# The University of Reading

# Unbiased Ensemble Square Root Filters

D.M. Livings, S.L. Dance, and N.K. Nichols

*Department of Mathematics*

*The University of Reading*

*Whiteknights, PO Box 220*

*Reading*

*Berkshire RG6 6AX*

# Department of Mathematics

**Abstract**

There is an established framework that describes a class of ensemble Kalman filter algorithms as square root filters (SRFs). These schemes carry out analyses by updating the ensemble mean and a square root of the ensemble covariance matrix. The matrix square root of the forecast covariance is post-multiplied by another matrix to give a matrix square root of the analysis covariance. The choice of post-multiplier is not unique, but can be multiplied by any orthogonal matrix to give another scheme that also fits into this framework.

In this work we re-examine the ensemble SRF framework. The key result is that not all filters of this type bear the desired relationship to the forecast ensemble: there can be a systematic bias in the analysis ensemble mean and consequently an accompanying shortfall in the spread of the analysis ensemble as expressed by the ensemble covariance matrix. This points to the need for a restricted version of the notion of an ensemble SRF, which we call an unbiased ensemble SRF. We have established a set of necessary and sufficient conditions for the scheme to be unbiased. Many (but not all) published ensemble SRF algorithms satisfy them. Whilst these conditions are not a cure-all and cannot deal with independent sources of bias such as model and observation errors, they should be useful to designers of ensemble SRFs in the future.

# 1 Introduction

Data assimilation seeks to solve the following problem: given an imperfect discrete model of the dynamics of a system and noisy observations of the system, find estimates of the state of the system. Sequential data assimilation techniques break this problem into a cycle of alternating forecast and analysis steps. In the forecast step the system dynamical model is used to evolve an earlier state estimate forward in time, giving a forecast state at the time of the latest observations. In the analysis step the observations are used to update the forecast state, giving an improved state estimate called the analysis. This analysis is used as the starting point for the next forecast.

Sequential data assimilation techniques include the optimal linear Kalman filter (KF) and its nonlinear generalisation, the extended Kalman filter (EKF) (Gelb, 1974; Jazwinski, 1970). As well as an estimate of the state of the system, these filters maintain an error covariance matrix that acts as a measure of the uncertainty in the estimate. This covariance matrix is updated with every analysis to reflect the new information provided by the observations, and is evolved along with the state estimate from the time of an analysis to the time of the next forecast. This maintenance of a flow-dependent covariance matrix is a significant advantage of Kalman-type filters over variational data assimilation techniques such as 4D-Var that are popular in operational numerical weather forecasting (NWP) systems. However, the EKF rose to prominence in aerospace applications where the dimension of the state space for the model is relatively small, typically nine or less. Directly extending the filter to NWP systems where the state space dimension may be $10^7$ is beyond the capabilities of current computer technology. The EKF also shares with 4D-Var the need to implement tangent linear operators (Jacobians) for the nonlinear forecast model and the model of how observations are related to the state of the system. This is a laborious activity for a large, complicated system such as an NWP model.

The ensemble Kalman filter (EnKF) is an attempt to overcome the drawbacks of the EKF. Two of its key ideas are to use an ensemble (statistical sample) of state estimates instead of a single state estimate and to calculate the error covariance matrix from this ensemble instead of maintaining a separate covariance matrix. If

the ensemble size is small, but not so small that it is statistically unrepresentative, then the extra work needed to maintain an ensemble of state estimates is more than offset by the work saved through not maintaining a separate covariance matrix. The EnKF also does not use tangent linear operators, which eases implementation and may lead to a better handling of nonlinearity. The KF aspect of the EnKF appears in the analysis step, which is designed so that the implied updates of the ensemble mean and ensemble covariance matrix mimic those of the state vector and covariance matrix in the standard KF.

The EnKF was originally presented in Evensen (1994). An important subsequent development was the recognition by Burgers et al. (1998) (and independently by Houtekamer and Mitchell (1998)) of the need to use an ensemble of pseudo-random observation perturbations to obtain the right statistics from the analysis ensemble. Deterministic methods for forming an analysis ensemble with the right statistics have also been presented. The former approach to the EnKF is comprehensively reviewed in Evensen (2003), whilst previously-published variants of the latter approach are placed in a uniform framework in Tippett et al. (2003). These variants include the ensemble transform Kalman filter (ETKF) of Bishop et al. (2001), the ensemble adjustment Kalman filter (EAKF) of Anderson (2001), and the filter of Whitaker and Hamill (2002). Filters that fit into the general framework are known as ensemble square root filters (SRFs).

This paper extends the results of Tippett et al. (2003) in two directions. Firstly, it explicitly shows that the ensemble SRF framework is sufficiently general to encompass all possible deterministic formulations of the analysis step of the EnKF, not just the specific cases considered by those authors. Secondly, and more importantly, it shows that the ensemble SRF framework also encompasses filters in which the analysis ensemble statistics do not bear the desired KF-like relationship to the forecast ensemble: there can be a systematic bias in the analysis ensemble mean and an accompanying shortfall in the spread of the analysis ensemble as expressed by the ensemble covariance matrix. This points to the need for a restricted version of the notion of an ensemble SRF, which we call an *unbiased ensemble SRF*. The unbiased ensemble SRF framework is still general enough to include all deterministic

3

formulations of the analysis step of the EnKF whilst excluding those filters that do not have the desired analysis ensemble statistics.

This paper adopts a formal style of presentation with explicitly stated definitions and theorems. The purpose of the definitions is to clarify key concepts, especially the distinctions between the different types of filter. The theorems are stated explicitly to help distinguish them from the rest of the text. Section 2 introduces some notation and defines what is meant by a deterministic analysis step for an EnKF. The definition is formulated in a way that makes sense for nonlinear observation operators as well as linear ones; care is taken to maintain this generality throughout the paper. Section 3 defines an ensemble SRF and shows that every EnKF with deterministic analysis step is an ensemble SRF. The bias issue is discussed in section 4, which shows how biased filters can arise within the ensemble SRF framework. The notion of an unbiased ensemble SRF is introduced in section 5, where it is shown that every EnKF with unbiased analysis step is an unbiased ensemble SRF and conversely. Section 6 presents conditions for an ensemble SRF to be unbiased. These conditions should be useful to workers devising new filters within the ensemble SRF framework. Filters within the ensemble SRF framework create a matrix of perturbations from the analysis ensemble mean by post-multiplying the corresponding matrix for the forecast ensemble. An alternative approach (used, for example, in the EAKF) is to pre-multiply the matrix of perturbations. This approach is discussed in section 7 and its relation to the unbiased ensemble SRF framework shown. Section 8 examines various published filter algorithms for bias, and shows that whilst many are unbiased some are not. Some concluding remarks are made in section 9.

## 2    Semi-deterministic EnKFs

This section introduces some notation and defines what is meant by a deterministic analysis step for an EnKF. An EnKF with a deterministic analysis step will be called *semi-deterministic* in this paper—'semi' because such a filter may still have stochastic elements in its forecast step.

Let $\{\mathbf{x}_i\}$ $(i = 1, \ldots, m)$ be an $m$-member ensemble in an $n$-dimensional state

space. The ensemble mean is the vector defined by

$$\overline{\mathbf{x}} = \frac{1}{m} \sum_{i=1}^{m} \mathbf{x}_i. \tag{1}$$

The ensemble perturbation matrix is the $n \times m$ matrix defined by

$$\mathbf{X} = \frac{1}{\sqrt{m-1}} \left( \begin{array}{cccc} \mathbf{x}_1 - \overline{\mathbf{x}} & \mathbf{x}_2 - \overline{\mathbf{x}} & \ldots & \mathbf{x}_m - \overline{\mathbf{x}} \end{array} \right). \tag{2}$$

The ensemble covariance matrix is the $n \times n$ matrix defined by

$$\mathbf{P} = \mathbf{X}\mathbf{X}^T = \frac{1}{m-1} \sum_{i=1}^{m} (\mathbf{x}_i - \overline{\mathbf{x}})(\mathbf{x}_i - \overline{\mathbf{x}})^T. \tag{3}$$

If the members of $\{\mathbf{x}_i\}$ are drawn independently from the same probability distribution, then $\overline{\mathbf{x}}$ is an unbiased estimate of the population mean and $\mathbf{P}$ is an unbiased estimate of the population covariance matrix. The first equality in (3) may be expressed by saying that $\mathbf{X}$ is a matrix square root of $\mathbf{P}$. Note that this use of the term 'square root' is inconsistent with its most common use in the mathematical literature, where a square root of the matrix $\mathbf{P}$ is often defined to be a matrix $\mathbf{X}$ such that $\mathbf{P} = \mathbf{X}^2$ (see, for example, Golub and Van Loan (1996, section 4.2.10)). However, the usage is well-established in the engineering literature (as in Andrews (1968) and Gelb (1974, section 8.4)) and has more recently been taken over into geophysical data assimilation (as in Tippett et al. (2003)).

Let $\mathbf{y}$ be an observation of the system in a $p$-dimensional observation space, let $H$ be the observation operator (possibly nonlinear) mapping state space to observation space, and let $\mathbf{R}$ be the $p \times p$ observation error covariance matrix. Let forecast quantities be denoted by the superscript $f$ and analysis quantities by the superscript $a$. There is no need for a notation to distinguish quantities at different times because this paper is concerned purely with the analysis step, which happens all at one time.

The analysis step of an ensemble filter updates the ensemble $\{\mathbf{x}_i^f\}$ resulting from the previous forecast to give a new ensemble $\{\mathbf{x}_i^a\}$ (the analysis ensemble) that will be used as the starting point for the next forecast. This update of the ensemble implies an update of the ensemble mean and the ensemble covariance matrix. In an EnKF the analysis step is designed so that these implied updates mimic the update of the state vector and covariance matrix in the KF. If the analysis algorithm includes stochastic elements, the resemblance to the KF may only be in some sort of mean

sense. This paper is confined to filters in which the resemblance is exact (in a sense to be made concrete in definition 1 below). In the case of a linear observation operator represented by the $p \times n$ matrix $\mathbf{H}$, the updates of the ensemble mean and ensemble covariance matrix are required to satisfy

$$\overline{\mathbf{x}^a} \;=\; \overline{\mathbf{x}^f} + \mathbf{K}(\mathbf{y} - \mathbf{H}\overline{\mathbf{x}^f}), \tag{4}$$

$$\mathbf{P}^a \;=\; (\mathbf{I} - \mathbf{KH})\mathbf{P}^f, \tag{5}$$

where $\mathbf{K}$ is the Kalman gain matrix defined by

$$\mathbf{K} = \mathbf{P}^f\mathbf{H}^T(\mathbf{H}\mathbf{P}^f\mathbf{H}^T + \mathbf{R})^{-1}. \tag{6}$$

These are just the update equations of the standard KF with the state vector replaced by the ensemble mean and the covariance matrix by the ensemble covariance matrix. The matrix inverted in the definition of the Kalman gain is of such frequent occurrence in what follows that it is convenient to introduce a special notation for it:

$$\mathbf{D} = \mathbf{H}\mathbf{P}^f\mathbf{H}^T + \mathbf{R}. \tag{7}$$

The above equations may be generalised to the case of a nonlinear observation operator by rewriting them in terms of a forecast observation ensemble $\{\mathbf{y}_i^f\}$ defined by

$$\mathbf{y}_i^f = H(\mathbf{x}_i^f). \tag{8}$$

The vector $\mathbf{y}_i^f$ is what the observation would be if the true state of the system were $\mathbf{x}_i^f$ and there were no observation noise. Like any other ensemble, the forecast observation ensemble has an ensemble mean $\overline{\mathbf{y}^f}$ and an ensemble perturbation matrix $\mathbf{Y}^f$. In the case of a linear observation operator, $\overline{\mathbf{y}^f} = \mathbf{H}\overline{\mathbf{x}^f}$ and $\mathbf{Y}^f = \mathbf{H}\mathbf{X}^f$. Using these relationships and the relationship $\mathbf{P}^f = \mathbf{X}^f(\mathbf{X}^f)^T$, equations (4)–(7) may be rewritten in a form that does not mention the linear operator $\mathbf{H}$. These rewritten equations form the basis of the following definition, which applies regardless of whether the observation operator is linear or nonlinear.

**Definition 1** *The analysis step of an EnKF is **deterministic** if the updates of the ensemble mean and ensemble covariance matrix exactly satisfy*

$$\overline{\mathbf{x}^a} \;=\; \overline{\mathbf{x}^f} + \mathbf{K}(\mathbf{y} - \overline{\mathbf{y}^f}), \tag{9}$$

$$P^a = (\mathbf{X}^f - \mathbf{K}\mathbf{Y}^f)(\mathbf{X}^f)^T, \tag{10}$$

$$\mathbf{K} = \mathbf{X}^f(\mathbf{Y}^f)^T\mathbf{D}^{-1}, \tag{11}$$

$$\mathbf{D} = \mathbf{Y}^f(\mathbf{Y}^f)^T + \mathbf{R}. \tag{12}$$

*An EnKF is* **semi-deterministic** *if its analysis step is deterministic.*

There is an alternative approach to extending an EnKF from linear to nonlinear observation operators, described in, for example, Evensen (2003, section 4.5). In this approach the state vector is augmented with a diagnostic variable that is the predicted observation vector:

$$\hat{\mathbf{x}} = \begin{pmatrix} \mathbf{x} \\ H(\mathbf{x}) \end{pmatrix} \tag{13}$$

and a linear observation operator is defined on augmented state space by

$$\widehat{\mathbf{H}} \begin{pmatrix} \mathbf{x} \\ \mathbf{y} \end{pmatrix} = \mathbf{y}. \tag{14}$$

The analysis step is then carried out in augmented state space using $\hat{\mathbf{x}}$ and $\widehat{\mathbf{H}}$ in place of $\mathbf{x}$ and $\mathbf{H}$. However, any filter that satisfies the analogues of (4)–(7) in augmented state space also satisfies the equations of definition 1 in unaugmented state space. The unaugmented state space approach will be used wherever possible in this paper. The only place where it cannot be used is in the formulation of some of the specific semi-deterministic EnKF algorithms in section 8, certain of which require a linear observation operator and thus an augmented state space must be used when applying them to nonlinear observation operators.

# 3   Ensemble square root filters

This section defines the notion of an ensemble SRF as introduced by Tippett et al. (2003). It goes beyond those authors by explicitly proving that every semi-deterministic EnKF is an ensemble SRF. It will be seen in section 4 that the converse to this statement is not generally true.

**Definition 2** *An **ensemble square root filter** is an ensemble filter in which the analysis ensemble is obtained by adding a column n-vector $\widetilde{\mathbf{x}}$ to the columns of a $n \times m$ matrix $\sqrt{m-1}\widetilde{\mathbf{X}}$ where $\widetilde{\mathbf{x}}$ and $\widetilde{\mathbf{X}}$ satisfy*

$$\widetilde{\mathbf{x}} = \overline{\mathbf{x}^f} + \mathbf{K}(\mathbf{y} - \overline{\mathbf{y}^f}), \tag{15}$$

$$\widetilde{\mathbf{X}} = \mathbf{X}^f \mathbf{T}, \tag{16}$$

$$\mathbf{T}\mathbf{T}^T = \mathbf{I} - (\mathbf{Y}^f)^T \mathbf{D}^{-1} \mathbf{Y}^f. \tag{17}$$

There is an obvious resemblance between (15) and the update (9) of the ensemble mean in a semi-deterministic EnKF. However, in an ensemble SRF the vector $\widetilde{\mathbf{x}}$ need not be the mean of the analysis ensemble. More will be said about this in section 4. A further point of resemblance between the semi-deterministic EnKF and the ensemble SRF is given by the following theorem, which should be compared to (10). It is the reason for the matrix square root condition (17) in the definition of an ensemble SRF.

**Theorem 1** *In an ensemble SRF, $\widetilde{\mathbf{X}}\widetilde{\mathbf{X}}^T = (\mathbf{X}^f - \mathbf{K}\mathbf{Y}^f)(\mathbf{X}^f)^T$.*

**Proof.**

$$\begin{aligned}
\widetilde{\mathbf{X}}\widetilde{\mathbf{X}}^f &= \mathbf{X}^f \mathbf{T}\mathbf{T}^T (\mathbf{X}^f)^T \\
&= \mathbf{X}^f (\mathbf{I} - (\mathbf{Y}^f)^T \mathbf{D}^{-1} \mathbf{Y}^f)(\mathbf{X}^f)^T \\
&= (\mathbf{X}^f - \mathbf{X}^f (\mathbf{Y}^f)^T \mathbf{D}^{-1} \mathbf{Y}^f)(\mathbf{X}^f)^T \\
&= (\mathbf{X}^f - \mathbf{K}\mathbf{Y}^f)(\mathbf{X}^f)^T. \tag{18}
\end{aligned}$$

$\square$

If $\widetilde{\mathbf{X}}$ were equal to $\mathbf{X}^a$, then theorem 1 would imply that the ensemble covariance matrix in the ensemble SRF updates as in the semi-deterministic EnKF. However, as $\widetilde{\mathbf{x}}$ need not equal $\overline{\mathbf{x}^a}$, so $\widetilde{\mathbf{X}}$ need not equal $\mathbf{X}^a$. Again, more will be said about this in section 4.

The following two theorems are necessary preliminaries for the main result of this section (theorem 7). The first is a simple consequence of the matrix square root condition (17) and is already present in Tippett et al. (2003).

**Theorem 2** *If* **T** *satisfies the matrix square root condition (17) and* **U** *is an* $m \times m$ *orthogonal matrix, then* **TU** *also satisfies (17).*

**Proof.** Recall that an orthogonal matrix satisfies $\mathbf{U}^T\mathbf{U} = \mathbf{U}\mathbf{U}^T = \mathbf{I}$. Thus $(\mathbf{T}\mathbf{U})(\mathbf{T}\mathbf{U})^T = \mathbf{T}\mathbf{U}\mathbf{U}^T\mathbf{T}^T = \mathbf{T}\mathbf{T}^T$ and so **TU** satisfies (17) if **T** does. $\square$

**Theorem 3** *If* $\mathbf{X}_1$ *and* $\mathbf{X}_2$ *are two* $n \times m$ *matrices such that* $\mathbf{X}_1\mathbf{X}_1^T = \mathbf{X}_2\mathbf{X}_2^T$, *then there exists an orthogonal matrix* **U** *such that* $\mathbf{X}_2 = \mathbf{X}_1\mathbf{U}$.

**Proof.** This proof makes use of the singular value decomposition (SVD) of a matrix; see, for example, Golub and Van Loan (1996, section 2.5). Start by taking the SVD of $\mathbf{X}_1$:

$$\mathbf{X}_1 = \mathbf{F}\mathbf{G}\mathbf{W}^T \tag{19}$$

where **G** is an $n \times m$ diagonal matrix (in the sense that $g_{ij} \neq 0$ if $i \neq j$) and **F** and **W** are orthogonal matrices of sizes $n \times n$ and $m \times m$ respectively. Without loss of generality it may be assumed that **G** can be expressed in the form

$$\mathbf{G} = \begin{pmatrix} \mathbf{G}_0 & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix} \tag{20}$$

where $\mathbf{G}_0$ is a nonsingular diagonal matrix of size $r \times r$ for some $r$. The orthogonal matrix **F** can then be expressed in the form

$$\mathbf{F} = \begin{pmatrix} \mathbf{F}_0 & \mathbf{F}_1 \end{pmatrix} \tag{21}$$

where $\mathbf{F}_0$ and $\mathbf{F}_1$ are column-orthogonal matrices of sizes $n \times r$ and $n \times (n-r)$ respectively (a column-orthogonal matrix is one in which the column vectors are orthonormal and thus the matrix satisfies $\mathbf{F}_0^T\mathbf{F}_0 = \mathbf{I}$). Similarly, **W** can be expressed in the form

$$\mathbf{W} = \begin{pmatrix} \mathbf{W}_0 & \mathbf{W}_1 \end{pmatrix} \tag{22}$$

where $\mathbf{W}_0$ and $\mathbf{W}_1$ are column-orthogonal matrices of sizes $m \times r$ and $m \times (m-r)$ respectively. It may be verified by substitution in (19) that

$$\mathbf{X}_1 = \mathbf{F}_0\mathbf{G}_0\mathbf{W}_0^T. \tag{23}$$

Let $\mathrm{ran}(\mathbf{X}_1)$ denote the range of $\mathbf{X}_1$ (the space of all $n$-vectors of the form $\mathbf{X}_1\mathbf{u}$ where $\mathbf{u}$ is an arbitrary $m$-vector). It follows from (23) that $\mathrm{ran}(\mathbf{X}_1) \subseteq \mathrm{ran}(\mathbf{F}_0\mathbf{G}_0)$. Furthermore, since any $r$-vector $\mathbf{v}$ can be written in the form $\mathbf{v} = \mathbf{W}_0^T\mathbf{u}$ by setting $\mathbf{u} = \mathbf{W}_0\mathbf{v}$, it follows that $\mathrm{ran}(\mathbf{X}_1) = \mathrm{ran}(\mathbf{F}_0\mathbf{G}_0)$.

Let $\mathbf{P} = \mathbf{X}_1\mathbf{X}_1^T$. Then $\mathbf{P} = \mathbf{F}_0\mathbf{G}_0\mathbf{W}_0^T\mathbf{W}_0\mathbf{G}\mathbf{F}_0^T = \mathbf{F}_0\mathbf{G}_0^2\mathbf{F}_0^T$. Thus $\mathrm{ran}(\mathbf{P}) \subseteq \mathrm{ran}(\mathbf{F}_0\mathbf{G}_0)$. Furthermore, since any $r$-vector $\mathbf{v}$ can be written in the form $\mathbf{v} = \mathbf{G}_0\mathbf{F}_0^T\mathbf{u}$ by setting $\mathbf{u} = \mathbf{F}_0\mathbf{G}_0^{-1}\mathbf{v}$, it follows that

$$\mathrm{ran}(\mathbf{P}) = \mathrm{ran}(\mathbf{F}_0\mathbf{G}_0) = \mathrm{ran}(\mathbf{X}_1). \tag{24}$$

By the hypothesis of the theorem $\mathbf{P} = \mathbf{X}_2\mathbf{X}_2^T$ as well. Therefore

$$\mathrm{ran}(\mathbf{X}_2) = \mathrm{ran}(\mathbf{P}) = \mathrm{ran}(\mathbf{F}_0\mathbf{G}_0). \tag{25}$$

It follows that every column of $\mathbf{X}_2$ can be expressed as a linear combination of the columns of $\mathbf{F}_0\mathbf{G}_0$. Thus there exists an $m \times r$ matrix $\widetilde{\mathbf{W}}_0$ such that

$$\mathbf{X}_2 = \mathbf{F}_0\mathbf{G}_0\widetilde{\mathbf{W}}_0^T. \tag{26}$$

Now

$$\mathbf{F}_0\mathbf{G}_0^2\mathbf{F}_0^T = \mathbf{P} = \mathbf{X}_2\mathbf{X}_2^T = \mathbf{F}_0\mathbf{G}_0\widetilde{\mathbf{W}}_0^T\widetilde{\mathbf{W}}_0\mathbf{G}_0\mathbf{F}_0^T. \tag{27}$$

Pre-multiplying the first and last terms of this chain of equations by $\mathbf{G}_0^{-1}\mathbf{F}_0^T$ and post-multiplying by $\mathbf{F}_0\mathbf{G}^{-1}$ gives $\mathbf{I} = \widetilde{\mathbf{W}}_0^T\widetilde{\mathbf{W}}_0$. Thus $\widetilde{\mathbf{W}}_0$ is a column-orthogonal matrix and may be extended to a full $m \times m$ orthogonal matrix

$$\widetilde{\mathbf{W}} = \left( \begin{array}{cc} \widetilde{\mathbf{W}}_0 & \widetilde{\mathbf{W}}_1 \end{array} \right). \tag{28}$$

It may be verified by substitution that

$$\mathbf{X}_2 = \mathbf{F}\mathbf{G}\widetilde{\mathbf{W}}^T \tag{29}$$

and thus

$$\mathbf{X}_2 = \mathbf{X}_1\mathbf{W}\widetilde{\mathbf{W}}^T. \tag{30}$$

Here $\mathbf{W}$ and $\widetilde{\mathbf{W}}$ are orthogonal matrices, and so therefore is $\mathbf{W}\widetilde{\mathbf{W}}^T$. Thus the theorem is proven by setting $\mathbf{U} = \mathbf{W}\widetilde{\mathbf{W}}^T$. $\square$

The following three theorems establish some useful facts about the matrix square root condition (17) that will be used later. In particular, theorem 6 establishes that there is always at least one solution $\mathbf{T}$ of (17), which fact will be used in the proof of the main result of this section (theorem 7).

**Theorem 4** *The RHS of (17) is symmetric positive definite.*

**Proof.** The RHS of (17) may be written

$$\mathbf{I} - (\mathbf{Y}^f)^T \mathbf{D}^{-1} \mathbf{Y}^f = (\mathbf{I} + (\mathbf{Y}^f)^T \mathbf{R}^{-1} \mathbf{Y}^f)^{-1}. \tag{31}$$

(This follows from Tippett et al. (2003, equation (15)) or may be verified by direct multiplication.) Since $\mathbf{I} + (\mathbf{Y}^f)^T \mathbf{R}^{-1} \mathbf{Y}^f$ is symmetric positive definite, so is $\mathbf{I} - (\mathbf{Y}^f)^T \mathbf{D}^{-1} \mathbf{Y}^f$. □

**Theorem 5** *Any solution of (17) is nonsingular.*

**Proof.** Suppose that $\mathbf{T}$ is singular. Then $\mathbf{T}\mathbf{T}^T$ is singular. But by theorem 4, $\mathbf{T}\mathbf{T}^T$ is symmetric positive definite and therefore nonsingular. This contradiction implies that $\mathbf{T}$ is nonsingular. □

**Theorem 6** *There is a unique symmetric positive definite solution to (17).*

**Proof.** This follows from theorem 4 and the theorem in linear algebra that a symmetric positive definite matrix has a unique symmetric positive definite square root (see, for example, Halmos (1974, section 82)). □

The following theorem is the main result of this section and shows the importance of the ensemble SRF framework to the EnKF.

**Theorem 7** *Every semi-deterministic EnKF is an ensemble SRF.*

**Proof.** It suffices to prove that $\mathbf{X}^a$ in a semi-deterministic EnKF can be expressed in the form $\mathbf{X}^a = \mathbf{X}^f \mathbf{T}$ where $\mathbf{T}$ satisfies (17). Let $\mathbf{T}_0$ be a solution of (17) (such as the positive definite solution that exists by theorem 6) and let $\widetilde{\mathbf{X}} = \mathbf{X}^f \mathbf{T}_0$. Then by theorem 1, $\widetilde{\mathbf{X}}\widetilde{\mathbf{X}}^T = \mathbf{X}^a (\mathbf{X}^a)^T$. Therefore by theorem 3 there exists an orthogonal

matrix $\mathbf{U}$ such that $\mathbf{X}^a = \widetilde{\mathbf{X}}\mathbf{U}$. Let $\mathbf{T} = \mathbf{T}_0\mathbf{U}$. Then $\mathbf{X}^a = \mathbf{X}^f\mathbf{T}$ and $\mathbf{T}$ satisfies (17) by theorem 2. $\square$

The results of this section make it possible to characterise the structure of the set of all ensemble SRF filters in terms of a well-known group of matrices; see appendix A for details.

# 4   Bias

A fact that appears to have been overlooked by Tippett et al. (2003) is that the ensemble SRF framework encompasses filters that are not EnKFs. To see this, suppose that an arbitrary ensemble SRF is a semi-deterministic EnKF. Then it follows from (9) and (15) that $\widetilde{\mathbf{x}}$ equals the analysis ensemble mean $\overline{\mathbf{x}^a}$ and that $\widetilde{\mathbf{X}}$ equals the analysis ensemble perturbation matrix $\mathbf{X}^a$. However, (2) implies that the sum of the columns of an ensemble perturbation matrix must be zero, and this does not necessarily follow from (16) and (17), which are the only constraints on $\widetilde{\mathbf{X}}$. To see this, let $\mathbf{T}$ be a particular solution of (17). Then by theorem 2 a more general solution is $\mathbf{T}\mathbf{U}$ where $\mathbf{U}$ is an arbitrary $m \times m$ orthogonal matrix. The corresponding general value of $\widetilde{\mathbf{X}}$ is

$$\widetilde{\mathbf{X}} = \mathbf{X}^f\mathbf{T}\mathbf{U}. \tag{32}$$

Now let $\mathbf{1}$ be a column $m$-vector in which every element is 1; that is

$$\mathbf{1} = \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix}. \tag{33}$$

The sum of the columns of $\widetilde{\mathbf{X}}$ is

$$\widetilde{\mathbf{X}}\mathbf{1} = \mathbf{X}^f\mathbf{T}\mathbf{U}\mathbf{1}. \tag{34}$$

Thus $\widetilde{\mathbf{X}}$ is a valid ensemble perturbation matrix if and only if $\mathbf{U}\mathbf{1}$ lies in the null space[1] of $\mathbf{X}^f\mathbf{T}$. The vector $\mathbf{U}\mathbf{1}$ is nonzero and can be made to point in any direction by an appropriate choice of $\mathbf{U}$. Therefore, unless $\mathbf{X}^f\mathbf{T} = \mathbf{0}$ (in which case the analysis

---

[1]The null space of an $n \times m$ matrix $\mathbf{M}$ is the set of all column $m$-vectors $\mathbf{u}$ such that $\mathbf{M}\mathbf{u} = \mathbf{0}$.

ensemble collapses to a point), there will be at least some choices of $\mathbf{U}$ that give $\widetilde{\mathbf{X}}\mathbf{1} \neq \mathbf{0}$ and hence an $\widetilde{\mathbf{X}}$ that is invalid as an ensemble perturbation matrix. In these cases the ensemble SRF cannot be a valid semi-deterministic EnKF.

To see the effect of treating an ensemble SRF that is not a valid semi-deterministic EnKF as though it were, let $\mathbf{x}_i'$ denote the $i$th column of $\sqrt{m-1}\widetilde{\mathbf{X}}$, and let $\overline{\mathbf{x}}^a$ and $\mathbf{P}^a$ denote the analysis ensemble mean and covariance matrix that result from a valid semi-deterministic EnKF. The members of the analysis ensemble from the ensemble SRF are

$$\mathbf{x}_i = \widetilde{\mathbf{x}} + \mathbf{x}_i' = \overline{\mathbf{x}}^a + \mathbf{x}_i'. \tag{35}$$

The mean of this ensemble is

$$\overline{\mathbf{x}} = \overline{\mathbf{x}}^a + \overline{\mathbf{x}'}. \tag{36}$$

But $\overline{\mathbf{x}'} = (\sqrt{m-1}/m)\widetilde{\mathbf{X}}\mathbf{1} \neq \mathbf{0}$ (because $\widetilde{\mathbf{X}}$ is invalid as an ensemble perturbation matrix), and so there is a bias in the ensemble mean. Furthermore, the ensemble covariance matrix of the ensemble $\{\mathbf{x}_i\}$ is

$$\mathbf{P} = \frac{1}{m-1}\sum_{i=1}^{m}(\mathbf{x}_i - \overline{\mathbf{x}})(\mathbf{x}_i - \overline{\mathbf{x}})^T = \mathbf{P}^a - \frac{m}{m-1}\overline{\mathbf{x}'}\,\overline{\mathbf{x}'}^T. \tag{37}$$

Therefore $\mathbf{P} \neq \mathbf{P}^a$ and in particular the ensemble standard deviation will be too small for any coordinate in which there is also a bias in the mean. Thus the analysis is not only biased but overconfident as well. This can lead to several problems that are discussed in section 9.

An example of output from a biased filter is shown in Fig. 1. The dynamical system is the two-dimensional swinging spring with nonlinear normal mode initialisation described in Lynch (2003) and summarised in appendix B. The coordinates are polar coordinates $r$, $\theta$ and the corresponding generalised momenta $p_r$, $p_\theta$. The filter is an ETKF (see Bishop et al. (2001) or section 8.1) with ensemble size $m = 10$. The model used in the forecast step is the same as the model used to generate the true trajectory. The model noise is taken to be zero. All coordinates are observed. The first observation is at time 0.1 and subsequent ones follow at intervals of 0.1. Although the actual observation errors are zero, the covariance matrix $\mathbf{R}$ passed to the filter is that of observations having uncorrelated errors with standard deviations in $\theta$, $p_\theta$, $r$, and $p_r$ of 0.1, 0.3, $7 \times 10^{-4}$, and $5 \times 10^{-3}$ respectively. These standard
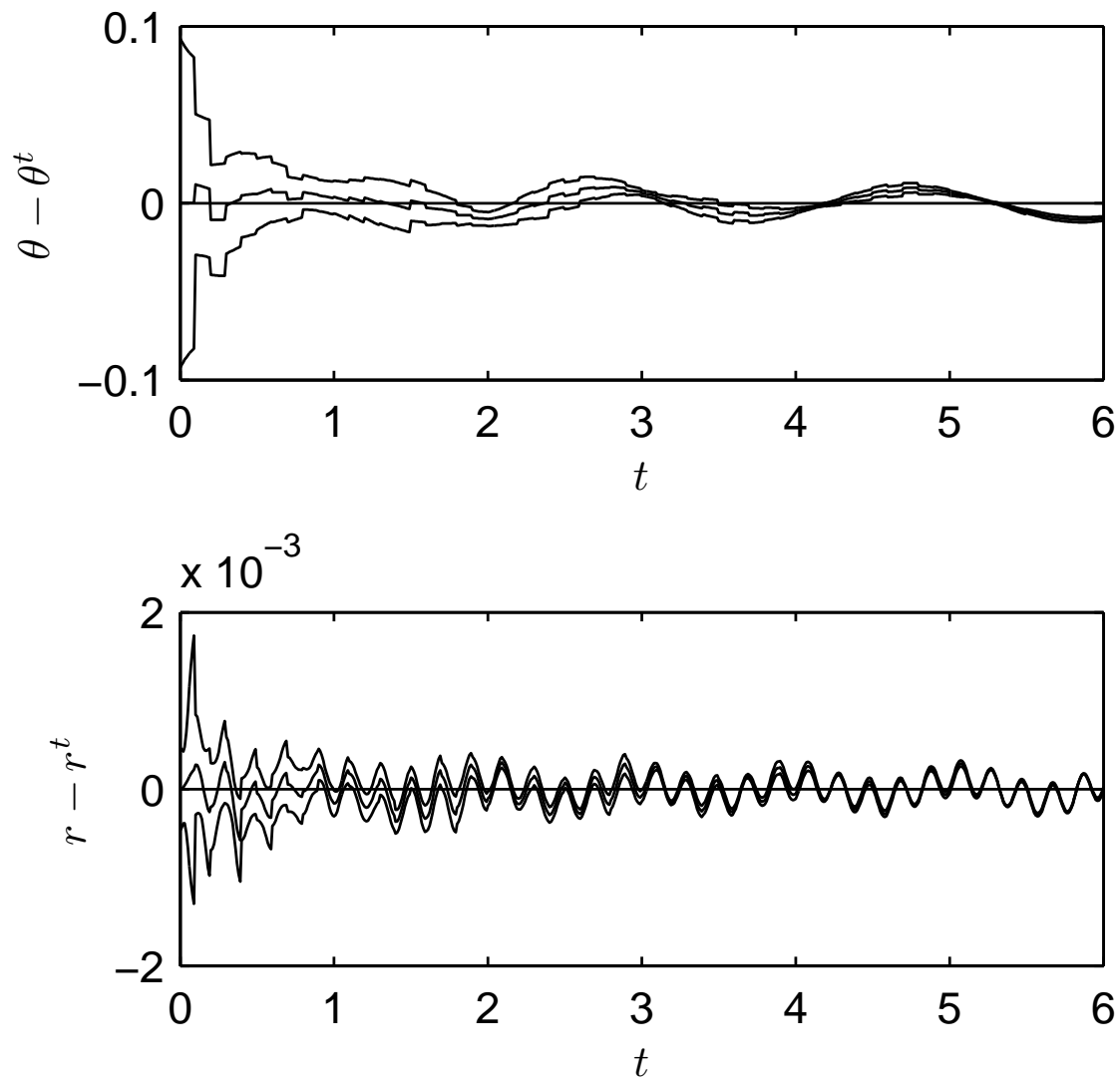
Figure 1: Output of ETKF with ensemble size $m = 10$ for swinging spring system with perfect observations. Coordinates are plotted relative to truth. Three lines show ensemble mean and ensemble mean $\pm$ ensemble standard deviation. Filter has a time-varying bias and is overconfident.

deviations are close to one-tenth of the amplitude of the oscillations in the truth. The same covariance matrix is used in generating a random initial ensemble, centred on the true initial state.

Fig. 1 shows the difference between the filter and the truth for $\theta$ and $r$. There are considerable intervals of time during which the true state of the system (represented by zero on the vertical axis) is outside the band defined by the ensemble mean $\pm$ ensemble standard deviation. This suggests that the ensemble statistics may be inconsistent with the actual error. This is confirmed by computing the fraction of analyses having an ensemble mean within one ensemble standard deviation of the truth for each coordinate. For unbiased, normally-distributed, analysis errors with standard deviation equal to the ensemble standard deviation, one would expect this fraction to be about 0.68. The actual errors need not be normally-distributed, but this is still a useful guide. In the case shown in Fig. 1 the actual fractions are 0.43 and 0.23 for $\theta$ and $r$ respectively. Further confirmation is provided by running the filter 100 times with different random initial ensembles and computing the same fractions. The results for $\theta$, $p_\theta$, $r$, and $p_r$ are 0.31, 0.32, 0.32, and 0.29 respectively.

# 5 Unbiased ensemble SRFs

Section 4 points to the need to supplement the ensemble SRF framework with an additional condition to rule out filters that do not have the desired KF-like analysis ensemble statistics. This leads to the following definition.

**Definition 3** *An **unbiased** ensemble SRF is an ensemble SRF in which* $\widetilde{\mathbf{X}}\mathbf{1} = \mathbf{0}$.

The following theorem shows that an unbiased ensemble SRF is the same thing as a semi-deterministic EnKF.

**Theorem 8** *Every semi-deterministic EnKF is an unbiased ensemble SRF, and conversely every unbiased ensemble SRF is a semi-deterministic EnKF.*

**Proof.** Suppose that a filter is a semi-deterministic EnKF. By theorem 7 the filter is an ensemble SRF as well. In this ensemble SRF, $\widetilde{\mathbf{X}} = \mathbf{X}^a$ where $\mathbf{X}^a$ is an ensemble

perturbation matrix and thus has columns that sum to zero. Therefore $\widetilde{\mathbf{X}}\mathbf{1} = \mathbf{0}$ and the ensemble SRF is unbiased.

Suppose conversely that a filter is an unbiased ensemble SRF. Unbiasedness implies that $\widetilde{\mathbf{x}}$ is the mean of the analysis ensemble and $\widetilde{\mathbf{X}}$ is the analysis ensemble perturbation matrix. Equation (15) implies that $\overline{\mathbf{x}^a}$ satisfies (9) in the definition of a semi-deterministic EnKF, whilst theorem 1 implies that $\mathbf{P}^a$ satisfies (10) in the same definition. Therefore the filter in a semi-deterministic EnKF. $\square$

# 6 Conditions for an unbiased ensemble SRF

Definitions 2 and 3 in conjunction with theorem 8 reduce the problem of constructing a semi-deterministic EnKF to finding a solution $\mathbf{T}$ of the matrix square root condition (17) and checking that the matrix $\widetilde{\mathbf{X}}$ defined by (16) satisfies $\widetilde{\mathbf{X}}\mathbf{1} = \mathbf{0}$. It would be useful to replace this condition on $\widetilde{\mathbf{X}}$ with one on $\mathbf{T}$, so that the problem of finding a semi-deterministic EnKF is reduced to one of finding $\mathbf{T}$ satisfying certain conditions. Such conditions for $\mathbf{T}$ are provided in this section. It is assumed throughout that $\mathbf{T}$ satisfies the matrix square root condition (17). The first theorem gives an additional sufficient condition for the resulting ensemble SRF to be unbiased.

**Theorem 9** *If* $\mathbf{1}$ *is an eigenvector of* $\mathbf{T}$*, then the ensemble SRF is unbiased.*

**Proof.** By hypothesis $\mathbf{T}\mathbf{1} = \lambda\mathbf{1}$ for some scalar $\lambda$. Therefore $\widetilde{\mathbf{X}}\mathbf{1} = \mathbf{X}^f\mathbf{T}\mathbf{1} = \lambda\mathbf{X}^f\mathbf{1} = \mathbf{0}$. Therefore the filter is unbiased. $\square$

An important special case is that of a symmetric $\mathbf{T}$. This is the subject of the following theorem and its corollary.

**Theorem 10** *If* $\mathbf{T}$ *is symmetric, then* $\mathbf{1}$ *is an eigenvector of* $\mathbf{T}$*.*

**Proof.** Since $\mathbf{Y}^f$ is an ensemble perturbation matrix, it satisfies $\mathbf{Y}^f\mathbf{1} = \mathbf{0}$. Therefore it follows from (17) that $\mathbf{T}^2\mathbf{1} = \mathbf{1}$ for symmetric $\mathbf{T}$. Thus $\mathbf{1}$ is an eigenvector of $\mathbf{T}^2$. But the eigenvectors of the square of a symmetric matrix are the same as those of the original matrix. Therefore $\mathbf{1}$ is an eigenvector of $\mathbf{T}$. $\square$

**Corollary 11** *If* $\mathbf{T}$ *is symmetric, then the ensemble SRF is unbiased.*

Although not as general as theorem 9, corollary 11 provides a particularly simple test for unbiasedness. However, it will be seen in section 8 that there exist unbiased filters for which $\mathbf{T}$ is not symmetric.

To state a partial converse to theorem 9 it is necessary to introduce the concept of a nondegenerate ensemble perturbation matrix. The columns of an ensemble perturbation matrix $\mathbf{X}$ are not linearly independent because there is at least one linear relation between them (they sum to zero). Thus the rank of $\mathbf{X}$ is at most $m - 1$.

**Definition 4** *An* $n \times m$ *ensemble perturbation matrix is* **nondegenerate** *if it has rank* $m - 1$.

Roughly speaking, $\mathbf{X}$ is nondegenerate if the ensemble perturbation vectors explore the state space to the maximum extent permitted by the ensemble size. Note that the nondegeneracy condition is equivalent to the null space of $\mathbf{X}$ being equal to the one-dimensional space spanned by $\mathbf{1}$. The following theorem states that the sufficient condition for unbiasedness in theorem 9 is also a necessary condition when $\mathbf{X}^f$ is nondegenerate.

**Theorem 12** *If* $\mathbf{X}^f$ *is nondegenerate and the ensemble SRF is unbiased, then* $\mathbf{1}$ *is an eigenvector of* $\mathbf{T}$.

**Proof.** Because the filter is unbiased, $\mathbf{0} = \widetilde{\mathbf{X}}\mathbf{1} = \mathbf{X}^f\mathbf{T}\mathbf{1}$. Therefore $\mathbf{T}\mathbf{1}$ is in the null space of $\mathbf{X}^f$. Because $\mathbf{X}^f$ is nondegenerate, this implies $\mathbf{T}\mathbf{1} = \lambda\mathbf{1}$ for some scalar $\lambda$. Therefore $\mathbf{1}$ is an eigenvector of $\mathbf{T}$. $\square$

The following theorems assume that the ensemble SRF with post-multiplier matrix $\mathbf{T}$ is unbiased and consider the ensemble SRF with post-multiplier matrix $\mathbf{T}\mathbf{U}$ where $\mathbf{U}$ is orthogonal.

**Theorem 13** *If the ensemble SRF with matrix* $\mathbf{T}$ *is unbiased and* $\mathbf{1}$ *is an eigenvector of* $\mathbf{U}$, *then the ensemble SRF with matrix* $\mathbf{T}\mathbf{U}$ *is unbiased.*

**Proof.** By hypothesis $\mathbf{X}^f\mathbf{T1} = \mathbf{0}$ and $\mathbf{U1} = \lambda\mathbf{1}$ for some scalar $\lambda$. Therefore $\mathbf{X}^f\mathbf{TU1} = \lambda\mathbf{X}^f\mathbf{T1} = \mathbf{0}$. Therefore the ensemble SRF with matrix $\mathbf{TU}$ is unbiased. $\square$

**Theorem 14** *If $\mathbf{X}^f$ is nondegenerate and the ensemble SRFs with matrices $\mathbf{T}$ and $\mathbf{TU}$ are unbiased, then $\mathbf{1}$ is an eigenvector of $\mathbf{U}$.*

**Proof.** Theorem 12 implies that $\mathbf{1}$ is an eigenvector of both $\mathbf{T}$ and $\mathbf{TU}$. Since $\mathbf{T}$ is invertible by corollary 5, $\mathbf{1}$ is also an eigenvector of $\mathbf{T}^{-1}\mathbf{TU} = \mathbf{U}$. $\square$

See appendix C for theorems on the structure of the set of all unbiased ensemble SRFs.

# 7  Pre-multiplier filters

Some semi-deterministic EnKFs (such as the EAKF of Anderson (2001)) have been presented in the pre-multiplier form $\mathbf{X}^a = \mathbf{AX}^f$ instead of the post-multiplier form $\mathbf{X}^a = \mathbf{X}^f\mathbf{T}$ used by the ensemble SRF framework. However, it follows from theorem 8 that such EnKFs can be written in post-multiplier form as well. The following two theorems show that the ability to write an ensemble SRF in pre-multiplier form provides an alternative to the tests for unbiasedness of section 6..

**Theorem 15** *If the analysis step of an ensemble SRF can be written in the form $\widetilde{\mathbf{X}} = \mathbf{AX}^f$, then the ensemble SRF is unbiased.*

**Proof.** $\widetilde{\mathbf{X}}\mathbf{1} = \mathbf{AX}^f\mathbf{1} = \mathbf{0}$. $\square$

**Theorem 16** *If $\mathbf{X}^f$ is nondegenerate and an ensemble SRF is unbiased, then the analysis step of the ensemble SRF can be written in the form $\widetilde{\mathbf{X}} = \mathbf{AX}^f$.*

**Proof.** Let $\mathbf{x}_i'^f$ and $\mathbf{x}_i'$ denote the $i$th columns of $\sqrt{m-1}\mathbf{X}^f$ and $\sqrt{m-1}\widetilde{\mathbf{X}}$ respectively. By the nondegeneracy of $\mathbf{X}^f$ there is a linearly independent set of $m-1$ members of $\{\mathbf{x}_i'^f\}$. Assume (without loss of generality) that these vectors correspond to the subscripts $i = 1, \ldots, m-1$. Because these vectors are linearly independent,

they can be mapped onto an arbitrary set of $m-1$ vectors by a linear transformation. In particular there exists a matrix $\mathbf{A}$ such that

$$\mathbf{x}'_i = \mathbf{A}\mathbf{x}'^f_i \quad \text{for } i = 1, \ldots, m-1. \tag{38}$$

Since $\mathbf{X}^f\mathbf{1} = \mathbf{0}$ and $\widetilde{\mathbf{X}}\mathbf{1} = \mathbf{0}$ it follows that

$$\mathbf{x}'^f_m = -\sum_{i=1}^{m-1} \mathbf{x}'^f_i, \tag{39}$$

$$\mathbf{x}'_m = -\sum_{i=1}^{m-1} \mathbf{x}'_i. \tag{40}$$

Therefore

$$\mathbf{A}\mathbf{x}'^f_m = -\sum_{i=1}^{m-1} \mathbf{A}\mathbf{x}'^f_i = -\sum_{i=1}^{m-1} \mathbf{x}'_i = \mathbf{x}'_m. \tag{41}$$

Equations (38) and (41) together imply $\widetilde{\mathbf{X}} = \mathbf{A}\mathbf{X}^f$. $\square$

In view of these theorems the reader may wonder whether a framework based on a pre-multiplier update of the ensemble perturbation matrix would be preferable to the ensemble SRF framework. Such a framework would be devoid of bias problems and theorem 16 implies that it would cover most semi-deterministic EnKFs of interest. However, an obstacle to formulating such as framework is the lack of a useful equivalent for $\mathbf{A}$ of the matrix square root condition (17) for $\mathbf{T}$. Also, in typical applications $n$ is much larger than $m$ and hence the $n \times n$ matrix $\mathbf{A}$ is very much larger than the $m \times m$ matrix $\mathbf{T}$.

# 8 Applications to specific filters

This section examines some published deterministic formulations of the analysis step of the EnKF from the point of view of bias. The filters discussed are the ensemble transform Kalman filter (ETKF) of Bishop et al. (2001), the ensemble adjustment Kalman filter (EAKF) of Anderson (2001), the filter of Whitaker and Hamill (2002), and the revised ETKF of Wang et al. (2004). The first three of these are placed in the ensemble SRF framework in Tippett et al. (2003), but without consideration of bias. This section, as well as containing a treatment of bias, also expands and extends the argument of that paper in a few places.

There is one further filter discussed in Tippett et al. (2003) that will not be considered in detail here. This is the *direct* method that consists of computing the RHS of (17) and then using some numerical method to find the matrix square root $\mathbf{T}$. The reason for not considering this method further is that its unbiasedness depends on the numerical method chosen, so that no general statement may be made. The filters discussed here all adopt a more indirect approach to finding $\mathbf{T}$, and indeed in many cases never find it at all in a practical implementation, its existence and relation to the filter merely being on a theoretical and analytical level.

## 8.1 The ensemble transform Kalman filter

The ensemble transform Kalman filter (ETKF) was originally introduced in Bishop et al. (2001), which describes its use to make rapid assessment of the future effect on error covariance of alternative strategies for deploying observational resources. The filter is placed within the ensemble SRF framework in Tippett et al. (2003). This original version of the ETKF starts by computing the $m \times m$ matrix $(\mathbf{Y}^f)^T \mathbf{R}^{-1} \mathbf{Y}^f$. It then computes the eigenvalue decomposition

$$(\mathbf{Y}^f)^T \mathbf{R}^{-1} \mathbf{Y}^f = \mathbf{C} \mathbf{\Gamma} \mathbf{C}^T \tag{42}$$

where $\mathbf{C}$ is orthogonal and $\mathbf{\Gamma}$ is diagonal. It follows from the identity (31) that

$$\mathbf{I} - (\mathbf{Y}^f)^T \mathbf{D}^{-1} \mathbf{Y}^f = \mathbf{C} (\mathbf{I} + \mathbf{\Gamma})^{-1} \mathbf{C}^T \tag{43}$$

and hence that a solution of (17) is

$$\mathbf{T} = \mathbf{C} (\mathbf{I} + \mathbf{\Gamma})^{-\frac{1}{2}}. \tag{44}$$

Using this $\mathbf{T}$ in (16) gives the original version of the ETKF. It is not obvious that such a $\mathbf{T}$ will always yield an unbiased filter. The experiment of Fig. 1 suggests that it may not and theorem 17 below shows that it rarely does. This issue of bias is addressed in the context of the ensemble generation problem in Wang et al. (2004), where two revised versions of the ETKF are proposed. The first such revision (Wang et al., 2004, appendix Ca) is equivalent to post-multiplying (44) by $\mathbf{C}^T$ to give

$$\mathbf{T} = \mathbf{C} (\mathbf{I} + \mathbf{\Gamma})^{-\frac{1}{2}} \mathbf{C}^T. \tag{45}$$

Since $\mathbf{C}^T$ is orthogonal, this is a solution of (17) by theorem 2, and since this $\mathbf{T}$ is symmetric, using it in (16) yields an unbiased filter by corollary 11.

The unbiasedness of this revised ETKF enables the conditions under which the original ETKF is unbiased to be clarified.

**Theorem 17** *If $\mathbf{X}^f$ is nondegenerate and the ETKF is unbiased, then $\mathbf{Y}^f = \mathbf{0}$.*

**Proof.** Since (44) may be obtained from (45) by post-multiplying by the orthogonal matrix $\mathbf{C}$, theorem 14 implies that $\mathbf{1}$ is an eigenvector of $\mathbf{C}$. Let $\mathbf{C1} = \lambda\mathbf{1}$. Then, by (42), $\mathbf{\Gamma 1} = \mathbf{C}^T(\mathbf{Y}^f)^T\mathbf{R}^{-1}\mathbf{Y}^f\mathbf{C1} = \lambda\mathbf{C}^T(\mathbf{Y}^f)^T\mathbf{R}^{-1}\mathbf{Y}^f\mathbf{1} = \mathbf{0}$. But $\mathbf{\Gamma 1}$ is the column vector of diagonal elements of the diagonal matrix $\mathbf{\Gamma}$. Thus $\mathbf{\Gamma} = \mathbf{0}$ and, by (42) again, $(\mathbf{Y}^f)^T\mathbf{R}^{-1}\mathbf{Y}^f = \mathbf{0}$. Since $\mathbf{R}^{-1}$ is positive definite, it follows that $\mathbf{Y}^f = \mathbf{0}$. $\square$

In words, the necessary condition $\mathbf{Y}^f = \mathbf{0}$ for unbiasedness in theorem 17 says that there is no observable difference between the members of the forecast ensemble. Thus an unbiased filter will be of rare occurrence for observation operators that supply useful information. When the necessary condition does occur, it follows from (42) that $\mathbf{\Gamma} = \mathbf{0}$ and $\mathbf{C}$ is arbitrary, and thus that the original ETKF reduces to $\mathbf{X}^a = \mathbf{X}^f\mathbf{C}$ where $\mathbf{C}$ is an arbitrary orthogonal matrix. Note that in this case $\mathbf{P}^a = \mathbf{P}^f$.

Repeating the experiment of Fig. 1 using the revised ETKF gives the results shown in Fig. 2. There is no sign of the bias present in the earlier experiment, the ensemble mean remaining within one ensemble standard deviation of the truth throughout. Running the filter 100 times with different random initial ensembles and computing the fraction of analyses with ensemble mean within one ensemble standard deviation of the truth gives 1.00 for each coordinate. The excess over the 0.68 expected for normally-distributed analysis errors is due to the perfect observations.

## 8.2 The ensemble adjustment Kalman filter

The ensemble adjustment Kalman filter (EAKF) was originally introduced in Anderson (2001). It is a pre-multiplier filter of the type discussed in section 7 and thus automatically unbiased. It also assumes a linear observation operator, so an
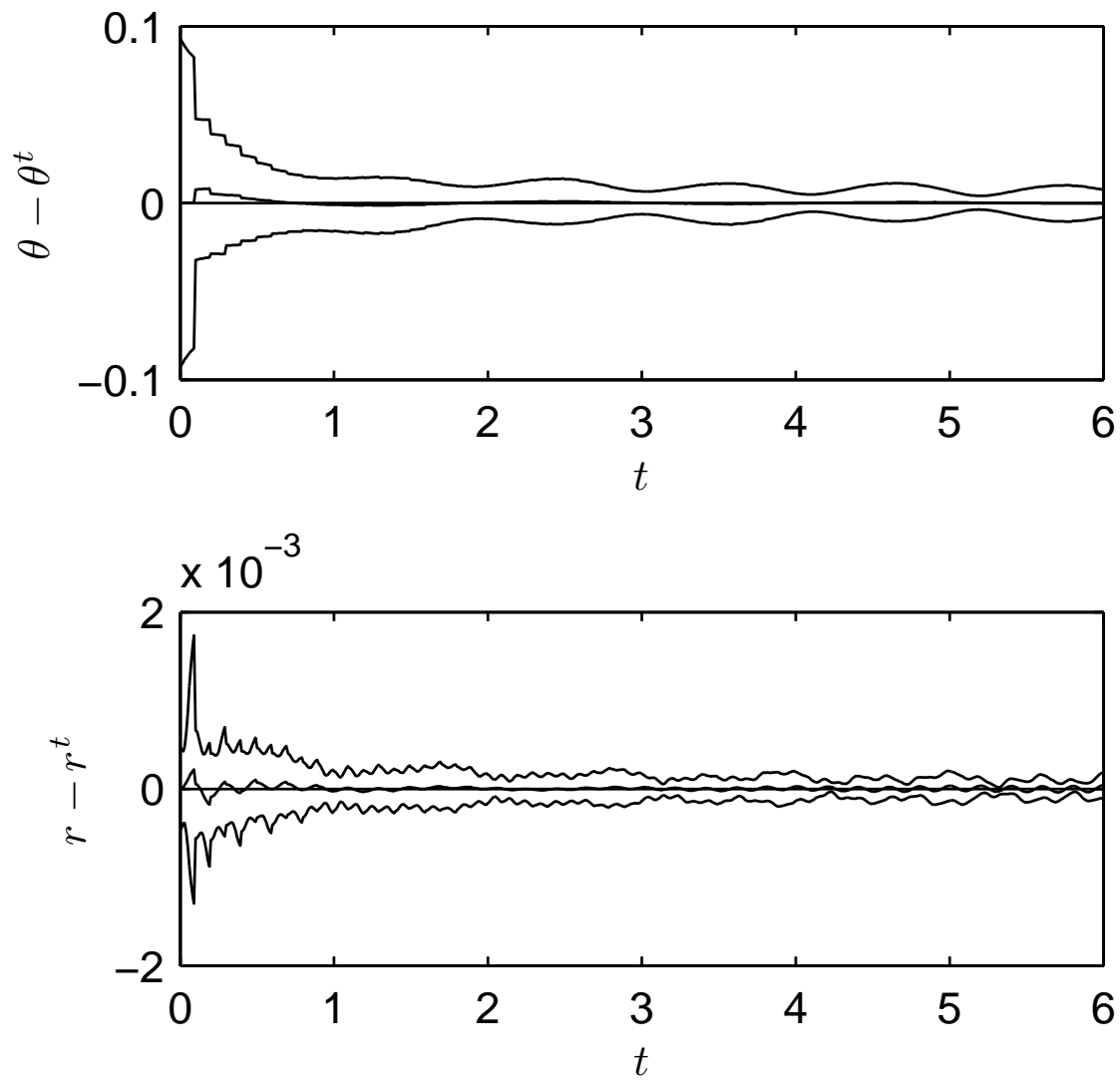
Figure 2: Output of revised ETKF with ensemble size $m = 10$ for swinging spring system with perfect observations. Plotting conventions as in Fig. 1. The bias evident in that figure is absent here.

augmented state space of the type discussed at the end of section 2 must be used to apply it to a nonlinear observation operator. Tippett et al. (2003, section 3a) outline how the EAKF may be expressed in post-multiplier form (16), but the demonstration glosses over a few details[2]. Therefore an alternative proof is presented here. This proof is based on the SVD of $\mathbf{X}^f$ instead of on the eigenvalue decomposition of $\mathbf{P}^f$ as in Tippett et al. (2003).

The first step is to construct the reduced SVD of $\mathbf{X}^f$ as in the proof of theorem 3:

$$\mathbf{X}^f = \mathbf{F}_0 \mathbf{G}_0 \mathbf{W}_0^T \tag{46}$$

where $\mathbf{G}_0$ is an $r \times r$ diagonal matrix of the nonzero singular values of $\mathbf{X}^f$ and $\mathbf{F}_0$ and $\mathbf{W}_0$ are column-orthogonal matrices of sizes $n \times r$ and $m \times r$ respectively. The next step is to compute the eigenvalue decomposition

$$(\mathbf{HF}_0\mathbf{G}_0)^T \mathbf{R}^{-1} \mathbf{HF}_0\mathbf{G}_0 = \widetilde{\mathbf{C}}_0 \widetilde{\mathbf{\Gamma}}_0 \widetilde{\mathbf{C}}_0^T \tag{47}$$

where $\widetilde{\mathbf{C}}_0$ is orthogonal, $\widetilde{\mathbf{\Gamma}}_0$ is diagonal, and both are $r \times r$. The pre-multiplier matrix for the EAKF is defined by

$$\mathbf{A} = \mathbf{F}_0 \mathbf{G}_0 \widetilde{\mathbf{C}}_0 (\mathbf{I} + \widetilde{\mathbf{\Gamma}}_0)^{-\frac{1}{2}} \mathbf{G}_0^{-1} \mathbf{F}_0^T, \tag{48}$$

which on substitution into $\mathbf{X}^a = \mathbf{A}\mathbf{X}^f$ gives

$$\mathbf{X}^a = \mathbf{F}_0 \mathbf{G}_0 \widetilde{\mathbf{C}}_0 (\mathbf{I} + \widetilde{\mathbf{\Gamma}}_0)^{-\frac{1}{2}} \mathbf{W}_0^T. \tag{49}$$

In typical applications where $n$ is much larger than $r$, it will be more efficient to directly calculate $\mathbf{X}^a$ in the form above rather than to calculate $\mathbf{A}$ and multiply by $\mathbf{X}^f$.

To reformulate the EAKF in post-multiplier form it is necessary to use the full SVD

$$\mathbf{X}^f = \mathbf{FGW}^T \tag{50}$$

where $\mathbf{G}$ is an $n \times m$ diagonal matrix and $\mathbf{F}$ and $\mathbf{W}$ are orthogonal matrices of sizes $n \times n$ and $m \times m$ respectively. As in the proof of theorem 3, $\mathbf{G}$, $\mathbf{F}$, and $\mathbf{W}$ are

---

[2]In particular, although it is mentioned that $\mathbf{G}_k$ (in the notation of that paper) may have to be of reduced size to ensure that $\mathbf{G}_k^{-1}$ exists, the consequences of this reduction are not followed through and the supposedly orthogonal matrix $\mathbf{G}_k^{-1}\mathbf{F}_k^T\mathbf{Z}_k^f$ finally obtained for converting the ETKF post-multiplier matrix into that for the EAKF need not be square.

related to $\mathbf{G}_0$, $\mathbf{F}_0$, and $\mathbf{W}_0$ by (20), (21), and (22). The eigenvalue decomposition (47) extends to

$$(\mathbf{HFG})^T \mathbf{R}^{-1} \mathbf{HFG} = \widetilde{\mathbf{C}} \widetilde{\mathbf{\Gamma}} \widetilde{\mathbf{C}}^T \tag{51}$$

where $\widetilde{\mathbf{C}}$ is orthogonal, $\widetilde{\mathbf{\Gamma}}$ is diagonal, and both are $m \times m$. This is achieved by setting

$$\widetilde{\mathbf{C}} = \begin{pmatrix} \widetilde{\mathbf{C}}_0 & \mathbf{0} \\ \mathbf{0} & \widetilde{\mathbf{C}}_1 \end{pmatrix}, \tag{52}$$

$$\widetilde{\mathbf{\Gamma}} = \begin{pmatrix} \widetilde{\mathbf{\Gamma}}_0 & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix}, \tag{53}$$

where $\widetilde{\mathbf{C}}_1$ is an $(m-r) \times (m-r)$ orthogonal matrix. It may be verified by substitution that (51) is satisfied. Similarly it follows by substitution that $\mathbf{X}^a$ defined by (49) also satisfies

$$\mathbf{X}^a = \mathbf{FG}\widetilde{\mathbf{C}}(\mathbf{I} + \widetilde{\mathbf{\Gamma}})^{-\frac{1}{2}} \mathbf{W}^T = \mathbf{X}^f \mathbf{T} \tag{54}$$

where

$$\mathbf{T} = \mathbf{W}\widetilde{\mathbf{C}}(\mathbf{I} + \widetilde{\mathbf{\Gamma}})^{-\frac{1}{2}} \mathbf{W}^T. \tag{55}$$

This gives the EAKF in post-multiplier form. To see that $\mathbf{T}$ satisfies (17) recall that $\mathbf{W}$ and $\widetilde{\mathbf{C}}$ are orthogonal, and so

$$
\begin{aligned}
\mathbf{TT}^T &= \mathbf{W}\widetilde{\mathbf{C}}(\mathbf{I} + \widetilde{\mathbf{\Gamma}})^{-1}\widetilde{\mathbf{C}}^T \mathbf{W}^T \\
&= (\mathbf{I} + \mathbf{W}\widetilde{\mathbf{C}}\widetilde{\mathbf{\Gamma}}\widetilde{\mathbf{C}}^T \mathbf{W}^T)^{-1} \\
&= (\mathbf{I} + \mathbf{W}(\mathbf{HFG})^T \mathbf{R}^{-1} \mathbf{HFGW}^T)^{-1} \\
&= (\mathbf{I} + (\mathbf{HZ}^f)^T \mathbf{R}^{-1} \mathbf{HZ}^f)^{-1} \\
&= (\mathbf{I} + (\mathbf{Y}^f)^T \mathbf{R}^{-1} \mathbf{Y}^f)^{-1} \\
&= \mathbf{I} - (\mathbf{Y}^f)^T \mathbf{D}^{-1} \mathbf{Y}^f.
\end{aligned}
\tag{56}
$$

where the last step uses the identity (31).

## 8.3 The filter of Whitaker and Hamill

Another filter given in pre-multiplier form, and thus unbiased, is the filter of Whitaker and Hamill (2002). The filter is formulated in terms that initially appear to

require a linear observation operator, but it will be shown below that this restriction can be eliminated. The filter is given by

$$\mathbf{X}^a = (\mathbf{I} - \widetilde{\mathbf{K}}\mathbf{H})\mathbf{X}^f \tag{57}$$

where $\widetilde{\mathbf{K}}$ is a solution of

$$(\mathbf{I} - \widetilde{\mathbf{K}}\mathbf{H})\mathbf{P}^f(\mathbf{I} - \widetilde{\mathbf{K}}\mathbf{H})^T = (\mathbf{I} - \mathbf{K}\mathbf{H})\mathbf{P}^f, \tag{58}$$

$\mathbf{K}$ being the standard Kalman gain defined by (6). This equation ensures that the ensemble covariance matrix updates as in (5). A solution of (58) is given in Whitaker and Hamill (2002), based on Andrews (1968). This solution is

$$\begin{aligned}
\widetilde{\mathbf{K}} &= \mathbf{P}^f\mathbf{H}^T \left[\left(\sqrt{\mathbf{H}\mathbf{P}^f\mathbf{H}^T + \mathbf{R}}\right)^T\right]^{-1} \left(\sqrt{\mathbf{H}\mathbf{P}^f\mathbf{H}^T + \mathbf{R}} + \sqrt{\mathbf{R}}\right)^{-1} \\
&= \mathbf{X}^f(\mathbf{Y}^f)^T \left(\sqrt{\mathbf{D}}^T\right)^{-1} \left(\sqrt{\mathbf{D}} + \sqrt{\mathbf{R}}\right)^{-1}
\end{aligned} \tag{59}$$

where, given a symmetric positive definite $p \times p$ matrix $\mathbf{V}$, the square root $\sqrt{\mathbf{V}}$ stands for a $p \times p$ matrix such that $\sqrt{\mathbf{V}}\sqrt{\mathbf{V}}^T = \mathbf{V}$. The general solution (59) is not considered in Tippett et al. (2003), which instead concentrates on the case of scalar observations where $p = 1$. However, it is not difficult to show that the more general form fits into the ensemble SRF framework. To do this, expand (57) and substitute (59) to obtain

$$\begin{aligned}
\mathbf{X}^a &= \mathbf{X}^f - \widetilde{\mathbf{K}}\mathbf{Y}^f \\
&= \mathbf{X}^f - \mathbf{X}^f(\mathbf{Y}^f)^T \left(\sqrt{\mathbf{D}}^T\right)^{-1} \left(\sqrt{\mathbf{D}} + \sqrt{\mathbf{R}}\right)^{-1}\mathbf{Y}^f \\
&= \mathbf{X}^f\mathbf{T}
\end{aligned} \tag{60}$$

where

$$\mathbf{T} = \mathbf{I} - (\mathbf{Y}^f)^T \left(\sqrt{\mathbf{D}}^T\right)^{-1} \left(\sqrt{\mathbf{D}} + \sqrt{\mathbf{R}}\right)^{-1}\mathbf{Y}^f. \tag{61}$$

Note that the linear operator $\mathbf{H}$ does not explicitly appear in this post-multiplier form of the filter, which may therefore be applied when the observation operator is nonlinear. It may be shown that $\mathbf{T}$ satisfies (17) as follows (which adapts a proof of Andrews (1968)):

$$\mathbf{T}\mathbf{T}^T = \left[\mathbf{I} - (\mathbf{Y}^f)^T \left(\sqrt{\mathbf{D}}^T\right)^{-1} \left(\sqrt{\mathbf{D}} + \sqrt{\mathbf{R}}\right)^{-1}\mathbf{Y}^f\right]$$

$$
\times \left[ \mathbf{I} - (\mathbf{Y}^f)^T \left( \sqrt{\mathbf{D}}^T \right)^{-1} \left( \sqrt{\mathbf{D}} + \sqrt{\mathbf{R}} \right)^{-1} \mathbf{Y}^f \right]^T
$$

$$
= \quad \mathbf{I} - (\mathbf{Y}^f)^T \left( \sqrt{\mathbf{D}}^T \right)^{-1} \left( \sqrt{\mathbf{D}} + \sqrt{\mathbf{R}} \right)^{-1}
$$

$$
\times \left[ \sqrt{\mathbf{D}} \left( \sqrt{\mathbf{D}} + \sqrt{\mathbf{R}} \right)^T + \left( \sqrt{\mathbf{D}} + \sqrt{\mathbf{R}} \right) \sqrt{\mathbf{D}}^T - \mathbf{Y}^f (\mathbf{Y}^f)^T \right]
$$

$$
\times \left[ \left( \sqrt{\mathbf{D}} + \sqrt{\mathbf{R}} \right)^T \right]^{-1} \sqrt{\mathbf{D}}^{-1} \mathbf{Y}^f
$$

$$
= \quad \mathbf{I} - (\mathbf{Y}^f)^T \left( \sqrt{\mathbf{D}}^T \right)^{-1} \left( \sqrt{\mathbf{D}} + \sqrt{\mathbf{R}} \right)^{-1}
$$

$$
\times \left[ \left( \sqrt{\mathbf{D}} + \sqrt{\mathbf{R}} \right) \left( \sqrt{\mathbf{D}} + \sqrt{\mathbf{R}} \right)^T \right]
$$

$$
\times \left[ \left( \sqrt{\mathbf{D}} + \sqrt{\mathbf{R}} \right)^T \right]^{-1} \sqrt{\mathbf{D}}^{-1} \mathbf{Y}^f
$$

$$
= \quad \mathbf{I} - (\mathbf{Y}^f)^T \left( \sqrt{\mathbf{D}} \sqrt{\mathbf{D}}^T \right)^{-1} \mathbf{Y}^f
$$

$$
= \quad \mathbf{I} - (\mathbf{Y}^f)^T \mathbf{D}^{-1} \mathbf{Y}^f. \tag{62}
$$

Although it was stated earlier that the filter is unbiased, this statement depended on the pre-multiplier form (57) in which the linear operator $\mathbf{H}$ appears. Unbiasedness follows in the general post-multiplier case by noting that $\mathbf{T1} = \mathbf{1}$ (because $\mathbf{Y}^f \mathbf{1} = \mathbf{0}$) and using theorem 9. Thus the general form of the filter of Whitaker and Hamill (2002) fits into the unbiased ensemble SRF framework.

# 9    Summary and discussion

Since its original presentation by Evensen (1994), several alternative formulations of the EnKF have been published. Some of these make use of an ensemble of pseudo-random observation perturbations in the analysis step, others do not. The ensemble SRF framework of Tippett et al. (2003) is a uniform framework encompassing several variants of the latter (deterministic) approach. This paper extends the results of Tippett et al. (2003) in two directions. Firstly, it explicitly shows that the ensemble SRF framework is sufficiently general to encompass all deterministic formulations of the analysis step of the EnKF (section 3). Secondly, it shows that the ensemble SRF framework also encompasses filters in which the analysis ensemble statistics do not bear the desired KF-like relationship to the forecast ensemble (section 4). In particular, there can be a systematic bias in the analysis ensemble mean.

The analysis ensemble statistics produced by a biased ensemble SRF are undesirable for a number of reasons beyond the simple fact that a biased mean tends to put the filter's best state estimate in the wrong place. Such a bias would not be too great a problem if it were accompanied by an increase in the size of the error estimate provided by the filter's covariance matrix. Users of the output would then be aware of the increased error, although they would remain unaware that part of the error is systematic rather than random. However, as is shown in section 4, there is actually a decrease in the size of the error estimate rather than an increase, and the worse the bias, the worse the overconfidence of the error estimate. A biased and overconfident analysis has the potential to create problems at later times in any Kalman-type filter. Such an analysis is likely to lead to a biased and overconfident forecast. The filter will then give more weight than it should to the forecast in the next analysis step and less to the observation. This will prevent the observation from properly correcting the bias in the forecast and the next analysis will be biased and overconfident as well. In extreme cases the filter may become increasingly overconfident until it is in effect a free-running forecast model diverging from the truth and taking no notice of observations.

To avoid the problems of a biased analysis ensemble, this paper introduces a restricted version of the ensemble SRF framework called the unbiased ensemble SRF framework (section 5). This is still sufficiently general to include all deterministic formulations of the analysis step of the EnKF, yet excludes filters with a biased analysis ensemble. Tests are provided in sections 6 and 7 that will enable designers of ensemble SRFs to test their algorithms for bias. Particularly simple tests are corollary 11 and theorem 15, which respectively state that a filter is unbiased if its post-multiplier matrix is symmetric or if it can be written in pre-multiplier form.

Of the filters discussed by Tippett et al. (2003), the EAKF and the filter of Whitaker and Hamill (2002) are unbiased, whilst the ETKF is biased. However, the more recent revised ETKF of Wang et al. (2004) is unbiased.

In this paper it has been assumed that the ensemble provides both a best state estimate (the ensemble mean) and a measure of the uncertainty in this estimate (the ensemble covariance matrix). However, an alternative approach is possible in which

a best state estimate is maintained separately from the ensemble, which still provides the measurement of estimation error. The forecast and analysis perturbation matrices are taken relative to the forecast and analysis best state estimates rather than the ensemble means. It is not necessary for the columns of these matrices to sum to zero and hence there is no need to impose an unbiasedness condition in the analysis step: the ensemble perturbation matrices may be updated using $\mathbf{X}^a = \mathbf{X}^f \mathbf{T}$ where $\mathbf{T}$ is any solution of the matrix square root condition (17). An example of such a filter is the maximum likelihood ensemble filter (MLEF) of Zupanski (2005), in which the analysis step updates the best state estimate using 3D-Var (with a cost function that uses the forecast ensemble covariance matrix instead of a static background error covariance matrix) and updates the ensemble perturbations using the ETKF. Although that paper references the original ETKF of Bishop et al. (2001), it is in fact the revised ETKF of Wang et al. (2004) that is used[3].

Finally, it must be stated the the type of bias discussed in this paper is not the only type of bias that may be encountered with an EnKF. Inconsistent ensemble statistics have been observed in formulations of the EnKF other than the original ETKF. Houtekamer and Mitchell (1998) present results showing problems with a stochastic EnKF and Anderson (2001) discusses the issue in the context of the EAKF. The causes of the inconsistencies in these cases must be different to that of the ETKF bias established in section 8.1. The authors attribute them to the use of small ensembles and to other approximations made in the course of deriving the filters. Various solutions to the problem have been proposed in the literature. Houtekamer and Mitchell (1998) use a pair of ensembles with the covariance calculated from each ensemble being used to assimilate observations into the other. The justification for such an approach is discussed further in van Leeuwen (1999) and Houtekamer and Mitchell (1999). Anderson (2001) uses a tunable scalar covariance inflation factor. The more fundamental problems of model and observation biases are not addressed here. Such biases may be estimated using data assimilation with an augmented state vector (e.g., Nichols (2003)). The incorporation of these

---

[3]By Zupanski (2005, (10) and (12)) the post-multiplier matrix is $\mathbf{V}(\mathbf{I}+\mathbf{\Lambda})^{-1/2}\mathbf{V}^T$ in the notation of that paper, which corresponds to (45) in this paper

techniques into the ensemble SRF framework is left for future work.

## Acknowledgements

## A    Structure of the set of ensemble SRFs

This appendix uses the results of section 3 to describe the set of all ensemble SRFs in terms of a well-known group of matrices.

**Theorem 18** *Let $\mathbf{T}_1$ be a solution of the matrix square root condition (17). Then any solution of (17) may be uniquely expressed in the form $\mathbf{T} = \mathbf{T}_1\mathbf{U}$ where $\mathbf{U}$ is orthogonal.*

**Proof.** Start with the special case $\mathbf{T}_1 = \mathbf{T}_s$ where $\mathbf{T}_s$ is the unique symmetric positive definite solution that exists by theorem 6. Recall that $\mathbf{T}$ is nonsingular by theorem 5. By a theorem in linear algebra (see, for example, Halmos (1974, section 83)) $\mathbf{T}$ has a unique polar decomposition $\mathbf{T} = \mathbf{T}_2\mathbf{U}$ where $\mathbf{T}_2$ is symmetric positive definite and $\mathbf{U}$ is orthogonal. Since $\mathbf{T}$ is a solution of (17),

$$\mathbf{I} - (\mathbf{Y}^f)^T\mathbf{D}^{-1}\mathbf{Y}^f = \mathbf{T}\mathbf{T}^T = \mathbf{T}_2\mathbf{T}_2^T. \tag{63}$$

Thus $\mathbf{T}_2$ is a symmetric positive definite solution of (17), which by theorem 6 implies $\mathbf{T}_2 = \mathbf{T}_s$. This establishes the theorem in the case $\mathbf{T}_1 = \mathbf{T}_s$.

Now consider the general case. By the special case above there exist orthogonal $\mathbf{U}_1$ and $\mathbf{U}_2$ such that $\mathbf{T}_1 = \mathbf{T}_s\mathbf{U}_1$ and $\mathbf{T} = \mathbf{T}_s\mathbf{U}_2$. Thus $\mathbf{T} = \mathbf{T}_s\mathbf{U}_2 = \mathbf{T}_1\mathbf{U}_1^{-1}\mathbf{U}_2$. Since $\mathbf{U}_1^{-1}\mathbf{U}_2$ is orthogonal this establishes the existence of $\mathbf{U}$ in the general case.

To show the uniqueness of $\mathbf{U}$ in the general case, suppose that there exists an orthogonal matrix $\mathbf{U}_3$ such that $\mathbf{T}_1\mathbf{U} = \mathbf{T}_1\mathbf{U}_3$. Since $\mathbf{T}_1$ is nonsingular it may be

cancelled from both sides of this equation to give $\mathsf{U} = \mathsf{U}_3$. Therefore $\mathsf{U}$ is unique.
□

Theorems 2 and 18 imply the following description of the set of all solutions $\mathsf{T}$ of (17) in terms of the group $O(m)$ of all $m \times m$ orthogonal matrices.

**Corollary 19** *Let $\mathsf{T}_1$ be a solution of (17). Then $\mathsf{U} \leftrightarrow \mathsf{T}_1 \mathsf{U}$ defines a one-to-one correspondence between $O(m)$ and the solutions $\mathsf{T}$ of (17).*

Thus the set of all ensemble SRFs is in one-to-one correspondence with $O(m)$.

# B  The Two-Dimensional Swinging Spring

The two-dimensional swinging spring (Lynch, 2003) consists of a heavy bob of mass $m$ suspended from a fixed point in a uniform gravitational field of acceleration $g$ by a light spring of unstretched length $\ell_0$ and elasticity $k$. The bob is constrained to move in a vertical plane. The system coordinates are polar coordinates $r$, $\theta$ ($r$ measured from the point of suspension, $\theta$ measured from the downward vertical) and the corresponding generalised momenta $p_r$, $p_\theta$. The equations of motion are

$$\dot{\theta} = \frac{p_\theta}{mr^2}, \tag{64}$$

$$\dot{p}_\theta = -mgr \sin\theta, \tag{65}$$

$$\dot{r} = \frac{p_r}{m}, \tag{66}$$

$$\dot{p}_r = \frac{p_\theta^2}{mr^3} - k(r - \ell_0) + mg \cos\theta. \tag{67}$$

The equilibrium length $\ell$ of the spring satisfies $k(\ell - \ell_0) = mg$. Following Lynch (2003) the parameter values used for the experiments in sections 4 and 8.1 are $m = 1$, $g = \pi^2$, $k = 100\pi^2$, and $\ell = 1$. The initial conditions are $(\theta, p_\theta, r, p_r) = (1, 0, 0.99540, 0)$, which is a case of nonlinear normal mode initialisation and largely suppresses the high frequency radial oscillations of the system.

# C  Structure of the set of unbiased ensemble SRFs

This appendix uses the results of section 6 and appendix A to describe the set of all unbiased ensemble SRFs in terms of a well-know group of matrices.

**Theorem 20** *Suppose that* $\mathbf{X}^f$ *is nondegenerate and that* $\mathbf{T}_1$ *is the post-multiplier matrix of an unbiased ensemble SRF. Then* $\mathbf{U} \leftrightarrow \mathbf{T}_1\mathbf{U}$ *defines a one-to-one correspondence between the subgroup of all matrices* $\mathbf{U}$ *in* $O(m)$ *that have* $\mathbf{1}$ *as an eigenvector and the set of post-multiplier matrices of unbiased ensemble SRFs.*

**Proof.** This follows from corollary 19 and theorems 13 and 14. □

The subgroup in theorem 20 is given more concrete form by the following theorem.

**Theorem 21** *There is a one-to-one correspondence between the subgroup of all matrices in* $O(m)$ *that have* $\mathbf{1}$ *as an eigenvector and the group* $O(1) \times O(m-1)$.

**Proof.** Let $\mathbf{W}$ be an orthogonal matrix in which the first column is a scalar multiple of $\mathbf{1}$. Then $\mathbf{U} \leftrightarrow \mathbf{W}^T\mathbf{U}\mathbf{W}$ is a one-to-one correspondence between $O(m)$ and itself. Under this correspondence, matrices $\mathbf{U}$ that have $\mathbf{1}$ as an eigenvector correspond to matrices that have the coordinate vector

$$\mathbf{e}_1 = \begin{pmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix} \tag{68}$$

as an eigenvector. The latter matrices are those of the form

$$\begin{pmatrix} \mathbf{U}_1 & \mathbf{0} \\ \mathbf{0} & \mathbf{U}_2 \end{pmatrix} \tag{69}$$

where $\mathbf{U}_1$ is an element of $O(1)$ (that is to say $\mathbf{U}_1 = \pm 1$) and $\mathbf{U}_2$ is an element of $O(m-1)$. This establishes the required correspondence. □

Thus, in the case of nondegenerate $\mathbf{X}^f$, the set of all unbiased ensemble SRFs is in one-to-one correspondence with $O(1) \times O(m-1)$.

# References

J. L. Anderson. An ensemble adjustment Kalman filter for data assimilation. *Mon. Wea. Rev.*, 129:2884–2903, 2001.

A. Andrews. A square root formulation of the Kalman covariance equations. *AIAA J.*, 6:1165–1166, 1968.

C. H. Bishop, B. J. Etherton, and S. J. Majumdar. Adaptive sampling with the ensemble transform Kalman filter. Part I: Theoretical aspects. *Mon. Wea. Rev.*, 129:420–436, 2001.

G. Burgers, P. J. van Leeuwen, and G. Evensen. Analysis scheme in the ensemble Kalman filter. *Mon. Wea. Rev.*, 126:1719–1724, 1998.

G. Evensen. Sequential data assimilation with a nonlinear quasi-geostrophic model using Monte Carlo methods to forecast error statistics. *J. Geophys. Res.*, 99(C5): 10143–10162, 1994.

G. Evensen. The Ensemble Kalman Filter: Theoretical formulation and practical implementation. *Ocean Dynamics*, 53:343–367, 2003.

A. Gelb, editor. *Applied Optimal Estimation.* The M.I.T. Press, 1974. 374 pp.

G. H. Golub and C. F. Van Loan. *Matrix Computations.* The Johns Hopkins University Press, 3rd edition, 1996. 694 pp.

P. R. Halmos. *Finite-Dimensional Vector Spaces.* Springer-Verlag, 1974. 200 pp.

P. L. Houtekamer and H. L. Mitchell. Data assimilation using an ensemble Kalman filter technique. *Mon. Wea. Rev.*, 126:796–811, 1998.

P. L. Houtekamer and H. L. Mitchell. Reply. *Mon. Wea. Rev.*, 127:1378–1379, 1999.

A. H. Jazwinski. *Stochastic Processes and Filtering Theory.* Academic Press, 1970. 376 pp.

P. Lynch. Introduction to initialization. In R. Swinbank, V. Shutyaev, and W. A. Lahoz, editors, *Data Assimilation for the Earth System*, pages 97–111. Kluwer Academic Publishers, 2003.

N. K. Nichols. Treating model error in 3-d and 4-d data assimilation. In R. Swinbank, V. Shutyaev, and W. A. Lahoz, editors, *Data Assimilation for the Earth System*, pages 127–135. Kluwer Academic Publishers, 2003.

M. K. Tippett, J. L. Anderson, C. H. Bishop, T. M. Hamill, and J. S. Whitaker. Ensemble square root filters. *Mon. Wea. Rev.*, 131:1485–1490, 2003.

P. J. van Leeuwen. Comment on "Data assimilation using an ensemble Kalman filter technique". *Mon. Wea. Rev.*, 127:1374–1377, 1999.

X. Wang, C. H. Bishop, and S. J. Julier. Which is better, an ensemble of positive-negative pairs or a centered spherical simplex ensemble? *Mon. Wea. Rev.*, 132:1590–1605, 2004.

J. S. Whitaker and T. M. Hamill. Ensemble data assimilation without perturbed observations. *Mon. Wea. Rev.*, 130:1913–1924, 2002.

M. Zupanski. Maximum likelihood ensemble filter: Theoretical aspects. *Mon. Wea. Rev.*, 133:1710–1726, 2005.