trajectories of the estimates of the unknown parameters in the value function generated using StaF kernel-based ADP. The ideal weights are unknown and time-varying; hence, the obtained weights can not be compared with their ideal weights.
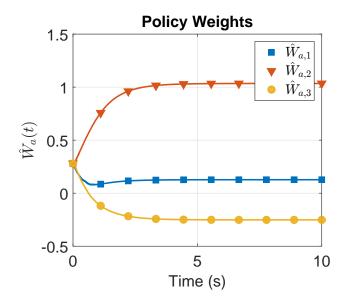
## Policy Weights



Figure 4. Trajectories of the estimates of the unknown parameters in the policy generated using StaF kernel-based ADP. The ideal weights are unknown and time-varying; hence, the obtained weights can not be compared with their ideal weights.
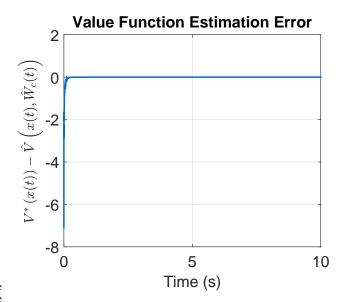
## Value Function Estimation Error



Figure 5. The error between the optimal and the estimated value function.

values of the weights are unknown, the weights can not directly be compared with their ideal values. However, since the optimal solution is known, the value function estimate corresponding to the weights in Figure 3 can be compared to the optimal value function at each time $t$. Figure 5 shows that the error between the optimal and the estimated value functions rapidly decays to zero.

### B. Optimal tracking problem with parametric uncertainties in the drift dynamics

*1) Simulation parameters:* Similar to [20], the developed StaF-based RL technique is extended to solve optimal tracking problems with parametric uncertainties in the drift dynamics.

The drift dynamics in the two-dimensional nonlinear dynamical system in (25) are assumed to be linearly parameterized as

$$f(x^o) = \underbrace{\begin{bmatrix} \theta_1 & \theta_2 & \theta_3 \\ \theta_4 & \theta_5 & \theta_6 \end{bmatrix}}_{\theta^T} \underbrace{\begin{bmatrix} x_1^o \\ x_2^o \\ x_2^o \left( cos \left( 2x_1^o \right) + 2 \right) \end{bmatrix}}_{\sigma_\theta(x^o)},$$

where $\theta \in \mathbb{R}^{3 \times 2}$ is the matrix of unknown parameters and $\sigma_\theta$ is the known vector of basis functions. The ideal values of the unknown parameters are $\theta_1 = -1$, $\theta_2 = 1$, $\theta_3 = 0$, $\theta_4 = -0.5$, $\theta_5 = 0$, and $\theta_6 = -0.5$. Let $\hat{\theta}$ denote an estimate of the unknown matrix $\theta$. The control objective is to drive the estimate $\hat{\theta}$ to the ideal matrix $\theta$, and to drive the state

$x$ to follow a desired trajectory $x_d$. The desired trajectory is selected to be solution of the initial value problem

$$\dot{x}_d(t) = \begin{bmatrix} -1 & 1 \\ -2 & 1 \end{bmatrix} x_d(t), \quad x_d(0) = \begin{bmatrix} 0 \\ 1 \end{bmatrix}, \qquad (27)$$

and the cost functional is selected to be $\int_0^\infty \left( e^T(t) \, \mathrm{diag}(10, 10) \, e(t) + (\mu(t))^2 \right) dt$, where $e(t) = x(t) - x_d(t)$, $\mu(t) = u(t) - g^+(x_d(t)) \left( \begin{bmatrix} -1 & 1 \\ -2 & 1 \end{bmatrix} x_d(t) - f(x_d(t)) \right)$, and $g^+(x^o)$ denotes the pseudoinverse of $g(x^o)$.

The value function is a function of the concatenated state $\zeta \triangleq \begin{bmatrix} e^T & x_d^T \end{bmatrix}^T \in \mathbb{R}^4$. The value function is approximated using five exponential StaF kernels given by $\sigma_i(\zeta^o, c_i^o)$, where the five centers are selected according to $c_i^o(\zeta^o) = \zeta^o + d_i(\zeta^o)$ to form a regular five dimensional simplex around the current state with $\nu^o(\zeta^o) \equiv 1$. Learning gains for system identification and value function approximation are selected as

$$\eta_{c1} = 0.001, \ \eta_{c2} = 2, \ \eta_{a1} = 2, \ \eta_{a2} = 0.001,$$
$$\beta = 0.01, \ \nu = 0.1, \ \nu_2 = 1, \ k = 500,$$
$$\Gamma_\theta = I_3, \ \Gamma(0) = 50 I_5, \ k_\theta = 20,$$

To implement BE extrapolation, a single state trajectory $\zeta_i$ is selected as $\zeta_i(\zeta^o, t) = \zeta^o + a_i(t)$, where $a_i(t)$ is sampled at each $t$ from a uniform distribution over the a $2.1 \times 2.1 \times 2.1 \times 2.1$ hypercube centered at the origin. The history stack required for CL contains ten points, and is recorded online using a singular value maximizing algorithm (cf. [17]), and the required state derivatives are computed using a fifth order Savitzky-Golay smoothing filter (cf. [30]).

The initial values for the state and the state estimate are selected to be $x(0) = [0,0]^T$ and $\hat{x}(0) = [0,0]^T$, respectively. The initial values for the NN weights for the value function, the policy, and the drift dynamics are selected to be $0.025 \times \mathbf{1}_5$, $0.025 \times \mathbf{1}_5$, and $\mathbf{0}_{3\times2}$, respectively, where $\mathbf{0}_{3\times2}$ denotes a $3 \times 2$ matrix of zeros. Since the system in (25) has no stable equilibria, the initial policy $\hat{\mu}(\zeta, \mathbf{0}_{3\times2})$ is not stabilizing. The stabilization demonstrated in Figure 6 is achieved via fast simultaneous learning of the system dynamics and the value function.

*2) Results:* Figures 6 and 7 demonstrate that the controller remains bounded and the tracking error is regulated to the origin. The NN weights are functions of the system state $\zeta$. Since $\zeta$ converges to a periodic orbit, the NN weights also converge to a periodic orbit (within the bounds of the excitation introduced by the BE extrapolation signal), as demonstrated in Figures 8 and 9. Figure 10 demonstrates that the unknown parameters in the drift dynamics, represented by solid lines, converge to their ideal values, represented by dashed lines.

*C. Comparison*

The developed technique is compared with the model-based RL method developed in [19] for regulation and [20] for tracking, respectively. Both the simulations are performed in MATLAB® SIMULINK® at 1000 Hz on the same machine. The regulation simulations run for 10 seconds of simulated
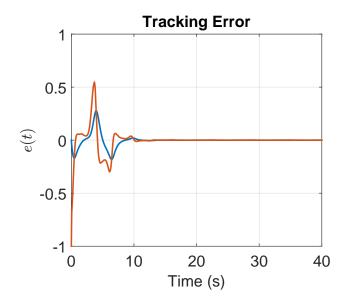


Figure 6. Tracking error trajectories generated using the proposed method for the nonlinear system.
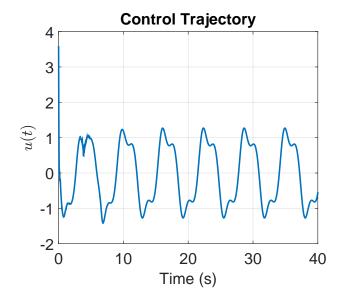


Figure 7. Control signal generated using the proposed method for the nonlinear system.

time, and the tracking simulations run for 40 seconds of simulated time. Tables I and II show that the developed controller requires significantly fewer computational resources than the controllers from [19] and [20].

Since the optimal solution for the regulation problem is known to be quadratic, the model-based RL method from [19] is implemented using three quadratic basis functions. Since the basis used is exact, the method from [19] yields a smaller steady-state error than the developed method, which uses three inexact, but generic StaF kernels. For the tracking problem, the method from [20] is implemented using ten polynomial basis functions selected based on a trial-and-error approach. The developed technique is implemented using five generic StaF kernels. In this case, since the optimal solution is unknown, both the methods use inexact basis functions, resulting in
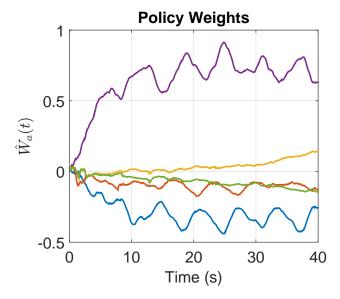
Figure 8. Policy weight trajectories generated using the proposed method for the nonlinear system. The weights do not converge to a steady-state value because the ideal weights are functions of the time-varying system state. Since an analytical solution of the optimal tracking problem is not available, weights cannot be compared against their ideal values
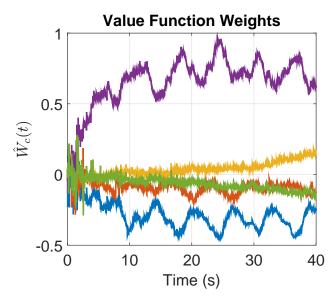


Figure 9. Value function weight trajectories generated using the proposed method for the nonlinear system. The weights do not converge to a steady-state value because the ideal weights are functions of the time-varying system state. Since an analytical solution of the optimal tracking problem is not available, weights cannot be compared against their ideal values

similar steady-state errors.

The two main advantages of StaF kernels are that they are universal, in the sense that they can be used to approximate a large class of value functions, and that they target local approximation, resulting in a smaller number of required basis functions. However, the StaF kernels trade optimality for universality and computational efficiency. The kernels are inexact, and the weight estimates need to be continually adjusted based on the system trajectory. Hence, as shown in Tables I and II, the developed technique results in a higher total cost than state-of-the-art model-based RL techniques.
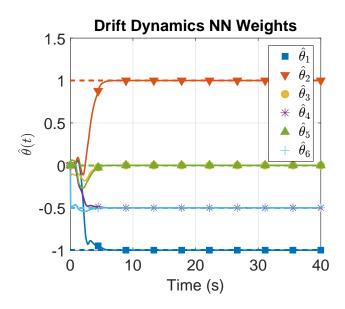


Figure 10. Trajectories of the unknown parameters in the system drift dynamics for the nonlinear system. The dotted lines represent the true values of the parameters.

| Method | Running time (s) | Total cost | Steady-state RMS error |
|---|---|---|---|
| StaF kernels with single moving extrapolation points | 0.95 | 2.82 | $2.5 \times 10^{-3}$ |
| Technique developed in [19] | 2 | 1.83 | $6.15 \times 10^{-6}$ |

Table I
REGULATION SIMULATION RUNNING TIMES FOR THE DEVELOPED
TECHNIQUE AND THE TECHNIQUE IN [19]

## VI. CONCLUSION

In this paper an infinite horizon optimal control problem is solved using a new approximation methodology called the StaF kernel method. Motivated by the fact that a smaller number of basis functions is required to approximate functions on smaller domains, the StaF kernel method aims to main-

| Method | Running time (s) | Total cost | Steady-state RMS error |
|---|---|---|---|
| StaF kernels with single moving extrapolation points | 15 | 6.38 | $2.13 \times 10^{-4}$ |
| Technique developed in [20] | 103 | 3.1 | $2.7 \times 10^{-4}$ |

Table II
TRACKING SIMULATION RUNNING TIMES FOR THE DEVELOPED
TECHNIQUE AND THE TECHNIQUE IN [20]

tain good approximation of the value function over a small neighborhood of the current state. Computational efficiency of model-based RL is improved by allowing selection of fewer time-varying extrapolation trajectories instead of a large number of autonomous extrapolation functions. Simulation results are presented that solve the infinite horizon optimal regulation and tracking problems online for a two state system using only three and five basis functions, respectively, via the StaF kernel method.

State-of-the-art solutions to solve infinite horizon optimal control problems online aim to approximate the value function over the entire operating domain. Since the approximate optimal policy is completely determined by the value function estimate, state-of-the-art solutions generate policies that are valid over the entire state space. Since the StaF kernel method aims at maintaining local approximation of the value function around the current system state, the StaF kernel method lacks memory, in the sense that the information about the ideal weights over a region of interest is lost when the state leaves the region of interest. Thus, unlike existing techniques, the StaF method generates a policy that is near-optimal only over a small neighborhood of the origin. A memory-based modification to the StaF technique that retains and reuses past information is a subject for future research.

## References

[1] R. S. Sutton and A. G. Barto, *Reinforcement Learning: An Introduction*. Cambridge, MA, USA: MIT Press, 1998.

[2] D. Bertsekas, *Dynamic Programming and Optimal Control*. Athena Scientific, 2007.

[3] P. Mehta and S. Meyn, "Q-learning and pontryagin's minimum principle," in *Proc. IEEE Conf. Decis. Control*, Dec. 2009, pp. 3598 –3605.

[4] K. Doya, "Reinforcement learning in continuous time and space," *Neural Comput.*, vol. 12, no. 1, pp. 219–245, 2000.

[5] R. Padhi, N. Unnikrishnan, X. Wang, and S. Balakrishnan, "A single network adaptive critic (SNAC) architecture for optimal control synthesis for a class of nonlinear systems," *Neural Netw.*, vol. 19, no. 10, pp. 1648–1660, 2006.

[6] A. Al-Tamimi, F. L. Lewis, and M. Abu-Khalaf, "Discrete-time nonlinear HJB solution using approximate dynamic programming: Convergence proof," *IEEE Trans. Syst. Man Cybern. Part B Cybern.*, vol. 38, pp. 943–949, 2008.

[7] F. L. Lewis and D. Vrabie, "Reinforcement learning and adaptive dynamic programming for feedback control," *IEEE Circuits Syst. Mag.*, vol. 9, no. 3, pp. 32–50, 2009.

[8] T. Dierks, B. Thumati, and S. Jagannathan, "Optimal control of unknown affine nonlinear discrete-time systems using offline-trained neural networks with proof of convergence," *Neural Netw.*, vol. 22, no. 5-6, pp. 851–860, 2009.

[9] K. Vamvoudakis and F. Lewis, "Online actor-critic algorithm to solve the continuous-time infinite horizon optimal control problem," *Automatica*, vol. 46, no. 5, pp. 878–888, 2010.

[10] H. Zhang, L. Cui, X. Zhang, and Y. Luo, "Data-driven robust approximate optimal tracking control for unknown general nonlinear systems using adaptive dynamic programming method," *IEEE Trans. Neural Netw.*, vol. 22, no. 12, pp. 2226–2236, 2011.

[11] S. Bhasin, R. Kamalapurkar, M. Johnson, K. Vamvoudakis, F. L. Lewis, and W. Dixon, "A novel actor-critic-identifier architecture for approximate optimal control of uncertain nonlinear systems," *Automatica*, vol. 49, no. 1, pp. 89–92, 2013.

[12] H. Zhang, L. Cui, and Y. Luo, "Near-optimal control for nonzero-sum differential games of continuous-time nonlinear systems using single-network adp," *IEEE Trans. Cybern.*, vol. 43, no. 1, pp. 206–216, 2013.

[13] H. Zhang, D. Liu, Y. Luo, and D. Wang, *Adaptive Dynamic Programming for Control Algorithms and Stability*, ser. Communications and Control Engineering. London: Springer-Verlag, 2013.

[14] K. Vamvoudakis and F. Lewis, "Online synchronous policy iteration method for optimal control," in *Recent Advances in Intelligent Control Systems*, W. Yu, Ed. Springer, 2009, pp. 357–374.

[15] G. Chowdhary, "Concurrent learning adaptive control for convergence without persistencey of excitation," Ph.D. dissertation, Georgia Institute of Technology, December 2010.

[16] G. Chowdhary and E. Johnson, "A singular value maximizing data recording algorithm for concurrent learning," in *Proc. American Control Conf.*, 2011, pp. 3547–3552.

[17] G. Chowdhary, T. Yucelen, M. Mühlegg, and E. N. Johnson, "Concurrent learning adaptive control of linear systems with exponentially convergent bounds," *Int. J. Adapt. Control Signal Process.*, vol. 27, no. 4, pp. 280–301, 2013.

[18] H. Modares, F. L. Lewis, and M.-B. Naghibi-Sistani, "Integral reinforcement learning and experience replay for adaptive optimal control of partially-unknown constrained-input continuous-time systems," *Automatica*, vol. 50, no. 1, pp. 193–202, 2014.

[19] R. Kamalapurkar, P. Walters, and W. E. Dixon, "Concurrent learning-based approximate optimal regulation," in *Proc. IEEE Conf. Decis. Control*, Florence, IT, Dec. 2013, pp. 6256–6261.

[20] R. Kamalapurkar, L. Andrews, P. Walters, and W. E. Dixon, "Model-based reinforcement learning for infinite-horizon approximate optimal tracking," in *Proc. IEEE Conf. Decis. Control*, 2014.

[21] R. Kamalapurkar, J. Klotz, and W. Dixon, "Concurrent learning-based online approximate feedback Nash equilibrium solution of N -player nonzero-sum differential games," *Acta Automatica Sinica*, to appear.

[22] B. Luo, H.-N. Wu, T. Huang, and D. Liu, "Data-based approximate policy iteration for affine nonlinear continuous-time optimal control design," *Automatica*, 2014.

[23] X. Yang, D. Liu, and Q. Wei, "Online approximate optimal control for affine non-linear systems with unknown internal dynamics using adaptive dynamic programming," *IET Control Theory Appl.*, vol. 8, no. 16, pp. 1676–1688, 2014.

[24] J. A. Rosenfeld, R. Kamalapurkar, and W. E. Dixon, "State following (StaF) kernel functions for function approximation part i: Theory and motivation," in *Proc. Am. Control Conf.*, 2015, to appear.

[25] G. G. Lorentz, *Bernstein polynomials*, 2nd ed. Chelsea Publishing Co., New York, 1986.

[26] D. Kirk, *Optimal Control Theory: An Introduction*. Dover, 2004.

[27] D. Liberzon, *Calculus of variations and optimal control theory: a concise introduction*. Princeton University Press, 2012.

[28] P. Ioannou and J. Sun, *Robust Adaptive Control*. Prentice Hall, 1996.

[29] H. K. Khalil, *Nonlinear Systems*, 3rd ed. Upper Saddle River, NJ, USA: Prentice Hall, 2002.

[30] A. Savitzky and M. J. E. Golay, "Smoothing and differentiation of data by simplified least squares procedures." *Anal. Chem.*, vol. 36, no. 8, pp. 1627–1639, 1964.