

Nonlinear Projection Filter Based on Galerkin Approximation

Randal Beard,* John Kenney,[†] Jacob Gunther,[‡] Jonathan Lawton,[†] and Wynn Stirling[§]
Brigham Young University, Provo, Utah 84602

The conditional probability density function of the state of a stochastic dynamic system represents the complete solution to the nonlinear filtering problem because, with the conditional density in hand, all estimates of the state, optimal or otherwise, can be computed. It is well known that, for systems with continuous dynamics, the conditional density evolves, between measurements, according to Kolmogorov's forward equation. At a measurement, it is updated according to Bayes formula. Therefore, these two equations can be viewed as the dynamic equations of the conditional density and, hence, the exact nonlinear filter. In this paper, Galerkin's method is used to approximate the nonlinear filter by solving for the entire conditional density. Using a discrete cosine transform to approximate the projections required in Galerkin's method leads to a computationally realizable nonlinear filter. The implementation details are given and performance is assessed through simulations.

I. Introduction

ESTIMATING the state of a stochastic dynamic system from noisy observations is an important problem in engineering. The extensive work on this problem for linear systems was initiated by Kalman and Bucy.^{1,2} In the decades since this early work, many important theoretical results for the linear problem have emerged and the linear filter has found wide practical application. However, since most systems are not truly linear, linear filtering theory does not apply directly to most physical systems.

In the more general setting of nonlinear systems, filter theory is less developed but has received attention from a number of researchers since the 1960s.^{3–6} For linear systems with Gaussian inputs, the probability density function of the state conditioned on the measurements is Gaussian. Hence, optimal linear filters need only propagate the conditional mean and covariance that completely describes the density function. However, for nonlinear systems, the conditional density may not have a finite parameterization. General nonlinear filters must therefore propagate the entire density function.

The extended Kalman filter (EKF), a heuristic filter based on the linearized dynamics of a system,^{7–10} has become the standard for state estimation of nonlinear systems. The EKF assumes that 1) deviations from the reference state trajectory are small, 2) the mathematical description of the system dynamics and observations is accurate, and 3) the conditional density function of the state is Gaussian. For large deviations from the reference trajectory, the EKF performs poorly or becomes unstable. Furthermore, if constructed from an erroneous model, the EKF state estimates can diverge from the true state. In addition, because of the nonlinear nature of the EKF, its estimate may depend on the initial conditions of the filter. In the final section of this paper, we show an example where the EKF gives erroneous results if it is not correctly initialized. These statements are not meant to discredit the EKF, the practical importance of which cannot be overstated, only to point out that the EKF is based on a linearization argument that, for severely nonlinear systems, can cause it to perform poorly. It should be noted there has been extensive work on alleviating some of the preceding problems.^{11–15}

In this paper, we construct a nonlinear filter that approximates the exact nonlinear filter for systems with continuous nonlinear dy-

namics and discrete nonlinear observations. Accordingly, the paper represents a departure from current research directions in nonlinear filtering. Rather than pursue enhancements or modifications of the EKF, or explore dual relationships to various nonoptimal control algorithms, we investigate a method of approximating the exact nonlinear filtering problem. Note that if the exact nonlinear filtering problem could be solved directly or approximated efficiently, it would be widely used today in various industrial applications. We hope that this paper is a step in that direction.

The exact nonlinear filter consists of two dynamic equations³:

1) A partial differential equation (Kolmogorov's forward equation) that describes how the conditional density evolves between measurements, and

2) A difference equation (Bayes formula) that describes how it is modified by information supplied by new measurements.

To solve these equations we employ Galerkin's method, a classic procedure for approximating solutions of partial differential equations (PDEs).^{16,17} In the context of the forward Kolmogorov equation, Galerkin's method was suggested by Risken¹⁸ as a possible way to approximate the solution. However, a detailed analysis was not pursued.

Galerkin's method assumes that the exact solution to a PDE can be expanded as an infinite sum of basis elements. An approximate solution is found by truncating this sum and projecting the resulting error onto the finite subspace spanned by the basis elements used to approximate the solution.¹⁶ To distinguish the approximate filter from the exact filter, and for lack of a better name, we will refer to the resulting filter as the nonlinear projection filter (NPF). Using a complex exponential basis as approximating elements, we show that the nonlinear filter can be implemented efficiently (for low-order systems) using discrete cosine transforms (DCT) resulting in a fast nonlinear filter that could be implemented in real time on a digital signal processor. Sinusoidal bases have been used before to implement Galerkin-based algorithms, but in much different contexts.^{19–21} This work is also a natural extension of the application of Galerkin's method to optimal and robust control.^{22–25} While there have been numerous studies that have applied numerical methods for solving PDEs,²⁶ we are not aware of a careful study of Galerkin's spectral method to the nonlinear filtering problem.

An important issue concerns the convergence of the Galerkin approximation. We shall prove that the approximation residual converges to zero as the dimension of the finite dimensional subspace used in the approximation tends to infinity. In other words, the NPF converges to the exact nonlinear filter. One limitation in our convergence result is that we do not obtain an explicit estimate of the approximation error. We only show that the approximation error can be made arbitrarily small by making the order of the approximation large enough.

Because the NPF is an approximation of the exact nonlinear filter, it can outperform the EKF for large variations in the state, for model

Received Dec. 8, 1997; revision received July 23, 1998; accepted for publication Oct. 5, 1998. Copyright © 1998 by the American Institute of Aeronautics and Astronautics, Inc. All rights reserved.

*Assistant Professor, Department of Electrical and Computer Engineering.

[†]Research Assistant, Department of Electrical and Computer Engineering.

[‡]Research Assistant, Department of Electrical and Computer Engineering; currently Technical Staff, Merisoft, Provo, UT 84602.

[§]Professor, Department of Electrical and Computer Engineering.

mismatches, and for arbitrary initial conditions. In addition, because we propagate the conditional density, we can compute optimal state estimates for any criterion. In contrast, the EKF is only capable of approximating the conditional mean.

The remainder of the paper is organized as follows. The nonlinear filtering problem and the exact nonlinear filter is given in Sec. II. In Sec. III, Galerkin's method is used to approximate the exact nonlinear filter resulting in the NPF. In Sec. IV we show that the NPF converges to the exact nonlinear filter as the order of approximation increases. Practical implementation of the filter and the use of exponential basis elements is discussed in Sec. V. Finally, Sec. VI contains simulations comparing the NPF and the EKF.

II. Exact Nonlinear Filter

Most physical systems evolve continuously in time while measurements may only be taken periodically at discrete time instants. Suppose the n -dimensional state x_t of a continuous nonlinear stochastic dynamic system satisfies

$$dx_t = f(t, x_t) dt + G(t, x_t) d\beta_t, \quad t \geq t_0 \quad (1)$$

where $\{\beta_t, t \geq t_0\}$ is a p -dimensional Brownian motion with covariance matrix $Q(t) dt$. Let m -dimensional noisy measurements be made at discrete times t_k

$$y_k = h(x_{t_k}, t_k) + v_k, \quad k = 1, 2, \dots \quad (2)$$

where $\{v_k, k \geq 1\}$ is an m -dimensional white Gaussian sequence independent of $d\beta_t$ with covariance matrix R_k . Define the collection of measurements taken up to and including time t as $Y_t = \{y_k : t_k \leq t\}$. We seek equations of evolution for the conditional density $p(t, x | Y_t)$ because it summarizes all the statistical information about the state contained in the measurements Y_t and the initial condition $p(t_0, x)$. From $p(t, x | Y_t)$, the conditional mean and variance can be computed, which for nonlinear systems generally depend on all of the higher order moments.

Between observations at t_k and t_{k+1} , $p \equiv p(t, x | Y_t)$ diffuses according to Kolmogorov's forward equation

$$\frac{\partial p}{\partial t} = - \sum_{i=1}^n \frac{\partial (p f_i)}{\partial x_i} + \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \frac{\partial^2 [p (G Q G^T)_{ij}]}{\partial x_i \partial x_j} \quad (3)$$

where either $p(t_0, x)$ or $p(t_k, x | Y_{t_k})$, the measurement update at t_k , is used as the initial condition, and where $(A)_{ij}$ is the (i, j) th element of the matrix A . At an observation, p satisfies the difference equation

$$p(t_{k+1}, x | Y_{t_{k+1}}) = \frac{p(y_{k+1} | x) p(t_{k+1}^-, x | Y_{t_k})}{\int p(y_{k+1} | \xi) p(\xi, t_{k+1}^- | Y_{t_k}) d\xi} \quad (4)$$

where

$$p(y_{k+1} | x) = \frac{\exp\{-\frac{1}{2}[y_k - h(x, t_k)]^T R_k^{-1}[y_k - h(x, t_k)]\}}{\sqrt{(2\pi)^m |R_k|}}$$

and t_{k+1}^- is the instant in time right before the $(k+1)$ sample. For a detailed derivation and discussion of these results, see Ref. 3.

Equations (3) and (4) represent dynamic equations for the exact nonlinear filter. Equation (3) is used to compute predictions between measurements, while measurements are used to update the information about the state via Eq. (4). Closed-form solutions to these equations are generally not available. One exception is the linear-Gaussian case: For a linear dynamic system driven by Gaussian noise, this filter reduces to the standard Kalman filter.³ In the next section we show that Galerkin's approximation method can be used to reduce Eq. (3) to a linear ordinary differential equation and Eq. (4) to an algebraic update equation.

Other nonlinear filters have also been derived. In particular Zakai has developed a nonlinear filter that has been studied by a number of researchers.^{27,28} The main reason that we have chosen to study the Kolmogorov equation instead of the Zakai equation is that the Zakai equation propagates a density that does not have unit mass. Therefore the solution must be normalized at each point in time before

useful data can be extracted. Also related to this paper is the work of Daum,²⁹ which derives a finite-dimensional nonlinear filter for continuous processes with discrete measurements. However, Daum makes several restrictive assumptions, most notably the use of a linear observer. In addition, Daum does not discuss the solution of the forward Kolmogorov equation, rather he assumes that the solution is known. The global Galerkin approximation method outlined in this paper could possibly be used in conjunction with Daum's filter. For continuous processes with continuous measurements, Brigo et al.³⁰ used a method related to our's, where the Kushner-Stratonovich equation is projected onto the tangent space of a finite dimensional manifold of probability densities to produce a finite dimensional approximation of the full nonlinear filter.

III. Approximate Filter

A. Prediction Equation

In this section we use the global or spectral Galerkin approximation method to reduce Eq. (3) from a function of space and time to a function of time only. The result is a linear ordinary differential equation that is easy to solve numerically. Galerkin's method is discussed in most textbooks on PDEs. Particularly good introductions to the topic can be found in Refs. 16, 17, and 31.

To apply Galerkin's method, we first assume that the solution of Eq. (3) satisfies

$$p(t, x | Y_t) = \sum_{\ell=0}^{\infty} b_{\ell}(t) \phi_{\ell}(x)$$

where equality is in the sense of the L_2 norm, and where $\{\phi_{\ell}\}_{\ell=0}^{\infty}$ is a complete set of basis functions for L_2 . We approximate p by truncating the sum:

$$p_N(t, x | Y_t) \triangleq \sum_{\ell=0}^{N-1} c_{\ell}(t) \phi_{\ell}(x) \quad (5)$$

where the coefficients c_{ℓ} satisfy the projection equation

$$\int_{\Omega} \left\{ \frac{\partial p_N}{\partial t} + \sum_{i=1}^n \frac{\partial (p_N f_i)}{\partial x_i} - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \frac{\partial^2 [p_N (G Q G^T)_{ij}]}{\partial x_i \partial x_j} \right\} \phi_q dx = 0$$

for $q = 0, \dots, N-1$. Because p does not have compact support, the set Ω should be \mathbb{R}^n . However, p will be almost zero over most of \mathbb{R}^n . Practically we must select Ω to be a closed and bounded subset of \mathbb{R}^n . The convergence results in Sec. IV assume that Ω is a sufficiently large compact set. Interchanging summation and differentiation, we obtain

$$\sum_{\ell=0}^{N-1} \dot{c}_{\ell} \int_{\Omega} \phi_{\ell} \phi_q dx = \sum_{\ell=0}^{N-1} c_{\ell} \left\{ - \sum_{i=1}^n \int_{\Omega} \frac{\partial [\phi_{\ell} f_i]}{\partial x_i} \phi_q dx + \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \int_{\Omega} \frac{\partial^2 [\phi_{\ell} (G Q G^T)_{ij}]}{\partial x_i \partial x_j} \phi_q dx \right\} \quad (6)$$

for $q = 0, \dots, N-1$. Equation (6) is a system of N linear ordinary differential equations that may be written in matrix notation by defining

$$c = [c_0, \dots, c_{N-1}]^T, \quad [M]_{i,j} = \int_{\Omega} \phi_i \phi_j dx \quad (7)$$

$$[A_1(t)]_{i,j} = - \sum_{k=1}^n \int_{\Omega} \frac{\partial [\phi_j f_k]}{\partial x_k} \phi_i dx \quad (8)$$

$$[A_2(t)]_{i,j} = \frac{1}{2} \sum_{k=1}^n \sum_{\ell=1}^n \int_{\Omega} \frac{\partial^2 [\phi_j (G Q G^T)_{k\ell}]}{\partial x_k \partial x_{\ell}} \phi_i dx \quad (9)$$

With these definitions, Eq. (6) can be rewritten as

$$\mathbf{M}\dot{\mathbf{c}} = [\mathbf{A}_1(t) + \mathbf{A}_2(t)]\mathbf{c}$$

Letting $\mathbf{A}_N(t) = \mathbf{M}^{-1}[\mathbf{A}_1(t) + \mathbf{A}_2(t)]$, we have

$$\dot{\mathbf{c}}(t) = \mathbf{A}_N(t)\mathbf{c}(t)$$

The initial condition for this ordinary differential equation can be obtained from either the initial density function $p(t_0, x)$ or from the measurement update $p(t_k, x | Y_k)$. In Sec. III.B we will show how to obtain the initial coefficients in the case of a measurement update. At time $t = t_0$, the initial coefficients are obtained by choosing $\mathbf{c}(t_0)$ such that

$$\begin{aligned} \int_{\Omega} p(t_0, x) \phi_q dx &= \int_{\Omega} p_N(x, t_0 | Y_{t_0}) \phi_q dx \\ &= \sum_{\ell=0}^{N-1} c_{\ell}(t_0) \int_{\Omega} \phi_{\ell} \phi_q dx \end{aligned}$$

for $q = 0, \dots, N-1$. To write this in matrix notation define

$$\mathbf{s} \triangleq \left[\int_{\Omega} p(t_0, x) \phi_0 dx, \dots, \int_{\Omega} p(t_0, x) \phi_{N-1} dx \right]^T$$

The initial conditions then become $\mathbf{M}\mathbf{c}(t_0) = \mathbf{s}$ or $\mathbf{c}(t_0) = \mathbf{M}^{-1}\mathbf{s}$.

We have therefore converted Eq. (3), describing the evolution of p into a system of ordinary differential equations

$$\dot{\mathbf{c}}_N(t) = \mathbf{A}_N(t)\mathbf{c}_N(t), \quad \mathbf{c}_N(t_0) = \mathbf{M}^{-1}\mathbf{s} \quad (10)$$

which describes the evolution of the coordinates \mathbf{c} of p_N in $\text{span}\{\phi_{\ell}\}_{\ell=0}^{N-1}$. Furthermore, if f , G , and Q are independent of time, then $\mathbf{A}_N(t) = \mathbf{A}_N$ and Eqs. (10) have the following simple solution that describes how \mathbf{c} changes between measurements: $\mathbf{c}(t) = e^{\mathbf{A}_N(t-t_k)}\mathbf{c}(t_k)$, $t \in [t_k, t_{k+1})$.

B. Measurement Update

At an observation, p satisfies Eq. (4). Again we apply Galerkin's method and replace p by p_N in Eq. (4) and define the updated approximate conditional density

$$p_N(x, t_{k+1} | Y_{t_{k+1}}) = \frac{p(y_{k+1} | x) p_N(x, t_{k+1}^- | Y_{t_k})}{\int_{\Omega} p(y_{k+1} | \xi) p_N(\xi, t_{k+1}^- | Y_{t_k}) d\xi}$$

Projecting this equation onto the space $\text{span}\{\phi_{\ell}\}_{\ell=0}^{N-1}$, we have

$$\begin{aligned} \sum_{\ell=0}^{N-1} c_{\ell}(t_{k+1}) \int_{\Omega} \phi_{\ell} \phi_q dx &= \frac{\sum_{\ell=0}^{N-1} c_{\ell}(t_{k+1}^-) \int_{\Omega} p(y_{k+1} | x) \phi_{\ell} \phi_q dx}{\sum_{\ell=0}^{N-1} c_{\ell}(t_{k+1}^-) \int_{\Omega} p(y_{k+1} | x) \phi_{\ell} dx} \end{aligned}$$

for $q = 0, \dots, N-1$. Or, in matrix notation,

$$\mathbf{c}(t_{k+1}) = \frac{\mathbf{M}^{-1} \mathbf{Y}(y_{k+1}) \mathbf{c}(t_{k+1}^-)}{\mathbf{v}(y_{k+1})^T \mathbf{c}(t_{k+1}^-)} \quad (11)$$

where

$$[\mathbf{Y}(y_{k+1})]_{q,\ell} = \int_{\Omega} p(y_{k+1} | x) \phi_{\ell} \phi_q dx \quad (12)$$

$$[\mathbf{v}(y_{k+1})]_{\ell} = \int_{\Omega} p(y_{k+1} | x) \phi_{\ell} dx \quad (13)$$

Therefore, at a measurement update, the solution to Eq. (10) is restarted with initial condition $\mathbf{c}(t_k)$ obtained from Eq. (11). It is important to note that the matrices $\mathbf{Y}(y_{k+1})$ and $\mathbf{v}(y_{k+1})$ depend on the measurement at time t_{k+1} , and hence cannot be computed off-line like the matrix \mathbf{A}_N in Eq. (10). To make this filter practical,

we need to be able to compute the elements of $\mathbf{Y}(y_k)$ and $\mathbf{v}(y_k)$ quickly. In Sec. V, we will use complex exponential basis functions to reduce implementation of the projections to computing DCTs and inverse discrete cosine transforms (IDCTs). Before doing so, we address the convergence of the NPF.

IV. Convergence

In this section, we state conditions that guarantee that the NPF given by Eqs. (10) and (11) converges (in the L_2 norm) to the exact nonlinear filter given by Eqs. (3) and (4). The essence of the proof is to show that the bound on the estimation error is given by a graph that looks like Fig. 1. The evolution of the error between the Galerkin approximation and the actual nonlinear filter is bounded by the curve, where ϵ can be made arbitrarily small by making the number of approximating terms N large. The bound on the error grows exponentially in between each sampling period and then jumps a finite amount at each measurement update. The bound derived in this section and depicted in Fig. 1 is conservative.

To make the convergence arguments transparent, we introduce some special notation that will only be used in this section. Assume that $\Omega \subset \mathbb{R}^n$ is a closed and bounded subset of \mathbb{R}^n and let \hat{p}_{Ω} be the truncation of p to Ω , i.e.,

$$\hat{p}_{\Omega}(t, x) = \begin{cases} p(t, x) & \text{if } x \in \Omega \\ 0 & \text{otherwise} \end{cases}$$

We assume that p is integrable, which implies that for every $\epsilon > 0$, there exists a closed and bounded set $\Omega \subset \mathbb{R}^n$ such that $\|p - \hat{p}_{\Omega}\|_{L_2} < \epsilon$, for all $t \in [t_0, t_1)$.

We write the Kolmogorov equation as

$$\frac{\partial p}{\partial t} + \mathcal{L}(p) = 0$$

where

$$\mathcal{L}(p) \triangleq \sum_{i=1}^n \frac{\partial(p f_i)}{\partial x_i} - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \frac{\partial^2 [p(G Q G^T)_{ij}]}{\partial x_i \partial x_j}$$

We assume that the basis functions $\{\phi_{\ell}\}_{\ell=0}^{\infty}$ form a complete basis for $L_2(\Omega)$ and that $\int_{\Omega} \mathcal{L}(\phi_j) \phi_k dx < \infty$ for all (j, k) . Using arguments similar to those outlined in Ref. 24, we can assume, without loss of generality, that the basis functions $\{\phi_{\ell}\}_{\ell=0}^{\infty}$ are orthonormal on the set Ω . This assumption is made to simplify the convergence proof, but does not impose a practical limitation on the selection of the basis functions. Practically, the basis functions are only required to be linearly independent.

Because $\hat{p}_{\Omega} \in L_1 \cap L_2$ it has a Fourier expansion, which we denote as

$$\hat{p}_{\Omega}(t, x | Y_t) = \sum_{\ell=0}^{\infty} b_{\ell}(t) \phi_{\ell}(x) \quad (14)$$

Because for each t in the interval of interest, \hat{p}_{Ω} satisfies Kolmogorov's equation on Ω , we get

$$\int_{\Omega} \left[\frac{\partial(\hat{p}_{\Omega} - p_N)}{\partial t} + \mathcal{L}(\hat{p}_{\Omega} - p_N) \right] \phi_q dx = 0$$

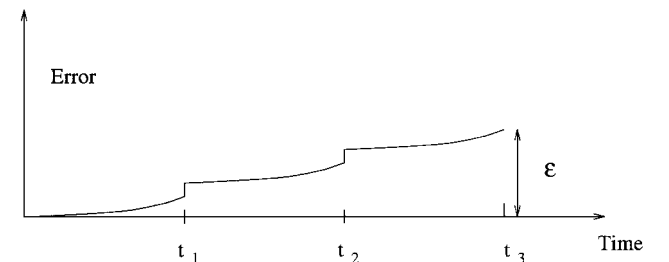


Fig. 1 Bound on the evolution of the approximation error. ϵ can be made arbitrarily small by making the order of the approximation N large enough.

$q = 0, \dots, N-1$. Using Eqs. (14) and (5) we get

$$\begin{aligned} & \sum_{\ell=0}^{N-1} (\dot{b}_\ell - \dot{c}_\ell) \int_{\Omega} \phi_\ell \phi_q \, dx + \sum_{\ell=0}^{N-1} (b_\ell - c_\ell) \int_{\Omega} \mathcal{L}(\phi_\ell) \phi_q \, dx \\ &= - \sum_{\ell=N}^{\infty} \dot{b}_\ell \int_{\Omega} \phi_\ell \phi_q \, dx - \sum_{\ell=N}^{\infty} b_\ell \int_{\Omega} \mathcal{L}(\phi_\ell) \phi_q \, dx \end{aligned}$$

$q = 0, \dots, N-1$. Define $\mathbf{b}_N = (b_0, \dots, b_{N-1})^T$ and $\Phi_N = (\phi_0, \dots, \phi_{N-1})^T$. Using the orthonormality of the basis functions and using the previously defined quantity A_N , we obtain

$$(\dot{\mathbf{b}}_N - \dot{\mathbf{c}}_N) + A_N(\mathbf{b}_N - \mathbf{c}_N) = - \sum_{\ell=N}^{\infty} b_\ell \int_{\Omega} \mathcal{L}(\phi_\ell)^{M^{-1}} \Phi_N \, dx \quad (15)$$

Define

$$\mu_N = \left\| \sum_{\ell=N}^{\infty} b_\ell \int_{\Omega} \mathcal{L}(\phi_\ell)^{M^{-1}} \Phi_N \, dx \right\|$$

The following theorem states the main result of this section.

Theorem 1. If Ω is compact, p is integrable, and the basis functions $\{\phi_\ell\}_{\ell=0}^{\infty}$ form a complete orthonormal basis for $\mathcal{L}_2(\Omega)$, and if in addition

$$\frac{\mu_N}{\|A_N\|} e^{\|A_N\|} \rightarrow 0 \quad \text{as } N \rightarrow \infty$$

then given a fixed time $\hat{t} > 0$, and any $\epsilon > 0$, there exists a K such that for $N > K$, $\|p(t, x | Y_t) - p_N(t, x | Y_t)\|_{L_2} < \epsilon$ for all $t < \hat{t}$.

The proof to this theorem is made transparent by the following three lemmas.

Lemma 1. Under the hypothesis of Theorem 1, for any $\epsilon > 0$, there exists a positive integer K and a $\delta_k > 0$ such that if $N > K$ and $\|\mathbf{b}_N(t_k) - \mathbf{c}(t_k)\| < \delta_k$, then $\|\mathbf{b}_N(t) - \mathbf{c}(t)\| < \epsilon$, for all $t \in [t_k, t_{k+1})$.

Proof. Since the differential equation $\dot{\xi} = -A_N \xi$ is globally Lipschitz with Lipschitz constant $\|A_N\|$, the theorem on the continuous dependence of ordinary differential equations on initial conditions and bounded disturbances,³² and Eq. (15), imply that

$$\|\mathbf{b}_N(t) - \mathbf{c}_N(t)\| \leq \delta_k e^{\|A_N\|(t-t_k)} + \frac{\mu_N}{\|A_N\|} (e^{\|A_N\|(t-t_k)} - 1) \quad (16)$$

The first term can be made less than $\epsilon/2$ by making

$$\delta_k < (\epsilon/2) \exp[-\|A_N\|(t_{k+1}^- - t_k)]$$

The hypothesis of Theorem 1 implies that there exists a K such that $N > K$ implies that the second term is less than $\epsilon/2$. QED

Lemma 2. Under the hypothesis of Theorem 1, for any $\epsilon > 0$, there exists a positive integer K and a $\delta_{k+1}^- > 0$ such that if $N > K$ and $\|\mathbf{b}_N(t_{k+1}^-) - \mathbf{c}_N(t_{k+1}^-)\| < \delta_{k+1}^-$, then $\|\mathbf{b}_N(t_{k+1}) - \mathbf{c}_N(t_{k+1})\| < \epsilon$.

Proof. Equation (11) shows that

$$\mathbf{c}(t_{k+1}) = \frac{M^{-1} \mathbf{Y}(y_{k+1}) \mathbf{c}(t_{k+1}^-)}{\mathbf{v}(y_{k+1})^T \mathbf{c}(t_{k+1}^-)} \quad (17)$$

If we truncate Eq. (4) to Ω and integrate, we obtain

$$\mathbf{b}(t_{k+1}) = \frac{M^{-1} \mathbf{Y}(y_{k+1}) \mathbf{b}(t_{k+1}^-) + \rho_{1N}}{\mathbf{v}(y_{k+1})^T \mathbf{b}(t_{k+1}^-) + \rho_{2N}} + \rho_{3N} \quad (18)$$

where

$$\rho_{1N} = \int_{\Omega} \sum_{\ell=N}^{\infty} b_\ell p(y_{k+1} | \xi) \phi_\ell(\xi) \Phi_N(\xi) \, d\xi$$

$$\rho_{2N} = \int_{\Omega} \sum_{\ell=N}^{\infty} b_\ell p(y_{k+1} | \xi) \phi_\ell(\xi) \, d\xi$$

$$\rho_{3N} = \int_{\Omega} \sum_{\ell=N}^{\infty} b_\ell \Phi(\xi) \, d\xi$$

Since the Fourier series converges, $\|\rho_{1N}\|$, $|\rho_{2N}|$, and $\|\rho_{3N}\|$ can be made less than any arbitrary positive number for N sufficiently large.

Subtracting Eq. (18) from Eq. (17) we obtain

$$\begin{aligned} & \|\mathbf{b}_N(t_{k+1}) - \mathbf{c}_N(t_{k+1})\| \\ & \leq \left\| \frac{M^{-1} \mathbf{Y} \mathbf{b}(t_{k+1}^-)}{\mathbf{v}^T \mathbf{b}(t_{k+1}^-) + \rho_{2N}} - \frac{M^{-1} \mathbf{Y} \mathbf{c}(t_{k+1}^-)}{\mathbf{v}^T \mathbf{c}(t_{k+1}^-)} \right\| \\ & \quad + \left\| \frac{\rho_{1N}}{\mathbf{v}^T \mathbf{b}(t_{k+1}^-) + \rho_{2N}} \right\| + \|\rho_{3N}\| \end{aligned}$$

Define

$$\alpha \triangleq \frac{1}{\mathbf{v}^T \mathbf{b}(t_{k+1}^-) + \rho_{2N}}, \quad \alpha_N \triangleq \frac{1}{\mathbf{v}^T \mathbf{c}(t_{k+1}^-)}$$

Using this notation we have the following sequence of inequalities:

$$\begin{aligned} & \left\| \frac{M^{-1} \mathbf{Y} \mathbf{b}(t_{k+1}^-)}{\mathbf{v}^T \mathbf{b}(t_{k+1}^-) + \rho_{2N}} - \frac{M^{-1} \mathbf{Y} \mathbf{c}(t_{k+1}^-)}{\mathbf{v}^T \mathbf{c}(t_{k+1}^-)} \right\| \\ &= \left\| \alpha M^{-1} \mathbf{Y} \mathbf{b}(t_{k+1}^-) - \alpha_N M^{-1} \mathbf{Y} \mathbf{c}(t_{k+1}^-) \right\| \\ &\leq \|M^{-1} \mathbf{Y}\| \|\alpha \mathbf{b}(t_{k+1}^-) - \alpha_N \mathbf{c}(t_{k+1}^-)\| \\ &= \|M^{-1} \mathbf{Y}\| \|\alpha (\mathbf{b}(t_{k+1}^-) - \mathbf{c}(t_{k+1}^-)) + (\alpha - \alpha_N) \mathbf{c}(t_{k+1}^-)\| \\ &\leq \|M^{-1} \mathbf{Y}\| [\|\alpha\| \|\mathbf{b}(t_{k+1}^-) - \mathbf{c}(t_{k+1}^-)\| + |\alpha - \alpha_N| \|\mathbf{c}(t_{k+1}^-)\|] \end{aligned}$$

Since $p(y_{k+1} | x)$ is a Gaussian density we know that $\|M^{-1} \mathbf{Y}(y_{k+1})\| = B_1 < \infty$ and $|\alpha| = B_2 < \infty$. We also know that $\|\mathbf{c}(t_{k+1}^-)\| = B_3 < \infty$, where B_1 , B_2 , and B_3 are finite numbers.

Also note that

$$\begin{aligned} |\alpha - \alpha_N| &= \left| \frac{1}{\mathbf{v}^T \mathbf{b}(t_{k+1}^-) + \rho_{2N}} - \frac{1}{\mathbf{v}^T \mathbf{c}(t_{k+1}^-)} \right| \\ &\leq \left\| \frac{\mathbf{v}}{[\mathbf{v}^T \mathbf{b}(t_{k+1}^-) + \rho_{2N}][\mathbf{v}^T \mathbf{c}(t_{k+1}^-)]} \right\| \\ &\quad \times \|\mathbf{b}_N(t_{k+1}^-) - \mathbf{c}_N(t_{k+1}^-)\| \\ &\quad + \left| \frac{\rho_{2N}}{[\mathbf{v}^T \mathbf{b}(t_{k+1}^-) + \rho_{2N}][\mathbf{v}^T \mathbf{c}(t_{k+1}^-)]} \right| \end{aligned}$$

The first quantity is bounded by a finite number B_4 , and the second quantity is a sequence, say ρ_{4N} that converges in N . Therefore

$$|\alpha - \alpha_N| \leq B_2 \|\mathbf{b}_N(t_{k+1}^-) - \mathbf{c}_N(t_{k+1}^-)\| + \rho_{4N}$$

Collecting the preceding results we get

$$\begin{aligned} \|\mathbf{b}_N(t_{k+1}) - \mathbf{c}_N(t_{k+1})\| &\leq B_1 [B_2 \|\mathbf{b}_N(t_{k+1}^-) - \mathbf{c}_N(t_{k+1}^-)\| \\ &\quad + (B_4 \|\mathbf{b}_N(t_{k+1}^-) - \mathbf{c}_N(t_{k+1}^-)\| + \rho_{4N}) B_2] + \rho_{3N} + \rho_{5N} \end{aligned}$$

where

$$\rho_{5N} = \left\| \frac{\rho_{1N}}{\mathbf{v}^T \mathbf{b}(t_{k+1}^-) + \rho_{2N}} \right\|$$

The lemma follows from this formula.

QED

Lemma 3. If at a particular instant of time t ,

$$\|\mathbf{b}_N(t) - \mathbf{c}_N(t)\| \rightarrow 0$$

as $N \rightarrow \infty$, then

$$\|\hat{p}_\Omega(t, x | Y_t) - p_N(t, x | Y_t)\|_{L_2} \rightarrow 0$$

Proof.

$$\begin{aligned} \|\hat{p}_\Omega - p_N\|_{L_2}^2 &= \int_\Omega |\hat{p}_\Omega - p_N|^2 dx \leq \int_\Omega |(\mathbf{b}_N - \mathbf{c}_N)\Phi_N|^2 dx \\ &+ \int_\Omega \left| \sum_{\ell=N}^\infty b_\ell \phi_\ell \right|^2 dx = (\mathbf{b}_N - \mathbf{c}_N)^T \mathbf{M} (\mathbf{b}_N - \mathbf{c}_N) \\ &+ \int_\Omega \left| \sum_{\ell=N}^\infty b_\ell \phi_\ell \right|^2 dx \end{aligned}$$

By the mean value theorem of calculus, there exists a $\xi \in \Omega$ such that

$$\|\hat{p}_\Omega - p_N\|_{L_2}^2 \leq \|\mathbf{b}_N - \mathbf{c}_N\|^2 + \lambda(\Omega) \left| \sum_{\ell=N}^\infty b_\ell \phi_\ell(\xi) \right|^2$$

where $\lambda(\Omega)$ is the Lebesgue measure of Ω . Because the second term on the right-hand side is the tail of a convergent Fourier series, it converges to zero. The lemma then follows from the hypothesis on $\|\mathbf{b}_N - \mathbf{c}_N\|$. QED

Proof of Theorem 1. First note that

$$\|p - p_N\|_{L_2} \leq \|p - \hat{p}_\Omega\|_{L_2} + \|\hat{p}_\Omega - p_N\|_{L_2}$$

Because an appropriate choice of Ω guarantees that $\|p - \hat{p}_\Omega\|_{L_2} < \epsilon/2$ for all $t < \hat{t}$, we need to show that $\|\hat{p}_\Omega - p_N\|_{L_2} < \epsilon/2$ for all $t < \hat{t}$.

Given $t_k < \hat{t} \leq t_{k+1}$ and ϵ , Lemma 3 indicates that there exists a δ such that if $\|\mathbf{b}_N(\hat{t}) - \mathbf{c}_N(\hat{t})\| < \delta$, then $\|p(\hat{t}, x | Y_{\hat{t}}) - p_N(\hat{t}, x | Y_{\hat{t}})\|_{L_2} < \epsilon/2$. From Lemma 1, there exists a K_k and a δ_k such that if $N > K_k$ and $\|\mathbf{b}_N(t_k) - \mathbf{c}_N(t_k)\| < \delta_k$, then $\|\mathbf{b}_N(\hat{t}) - \mathbf{c}_N(\hat{t})\| < \delta$. From Lemma 2, there exists a K_k^- and a δ_k^- such that if $N > K_k^-$ and $\|\mathbf{b}_N(t_k^-) - \mathbf{c}_N(t_k^-)\| < \delta_k^-$, then $\|\mathbf{b}_N(t_k) - \mathbf{c}_N(t_k)\| < \delta_k$. Another application of Lemma 1 implies that there exists a K_{k-1} and a δ_{k-1} such that if $N > K_{k-1}$ and if $\|\mathbf{b}_N(t_{k-1}) - \mathbf{c}_N(t_{k-1})\| < \delta_{k-1}$, then $\|\mathbf{b}_N(t_k^-) - \mathbf{c}_N(t_k^-)\| < \delta_k^-$.

The preceding argument is repeated k times, resulting in the requirement that $\|\mathbf{b}_N(0) - \mathbf{c}_N(0)\| < \delta_0$. If we assume that the basis functions are orthonormal, then at time $t = 0$ there is zero approximation error because

$$\hat{p}_\Omega(0, x | Y_0) = p_0(x)$$

and

$$\int_\Omega p_N(0, t | Y_0) \Phi_N dx = \int_\Omega p_0(x) \Phi_N dx$$

implies that

$$\|\mathbf{b}_N(0) - \mathbf{c}_N(0)\| = 0$$

Therefore letting

$$K = \max_{0 \leq \ell \leq k} \{K_\ell, K_\ell^-\}$$

proves the theorem. QED

Remark 1. The bound derived in the preceding proof is given qualitatively by Fig. 1.

Remark 2. The requirement in the hypothesis of Theorem 1 that $(\mu_N / \|\mathbf{A}_N\|) e^{\|\mathbf{A}\|} \rightarrow 0$ is somewhat unsatisfactory because it is not clear at this stage how to guarantee that this is true a priori. The requirement, however, can be tested for any given system by computing \mathbf{A}_N and μ_N for various values of N and testing to see if the expression is converging to zero. Further research could focus on deriving conditions under which this hypothesis holds.

Remark 3. The convergence proof given in this section does not give explicit bounds on the approximation error. In other words, given a desired ϵ , we cannot say what K must be to guarantee that

an approximation error of ϵ is achieved. We can only say that such a K does in fact exist.

Remark 4. The convergence proof guarantees a small approximation error for any finite time. The arguments given in this section do not say anything about the steady-state error.

V. Cosine Basis

In this section we show how to implement Eqs. (10) and (11), using a cosine basis and the DCT fast transform algorithm. The reasons for choosing a cosine series as the basis for the Hilbert space are 1) fast algorithms for implementing the DCT exist³³ and may be used to approximate the integrals of Sec. III, and 2) each function in this set satisfies boundary conditions that adequately approximates the boundary conditions that should be met for this problem. The boundary condition that should be imposed is that of a totally absorbing boundary similar to matching the impedance at the terminal end of a wave guide. Because systematic methods for determining such boundary conditions for PDEs of the form of Eq. (3) are elusive at this point, we impose the condition that the spatial derivatives be zero at the boundaries. This is accomplished with a cosine basis because each function in the set satisfies this condition. Simulations have shown that this is an adequate approximation. We now proceed to show how the DCT may be used to implement the NPF.

Because the general n -dimensional case is notationally messy and the basic ideas are contained in one dimension, we will restrict ourselves to this case. Suppose we choose $\Omega = [a, b]$, then the cosine basis set is

$$\{\phi_\ell(x)\}_{\ell=0}^{N-1} = \begin{cases} \frac{1}{\sqrt{b-a}} & \ell = 0 \\ \sqrt{\frac{2}{b-a}} \cos\left(\frac{2\pi\ell}{b-a}(x-a)\right) & 1 \leq \ell \leq N-1 \end{cases} \quad (19)$$

This set of basis functions for S_N are the first N functions in a complete orthonormal sequence that spans the Hilbert space $L_2[a, b]$.

The DCT algorithm can be used to approximate the sequence of inner products of the real valued function $\eta(x)$ with each of the basis functions (19) simultaneously as follows:

$$\int_a^b \eta(x) \phi_k(x) dx \approx \sqrt{\frac{b-a}{N}} \text{DCT}_k[\eta(\zeta)] \quad (20)$$

where the input to the DCT algorithm is the N samples of $\eta(\cdot)$ over the interval $[a, b]$ given by

$$\zeta = [(2\zeta + 1)/2N](b-a) + a, \quad \zeta = 0, \dots, N-1 \quad (21)$$

Therefore, the DCT returns N -weighted integrals of η that would normally require $\mathcal{O}(N^2)$ operations; because the DCT is a fast transform [$\mathcal{O}(N \log N)$ similar to the fast Fourier transform], this results in substantial reduction in computational burden for large N .

With Eq. (20), we can use the DCT to approximate Eqs. (7-9) and Eqs. (12-13) as

$$\begin{aligned} \mathbf{M} &= \mathbf{I} \\ [\mathbf{A}_1(t)]_\ell &= \sqrt{\frac{b-a}{N}} \text{DCT} \left[\left(\frac{\partial \phi_\ell f}{\partial x} \right) (\zeta) \right] \\ [\mathbf{A}_2(t)]_\ell &= \sqrt{\frac{b-a}{N}} \text{DCT} \left\{ \frac{Q}{2} \left[\frac{\partial^2 (\phi_\ell G^2)}{\partial x^2} \right] (\zeta) \right\} \\ s &= \frac{\sqrt{b-a}}{N} \text{DCT}[p(\zeta, t_0)] \\ [\mathbf{Y}_k]_\ell &= \frac{\sqrt{b-a}}{N} \text{DCT}[p(y_k | \zeta) \phi_\ell(\zeta)] \\ v(y_k) &= \frac{\sqrt{b-a}}{N} \text{DCT}[p(y_k | \zeta)] \end{aligned}$$

These equations show that all of the quantities needed to propagate $\mathbf{c}(t)$ between measurements and to update it at a measurement can be computed quickly using the DCT.

In particular, the procedure for implementing the filter is as follows. First note that \mathbf{M} is the identity for this choice of basis functions and thus may be ignored. Second we compute $\mathbf{c}(t_0)$ by arraying the samples of $p(\zeta, t_0)$ [where ζ takes on the values given by Eq. (21)] in a vector, applying the DCT algorithm, and scaling the output by $\sqrt{[(b-a)/N]}$. Each column of $\mathbf{A}_1(t_0)$ and $\mathbf{A}_2(t_0)$ is computed in a similar manner and their sum becomes $\mathbf{A}_N(t_0)$ (actually the integrands associated with corresponding elements of $\mathbf{A}_1(t_0)$ and $\mathbf{A}_2(t_0)$ may be summed prior to taking the DCT). As previously stated, if f , G , and Q are independent of time, then so is $\mathbf{A}_N(t_0)$, so that \mathbf{A}_N and its associate state transition matrix $e^{\mathbf{A}_N(t_1-t_0)}$ may be computed off-line; otherwise, they must be computed at each time step. The parameter vector $\mathbf{c}(t_1^-)$ for the conditional density just prior to the first measurement is given by

$$\mathbf{c}(t_1^-) = e^{\mathbf{A}_N(t_1-t_0)} \mathbf{c}(t_0) \quad (22)$$

At the observation time t_1 , the measurement update is performed as follows. First, the vector $\mathbf{v}(y_1)$ and the ℓ th column of $\mathbf{Y}(y_1)$ for each $\ell = 0, \dots, N-1$ are calculated in a similar manner as $\mathbf{c}(t_0)$ and the columns of $\mathbf{A}_N(t_0)$. These quantities cannot be computed off-line because they depend on the measurement. The variables $\mathbf{Y}(y_1)$ and $\mathbf{v}(y_1)$ are then used to compute the updated parameter vector $\mathbf{c}(t_1)$ according to Eq. (11). This vector is then used as the initial condition to Eq. (22) for the next time interval, and the cycle repeats.

As the output of our filter we may wish to compute the conditional mean, the conditional covariance or the density function itself. These quantities can also be computed efficiently using the DCT as shown next.

Conditional Mean. By definition the conditional mean of the state at time t based on observations up to and including time $\tau \leq t$ is

$$\bar{x}_t^\tau = \int \xi p(\xi, t | Y_\tau) d\xi$$

We may compute the mean of our approximate conditional density as follows:

$$\begin{aligned} \hat{x}_t^\tau &= \int \xi p_N(\xi, t | Y_\tau) d\xi = \int \xi \sum_{\ell=0}^{N-1} c_\ell(t) \phi_\ell(\xi) d\xi \\ &= \sum_{\ell=0}^{N-1} c_\ell(t) \int_{\Omega} x \phi_\ell dx = \gamma^T \mathbf{c} \end{aligned}$$

where $\gamma = \sqrt{[(b-a)/N]} \text{DCT}[\zeta]$.

Conditional Covariance. The covariance of the approximate conditional density function is

$$\begin{aligned} \hat{P}_t^\tau &= \int \xi^2 \hat{p}_N(\xi, t | Y_\tau) d\xi \\ &- (\hat{x}_t^\tau)^2 = \sum_{\ell=0}^{N-1} c_\ell(t) \int_{\Omega} x^2 \phi_\ell dx - (\gamma^T \mathbf{c})^2 \\ &= \Gamma^T \mathbf{c} - (\gamma^T \mathbf{c})^2 \end{aligned}$$

where $\Gamma = \sqrt{[(b-a)/N]} \text{DCT}[\zeta^2]$.

Density Function. To recover the approximate density function itself we have

$$p_N(x, t | Y_t) = \sum_{\ell=0}^{N-1} c_\ell(t) \phi_\ell(x) = \sqrt{\frac{N}{b-a}} \text{IDCT}[\mathbf{c}(t)]$$

VI. Examples

In this section we give simulation results comparing the NPF to the EKF for three continuous discrete filtering scenarios: 1) a standard nonlinear filtering problem, 2) an unobservable bimodal system, and 3) a two-dimensional system with an extremely poor sensor. The objective of the first example is to show that for a straightforward nonlinear system, the NPF and the EKF give similar results. The purpose of the second two examples is not to discredit the EKF, but rather to show that there are some simple examples whose nonlinearities are too difficult for the EKF, and yet can be handled by the NPF. In all three scenarios the EKF and the NPF have the same initial conditions and parameters. For comparison, state estimates are plotted with the true state trajectory.

A. Nominal Example

This scenario compares the performance of the EKF and NPF for a standard filtering problem. The system dynamics and output equation are given by

$$dx_t = \sin(x_t) dt + d\beta_t, \quad y_k = x(t_k) + v_k$$

where $Q = 0.5$ and $R = 0.5$. The sampling period is $T = 0.1$ s. The initial state and the initial state estimate are both set to zero. The initial covariance estimate is one. The approximation interval is $\Omega = [-5\pi/2, 5\pi/2]$ and the number of Fourier basis elements is $N = 128$. Figure 2 shows the true state trajectory as well as the estimated state trajectories produced by the NPF and EKF. As a figure of merit, the rms errors between the true and estimated states is 0.275 for NPF and 0.235 for the EKF. Figure 2 shows that both the NPF and the EKF give very good state estimates.

B. Unobservable System

Consider the system

$$dx_t = \sin(x_t) dt + d\beta_t \quad (23)$$

$$y_k = |x_{t_k}| + v_k \quad (24)$$

which is nonlinear unobservable at the origin. In the absence of noise, system (23) has stable equilibria at the points $\{\pi \pm 2\pi\ell, \ell = 0, 1, \dots\}$ and unstable equilibria at the points $\{\pm 2\pi\ell, \ell = 0, 1, \dots\}$. In the presence of noise the state x_t will tend to the closest stable equilibrium. The interesting thing is that the output cannot be used to distinguish between positive and negative values of the state variable. This means that there is very little information available to distinguish between positive equilibria and the corresponding negative ones. Because we cannot distinguish between positive and negative values of x , we expect that the conditional density function for the state will be multimodal with major lobes centered at approximately $\pi \pm 2\pi\ell, \ell = 0, 1, \dots$. It is interesting to note that the probability mass that is initially located to the right of the origin will migrate to an equilibrium point with a positive coordinate while the mass that starts out to the left of the origin will tend to migrate to a negative equilibrium. The conditional mean at each time instance will depend largely on the prior distribution that is “seldom precisely known and generally rather arbitrarily specified” (Ref. 3, p. 244). The EKF for this problem is given by

$$\frac{d\hat{x}_t}{dt} = \sin(\hat{x}_t), \quad \dot{P}_t = 2P_t \cos(\hat{x}_t) + Q$$

$$K_k = \begin{cases} \frac{P_{t_k}^-}{P_{t_k}^- + R} & \hat{x}_{t_k}^- \geq 0 \\ \frac{-P_{t_k}^-}{P_{t_k}^- + R} & \hat{x}_{t_k}^- < 0 \end{cases}$$

$$P_{t_k} = \begin{cases} (1 - K_k) P_{t_k}^- & \hat{x}_{t_k}^- \geq 0 \\ (1 + K_k) P_{t_k}^- & \hat{x}_{t_k}^- < 0 \end{cases}$$

$$\hat{x}_{t_k} = \hat{x}_{t_k}^- + K_k (y_k - |\hat{x}_{t_k}^-|)$$

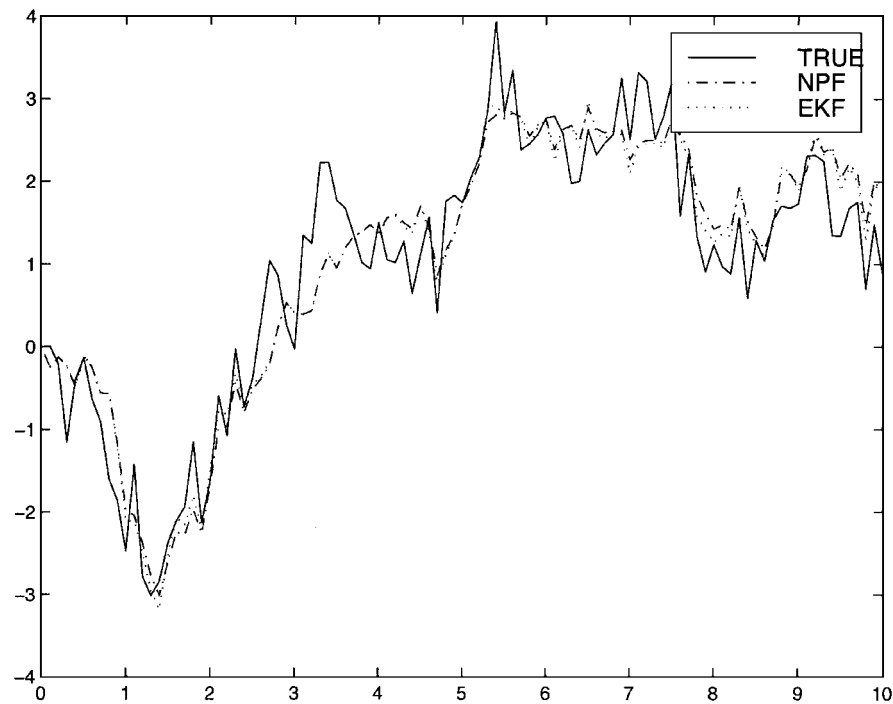


Fig. 2 Comparison of state trajectory estimates of an EKF and the NPF. Points are plotted at observation times.

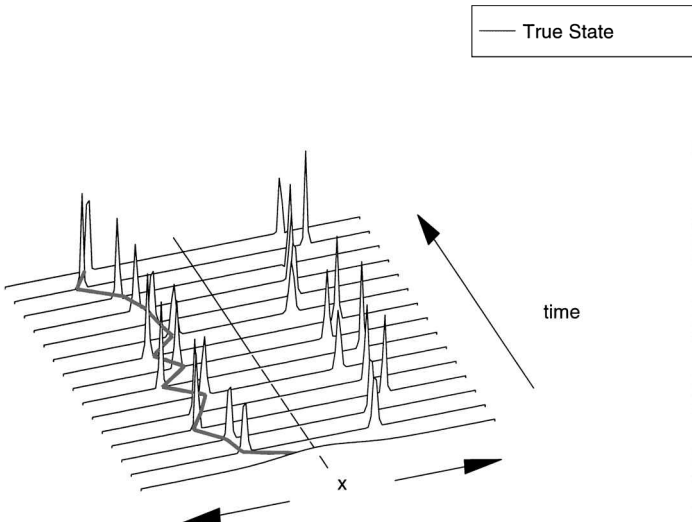


Fig. 3 Approximate density function and the true state at the measurement updates.

From the prediction equation we can see that the EKF has multiple equilibria as well. Which equilibrium the filter predicts as the state will depend on the initial condition $\hat{x}(0)$, and it will not approximate the correct conditional mean unless the prior distribution is accurate and all of the prior probability density mass is located on one side of the origin. If $\hat{x}(0) < 0$, then $\hat{x}(t)$ goes to an equilibrium to the left and the covariance converges to an inaccurately low value unless the preceding condition is satisfied. An analogous situation holds for $\hat{x}(0) > 0$. If the initial state of the system is zero, then the noise during the initial transients of the system will determine whether the system goes to the left or to the right. It is therefore impossible to predict a priori the sign with which we should initialize the EKF.

The output of the NPF is shown in Fig. 3. The conditional density functions are plotted at discrete times. The thick line shows the true state of the system. After a few measurements the predicted density function is bimodal with modes approximately centered at plus and minus the absolute value of the actual state, which is what we expect.

C. Second-Order System

In this section we consider the NPF for a system with a two-dimensional state space and a bounded region of attraction. The system dynamics are given by

$$dx = \begin{bmatrix} -x_2 \\ 0.2(x_1^2 - 1)x_2 + x_1 \end{bmatrix} + d\beta \tag{25}$$

The covariance of $d\beta$ used for this example is $0.5I$ for all time. This system has an unstable limit cycle that bounds the region of attraction to the equilibrium point at the origin. In the absence of noise, state trajectories with initial points inside of the region of convergence spiral inward to the origin, while those starting outside of this region diverge.

The measurement function for this system is the scalar function shown in Fig. 4a as a contour plot. This form of measurement function corresponds to a sensor that gives no directional information. The inner contour bounds a region where the sensor is saturated and the signal is nearly one whenever the system state enters this region. Beyond the outer contour the signal is nearly zero. The region between the two contours corresponds to an approximately linear region. The covariance for the measurement noise used in this example is $R_k = 0.17$ for all k .

The simulated results for this example are shown in Fig. 4b. The true state trajectory and the trajectories of the conditional means produced by the NPF and the EKF are plotted. The prior density for both filters is Gaussian with covariance $0.1I$. After a few measurement updates we see that the EKF has made a fatal error. The EKF measurement function at each measurement is linearized about the current mean estimate, but the actual state is outside of the region where the linearization is valid. The result is that the EKF estimate jumps outside of the region of attraction. The linearization point is now in a region where the measurement function has nearly zero gradient, so that the linearized system is almost completely unobservable. The EKF state estimate now evolves open loop, and because it is outside of the region of attraction, it ultimately diverges even though the true state remains inside the region of attraction. As seen in Fig. 4b the NPF performs reasonably well given the limited information available from the measurements and does not fail because of the nonlinearities in the system and measurement characteristics.

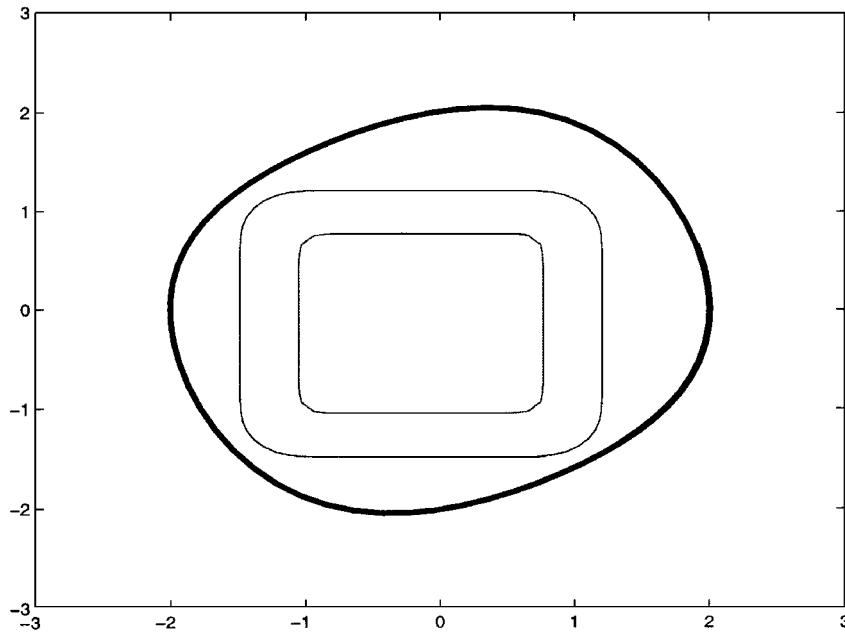


Fig. 4a Bound for the region of attraction for the dynamic system (25) and the contour plot of the measurement function.

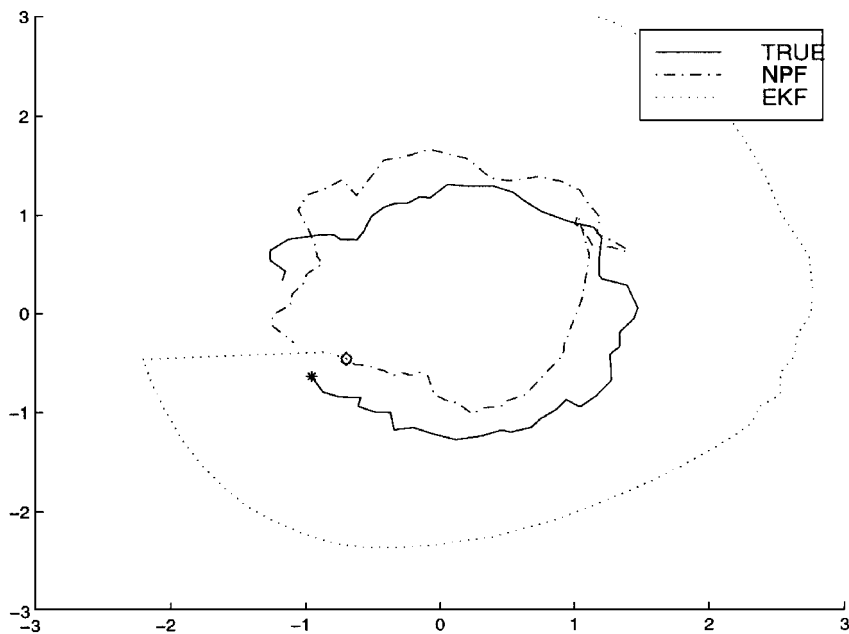


Fig. 4b Simulated state and estimated state trajectories.

VII. Conclusions

In this paper we have used the Galerkin spectral method to derive a nonlinear filter that approximates the exact evolution of the conditional density function of a nonlinear system. We have shown that our filter converges to the true density function as the order of the approximation is increased. We have also shown how to use the DCT algorithm efficiently to implement the algorithm and demonstrated several simple examples where our filter outperforms the EKF. The results contribute to the state of the art in nonlinear filtering by deriving an approximate filter that returns the entire density function and is not based on the linearization of the system. Various filters other than the minimum variance filter could be derived from our approximation because we approximate the conditional density on the state and not its moments. The major limitation of the NPF filter is the curse of dimensionality in that the projection equations must be implemented over an n -dimensional state space. The computational complexity and the memory requirements grow exponentially with n . This currently limits the practical applicabil-

ity of this filter to low-order systems. It is envisioned that the filter could be used to process data from sensors with severely nonlinear output characteristics.

Acknowledgment

This paper was presented at the 1997 American Control Conference.

References

- ¹Kalman, R. E., "A New Approach to Linear Filtering and Prediction Problems," *Transactions of the ASME, Journal of Basic Engineering*, Vol. 82, No. 1, 1960, pp. 34, 35.
- ²Kalman R. E., and Bucy, R. S., "New Results in Linear Filtering and Prediction Theory," *Transaction of the ASME, Journal of Basic Engineering*, Vol. 83, 1961, 95–108.
- ³Jazwinski, A. H., *Stochastic Processes and Filtering Theory*, Vol. 64, *Mathematics in Science and Engineering*, Academic Press, New York, 1970.
- ⁴Kushner, H. J., "Dynamical Equations for Optimal Nonlinear Filter," *Journal of Differential Equations*, Vol. 3, No. 2, 1967, pp. 179–190.

- ⁵Kushner, H. J., "Nonlinear Filtering: The Exact Dynamical Equations Satisfied by the Conditional Mode," *IEEE Transactions on Automatic Control*, Vol. 12, 1967, pp. 262-267.
- ⁶Stratonovich, R. L., "Conditional Markov Processes," *Theory of Probability and Its Applications*, Vol. 5, 1960, pp. 156-178.
- ⁷Bryson, A. E., and Ho, Y. C., *Applied Optimal Control*, Hemisphere, New York, 1975.
- ⁸Lewis, F. L., *Optimal Estimation: With an Introduction to Stochastic Control Theory*, Wiley, New York, 1986.
- ⁹Ohap, R. F., and Stubberud, A. R., "A Technique for Estimating the State of a Nonlinear System," *IEEE Transactions on Automatic Control*, Vol. 10, No. 4, 1965, pp. 150-155.
- ¹⁰Gelb, A. (ed.), *Applied Optimal Estimation*, MIT Press, Cambridge, MA, 1974.
- ¹¹Ljung, L., "Asymptotic Behavior of the Extended Kalman Filter Using Filtered State Estimates," *IEEE Transactions on Automatic Control*, Vol. 24, No. 1, 1979, pp. 36-50.
- ¹²Ursin, B., "Asymptotic Convergence Properties of the Extended Kalman Filter Using Filtered State Estimates," *IEEE Transactions on Automatic Control*, Vol. 25, No. 12, 1980, pp. 1207-1211.
- ¹³Galkowski, P. J., and Islam, M. A., "An Alternative Derivation of the Modified Gain Function of Song and Speyer," *IEEE Transactions on Automatic Control*, Vol. 36, No. 11, 1991, pp. 1323-1326.
- ¹⁴Bell, B. M., and Cathey, F. W., "The Iterated Kalman Filter Update as a Gauss-Newton Method," *IEEE Transactions on Automatic Control*, Vol. 38, No. 2, 1993, pp. 294-297.
- ¹⁵Bell, B. M., "The Iterated Kalman Smoother as a Gauss-Newton Method," *SIAM Journal of Optimization*, Vol. 4, No. 8, 1994, pp. 626-636.
- ¹⁶Mikhlin, S. G., *Variational Methods in Mathematical Physics*, MacMillan, New York, 1964.
- ¹⁷Kantorovich, L. V., and Krylov, V. I., *Approximate Methods of Higher Analysis*, Interscience Publishers, New York, 1958.
- ¹⁸Risken, H., *The Fokker-Planck Equation: Methods of Solution and Application*, Vol. 18, *Springer Series in Synergetics*, Springer-Verlag, Berlin, 1989.
- ¹⁹Ling, F. H., and Wu, X. X., "Fast Galerkin Method and Its Application to Determine Periodic Solutions of Non-Linear Oscillators," *International Journal of Non-Linear Mechanics*, Vol. 22, No. 2, 1987, pp. 89-98.
- ²⁰Bouc, R., and Defilippi, M., "A Galerkin Multiharmonic Procedure for Nonlinear Multidimensional Random Vibration," *International Journal of Engineering Science*, Vol. 25, No. 6, 1987, pp. 723-733.
- ²¹Weber, H., "Nonlinear Solution of a Class of Nonlinear Boundary Value Problems for Analytic Functions," *Journal of Applied Mathematics and Physics*, Vol. 33, May 1982, pp. 301-314.
- ²²Beard, R., "Improving the Closed-Loop Performance of Nonlinear Systems," Ph.D. Thesis, Dept. of Electrical, Computer, and Systems Engineering, Rensselaer Polytechnic Inst., Troy, NY, 1995.
- ²³Beard, R., Saridis, G., and Wen, J., "Improving the Performance of Stabilizing Control for Nonlinear Systems," *Control Systems Magazine*, Vol. 16, No. 10, 1996, pp. 27-35.
- ²⁴Beard, R., Saridis, G., and Wen, J., "Galerkin Approximation of the Generalized Hamilton-Jacobi-Bellman Equation," *Automatica*, Vol. 33, No. 12, 1997, pp. 2159-2177.
- ²⁵Beard, R., Saridis, G., and Wen, J., "Approximate Solution to the Time-Invariant Hamilton-Jacobi-Bellman Equation," *Journal of Optimization Theory and Application*, Vol. 96, March 1998, pp. 589-626.
- ²⁶Baras, J. S., "Real Time Architectures for the Zakai Equation and Application," *Stochastic Analysis: Liber Amicorum for Moshe Zakai*, edited by E. Mayer-Wolf, E. Merzbach, and A. Shwartz, Academic Press, Boston, MA, 1991, pp. 15-38.
- ²⁷Zakai, M., "On the Optimal Filtering of Diffusion Processes," *Zeitschrift für Wahrscheinlichkeitstheorie und verw. Geb.*, Vol. 2, 1969, pp. 230-243.
- ²⁸Daum, F. E., "Solution of the Zakai, Equation by Separation of Variables," *IEEE Transactions on Automatic Control*, Vol. 32, No. 10, 1987, pp. 941-943.
- ²⁹Daum, F. E., "Exact Finite-Dimensional Nonlinear Filters," *IEEE Transactions on Automatic Control*, Vol. 31, No. 7, 1986, pp. 616-622.
- ³⁰Brigo, D., Hanzon, B., and LeGlanc, F., "A Differential Geometric Approach to Nonlinear Filtering: The Projection Filter," *IEEE Transactions on Automatic Control*, Vol. 43, No. 2, 1998, pp. 247-252.
- ³¹Zeidler, E., *Nonlinear Functional Analysis and Its Applications, II/A: Linear Monotone Operators*, Springer-Verlag, Berlin, 1990.
- ³²H. K., Khalil, *Nonlinear Systems*, Macmillan, New York, 1992.
- ³³Jain, A. K., *Fundamentals of Digital Image Processing*, Prentice-Hall, Englewood Cliffs, NJ, 1989.