



Injection-Constrained State Estimation

Ankit Goel*

University of Maryland, Baltimore County, Maryland 21250

and

Dennis S. Bernstein†

University of Michigan, Ann Arbor, Michigan 48109

<https://doi.org/10.2514/1.G006108>

In applications of state estimation involving data assimilation over a spatial region, it is often convenient, and sometimes necessary, to confine the state correction to a prescribed subspace of the state space that corresponds to the measurement location. This is the injection-constrained state-estimation problem, where the injection of the output error is constrained to a specified subspace of the state space. Unlike full-state output-error injection, which is the dual of static full-state feedback, constrained output-error injection is the dual of static output feedback. To address the injection-constrained state-estimation problem, this paper develops the injection-constrained unscented Kalman filter (IC-UKF) and the injection-constrained retrospective cost filter (IC-RCF). The performance of these filters is evaluated numerically for linear and nonlinear state-estimation problems in order to compare their accuracy and determine their suboptimality relative to full-state output-error injection. As a benchmark test case, IC-UKF and IC-RCF are applied to the viscous Burgers equation for state and parameter estimation.

Nomenclature

e_i	=	i th column of the $n \times n$ identity matrix
$e_{k+1 k}$	=	prior error at step $k + 1$
$e_{k+1 k+1}$	=	posterior error at step $k + 1$
G_f	=	filter
I_n	=	identity matrix of size $n \times n$
K_k	=	filter gain
N_i	=	filter coefficient
$P_{k+1 k}$	=	prior error covariance at step $k + 1$
$P_{k+1 k+1}$	=	posterior error covariance at step $k + 1$
q	=	forward-shift operator
u_k	=	measured input
v_k	=	measurement noise
w_k	=	process noise
x_k	=	state
$\hat{x}_{k+1 k}$	=	prior state estimate at step $k + 1$
$\hat{x}_{k+1 k+1}$	=	posterior state estimate at step $k + 1$
y_k	=	measured output
$z_{k+1 k}$	=	prior output error at step $k + 1$
$z_{k+1 k+1}$	=	posterior output error at step $k + 1$
Γ_k	=	injection-constraint matrix
η_k	=	injection signal
λ	=	forgetting factor
$1_{n \times m}$	=	$n \times m$ matrix of ones

Subscripts

f	=	filter
k	=	iteration step

I. Introduction

THE classical Kalman filter and its variants, such as the extended Kalman filter (EKF) [1], unscented Kalman filter (UKF) [2,3],

and ensemble Kalman filter (EnKF) [4], construct state estimates by injecting the output error into a model of the system dynamics. In EKF, the state estimates are first propagated using the nonlinear dynamics of the system, and the propagated state is then corrected using the measurement data. The first step is referred to as the *physics update* or the *prediction step*, whereas the second step is referred to as *data assimilation* or the *correction step*. The Kalman gain used in the data-assimilation step is computed using the Jacobian of the dynamics along the estimated trajectory. In contrast, ensemble-based estimation methods, such as EnKF, UKF, and particle filters, propagate an ensemble of estimation models to compute the Kalman gain and the state estimate.

As the complexity of the system increases, both of these approaches become intractable. An example of such a system is the upper atmosphere of a planet. The global ionosphere–thermosphere model (GITM) is a model of the upper atmosphere that propagates the state of the atmosphere by solving coupled continuity, momentum, and energy equations in the computational domain comprising the atmosphere between the altitude of 100 and 600 km [5]. In a typical simulation, GITM propagates approximately 10 million states. EKF, EnKF, and UKF are thus not practical in this application.

Furthermore, large-scale complex models such as GITM often depend on parallel computing for simulation. In such cases, the computational domain is divided into blocks, and each block is solved on one processor. The processors propagate the states in each block independently of the other blocks except at the boundaries where constraints are imposed to maintain continuity. The application of the Kalman filter and its variants requires the collection of all of the states in the estimation model at a single processor to compute the Kalman gain and facilitate communication of the correction term to all of the processors. In addition to time and memory requirements, such an implementation also requires considerable programming expertise and effort. Thus, for high-dimensional models that depend on parallel computing for simulation, it is convenient to restrict the data injection to a subset of the processors, thus reducing the computational cost and programming effort.

Furthermore, in applications that encompass large spatial regions, such as in weather forecasting, the measurement data may be correlated with states within only a localized region. In such cases, it may be sufficient to confine the output-error injection to a subspace of the state space [6–10]. Finally, for ensemble-based estimation methods, localized subspace injection can also reduce the size of the ensemble.

For applications with linear dynamics, an injection-constrained Kalman filter was derived in [11]. This estimator uses a modified Riccati difference equation to update the error covariance, where an

Received 5 April 2021; revision received 5 December 2021; accepted for publication 6 December 2021; published online 8 January 2022. Copyright © 2021 by Ankit Goel. Published by the American Institute of Aeronautics and Astronautics, Inc., with permission. All requests for copying and permission to reprint should be submitted to CCC at www.copyright.com; employ the eISSN 1533-3884 to initiate your request. See also AIAA Rights and Permissions www.aiaa.org/randp.

*Assistant Professor, Department of Mechanical Engineering, UMBC, 1000 Hilltop Circle; ankgoel@umbc.edu.

†Professor, Department of Aerospace Engineering, UM Ann Arbor, 1320 Beal Avenue; dsbaero@umich.edu.

additional term involving an oblique projector accounts for the injection constraint.

Although [11] can be applied to the linearized dynamics as in the case of the EKF, the present paper focuses on alternative injection-constrained state-estimation techniques that are applicable to nonlinear systems without requiring linearized dynamics. To address this problem, the present paper develops two injection-constrained state estimators, namely, the injection-constrained unscented Kalman filter (IC-UKF) and the injection-constrained retrospective cost filter (IC-RCF). IC-UKF is an extension of UKF, where the data injection is constrained to a specified subspace. This constraint allows the filter gain to be computed using a smaller ensemble size, thus reducing the computational cost. In particular, IC-UKF requires propagation of an ensemble of $2l_\eta + 1$ copies of the model instead of the $2l_x + 1$ copies propagated by UKF, where l_η and l_x are the dimensions of the subspace used for output-error injection and the state, respectively. Like IC-UKF, IC-RCF constrains the data injection to a specified subspace; however, the filter gain is computed using retrospective cost optimization. A similar technique was used for state estimation in [12]. The goal of the present paper is thus to assess the performance of IC-UKF and IC-RCF for injection-constrained state estimation (ICSE).

A special case of ICSE is addressed by the Schmidt-Kalman filter (SKF) [13–15]. In its original form, SKF distinguishes between uncertain parameters and dynamic states, where the covariance of the parameters and states is propagated, but the output injection is confined to the dynamic states. A UKF extension of SKF is presented in [16].

The local ensemble Kalman filter (LEKF), which was motivated by atmospheric data assimilation, also localizes the effect of the measured data by restricting the data assimilation step to a subset of states [7, 17, 18]. In LEKF, the word local implies that the states to be updated are selected as the physical variables in the spatial vicinity surrounding the observation location. In this sense, LEKF can be interpreted as a special case of ICSE, where the injection constraint is motivated by the physical proximity of observation locations and states.

A key distinction between IC-UKF and IC-RCF is the fact that IC-UKF requires propagation of an ensemble of $2l_\eta + 1$ copies of the model. In contrast, IC-RCF requires only the output error, which can be computed using the propagation of only a single copy of the system model. Furthermore, implementing IC-UKF requires that the entire state be collected and assembled in order to propagate the ensemble and compute the filter gain. The programming effort required to implement IC-UKF can thus be prohibitively large for high-dimensional models consisting of numerous submodels. In contrast, IC-RCF is a modular scheme, which uses the past computed error but does not require the state of the model. Consequently, the programming effort required to implement IC-RCF is independent of the complexity of the application.

As a special case of ICSE for nonlinear systems, this paper considers the problem of estimating unknown parameters. By viewing the unknown parameters as constant states, the data injection is confined to the subspace corresponding to the unknown parameters. For the case of linear systems with unknown coefficients, this problem is typically addressed by means of the EKF [19].

The main contribution of the present paper is the development of injection-constrained state estimators for nonlinear systems, namely, the injection-constrained UKF and the IC-RCF. The extension of UKF to ICSE facilitates implementation of UKF and reduces the required ensemble size, whereas the IC-RCF provides an ensemble-free technique for ICSE in nonlinear systems. A numerical comparison of the performance of these filters is presented. A preliminary version of some of the results in the present paper appeared in [20].

This paper is organized as follows. The ICSE problem is described in Sec. II and a general form of injection-constrained filter is introduced. Section III presents the optimal injection-constrained filter for linear systems, and Secs. IV and V present the injection-constrained UKF and retrospective cost filter for nonlinear systems. Numerical examples comparing the performance of the various filters are presented in Sec. VI. Finally, conclusions are discussed in Sec. VIII.

II. Injection-Constrained State Estimation Problem

Consider the system

$$x_{k+1} = f_k(x_k) + w_k \quad (1)$$

$$y_k = g_k(x_k) + v_k \quad (2)$$

where, for all $k \geq 0$, $x_k \in \mathbb{R}^{l_x}$, $y_k \in \mathbb{R}^{l_y}$, $f_k: \mathbb{R}^{l_x} \rightarrow \mathbb{R}^{l_x}$, and $g_k: \mathbb{R}^{l_x} \rightarrow \mathbb{R}^{l_y}$. Furthermore, let $x_0 \sim \mathcal{N}(\bar{x}_0, P_{0|0})$, and, for all $k \geq 0$, let $w_k \sim \mathcal{N}(0, Q_k)$ and $v_k \sim \mathcal{N}(0, R_k)$ denote the disturbance and sensor noise, respectively. The goal is to estimate x_k using knowledge of the functions f and g as well as the measurement y_k ; this is the ICSE problem. A special case of Eqs. (1) and (2) is the linear system

$$x_{k+1} = A_k x_k + w_k \quad (3)$$

$$y_k = C_k x_k + v_k \quad (4)$$

where, for all $k \geq 0$, A_k, C_k are real matrices.

Now consider the *injection-constrained state estimator*

$$\hat{x}_{k+1|k} = \bar{f}_k(\hat{x}_{k|k}) \quad (5)$$

$$\hat{x}_{k+1|k+1} = \hat{x}_{k+1|k} + \Gamma_k \eta_{k+1} \quad (6)$$

where $\hat{x}_{k+1|k}$ is the prior estimate of x_{k+1} ; $\hat{x}_{k+1|k+1}$ is the posterior estimate of x_{k+1} ; $\eta_{k+1} \in \mathbb{R}^{l_\eta}$, where $l_\eta < l_x$, is the *injection signal*; and, for all $k \geq 0$, $\Gamma_k \in \mathbb{R}^{l_x \times l_\eta}$ is the *injection-constraint matrix* that constrains the injection of η_{k+1} to a specified subspace. In ensemble-based filters,

$$\bar{f}_k(x) = \sum_{i=1}^n \sigma_i f_k(x + \varepsilon_i) \quad (7)$$

where n is the ensemble size, σ_i is a scalar, and ε_i is the state perturbation; otherwise, $\bar{f}_k = f_k$. For all $k \geq 0$, Γ_k is assumed to have full-column rank. For example, the $l_x \times l_\eta$ injection-constraint matrix

$$\Gamma_k = \begin{bmatrix} I_{l_\eta} \\ 0_{l_x \times l_\eta} \end{bmatrix} \quad (8)$$

where $l_\rho \triangleq l_x - l_\eta$, constrains the injection of η_{k+1} to the first l_η components of x_k . Note that the injection-constraint matrix can always be expressed in the form (8) by permuting the components of state x_k . In this paper, the columns of Γ_k are assumed to be selected from the columns of the $l_x \times l_x$ identity matrix. Consequently, with appropriate permutation of the state in Eq. (1), Γ_k can always be transformed in to the form given by Eq. (8). Finally, when Γ_k is constant, it is written as Γ for convenience.

The injection signal η_{k+1} is given by

$$\eta_{k+1} = \hat{K}[y_{k+1} - g_k(\hat{x}_{k+1|k})] \quad (9)$$

where \hat{K} is determined by optimization below and $g_k(\hat{x}_{k+1|k})$ is the predicted output based on the prior estimate $\hat{x}_{k+1|k}$. Note that the notation η_{k+1} for the injection signal reflects the fact that the signal is injected at step $k + 1$.

III. Optimal Injection-Constrained Filter for Linear Systems

This section reviews the optimal injection-constrained filter (OICF) presented in [11] to establish notation for the development of injection-constrained state estimators. For the linear system described by (3), (4), consider the injection-constrained filter

$$\hat{x}_{k+1|k} = A_k \hat{x}_{k|k} \quad (10)$$

$$\hat{x}_{k+1|k+1} = \hat{x}_{k+1|k} + \Gamma_k \hat{K} (y_{k+1} - C_{k+1} \hat{x}_{k+1|k}) \quad (11)$$

where the gain matrix $\hat{K} \in \mathbb{R}^{l_y \times l_x}$ is determined by optimization below. For all $k \geq 0$, define the prior error $e_{k+1|k}$ and the posterior error $e_{k+1|k+1}$ by

$$e_{k+1|k} \triangleq x_{k+1} - \hat{x}_{k+1|k} \quad (12)$$

$$e_{k+1|k+1} \triangleq x_{k+1} - \hat{x}_{k+1|k+1} \quad (13)$$

and the corresponding covariances of $e_{k+1|k}$ and $e_{k+1|k+1}$ by

$$P_{k+1|k} \triangleq \mathbb{E} \left[e_{k+1|k} e_{k+1|k}^T \right] \quad (14)$$

$$P_{k+1|k+1} \triangleq \mathbb{E} \left[e_{k+1|k+1} e_{k+1|k+1}^T \right] \quad (15)$$

The following result is given in [11]. For completeness, an abbreviated proof is given.

Proposition III.1: For all $k \geq 0$, let the prior covariance $P_{k+1|k}$ and the posterior covariance $P_{k|k}$ be given by Eqs. (14) and (15), respectively. Let K_{k+1} denote the minimizer of $\text{tr} P_{k+1|k+1}$. Then, for all $k \geq 0$, the optimal injection-constrained filter gain K_{k+1} is given by

$$K_{k+1} = (\Gamma_k^T \Gamma_k)^{-1} \Gamma_k^T P_{k+1|k} C_{k+1}^T \bar{R}_{k+1}^{-1} \quad (16)$$

where

$$\bar{R}_{k+1} \triangleq C_{k+1} P_{k+1|k} C_{k+1}^T + R_{k+1} \quad (17)$$

and the corresponding posterior covariance at step $k+1$ is given by

$$P_{k+1|k+1} = P_{k+1|k} - P_{k+1|k} C_{k+1}^T \bar{R}_{k+1}^{-1} C_{k+1} P_{k+1|k} + \pi_k P_{k+1|k} C_{k+1}^T \bar{R}_{k+1}^{-1} C_{k+1} P_{k+1|k} \pi_k \quad (18)$$

where

$$P_{k+1|k} = A_k P_{k|k} A_k^T + Q_k \quad (19)$$

and

$$\pi_k \triangleq I - \Gamma_k (\Gamma_k^T \Gamma_k)^{-1} \Gamma_k^T \quad (20)$$

Proof: Note that the prior and posterior errors satisfy

$$e_{k+1|k} = A_k e_{k|k} + w_k, \\ e_{k+1|k+1} = (I - \Gamma_k \hat{K} C_{k+1}) e_{k+1|k} - \Gamma_k \hat{K} v_{k+1}$$

and thus the prior and posterior covariances satisfy

$$P_{k+1|k} = A_k P_{k|k} A_k^T + Q_k, \\ P_{k+1|k+1} = P_{k+1|k} + \Gamma_k \hat{K} \bar{R}_{k+1}^{-1} \hat{K}^T \Gamma_k^T - \Gamma_k \hat{K} C_{k+1} P_{k+1|k} - P_{k+1|k} C_{k+1}^T \hat{K}^T \Gamma_k^T \quad (21)$$

To minimize $\text{tr} P_{k+1|k+1}$, note that

$$\frac{d}{d\hat{K}} \text{tr} P_{k+1|k+1} = 2\Gamma_k^T \Gamma_k \hat{K} \bar{R}_{k+1}^{-1} - 2\Gamma_k^T P_{k+1|k} C_{k+1}^T \quad (22)$$

Setting Eq. (22) to zero yields the optimal injection-constrained filter gain (16), and setting $\hat{K} = K_{k+1}$ in Eq. (21) yields Eq. (18). \square

In the case where $\Gamma_k \equiv I_{l_x}$, note that $\pi_k \equiv 0$ and Proposition III.1 thus yields the classical Kalman filter. This section reviews the optimal injection-constrained filter (OICF) presented in [11] to establish notation for the development of injection constrained state estimators.

Proposition III.2: Let $k \geq 0$. Let $P_{k|k}$ be the posterior covariance of $e_{k|k}$. Let $\Gamma_k \neq I_{l_x}$ and let $P_{k+1|k+1}^{\text{ICF}}$ denote the posterior covariance at step $k+1$, which is given by Eq. (18). Let $P_{k+1|k+1}^{\text{KF}}$ denote the posterior covariance at step $k+1$ given by the Kalman filter, which is obtained by setting $\pi_k = 0$ in Eq. (18). Then,

$$\text{tr} P_{k+1|k+1}^{\text{KF}} \leq \text{tr} P_{k+1|k+1}^{\text{ICF}} \quad (23)$$

Proof: Note that

$$P_{k+1|k+1}^{\text{KF}} - P_{k+1|k+1}^{\text{ICF}} = -\pi_k P_{k+1|k} C_{k+1}^T \bar{R}_{k+1}^{-1} C_{k+1} P_{k+1|k} \pi_k \leq 0$$

The fact that trace is a linear operator yields Eq. (23). \square

If A_k, C_k , and Γ_k are constant, then the posterior error satisfies

$$e_{k+1|k+1} = (I - \Gamma \hat{K} C) A e_{k|k} + (I - \Gamma \hat{K} C) w_k - \Gamma \hat{K} v_{k+1} \quad (24)$$

The stability of Eq. (24) depends on the existence of a gain \hat{K} such that $A_f \triangleq A - \Gamma \hat{K} C A$ is asymptotically stable; when such a gain exists, the triple (A, Γ, CA) is called *static-output-feedback stabilizable*. Various necessary and/or sufficient conditions for static-output-feedback stabilizability are given in [21–23]. It is easy to see that, if (A, Γ, CA) is static-output-feedback stabilizable, then (A, Γ) is stabilizable and (A, CA) is detectable. The converse is not true; that is, if (A, Γ) is stabilizable and (A, CA) is detectable, then (A, Γ, CA) may or may not be static-output-feedback stabilizable. However, if (A, Γ) is not stabilizable, then (A, Γ, CA) is not static-output-feedback stabilizable.

Furthermore, as shown by the next example, static-output-feedback stabilizability of (A, Γ, CA) is not a sufficient condition for the OICF to be asymptotically stable.

Example III.1: Consider Eqs. (3) and (4) with the Lyapunov-stable LTI dynamics

$$A = \begin{bmatrix} 1.1 & -0.2 & -0.5 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix}, \quad C = [0.9 \quad 0.1 \quad -0.9] \quad (25)$$

For $\Gamma = e_1$ and $K = 0.2$, the spectral radius of A_f is 0.87, and, for $\Gamma = e_2$ and $K = -0.2$, the spectral radius of A_f is 0.95. The triple (A, Γ, CA) is thus static-output-feedback stabilizable for both choices of Γ .

For all $k \geq 0$, let $Q_k = 10^{-4} I_3$ and $R_k = 10^{-3}$. Let $\bar{x}_0 = [1 \quad 1 \quad 1]^T$ and $P_{0|0} = 10I_3$. Figure 1 shows the norm of the posterior error $e_{k|k}$, the trace of the posterior covariance $P_{k|k}$, and the filter gain K_k for two choices of Γ . For $\Gamma = e_1$, OICF is stable as shown in the subplots on the left by the decrease in the posterior error and the convergence of the posterior covariance. However, for $\Gamma = e_2$, OICF is unstable as shown in the subplots on the right by the diverging posterior error and posterior covariance. Note that OICF gain converges for both choice of Γ . \diamond

Example III.1 shows that the stability of the error system (24) is not guaranteed for the OICF gain given by Eq. (16) and, unlike the Kalman filter, OICF is not necessarily asymptotically stable. However, as shown in Example VI.2, IC-RCF converges to a stabilizing gain in this particular case.

Next, to facilitate the development of IC-UKF, the OICF gain and the posterior covariance are reformulated in terms of covariance matrices. For all $k \geq 0$, define the *prior output error* $z_{k+1|k}$ and the *posterior output error* $z_{k+1|k+1}$ by

$$z_{k+1|k} \triangleq C_{k+1} e_{k+1|k} \quad (26)$$

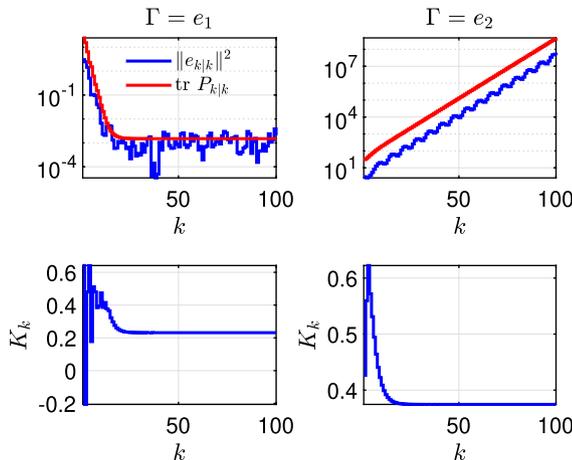


Fig. 1 Example III.1. OICF posterior error, posterior covariance, and the filter gain for two choices of the injection-constraint matrix Γ .

$$z_{k+1|k+1} \triangleq C_{k+1}e_{k+1|k} + v_{k+1} \quad (27)$$

and define the covariance of $z_{k+1|k+1}$ and the cross-covariance of $e_{k+1|k}$ and $z_{k+1|k}$ by

$$P_{z_{k+1|k+1}} \triangleq \mathbb{E}\left[z_{k+1|k+1}z_{k+1|k+1}^T\right] \quad (28)$$

$$P_{e,z_{k+1|k}} \triangleq \mathbb{E}\left[e_{k+1|k}z_{k+1|k}^T\right] \quad (29)$$

which, for all $k \geq 0$, satisfy

$$P_{z_{k+1|k+1}} = C_{k+1}P_{k+1|k}C_{k+1}^T + R_{k+1} \quad (30)$$

$$P_{e,z_{k+1|k}} = P_{k+1|k}C_{k+1}^T \quad (31)$$

Note that $P_{z_{k+1|k+1}} = \bar{R}_{k+1}$.

Next, substituting Eqs. (30) and (31) in Eqs. (21), (16), and (18), the posterior covariance for any value of \hat{K} can be written as

$$P_{k+1|k+1} = P_{k+1|k} + \Gamma_k \hat{K} P_{z_{k+1|k+1}} \hat{K}^T \Gamma_k^T - \Gamma_k \hat{K} P_{e,z_{k+1|k}}^T - P_{e,z_{k+1|k}} \hat{K}^T \Gamma_k^T \quad (32)$$

the optimal gain can be written as

$$K_{k+1} = (\Gamma_k^T \Gamma_k)^{-1} \Gamma_k^T P_{e,z_{k+1|k}} P_{z_{k+1|k+1}}^{-1} \quad (33)$$

and the corresponding posterior covariance can be written as

$$P_{k+1|k+1} = P_{k+1|k} - P_{e,z_{k+1|k}} P_{z_{k+1|k+1}}^{-1} P_{e,z_{k+1|k}}^T + \pi_k P_{e,z_{k+1|k}} P_{z_{k+1|k+1}}^{-1} P_{e,z_{k+1|k}}^T \pi_k^T \quad (34)$$

The posterior covariance and the optimal filter gain in Eqs. (32)–(34) are written in terms of covariance matrices instead of the matrices A_k and C_k to facilitate presentation in the later sections. As shown in the next section, UKF approximates these covariance matrices using ensembles.

IV. Injection-Constrained Unscented Kalman Filter

This section briefly reviews the UKF to establish notation and terminology used in the development of the injection-constrained filters. The UKF algorithm is formulated using a compact matrix-based notation and is based on the algorithm presented in [2].

Let $x \in \mathbb{R}^{l_x}$ and $P \in \mathbb{R}^{l_x \times l_x}$ be positive definite. The ensemble $X(x, P) \in \mathbb{R}^{l_x \times (2l_x+1)}$ is the matrix

$$X(x, P) \triangleq [xx + p_1 \cdots x + p_{l_x}x - p_1 \cdots x - p_{l_x}] \quad (35)$$

Let $\alpha > 0$. Define

$$W \triangleq \frac{1}{2\alpha^2 l_x} \begin{bmatrix} 2(\alpha^2 - 1)l_x \\ I_{2l_x \times 1} \end{bmatrix} \in \mathbb{R}^{2l_x+1} \quad (36)$$

The weighted mean of the ensemble X is $\bar{x} \triangleq XW$, and the ensemble perturbation is

$$\tilde{X} \triangleq X - H(\bar{x}) \quad (37)$$

where, for $v \in \mathbb{R}^n$,

$$H(v) \triangleq 1_{1 \times 2l_x+1} \otimes v \in \mathbb{R}^{n \times 2l_x+1} \quad (38)$$

Note that \otimes is the Kronecker product [24].

A. Summary of the Unscented Kalman Filter

To compute the filter gain K_{k+1} , UKF approximates the covariance matrices $P_{k+1|k}$, $P_{z_{k+1|k+1}}$, and $P_{e,z_{k+1|k}}$ in Eqs. (33) and (34) by propagating an ensemble of $2l_x + 1$ sigma points. For all $k \geq 0$, the i th sigma point $\hat{x}_{\sigma_i,k}$ is defined as the i th column of the $l_x \times (2l_x + 1)$ ensemble matrix

$$X_{k|k} \triangleq X \left(\hat{x}_{k|k}, \alpha \sqrt{l_x P_{k|k}} \right) \quad (39)$$

Then, for all $i = 1, \dots, 2l_x + 1$, the sigma points are propagated as

$$\hat{x}_{\sigma_i,k+1} = f_k(\hat{x}_{\sigma_i,k}) \quad (40)$$

The prior estimate and the prior covariance at step $k + 1$ are given by

$$\hat{x}_{k+1|k} = X_{k+1|k} W \quad (41)$$

$$P_{k+1|k} = \tilde{X}_{k+1|k} W_d \tilde{X}_{k+1|k}^T + Q_k \quad (42)$$

where

$$X_{k+1|k} \triangleq \begin{bmatrix} \hat{x}_{\sigma_1,k+1} & \cdots & \hat{x}_{\sigma_{2l_x+1},k+1} \end{bmatrix} \in \mathbb{R}^{l_x \times 2l_x+1} \quad (43)$$

Note that the prior estimate is the weighted sum of the columns of the propagated ensemble $X_{k+1|k}$. The prior ensemble $X'_{k+1|k}$, generated using the prior estimate and the prior covariance, is given by

$$X'_{k+1|k} \triangleq X \left(\hat{x}_{k+1|k}, \alpha \sqrt{l_x P_{k+1|k}} \right) \quad (44)$$

and, for $i = 1, \dots, 2l_x + 1$, the corresponding outputs are given by

$$\hat{y}_{\sigma_i,k+1} = g_{k+1} \left(X'_{k+1|k} e_i \right) \quad (45)$$

where e_i is the i th column of I_{2l_x+1} . The covariance matrices $P_{z_{k+1|k+1}}$ and $P_{e,z_{k+1|k}}$ are then given by

$$P_{z_{k+1|k+1}} = \tilde{Y}_{k+1} W_d \tilde{Y}_{k+1}^T + R_{k+1} \quad (46)$$

$$P_{e,z_{k+1|k}} = \tilde{X}'_{k+1|k} W_d \tilde{Y}_{k+1}^T \quad (47)$$

where

$$Y_{k+1} \triangleq \begin{bmatrix} \hat{y}_{\sigma_1,k+1} & \cdots & \hat{y}_{\sigma_{2l_x+1},k+1} \end{bmatrix} \in \mathbb{R}^{l_y \times 2l_x+1} \quad (48)$$

Algorithm 1: Unscented Kalman filter

Input : y_{k+1}, α
Output : $P_{k+1|k+1}, \hat{x}_{k+1|k+1}, K_{k+1}$
Data : $P_{k|k}, \hat{x}_{k|k}$
for $k > 0$ **do**
 Build $2l_x + 1$ -member ensemble using Eq. (39);
 Propagate ensemble using Eq. (40);
 Compute prior estimate and prior covariance using Eqs. (41) and (42);
 Build $2l_x + 1$ -member ensemble using Eq. (44) and compute ensemble output using Eq. (45);
 Compute output covariance using Eq. (46) and cross-covariance using Eq. (47);
 Compute UKF gain using Eq. (51), the posterior covariance using Eq. (50), and the posterior update using Eq. (49);
end

Finally, the posterior estimate at step $k + 1$ is

$$\hat{x}_{k+1|k+1} = \hat{x}_{k+1|k} + K_{k+1}(y_{k+1} - Y_{k+1}W) \quad (49)$$

and the posterior covariance at step $k + 1$ is

$$P_{k+1|k+1} = P_{k+1|k} - K_{k+1}P_{e,z_{k+1|k}}^T \quad (50)$$

where

$$K_{k+1} = P_{e,z_{k+1|k}} P_{z_{k+1|k}}^{-1} \quad (51)$$

UKF is summarized in Algorithm 1. Note that, as in the case of the Kalman filter, the output-error injection in Eq. (49) is unconstrained, that is, $\Gamma = I_{l_x}$.

B. Injection-Constrained Unscented Kalman Filter

To constrain the injection of η_k to a subspace, $\hat{x}_{k|k}$ is partitioned as

$$\hat{x}_{k|k} = \begin{bmatrix} \hat{x}_{c,k|k} \\ \hat{x}_{c\perp,k|k} \end{bmatrix} \quad (52)$$

where $\hat{x}_{c,k|k} \in \mathbb{R}^{l_\eta}$ is the portion of $\hat{x}_{k|k}$ in the subspace specified by Γ_k given by Eq. (8). The form of Eq. (52) implies that Γ_k has the form given by Eq. (8). This assumption entails no loss of generality as long as it is feasible to change the basis of Eq. (1) if needed.

Since the output error is injected into the subspace corresponding to $\hat{x}_{c,k}$, IC-UKF uses $\hat{x}_{c,k}$ to construct a $2l_\eta + 1$ -member ensemble to compute the filter gain and the corresponding posterior covariance. This approach is analogous to the unconstrained UKF, which uses a $2l_x + 1$ -member ensemble to obtain a filter gain for output-error injection. Note that each member of the reduced-order ensemble is of dimension l_x , and thus includes information from the full state of the system.

The IC-UKF is computed as follows. Let the i th sigma point $\hat{x}_{c\sigma_i,k}$ be given by the i th column of the $l_x \times (2l_\eta + 1)$ matrix

$$X_{k|k} \triangleq \begin{bmatrix} X \left(\Gamma_k^T \hat{x}_{c,k|k}, \alpha \sqrt{l_\eta P_{c,k|k}} \right) \\ H(\hat{x}_{c\perp,k|k}) \end{bmatrix} \quad (53)$$

where $P_{c,k|k} \triangleq \Gamma_k^T P_{k|k} \Gamma_k$. Note that $P_{c,k|k}$ is the covariance of $\hat{x}_{c,k|k}$. For all $i = 1, \dots, 2l_\eta + 1$, the sigma points are propagated as

$$\hat{x}_{\sigma_i,k+1} = f_k(\hat{x}_{\sigma_i,k}) \quad (54)$$

The prior estimate and the prior covariance at step $k + 1$ are given by

$$\hat{x}_{k+1|k} = X_{k+1|k}W \quad (55)$$

$$P_{c,k+1|k} = \Gamma_k^T \tilde{X}_{k+1|k} W_d \tilde{X}_{c,k+1|k}^T \Gamma_k + \Gamma_k^T Q_k \Gamma_k \quad (56)$$

where

$$X_{k+1|k} \triangleq \begin{bmatrix} \hat{x}_{\sigma_1,k+1} & \cdots & \hat{x}_{\sigma_{2l_\eta+1},k+1} \end{bmatrix} \in \mathbb{R}^{l_x \times (2l_\eta+1)} \quad (57)$$

Note that $P_{c,k+1|k} \in \mathbb{R}^{l_\eta \times l_\eta}$ is the covariance of the partitioned state. The prior ensemble $X'_{k+1|k}$, generated using the prior estimate and the prior covariance, is given by

$$X'_{k+1|k} \triangleq \begin{bmatrix} X \left(\Gamma_k^T \hat{x}_{k+1|k}, \alpha \sqrt{l_\eta P_{c,k+1|k}} \right) \\ H(\hat{x}_{c\perp,k+1|k}) \end{bmatrix} \quad (58)$$

and, for $i = 1, \dots, 2l_x + 1$, the output corresponding to the propagated sigma points are given by

$$\hat{y}_{\sigma_i,k+1} = g_{k+1}(X'_{k+1|k} e_i) \quad (59)$$

where e_i is the i th column of I_{2l_x+1} . The covariance matrices $P_{z_{k+1|k+1}}$ and $P_{e,z_{k+1|k}}$ are then given by

$$P_{z_{k+1|k+1}} = \tilde{Y}_{k+1} W_d \tilde{Y}_{k+1}^T + R_{k+1} \quad (60)$$

$$P_{e,z_{k+1|k}} = \Gamma_k^T \tilde{X}'_{k+1|k} W_d \tilde{Y}_{k+1}^T \quad (61)$$

where

$$Y_{k+1} \triangleq \begin{bmatrix} \hat{y}_{\sigma_1,k+1} & \cdots & \hat{y}_{\sigma_{2l_\eta+1},k+1} \end{bmatrix} \in \mathbb{R}^{l_y \times (2l_\eta+1)} \quad (62)$$

Finally, the posterior estimate at step $k + 1$ is

$$\hat{x}_{k+1|k+1} = \hat{x}_{k+1|k} + \Gamma_k K_{k+1}(y_{k+1} - Y_{k+1}W) \quad (63)$$

and the posterior covariance of the partitioned state at step $k + 1$ is given by

$$P_{c,k+1|k+1} = P_{c,k+1|k} - K_{k+1}P_{e,z_{k+1|k}}^T \quad (64)$$

where the IC-UKF gain is

$$K_{k+1} = P_{e,z_{k+1|k}} P_{z_{k+1|k+1}}^{-1} \quad (65)$$

Note that, since $X_{k+1|k}$, $X'_{k+1|k}$, and Y_{k+1} have $2l_\eta + 1$ in contrast to $2l_x + 1$ columns, the computation of $P_{c,k+1|k+1}$ requires a smaller ensemble in contrast to UKF, but each ensemble member is still the full state.

The full-state posterior covariance at step $k + 1$ is

$$P_{k+1|k+1} = P_{k+1|k} + \Gamma_k K_{k+1} P_{z_{k+1|k+1}} K_{k+1}^T \Gamma_k^T - \Gamma_k K_{k+1} P_{e,z_{k+1|k}}^T - P_{e,z_{k+1|k}} K_{k+1}^T \Gamma_k^T \quad (66)$$

where K_{k+1} is given by Eq. (65) and the matrices $P_{k+1|k}$, $P_{z_{k+1|k+1}}$, $P_{e,z_{k+1|k}}$ can be computed as shown in Sec. IV.A using an ensemble of arbitrary size. As expected, accuracy of $P_{k+1|k+1}$ improves as the ensemble size increases. Note that Eq. (66) is motivated by Eq. (32). Finally, since IC-UKF requires the covariance of $\hat{x}_{c,k|k}$ to compute the sigma-points, the full-state posterior covariance $P_{k+1|k+1}$ does not need to be computed to implement IC-UKF. IC-UKF is summarized in Algorithm 2.

V. Injection-Constrained Retrospective-Cost Filter

This section presents the injection-constrained retrospective-cost filter (IC-RCF), which obviates the need to propagate an ensemble in

Algorithm 2: Injection-constrained nscented Kalman filter (IC-UKF)

Input : $y_{k+1}, \alpha, \Gamma_k$
Output : $P_{c,k+1|k+1}, \hat{x}_{k+1|k+1}, K_{k+1}$
Data : $P_{c,k|k}, \hat{x}_{k|k}$
for $k > 0$ **do**
 Build $2l_\eta + 1$ -member ensemble using Eq. (53);
 Propagate ensemble using Eq. (54);
 Compute prior estimate and prior covariance using Eqs. (55) and (56);
 Build $2l_\eta + 1$ -member ensemble using Eq. (58) and compute ensemble output using Eq. (59);
 Compute output covariance using Eq. (60) and cross-covariance using Eq. (61);
 Compute UKF gain using Eq. (65), and the posterior covariance using Eq. (66), and the posterior update using Eq. (63);
end

order to compute the filter gain. For the system described by (1), (2), IC-RCF has the form (5), (6), and (9), where the IC-RCF gain \hat{K} is determined by minimizing the retrospective cost as shown below. To obtain a computationally efficient, recursive formula for the filter gain, the dimension of injection-matrix Γ_k is assumed to be a constant. It thus follows that the size of η_k and K_k is a constant.

For all $k \geq 0$, the injection signal η_{k+1} (9) is factored as

$$\eta_{k+1} = \Phi_{k+1} \theta_{k+1} \quad (67)$$

where

$$\Phi_{k+1} \triangleq I_{l_\eta} \otimes z_k^T \in \mathbb{R}^{l_\eta \times l_\eta l_y} \quad (68)$$

$$\theta_{k+1} \triangleq \text{vec } K_{k+1}^T \in \mathbb{R}^{l_\eta l_y} \quad (69)$$

where the *vec* operator stacks the columns of a matrix in a vector, and

$$z_{k+1} \triangleq y_{k+1} - g_{k+1}(\hat{x}_{k+1|k}) \quad (70)$$

Note that Eq. (67) shows that η_{k+1} is the product of the regressor matrix Φ_{k+1} , whose nonzero entries are measured data, and the vectorized IC-RCF gain K_{k+1} . The filter gain K_{k+1} is obtained by optimization of the *retrospective cost function* is defined by

$$J_k(\hat{\theta}) \triangleq \sum_{i=0}^k \lambda^{k-i} \left[z_i + G_f(\mathbf{q})(\Phi_i \hat{\theta} - \eta_i) \right]^T \left[z_i + G_f(\mathbf{q})(\Phi_i \hat{\theta} - \eta_i) \right] + \lambda^k \hat{\theta}^T R_\theta \hat{\theta} \quad (71)$$

where $\lambda \in (0, 1]$ is the forgetting factor, $R_\theta \in \mathbb{R}^{l_\eta l_y \times l_\eta l_y}$ is positive definite, \mathbf{q} is the forward-shift operator [25], and G_f is an $l_y \times l_\eta$ FIR (finite-impulse-response) filter. In particular, G_f has the form

$$G_f(\mathbf{q}) = \sum_{i=1}^{n_f} \frac{1}{\mathbf{q}^i} N_i \quad (72)$$

where the length n_f of the filter window and the filter coefficients $N_1, \dots, N_{n_f} \in \mathbb{R}^{l_y \times l_\eta}$ depend on the measurements l_y and the dimension of the injection signal l_η as shown later in this section.

Next, for all $k \geq 0$, the optimal gain

$$\theta_{k+1} \triangleq \underset{\hat{\theta} \in \mathbb{R}^{l_\eta l_y}}{\text{argmin}} J_k(\hat{\theta}) \quad (73)$$

is computed by recursive least squares as

$$\mathcal{P}_{k+1} = \frac{1}{\lambda} \mathcal{P}_k - \frac{1}{\lambda} \mathcal{P}_k \Phi_{f,k}^T (\lambda I_{l_y} + \Phi_{f,k} \mathcal{P}_k \Phi_{f,k}^T)^{-1} \Phi_{f,k} \mathcal{P}_k \quad (74)$$

Algorithm 3: Injection-constrained retrospective-cost filter (IC-RCF)

Input : z_k, η_k
Output : η_{k+1}, K_{k+1}
Data : \mathcal{P}_k, K_k
for $k > 0$ **do**
 Build regressor using Eq. (68);
 Filter data using Eqs. (76) and (77);
 Compute \mathcal{P}_{k+1} using Eq. (74);
 Compute θ_{k+1} using Eq. (75);
 Compute η_{k+1} using Eq. (67);
end

$$\theta_{k+1} = \theta_k - \mathcal{P}_{k+1} \Phi_{f,k}^T (z_k + \Phi_{f,k} \theta_k - \eta_{f,k}) \quad (75)$$

where the filtered regressor $\Phi_{f,k} \triangleq G_f(\mathbf{q}) \Phi_k$, the filtered injection signal $\eta_{f,k} \triangleq G_f(\mathbf{q}) \eta_k$, $\mathcal{P}_0 \triangleq R_\theta^{-1}$, and $\theta_0 = 0 \in \mathbb{R}^{l_\eta l_y}$. Note that

$$\Phi_{f,k} = N \bar{\Phi}_k \quad (76)$$

$$\eta_{f,k} = N \bar{\eta}_k \quad (77)$$

where

$$N \triangleq \begin{bmatrix} N_1 & \cdots & N_{n_f} \end{bmatrix}, \quad \bar{\Phi}_k \triangleq \begin{bmatrix} \Phi_{k-1} \\ \vdots \\ \Phi_{k-n_f} \end{bmatrix}, \quad \bar{\eta}_k \triangleq \begin{bmatrix} \eta_{k-1} \\ \vdots \\ \eta_{k-n_f} \end{bmatrix} \quad (78)$$

Finally, at step $k + 1$, the injection signal η_{k+1} is given by Eq. (67), where the filter gain θ_{k+1} is given by Eq. (75). IC-RCF is summarized in Algorithm 3. Note that IC-RCF does not require $\mathcal{P}_{k|k}$ to compute the filter gain. The posterior covariance, however, can be propagated using the method described at the end of Sec. IV.B, where K_{k+1} is given by inverting the *vec* operator in Eq. (69).

The next result shows that the injection signal η_k is constrained to lie in a subspace determined by the coefficients of G_f .

Lemma V.1: Let $\beta > 0$; $R_\theta = \beta I_\theta$; η_k be given by Eq. (67); Φ_k be given by Eq. (68); $\bar{\Phi}_{f,k}, \bar{\eta}_{f,k}, N$, be given by Eqs. (76–78); and θ_{k+1} be defined by Eq. (73). Then, for all $k \geq 1$,

$$\eta_{k+1} = -\frac{1}{\beta} [N_1^T \cdots N_{n_f}^T] \sum_{i=1}^k \lambda^{-i} \Psi_{k,i} [z_i + N \bar{\Phi}_i \theta_{k+1} - N \bar{\eta}_i] \in \mathcal{R}([N_1^T \cdots N_{n_f}^T]) \quad (79)$$

where

$$\Psi_{k,i} \triangleq \begin{bmatrix} z_{k+1}^T z_{i-1} \otimes I_{l_y} \\ \vdots \\ z_{k+1}^T z_{i-n_f} \otimes I_{l_y} \end{bmatrix} \quad (80)$$

Proof: Note that the cost function (71) can be rewritten as

$$J_k(\hat{\theta}) = \hat{\theta}^T \mathcal{A}_k \hat{\theta} + 2b_k^T \hat{\theta} + c_k \quad (81)$$

where

$$\mathcal{A}_k \triangleq \sum_{i=1}^k \lambda^{k-i} \bar{\Phi}_i^T N^T N \bar{\Phi}_i + \lambda^k R_\theta \quad (82)$$

$$b_k \triangleq \sum_{i=1}^k \lambda^{k-i} \bar{\Phi}_i^T N^T (z_i - N \bar{\eta}_i) \quad (83)$$

$$c_k \triangleq \sum_{i=1}^k \lambda^{k-i} (z_i - N\bar{\eta}_i)^T (z_i - N\bar{\eta}_i) \quad (84)$$

At step k , the batch least-squares minimizer θ_{k+1} of Eq. (81) is given by

$$\theta_{k+1} = -\mathcal{A}_k^{-1} b_k \quad (85)$$

which is equal to θ_{k+1} given by Eq. (75). Note that

$$\begin{aligned} \Phi_{k+1} \mathcal{A}_k \theta_{k+1} &= \sum_{i=1}^k \lambda^{k-i} \Phi_{k+1} (\bar{\Phi}_i^T N^T N \bar{\Phi}_i) \theta_{k+1} + \lambda^k \beta \Phi_{k+1} \theta_{k+1} \\ &= \sum_{i=1}^k \left(\lambda^{k-i} \sum_{j=1}^{n_f} (I_{l_u} \otimes z_{k+1}^T) (I_{l_u} \otimes z_{i-j}^T) N_j^T \right) N \bar{\Phi}_i \theta_{k+1} \\ &\quad + \lambda^k \beta \Phi_{k+1} \theta_{k+1} \\ &= \sum_{i=1}^k \left(\lambda^{k-i} \sum_{j=1}^{n_f} N_j^T z_{k+1}^T z_{i-j} \right) N \bar{\Phi}_i \theta_{k+1} + \lambda^k \beta \Phi_{k+1} \theta_{k+1} \\ &= [N_1^T \dots N_{n_f}^T] \sum_{i=1}^k \lambda^{k-i} \Psi_{k,i} N \bar{\Phi}_i \theta_{k+1} + \lambda^k \beta \Phi_{k+1} \theta_{k+1} \end{aligned} \quad (86)$$

and

$$\begin{aligned} \Phi_{k+1} b_k &= \Phi_{k+1} \sum_{i=1}^k \lambda^{k-i} \bar{\Phi}_i^T N^T (z_i - N\bar{\eta}_i) \\ &= [N_1^T \dots N_{n_f}^T] \sum_{i=1}^k \lambda^{k-i} \Psi_{k,i} (z_i - N\bar{\eta}_i) \end{aligned} \quad (87)$$

Writing Eq. (85) as $\mathcal{A}_k \theta_{k+1} = -b_k$, multiplying by Φ_{k+1} , and using Eqs. (86) and (87) yields Eq. (79). \square

It follows from Lemma V.1 that the injection signal η_k is constrained to lie in the subspace of \mathbb{R}^{l_η} spanned by the coefficients of the filter G_f . To allow the injection signal to be unconstrained in \mathbb{R}^{l_η} , the filter coefficients are chosen such that $\mathcal{R}([N_1^T \dots N_{n_f}^T]) = \mathbb{R}^{l_\eta}$. In view of Lemma V.1, for all examples in this paper, N_1, \dots, N_{n_f} are constructed using blocks of I_{l_η} as shown in Sec. VI.

If the filter coefficients are chosen such that $\mathcal{R}([N_1^T \dots N_{n_f}^T]) \subset \mathbb{R}^{l_\eta}$, then η_k is constrained to a subspace whose dimension is strictly less than l_η . Thus, the injection signal can be written as

$$\eta_k = G \zeta_k \quad (88)$$

where $l_\zeta < l_\eta$, $\zeta \in \mathbb{R}^{l_\zeta}$, and $G \in \mathbb{R}^{l_\eta \times l_\zeta}$ has full column rank. In this case, the injection-matrix Γ_k can be redefined as $\Gamma_k G$ and the injection signal η_k can be redefined as ζ_k . Thus, such a choice of filter coefficients is equivalent to choosing a smaller injection subspace.

VI. Numerical Examples

This section presents numerical examples that illustrate and compare the injection-constrained state estimators presented in this paper. In particular, OICF, IC-UKF, and IC-RCF are used to estimate the state in a Lyapunov-stable linear system, and IC-UKF and IC-RCF are used to estimate the state in the chaotic Lorenz system and the inviscid Burgers equation.

Example VI.1: Injection-constrained estimation in a simple harmonic oscillator

Consider a simple harmonic oscillator

$$\ddot{y} + \omega^2 y = 0 \quad (89)$$

which can be written in state-space form as

$$\dot{x} = Ax \quad (90)$$

$$y = Cx \quad (91)$$

where

$$x \triangleq \begin{bmatrix} y \\ \dot{y} \end{bmatrix}, \quad A \triangleq \begin{bmatrix} 0 & 1 \\ -\omega^2 & 0 \end{bmatrix}, \quad C \triangleq [1 \quad 0] \quad (92)$$

Defining $x_k = x(kT)$, where $T > 0$ is the discretization time step, it follows that

$$x_{k+1} = A_d x_k \quad (93)$$

where $A_d \triangleq e^{AT}$. Let $\omega = 5$ rad/s and $T = 0.01$ s. Let $x_0 = [1 \quad 1]$, $Q = 10^{-4}$, and $R = 10^{-4}$. Figures 2a and 2b show the state error and the error covariance for $\Gamma = e_1$ and $\Gamma = I_2$. Note that $P_{k|k}$ converges to similar values for both choice of Γ , although the convergence rate is slower for $\Gamma = e_1$. Next, the posterior covariance is propagated for several values of the process noise. Figure 2c shows $P_{k|k}$ at $k = 3000$ for various values of Q .

This numerical example shows that the full state can be estimated by restricting the filter correction to only the position estimate. Furthermore, in the case where process noise is smaller than the measurement noise, the trace of asymptotic error covariance for both injection-constrained and the unconstrained filter is similar.

Example VI.2: Injection-constrained state estimation for a linear system

Consider Eqs. (3) and (4) with the Lyapunov-stable LTI dynamics given by Eq. (25). For all $k \geq 0$, let $Q_k = 10^{-4} I_3$, $R_k = 10^{-3}$, and $u_k = 0$, and let $\bar{x}_0 = [1 \quad 1 \quad 1]^T$ and $P_{0|0} = 10I_3$. The state x_k is estimated using OICF, IC-UKF, and IC-RCF with the values of Γ listed in Table 1. In each estimator, same tuning parameters are used for all choices of Γ . In particular, IC-UKF uses $\alpha = 1.2$ and IC-RCF uses $\lambda = 1$, $R_\theta = I_{l_\eta}$, and the filters $G_f(\mathbf{q})$ given in Table 1.

For all three estimators and all values of Γ in Table 1, Fig. 3 shows the norm of the posterior error $e_{k|k}$ and the trace of the posterior covariance $P_{k|k}$. Note that, in all cases, the posterior error $e_{k|k}$ with IC-RCF decreases, whereas, as shown in the second row of Fig. 3,

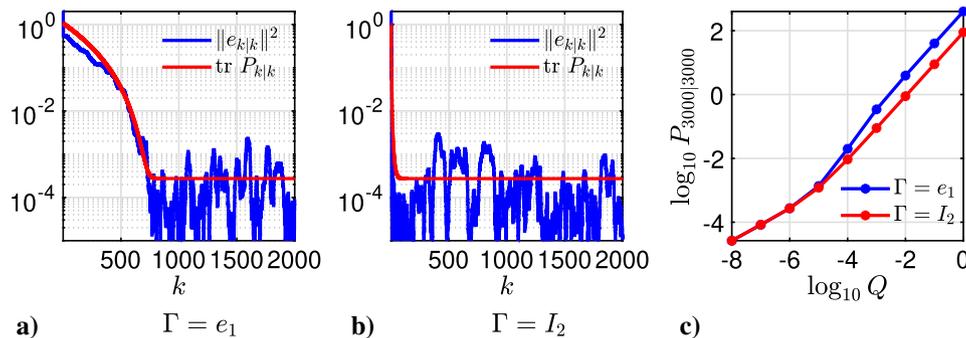


Fig. 2 Example VI.1. Injection-constrained state estimation in a simple harmonic oscillator.

Table 1 Example VI.2: filter used for IC-RCF

Γ	l_η	$G_\Gamma(q)$
e_1	1	$\frac{-1}{q}$
e_2	1	$\frac{1}{q}$
e_3	1	$\frac{1}{q}$
$[e_1 \ e_2]$	2	$\frac{[-1 \ 0]}{q} + \frac{[0 \ -1]}{q^2}$
$[e_2 \ e_3]$	2	$\frac{[-1 \ 0]}{q} + \frac{[0 \ -1]}{q^2}$
$[e_1 \ e_3]$	2	$\frac{[-1 \ 0]}{q} + \frac{[0 \ -1]}{q^2}$
I_3	3	$\frac{[-1 \ 0 \ 0]}{q} + \frac{[0 \ -1 \ 0]}{q^2} + \frac{[0 \ 0 \ -1]}{q^3}$

the posterior error $e_{k|k}$ in OICF diverges, and as shown in the second and third row of Fig. 3, the posterior error $e_{k|k}$ for IC-UKF does not decrease. This example suggests that IC-RCF can improve the state estimation accuracy in the cases where OICF and IC-UKF are

ineffective. Furthermore, the asymptotic performance of IC-RCF is similar to that of OICF. Note that, for $\Gamma = I_3$, OICF and IC-UKF estimates are nearly identical and thus the blue trace is completely overlaid by the red trace.

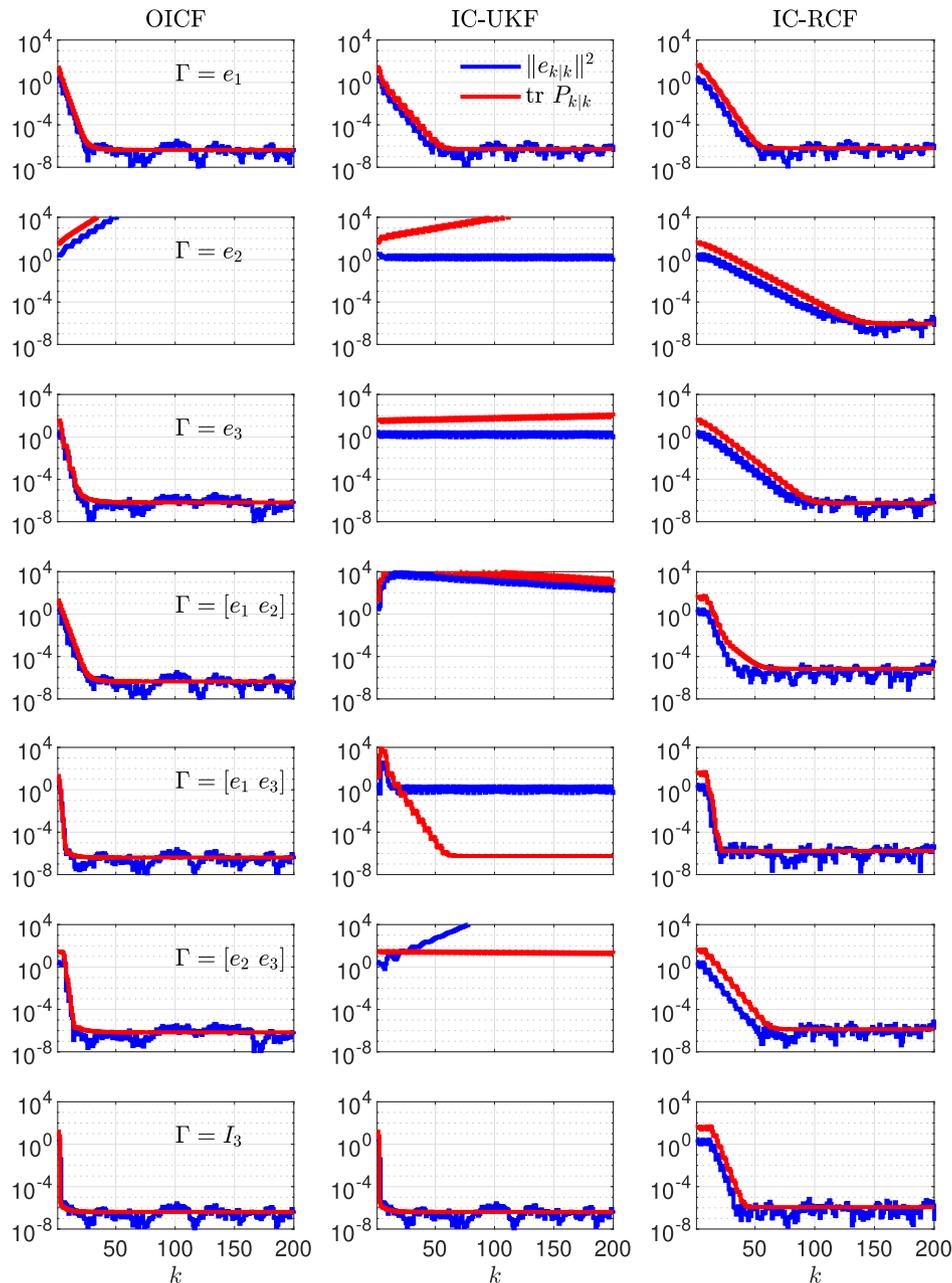
For $\Gamma = e_2$, IC-UKF is ineffective as shown in the second row of Fig. 3. In this case, for all $k \geq 0$,

$$X_k = \begin{bmatrix} \hat{x}_{2,k|k} & \hat{x}_{2,k|k} + p & \hat{x}_{2,k|k} - p \end{bmatrix} \quad (94)$$

where $p = \alpha \sqrt{P_{(2,2)k|k}}$ and thus the propagated ensemble

$$X_{k+1|k} = \begin{bmatrix} \hat{x}_{1,k|k} & \hat{x}_{1,k|k} & \hat{x}_{1,k|k} \end{bmatrix} \quad (95)$$

It thus follows from Eq. (37) that $\tilde{X}_{k+1} = 0$, and therefore $K_k = 0$. In this example, the collapse of the propagated sigma points to a single point is due to the structure of A and the choice of Γ . Similarly, with $\Gamma = e_3$ it follows that, for all $k \geq 0$, $K_k = 0$, and thus IC-UKF is ineffective. In contrast, for the chosen tuning parameters and for all

**Fig. 3** Example VI.2. Norm of posterior error and trace of posterior covariance for all three estimators.

values of Γ considered in this example, the IC-RCF error system is asymptotically stable.

Note that, in Table 1, the sign of the filter coefficient for second and third choice of Γ is negative, whereas it is positive for first choice of Γ . This example shows that both signs of the filter coefficients need to be tested in order to determine the sign that yields a stable IC-RCF. Finally, note that as the rank of the injection-constraint matrix increases, the convergence rate improves, although the asymptotic error is of similar magnitude for all choices of Γ . \diamond

Example VI.3: Injection-constrained state estimation for the Lorenz system

Consider the chaotic Lorenz system

$$\dot{x}_1 = \sigma(x_2 - x_1) \tag{96}$$

$$\dot{x}_2 = x_1(\rho - x_3) - x_2 \tag{97}$$

$$\dot{x}_3 = x_1x_2 - \beta x_3 \tag{98}$$

where $\sigma = 10$, $\rho = 28$, and $\beta = 8/3$. Lorenz system exhibits chaotic behavior for these parameter values. Let $x(0) = [10 \ 10 \ 10]^T$. For all $k \geq 0$, let $y_k \triangleq x_2(kT_s) + v_k$, where $T_s = 0.01$ s and $v_k \sim \mathcal{N}(0, 10^{-6})$. The Lorenz system is propagated using the explicit Runge–Kutta (4,5) method implemented in the MATLAB function ode45 in between the measurements. Process noise $w_k \sim \mathcal{N}(0, 10^{-6}I_3)$ is added to the state at each $t = kT_s$.

Letting $P_{0|0} = I_3$, the state is estimated using IC-UKF and IC-RCF with the values of Γ listed in Table 2. For each estimator, the same tuning parameters are used for all values of Γ . In particular, IC-UKF uses $\alpha = 1.2$ and IC-RCF uses $\lambda = 1$, $R_\theta = I_{l_\eta}$, and the filters $G_f(q)$ given in Table 2. Note that filter gain K_k and the injection signal η_k are l_η -dimensional vectors at each step k .

Table 2 Example VI.3: filters used for IC-RCF

Γ	l_η	$G_f(q)$
e_1	1	$\frac{-100}{q}$
$[e_1 \ e_2]$	2	$\frac{[-100 \ 0]}{q} + \frac{[0 \ -100]}{q^2}$
I_3	3	$\frac{[-100 \ 0 \ 0]}{q} + \frac{[0 \ -100 \ 0]}{q^2} + \frac{[0 \ 0 \ -100]}{q^3}$

Each subplot in Fig. 4 shows the norm of the posterior error and the trace of posterior covariance obtained with for the corresponding choice of the filter and the injection-constraint matrix Γ . In particular, the subplots in the first column are obtained with IC-RCF and the subplots in the first column are obtained with IC-UKF with the injection-constraint matrix specified below the x axis.

In this example, IC-UKF outperforms IC-RCF, in terms of both accuracy and convergence rate, for all choices of Γ . However, IC-UKF propagates $2l_\eta + 1$ -member ensemble, whereas IC-RCF requires only the predicted output in order to compute the injection signal. Consequently, as the dimension of the system increases, the computational cost of IC-UKF will increase proportionately; however, the computational cost of IC-RCF is independent of the system dimension. \diamond

Example VI.4: Injection-constrained state estimation for the inviscid Burgers equation

Consider the one-dimensional inviscid Burgers equation

$$\frac{\partial u}{\partial t} + \frac{\partial u^2}{\partial x} = 0 \tag{99}$$

where $u: [0, \infty) \times [0, 1] \rightarrow \mathbb{R}$. To simulate the flow on an infinite domain, the boundary conditions are set as $u(t, 0) = u(t, 1)$. The spatial domain $[0, 1]$ is discretized using N equally spaced grid points, so that $\Delta x \triangleq (1/N - 1)$. The time step Δt is chosen such that the Courant–Friedrichs–Lewy condition given in [26] is satisfied. For all $j \in \{1, \dots, N\}$, let $u_{j,k} \triangleq u(k\Delta t, j\Delta x)$. In this example, $N = 100$ and $\Delta t = 10^{-4}$. The inviscid Burgers equation (99) is discretized as shown in [27]. The initial condition is $u_{j,0} = 2 + \sin(2\pi j/25)$, and $P_{0|0} = 10^{-2}I_N$. Defining

$$U_k = [u_{1,k} \ \dots \ u_{N,k}]^T \in \mathbb{R}^N \tag{100}$$

the discretized inviscid Burgers equation is written as

$$U_{k+1} = F(U_k) \tag{101}$$

where F is a vector-valued function defined in [27]. For all $k \geq 0$, let

$$y_k = u_{87,k} + v_k \tag{102}$$

where $v_k \sim \mathcal{N}(0, 10^{-3})$. Process noise $w_k \sim \mathcal{N}(0, 10^{-3}I_N)$ is added to the state at each $t = k\Delta t$.

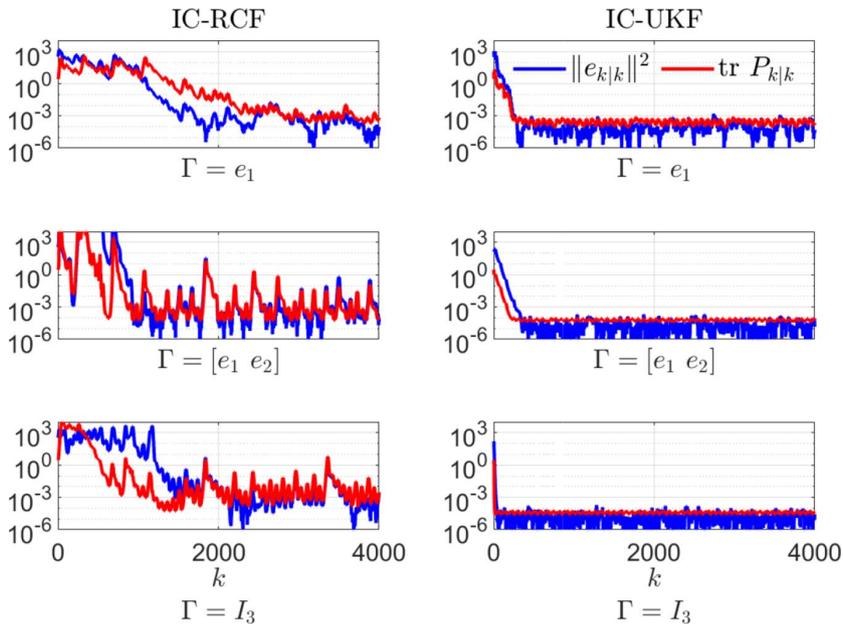


Fig. 4 Example VI.3. Norm of posterior error and trace of posterior covariance for IC-RCF and IC-UKF.

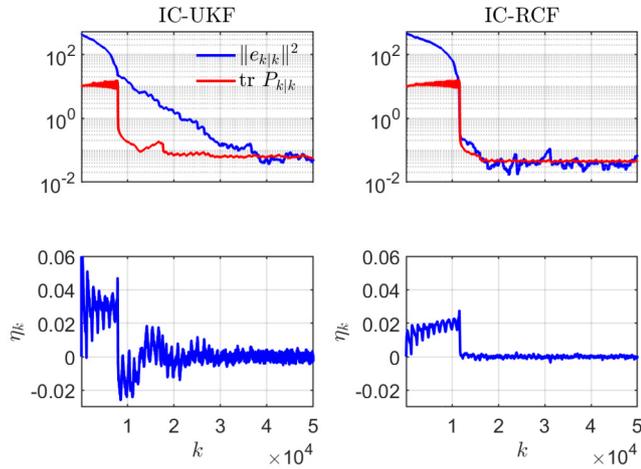


Fig. 5 Example VI.4. Norm of the posterior error, the trace of the posterior covariance, and the injection signal for IC-UKF and IC-RCF.

The state U_k is estimated using IC-UKF and IC-RCF with $\Gamma = e_{87}$. In particular, IC-UKF uses $\alpha = 1.2$ and IC-RCF uses $R_\theta = 10^5$, $\lambda = 0.9999$, and $G_f(\mathbf{q}) = (-1/\mathbf{q})$. Figure 5 shows the norm of the posterior error, the trace of the posterior covariance, and the injection signal for IC-UKF and IC-RCF. Note that, in this example, IC-RCF outperforms IC-UKF in terms of convergence rate and accuracy.

Finally, to investigate the effect of R_θ and Γ in IC-RCF, the state U_k is estimated with various values of R_θ , where, in all cases, $\Gamma = e_{87}$, $\lambda = 0.9999$, and $G_f(\mathbf{q}) = (-1/\mathbf{q})$. Figure 6 shows the norm of the posterior error and the trace of the posterior covariance for various

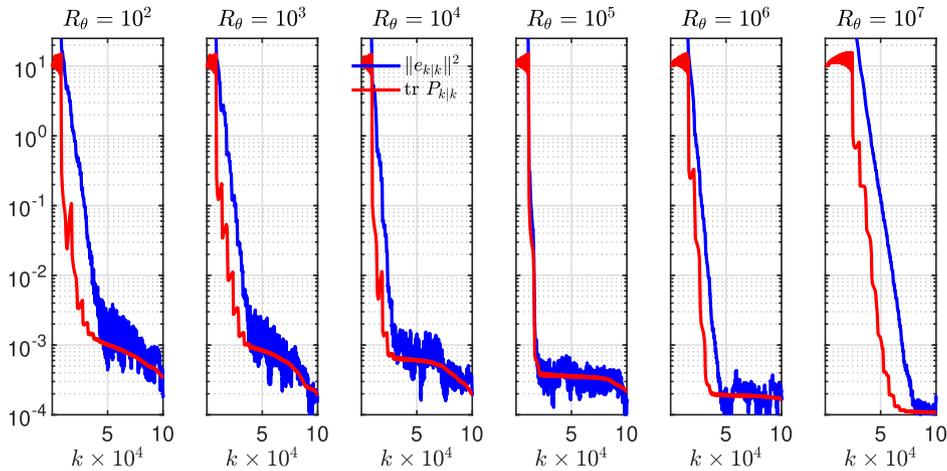


Fig. 6 Example VI.4. Effect of R_θ on the posterior error $e_{k|k}$ for IC-RCF.

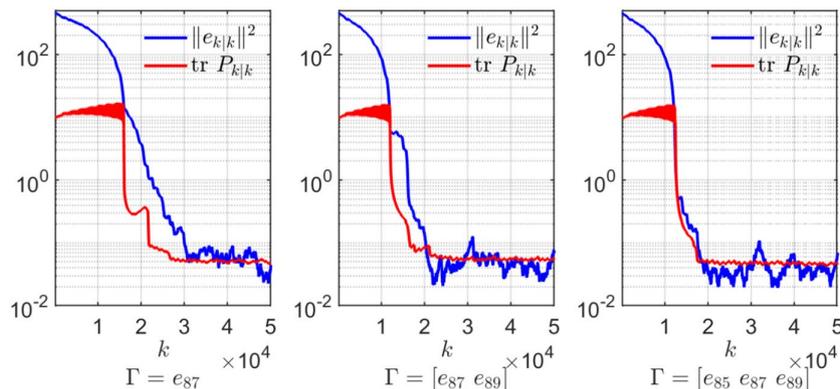


Fig. 7 Example VI.4. Effect of Γ on the posterior error $e_{k|k}$ for IC-RCF.

choice of R_θ in IC-RCF. Note that, as R_θ increases, the accuracy improves but convergence rate worsens. Next, the state U_k is estimated with various values of Γ , where, in all cases, $R_\theta = 10^6 I_{l_\theta}$, $\lambda = 0.9999$, and $G_f(\mathbf{q}) = (-1/\mathbf{q})$. Figure 7 shows the norm of the posterior error and the trace of the posterior covariance for various choices of Γ in IC-RCF. Note that asymptotic error is similar for all choices of Γ ; however, the convergence rate improves as the rank of Γ increases.

Despite retuning attempts with several values of α , the UKF estimate diverges. This observation suggests that constructing sigma points over the entire state space may result in state estimates that are physically unrealistic and result in divergence when propagated using the dynamics map. Both IC-UKF and IC-RCF can alleviate this problem since the perturbations in the state estimates can be restricted to physically realistic scenarios by appropriately choosing Γ_k \diamond

VII. Application of ICSE to Parameter Estimation in the Viscous Burgers Equation

This section shows that the parameter estimation problem is a special case of ICSE problem described in Sec. II. Consider the system

$$x_{k+1} = F_k(x_k, \mu) + w_k \quad (103)$$

$$y_k = G_k(x_k, \mu) + v_k \quad (104)$$

where $\mu \in \mathbb{R}^{l_\mu}$ is the vector of unknown constant parameters that parameterizes the dynamics and the measurement map. The goal is to compute an estimate $\hat{\mu}$ of the parameter μ using the measurements y_k . Defining the augmented dynamics

$$\begin{bmatrix} \mu_{k+1} \\ x_{k+1} \end{bmatrix} = \begin{bmatrix} \mu_k \\ F_k(x_k, \mu_k) + w_k \end{bmatrix} \quad (105)$$

$$y_k = G_k(x_k, \mu_k) + v_k \quad (106)$$

where the constant state μ_k reflects the constant “dynamics” of μ , and using the structure of Eq. (105), the *parameter estimator* has the form

$$\begin{bmatrix} \hat{\mu}_{k+1} \\ \hat{x}_{k+1} \end{bmatrix} = \begin{bmatrix} \hat{\mu}_k \\ F_k(\hat{x}_k, \hat{\mu}_k) \end{bmatrix} + \begin{bmatrix} \eta_k \\ 0 \end{bmatrix} \quad (107)$$

Note that Eq. (107) has the form of Eq. (5), where the injection-constraint matrix $\Gamma = [I_{l_\eta} \ 0_{l_\eta \times l_x}]^T$ constrains the injection signal η_k to the subspace corresponding to the vector μ of unknown parameters. ICSE can thus be applied to parameter estimation by using the output error to update the parameter estimate but not the state estimate. The following example uses ICSE to estimate an unknown parameter in the viscous Burgers equation using both IC-UKF and IC-RCF.

Example VII.1: Parameter estimation for the viscous Burgers equation

Consider the forced one-dimensional viscous Burgers equation

$$\frac{\partial u}{\partial t} + \frac{\partial u^2}{\partial x} \frac{1}{2} = \frac{\partial}{\partial x} \left(\mu \frac{\partial u}{\partial x} \right) + q(x, t) \quad (108)$$

where $u: [0, \infty) \times [0, 1] \rightarrow \mathbb{R}$, μ is the viscosity, and $q(x, t)$ is the external forcing given by

$$q(x, t) = \begin{cases} 0, & x \neq 1, \\ \sin(0.005t) + 0.25 \sin(0.01t), & x = 1 \end{cases} \quad (109)$$

In this example, $\mu = 0.3$. The spatial and temporal discretization is similar to Example VI.4. The initial condition is $u_{j,0} = \sin(2\pi j/25)$. Defining U_k by Eq. (100), the discretized Burgers equation is augmented with the unknown parameter dynamics as

$$\begin{bmatrix} \mu_{k+1} \\ U_{k+1} \end{bmatrix} = \begin{bmatrix} \mu_k \\ \bar{F}(U_k, \mu_k) \end{bmatrix} \quad (110)$$

where \bar{F} is the corresponding function. For all $k \geq 0$, let

$$y_k = u_{87,k} + v_k \quad (111)$$

where $v_k \sim \mathcal{N}(0, 10^{-5})$. Process noise $w_k \sim \mathcal{N}(0, 10^{-5}I_N)$ is added to the state at each $t = k\Delta t$. The augmented state $[\mu_k \ U_k^T]$ in Eq. (108) is estimated using injection-constrained filters with $\Gamma = e_1$. This value of Γ implies that only the first component of the

augmented state, that is, the parameter μ_k , is affected by the injection signal η_k , whereas U_k is propagated using the function \bar{F} . In particular, IC-UKF uses $\alpha = 1.2$ and $P_{0|0} = 10^{-3}I_{102}$ and IC-RCF uses $R_\theta = 10^5$, $\lambda = 0.9999$, and $G_f(q) = -1/q$. Note that $l_\eta = 1$ and thus the filter gain K_k is a scalar.

Figure 8 shows the state-error norm, the posterior covariance, the injection signal, and the parameter estimate obtained with IC-UKF and IC-RCF. Note that the parameter estimate does not converge to the true value of the unknown parameter. Despite several tuning efforts by varying α and $P_{0|0}$, IC-UKF did not yield convergence of the parameter estimate. Unlike IC-UKF, the parameter estimate converges to the true value of the unknown parameter and the posterior error decreases. \diamond

VIII. Conclusions

IC-UKF and IC-RCF were applied to state estimation in both linear and nonlinear systems. For the linear example, IC-UKF was ineffective for some injection-constraint matrices because the sigma points collapsed to a single point. In contrast, IC-RCF successfully estimated the states for all choices of the injection-constraint matrix. In the nonlinear example, both IC-UKF and IC-RCF successfully estimated the states for all choices of the injection-constraint matrices. IC-UKF and IC-RCF were also applied to the problem of parameter estimation in the Burgers equation. The IC-UKF parameter estimate did not converge to the true value, whereas IC-RCF successfully estimated the unknown parameter.

Numerical examples further revealed that the convergence rate of the injection-constrained filters depends on the rank of the injection-constraint matrix. As the dimension of the injection subspace increases, the convergence rate of the filter improves. Asymptotically, however, the output error and posterior covariance reached similar orders of magnitude for all of the injection-constraint matrices that were considered. This observation suggests that the degradation of the accuracy of the state estimates due to the output-error injection constraint is not as severe as might be expected. Consequently, the same asymptotic performance as the unconstrained state estimator may be achieved at a reduced computational cost with IC-UKF and IC-RCF. This trend is beneficial for IC-RCF, which requires specification of as many filter coefficients as the rank of the injection-constraint matrix, as well as IC-UKF, which requires a smaller ensemble and thus has a lower computational cost than the classical UKF.

Acknowledgments

This research was supported by Air Force Office of Scientific Research under Dynamic Data-Driven Applications Systems (DDDDAS; <http://www.1dddas.org/>) grant FA9550-16-1-0071. The authors thank the reviewers and associate editor for numerous comments and suggestions that helped improve this paper.

References

- [1] Simon, D., *Optimal State Estimation: Kalman, H-Infinity, and Non-linear Approaches*, Wiley, New York, 2006.
- [2] Wan, E. A., and Van Der Merwe, R., “The Unscented Kalman Filter for Nonlinear Estimation,” *Adaptive Systems for Signal Processing, Communications, and Control Symposium*, IEEE Publ., Piscataway, NJ, 2000, pp. 153–158. <https://doi.org/10.1109/ASSPCC.2000.882463>
- [3] Julier, S. J., and Uhlmann, J. K., “Unscented Filtering and Nonlinear Estimation,” *Proceedings of the IEEE*, Vol. 92, No. 3, 2004, pp. 401–422. <https://doi.org/10.1109/JPROC.2003.823141>
- [4] Evensen, G., *Data Assimilation: The Ensemble Kalman Filter*, Springer, New York, 2009.
- [5] Ridley, A. J., Deng, Y., and Toth, G., “The Global Ionosphere–Thermosphere Model,” *Journal of Atmospheric and Solar-Terrestrial Physics*, Vol. 68, No. 8, 2006, pp. 839–864. <https://doi.org/10.1016/j.jastp.2006.01.008>
- [6] Simon, D., “Reduced Order Kalman Filtering Without Model Reduction,” *Control & Intelligent Systems*, Vol. 35, No. 2, 2007, pp. 169–174. <https://doi.org/10.5555/1722239.1722249>

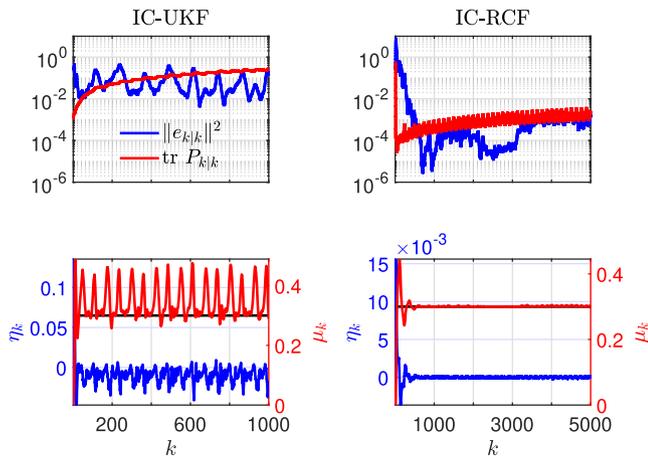


Fig. 8 Example VII.1. Norm of the posterior error, the trace of the posterior covariance, injection signal, and the parameter estimate for IC-UKF and IC-RCF.

- [7] Ott, E., Hunt, B. R., Szunyogh, I., Zimin, A. V., Kostelich, E. J., Corazza, M., Kalnay, E., Patil, D., and Yorke, J. A., "A Local Ensemble Kalman Filter for Atmospheric Data Assimilation," *Tellus A: Dynamic Meteorology and Oceanography*, Vol. 56, No. 5, 2004, pp. 415–428.
<https://doi.org/10.3402/tellusa.v56i5.14462>
- [8] Furrer, R., Genton, M. G., and Nychka, D., "Covariance Tapering for Interpolation of Large Spatial Datasets," *Journal of Computational and Graphical Statistics*, Vol. 15, No. 3, 2006, pp. 502–523.
<https://doi.org/10.1198/106186006X132178>
- [9] Kaufman, C. G., Schervish, M. J., and Nychka, D. W., "Covariance Tapering for Likelihood-Based Estimation in Large Spatial Data Sets," *Journal of the American Statistical Association*, Vol. 103, No. 484, 2008, pp. 1545–1555.
<https://doi.org/10.1198/016214508000000959>
- [10] Cressie, N., and Johannesson, G., "Fixed Rank Kriging for Very Large Spatial Data Sets," *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, Vol. 70, No. 1, 2008, pp. 209–226.
<https://doi.org/10.1111/j.1467-9868.2007.00633.x>
- [11] Chandrasekar, J., Bernstein, D. S., Barrero, O. O., and De Moor, B. L. R., "Kalman Filtering with Constrained Output Injection," *International Journal of Control*, Vol. 80, No. 12, 2007, pp. 1863–1879.
<https://doi.org/10.1080/00207170701373633>
- [12] Ansari, A., and Bernstein, D. S., "Adaptive Non-Bayesian State Estimation," *Proceedings of American Control Conference*, IEEE Publ., Piscataway, NJ, 2016, pp. 6977–6982.
<https://doi.org/10.1109/ACC.2016.7526772>
- [13] Schmidt, S. F., "Application of State-Space Methods to Navigation Problems," *Advances in Control Systems*, Vol. 3, edited by C. T. Leondes, Elsevier, New York, 1966, pp. 293–340.
- [14] Zanetti, R., and D'Souza, C., "Recursive Implementations of the Schmidt-Kalman 'Consider' Filter," *The Journal of the Astronautical Sciences*, Vol. 60, No. 3, 2013, pp. 672–685.
<https://doi.org/10.1007/s40295-015-0068-7>
- [15] Woodbury, D., and Junkins, J., "On the Consider Kalman Filter," *AIAA Guidance, Navigation, and Control Conference*, AIAA Paper 2010-7752, 2010.
<https://doi.org/10.2514/6.2010-7752>
- [16] Stauch, J., and Jah, M., "Unscented Schmidt-Kalman Filter Algorithm," *Journal of Guidance, Control, and Dynamics*, Vol. 38, No. 1, 2015, pp. 117–123.
<https://doi.org/10.2514/1.G000467>
- [17] Hunt, B. R., Kostelich, E. J., and Szunyogh, I., "Efficient Data Assimilation for Spatiotemporal Chaos: A Local Ensemble Transform Kalman Filter," *Physica D: Nonlinear Phenomena*, Vol. 230, Nos. 1–2, 2007, pp. 112–126.
<https://doi.org/10.1016/j.physd.2006.11.008>
- [18] Miyoshi, T., Yamane, S., and Enomoto, T., "Localizing the Error Covariance by Physical Distances Within a Local Ensemble Transform Kalman Filter (LETKF)," *Sola*, Vol. 3, Aug. 2007, pp. 89–92.
<https://doi.org/10.2151/sola.2007-023>
- [19] Ljung, L., "Asymptotic Behavior of the Extended Kalman Filter as a Parameter Estimator for Linear Systems," *IEEE Transactions on Automatic Control*, Vol. 24, No. 1, 1979, pp. 36–50.
<https://doi.org/10.1109/TAC.1979.1101943>
- [20] Goel, A., and Bernstein, D. S., "Adaptive State Estimation with Subspace-Constrained State Correction," *Proceedings of American Control Conference*, IEEE Publ., Piscataway, NJ, 2020, pp. 719–724.
<https://doi.org/10.23919/ACC45564.2020.9147916>
- [21] Syrmos, V. L., Abdallah, C. T., Dorato, P., and Grigoriadis, K., "Static Output Feedback—A Survey," *Automatica*, Vol. 33, No. 2, 1997, pp. 125–137.
[https://doi.org/10.1016/S0005-1098\(96\)00141-0](https://doi.org/10.1016/S0005-1098(96)00141-0)
- [22] Van der Woude, J., "A Note on Pole Placement by Static Output Feedback for Single-Input Systems," *Systems & Control Letters*, Vol. 11, No. 4, 1988, pp. 285–287.
- [23] Byrnes, C., and Anderson, B., "Output Feedback and Generic Stabilizability," *SIAM Journal on Control and Optimization*, Vol. 22, No. 3, 1984, pp. 362–380.
<https://doi.org/10.1137/0322024>
- [24] Van Loan, C. F., "The Ubiquitous Kronecker Product," *Journal of Computational and Applied Mathematics*, Vol. 123, Nos. 1–2, 2000, pp. 85–100.
[https://doi.org/10.1016/S0377-0427\(00\)00393-9](https://doi.org/10.1016/S0377-0427(00)00393-9)
- [25] Middleton, R. H., and Goodwin, G. C., *Digital Control and Estimation: A Unified Approach*, Prentice-Hall, Upper Saddle River, NJ, 1990.
- [26] Courant, R., Friedrichs, K., and Lewy, H., "On the Partial Difference Equations of Mathematical Physics," *IBM Journal of Research and Development*, Vol. 11, No. 2, 1967, pp. 215–234.
<https://doi.org/10.1147/rd.112.0215>
- [27] Goel, A., Duraisamy, K., and Bernstein, D. S., "Parameter Estimation in the Burgers Equation Using Retrospective-Cost Model Refinement," *Proceedings of American Control Conference*, IEEE Publ., Piscataway, NJ, 2016, pp. 6983–6988.
<https://doi.org/10.1109/ACC.2016.7526773>