

Comparison of Ensemble Kalman Filters Under Non-Gaussianity

JING LEI * AND PETER BICKEL

Department of Statistics, University of California, Berkeley

CHRIS SNYDER

National Center for Atmospheric Research, Boulder, CO USA

* *Corresponding author address:* Jing Lei, Department of Statistics, UC Berkeley, 367 Evans Hall, Berkeley, CA 94720.

E-mail: jinglei@stat.berkeley.edu

ABSTRACT

Recently various versions of ensemble Kalman filters (EnKF) has been proposed and studied. This work concerns, in a mathematically rigorous manner, the relative performance of two major versions of EnKF when the forecast ensemble is non-Gaussian. The approach is based on the stability of the filtering methods against small model violations, using the expected squared L_2 distance as a measure of the deviation between the updated distributions. Analytical and experimental results suggest that both stochastic and deterministic EnKFs are sensitive to the violation of the Gaussianity assumption, while the stochastic filter is relatively more stable than the deterministic filter under certain circumstances, especially when there are wild outliers. These results not only agree with previous empirical studies, but also suggest a natural choice of a free parameter in the square-root Kalman filter algorithm.

1. Introduction

The ensemble Kalman filter (EnKF, (Evensen 1994, 2003, 2007)) has become a popular tool for data assimilation because of its computational efficiency and flexibility (Anderson 2001; Whitaker and Hamill 2002; Ott et al. 2004; Bengtsson et al. 2003; Evensen 2007). In various versions of EnKFs, one major difference is how to get the updated ensemble after obtaining the updated mean and variance. Stochastic methods (Houtekamer and Mitchell 1998; Evensen 2003) directly use the Kalman gain together with random perturbations. On the other hand, deterministic methods (Anderson 2001; Bishop et al. 2001) use a non-random transformation on the forecast ensemble, which is also known as a special case of the Kalman square-root filter (Tippett et al. 2003).

The analysis error of EnKF consists of two parts: the use of a linear analysis algorithm that is suboptimal for all except Gaussian distributions; and the variance caused by using only a finite sample. The latter is studied for the stochastic filter by Sacher and Bartello (2008, 2009). In this paper we study the first part of error, that is, the error caused by non-Gaussianity.

Following the direction of Lawson and Hansen (2004), who did empirical comparison of the stochastic and deterministic filters, in this work we attempt to quantify the difference between these two methods under non-Gaussianity, through the perspective of *robustness*. It is known that in a Gaussian linear model both methods are consistent (Furrer and Bengtsson 2007). However, when the forecasting distribution is non-Gaussian both methods are biased even asymptotically, where the bias refers to the deviation from the true conditional distribution or equivalently the distribution given by the Bayes rules. Suppose the previous

updated ensemble is approximately Gaussian. After propagation through the non-linear dynamics, the resulting forecast ensemble will be slightly non-Gaussian if the time interval is short. Figure 1 gives such an example by looking at the first two coordinates of the Lorenz 63 3-dimensional system¹, where the previous update ensemble is Gaussian but the forecasting ensemble has some outliers. Therefore one would expect some bias in EnKF update due to the non-Gaussianity, and the bias could be different for different implementation of EnKF. Our question is: which method is more stable against non-Gaussianity? Here “stability” is a statistical notion which refers to the analysis being not seriously biased when the forecast distribution is slightly non-Gaussian. Another notion of “stability” is introduced by Sacher and Bartello (2009) which refers to the size of analysis error covariance being large enough to cover the true analysis center. We give a rigorous analysis of the sensitivity of the two EnKFs to non-Gaussianity of the forecasting ensemble based on the notion of *robustness* in statistics.

We show that the stochastic filter is more robust than the deterministic filter especially when the position of outliers is wild and/or the observation is accurate. Simulation results support our calculation not only for the L_2 distance but also for other quantities such as the third moment. These findings are consistent with those in Lawson and Hansen (2004). Moreover, such a comparison can be extended to many other types of model violations, such as the modeling error in the observation and the observation model. On the other hand, we also show that such a stability criterion leads to a natural choice of the orthogonal matrix

¹The Lorenz 63 system (Lorenz 1963) is a three dimensional continuous chaotic system, which is very sensitive to initial conditions in the discrete-step form. It has been used to test filtering methods in many data assimilation research works (see Anderson and Anderson 1999; Bengtsson, Snyder, and Nychka 2003).

in the unbiased ensemble square root filter Sakov and Oke (2007); Livings et al. (2008).

In Section 2 we introduce the ensemble Kalman filters, with a brief discussion on the large-ensemble behavior of the EnKF. Section 3 contains the main part of our comparison, beginning with some intuition in Section 3a; The basic concepts of asymptotic robustness can be found in Hampel et al. (1986), and we give a brief summary in Section 3b; In Section 3c we state our analytical results. Finally, in Section 4, we present various numerical experiments.

2. Ensemble Kalman filters

a. The Kalman filter

Consider a Gaussian linear model:

$$\mathbf{y} = \mathbf{H}\mathbf{x} + \epsilon,$$

where $\mathbf{x} \in \mathbb{R}^p$ is the hidden state variable, $\mathbf{y} \in \mathbb{R}^q$ the observation, $\epsilon \in \mathbb{R}^q$ an independent random noise, and $\mathbf{H} \in \mathbb{R}^{q \times p}$ the observation matrix. Assuming all the variables are Gaussian:

$$\mathbf{x} \sim N(\mu^f, \mathbf{P}^f), \quad \epsilon \sim N(0, \mathbf{R}),$$

then the updated state variable $\mathbf{x}|\mathbf{y}^o$ given a specific observation \mathbf{y}^o is still Gaussian²:

$$\mathbf{x}|\mathbf{y}^o \sim N(\mu^a, \mathbf{P}^a),$$

with

$$\mu^a = (\mathbf{I} - \mathbf{KH})\mu^f + \mathbf{K}\mathbf{y}^o, \quad \mathbf{P}^a = (\mathbf{I} - \mathbf{KH})\mathbf{P}^f, \quad (1)$$

²Throughout this paper we use superscript “f” and “a” to denote “forecast” and “analysis (update)” respectively.

where $\mathbf{K} = \mathbf{P}^f \mathbf{H}^T (\mathbf{H} \mathbf{P}^f \mathbf{H}^T + \mathbf{R})^{-1}$ is the *Kalman Gain*. Throughout this paper we always assume that \mathbf{P}^f and \mathbf{R} are positive definite.

Several practical issues arise in geophysics. First, the state variable is driven by non-linear geophysical dynamics, so its exact distribution is unknown and certainly is non-Gaussian. Usually only a random sample from the distribution is available. Second, the linear form of the observation is, again, only an approximation. The true observation model $\mathbf{y} = h(\mathbf{x}) + \varepsilon$ might involve a nonlinear $h(\cdot)$, or $h(\cdot)$ might even have no explicit functional form (e.g., a black-box function). These problems are partially addressed, as described below, by the ensemble Kalman filter.

b. The ensemble Kalman filter

Suppose $(\mathbf{x}^{f(i)})_{i=1}^n$ is an i.i.d (independent, identically distributed) sample from the forecast distribution of the state variable \mathbf{x}^f . The ensemble Kalman filter update consists of the following steps:

- i. Let $\hat{\boldsymbol{\mu}}^f$ and $\hat{\mathbf{P}}^f$ be the sample mean and covariance.
- ii. Estimate the Kalman gain: $\hat{\mathbf{K}} = \hat{\mathbf{P}}^f \mathbf{H}^T (\mathbf{H} \hat{\mathbf{P}}^f \mathbf{H}^T + \mathbf{R})^{-1}$.
- iii. Update the mean and covariance according to the Kalman filter:

$$\langle \hat{\boldsymbol{\mu}}^a \rangle = (\mathbf{I} - \hat{\mathbf{K}} \mathbf{H}) \hat{\boldsymbol{\mu}}^f + \hat{\mathbf{K}} \mathbf{y}^o, \quad \langle \hat{\mathbf{P}}^a \rangle = (\mathbf{I} - \hat{\mathbf{K}} \mathbf{H}) \hat{\mathbf{P}}^f,$$

where $\langle \cdot \rangle$ denotes the expectation over the randomness of the update procedure. If the update is deterministic, then $\langle \hat{\boldsymbol{\mu}}^a \rangle = \hat{\boldsymbol{\mu}}^a$ and $\langle \hat{\mathbf{P}}^a \rangle = \hat{\mathbf{P}}^a$.

iv. Update the ensemble $(\mathbf{x}^{f(i)})_1^n \rightarrow (\mathbf{x}^{a(i)})_1^n$, so that

$$\frac{1}{n} \sum_{i=1}^n \mathbf{x}^{a(i)} = \hat{\boldsymbol{\mu}}^a, \quad \frac{1}{n-1} \sum_{i=1}^n (\mathbf{x}^{a(i)} - \hat{\boldsymbol{\mu}}^a)(\mathbf{x}^{a(i)} - \hat{\boldsymbol{\mu}}^a)^T = \hat{\mathbf{P}}^a. \quad (2)$$

It is worth noting that in practice, the sample covariance matrix $\hat{\mathbf{P}}^f$ is not computed explicitly. Instead, it is sufficient to compute $\hat{\mathbf{P}}^f \mathbf{H}^T = \frac{1}{n-1} \sum (x^{f(i)}) (\mathbf{H}x^{f(i)})^T$, which is computationally more efficient if p is much larger than q .

The stochastic and the deterministic filters differ in step 4. In the stochastic filter,

$$\mathbf{x}_s^{a(i)} = \mathbf{x}^{f(i)} + \hat{\mathbf{K}}(\mathbf{y}^o - \mathbf{H}\mathbf{x}^{f(i)} + \epsilon^{(i)}), \quad \forall 1 \leq i \leq n, \quad (\text{STO.})$$

where $\epsilon^{(i)} \stackrel{iid}{\sim} N(0, R)$. The intuition is to use directly the Kalman gain to combine the forecast ensemble member $\mathbf{x}^{f(i)}$ and the observation \mathbf{y}^o , using additive noise $\epsilon^{(i)}$ to adjust the total variance of the updated ensemble, as if the perturbed observation associated with $\mathbf{x}^{f(i)}$ is another possible value of random variable \mathbf{y} . In some applications in order to reduce the sampling error of the noise, $\epsilon^{(i)}$'s are adjusted by a shifting and rescaling to ensure one of the following:

- $\epsilon^{(i)}$'s have zero mean.
- $\epsilon^{(i)}$'s have zero mean and covariance \mathbf{R} .
- $\epsilon^{(i)}$'s have zero mean, covariance \mathbf{R} and zero covariance with $X_f^{(i)}$'s.

When the ensemble size n is large, such a shifting and rescaling is negligible and all these variants are equivalent to the update given by (STO.). Therefore the analysis in this paper is applicable to these variants too.

The deterministic filter works in a different way:

$$\mathbf{x}_d^{a(i)} = \hat{\mu}^a + \hat{\mathbf{A}}(\mathbf{x}^{f(i)} - \hat{\mu}^f), \quad \forall 1 \leq i \leq n, \quad (\text{DET.})$$

where $\hat{\mathbf{A}}$ satisfies $\hat{\mathbf{A}}\hat{\mathbf{P}}^f\hat{\mathbf{A}}^T = \hat{\mathbf{P}}^a$. Loosely speaking, the matrix $\hat{\mathbf{A}}$ can be viewed as the square root of the difference between $\hat{\mathbf{P}}^a$ and $\hat{\mathbf{P}}^f$. The matrix $\hat{\mathbf{A}}$ is not unique in the multivariate case. Suppose $n > p$ and $\hat{\mathbf{P}}^f$ is full rank, then $\hat{\mathbf{A}}$ has the general form:

$$\hat{\mathbf{A}} = (\hat{\mathbf{P}}^a)^{\frac{1}{2}}\mathbf{U}(\hat{\mathbf{P}}^f)^{-\frac{1}{2}}, \quad (3)$$

where \mathbf{U} is any $p \times p$ orthogonal matrix chosen by the user. See Tippett et al. (2003); Sakov and Oke (2007) for further discussion on the choice of \mathbf{U} . If $n \leq p$ and $\hat{\mathbf{P}}^f$ is not full rank, (3) no longer holds but one can work on the principal components of the state space instead of the whole state space as described in Ott et al. (2004).

There is another formula for the update step of the deterministic filter using the right-multiplication:

$$\mathbf{x}_d^{a(i)} = \hat{\mu}^a + \sum_{j=1}^n \hat{a}'_{ij}(\mathbf{x}^{f(j)} - \hat{\mu}^f). \quad (4)$$

This formula can be shown to be closely related to (DET.) when the filter is unbiased, i.e., $\frac{1}{n} \sum_{i=1}^n \mathbf{x}_d^{a(i)} = \hat{\mu}^a$ (Tippett et al. 2003; Livings et al. 2008). We will use the left-multiplication throughout this paper because: 1) it has a clear geometrical interpretation; 2) we assume that n is large.

In practical applications, good performance of the EnKFs defined by (STO.) and (DET.) depends on a sufficiently large ensemble and on system dynamics and observation models that are sufficiently close to linear. For example, the EnKF will dramatically underestimate \mathbf{P}^a with small ensembles as it is analytically described by Sacher and Bartello (2008). As

a result, covariance localization and covariance inflation have been widely used to overcome such practical difficulties (Whitaker and Hamill 2002; Ott et al. 2004; Anderson 2003, 2007).

c. The large-ensemble behavior of the EnKF

If $n \rightarrow \infty$, then by law of large numbers, everything converges to its population counterpart. That is, $\hat{\mu}^f \xrightarrow{P} \mu^f$, $\hat{\mathbf{P}}^f \xrightarrow{P} \mathbf{P}^f$, $\hat{\mathbf{K}} \xrightarrow{P} \mathbf{K}$, $\hat{\mu}^a \xrightarrow{P} \mu^a$, $\hat{\mathbf{P}}^a \xrightarrow{P} \mathbf{P}^a$, and $\hat{\mathbf{A}} \xrightarrow{P} \mathbf{A}$ where $\mathbf{A} = (\mathbf{P}^a)^{\frac{1}{2}} \mathbf{U}(\mathbf{P}^f)^{-\frac{1}{2}}$ is the population counterpart of $\hat{\mathbf{A}}$. Here \xrightarrow{P} denotes convergence in probability³. Let $\delta_{\mathbf{x}}$ denote the point mass at \mathbf{x} (i.e., a probability distribution that puts all its mass at \mathbf{x}), then intuitively the empirical updated distributions $\hat{F}_s = \frac{1}{n} \sum \delta_{\mathbf{x}_s^{a(i)}}$ and $\hat{F}_d = \frac{1}{n} \sum \delta_{\mathbf{x}_d^{a(i)}}$ should converge weakly to the distribution of the random variables $(\mathbf{I} - \mathbf{KH})\mathbf{x} + \mathbf{K}(\mathbf{y} + \epsilon)$ and $\mu^a + \mathbf{A}(\mathbf{x} - \mu^f)$, respectively. In fact it can be shown that the above intuition is true (Appendix A, Proposition 6). As a result, our comparison between the stochastic filter and the deterministic filter will be based on the comparison between these two limiting distributions.

³For a sequence of random variables α_n , $n \geq 1$, and constant β , $\alpha_n \xrightarrow{P} \beta$ means that for any $\delta > 0$, $\lim_{n \rightarrow \infty} P(|\alpha_n - \beta| > \delta) = 0$.

3. Comparing the stochastic and the deterministic filters

a. Intuition and the contaminated Gaussian model

A simple and natural deviation from Gaussianity is a contaminated Gaussian model:

$$\mathbf{x}^f \sim F_r = (1 - r)F + rG, \quad (5)$$

where, without loss of generality, $F = N(0, \mathbf{P})$, $G = N(t, \mathbf{S})$, where \mathbf{P} and \mathbf{S} are positive definite, and $0 \leq r < 1$ is the amount of contamination. The interpretation of model (5) is that we assume a proportion of $(1 - r)$ of the forecast ensemble are drawn from a Gaussian distribution centered at 0, with covariance \mathbf{P} , while the rest are outliers coming from another Gaussian distribution centered at t with covariance \mathbf{S} . Since we use the Gaussian distribution $G = N(t, \mathbf{S})$ to model the outliers, we would expect G to be much different from $F = N(0, \mathbf{P})$, the majority of the forecast ensemble. That is, we expect (t, S) to be somewhat extreme: $\|t\|_2 \gg 0$ and/or $\|S\|_2 \gg \|\mathbf{P}\|_2$. For example, a large⁴ t and small \mathbf{S} mean that the outliers forms a small cluster far away from the majority, while a small t and a large \mathbf{S} mean that the outliers are widely dispersed. Also, denote $F_{o,r}(\cdot|\mathbf{y})$ the true distribution of \mathbf{x}^a , here the subindex “o” stands for “optimal”. Again, the optimal updated distribution refers to the one given by the Bayes rule. Similarly, the corresponding limiting updated distributions of EnKFs are denoted by $F_{s,r}(\cdot|\mathbf{y})$ and $F_{d,r}(\cdot|\mathbf{y})$, respectively. Here we keep in mind that t and \mathbf{S} are fixed. For simplicity, we focus on the case $q = p$ and $\mathbf{P} = \mathbf{I}_p$.

The merit of a filter can be characterized naturally in terms of the distance between

⁴Here and throughout this paper, by saying a vector or matrix is large we mean its L_2 norm is large.

the updated density and the optimal density $f_{o,r}$. Recall that if \mathbf{x}^f is Gaussian, i.e., $r = 0$, then $F_{s,0}$ and $F_{d,0}$ are both Gaussian, with the same mean and covariance agreeing with the optimal conditional distribution: $F_{s,0} = F_{d,0} = F_{o,0} = N(\mu_o^a, \mathbf{P}_o^a)$. Now the question is, when $r \neq 0$, i.e., \mathbf{x}^f is non-Gaussian, which one is closer to $F_{o,r}$?

We take a quick look at the densities of $F_{o,r}$, $F_{s,r}$ and $F_{d,r}$ in a simple one-dimensional setup similar to Lawson and Hansen (2004), but with $r = 0.05$ (right column of Figure 2). The original figure in Lawson and Hansen (2004) with $r = 0.5$ are included in the left column for comparison. We choose $t = 8$, $\mathbf{S} = 1$, and $\mathbf{y} = 0.5$, which makes \mathbf{y} a plausible observation from F_r . We consider three values of \mathbf{R} : In the top row, $\mathbf{R} = \mathbf{P}_r^f/4$, where \mathbf{P}_r^f is the variance of F_r . In this case the observation is accurate, which indicates that the likelihood function is highly unimodal (with a single high peak). As a result, the stochastic filter approximates the true density better because adding Gaussian perturbations to the bimodal ensemble will make the distribution more unimodal. In the middle row $\mathbf{R} = \mathbf{P}_r^f$, where the accuracy is modest and it is hard to tell which filter gives better approximation to the truth. Finally, in the bottom row we have $\mathbf{R} = 4\mathbf{P}_r^f$, a relatively inaccurate observation. Now when the two components are equally weighted (left column), the stochastic incorrectly populates the middle part because of the random perturbation while the deterministic retains the bimodal structure. In the right column, when the weights of two the components are very unbalanced, the deterministic update is closer to the optimal for a wide range of \mathbf{x} near the origin. However, it carries more outliers due to the small bump at $+7$, which might cause a larger bias in the higher moments.

Remark 1. In model (5) the assumption that G is Gaussian is only for mathematical conve-

nience. It can be an arbitrary distribution, since any distribution can be approximated by a mixture of Gaussian, and as we will see in the next section (eq. (7)), the effect of mixture contamination is approximately additive when the total amount of contamination r is small.

b. The robustness perspective

Robustness (Hampel et al. 1986) is a natural notion of the stability of an inference method against small model violation. Intuitively, a “good” method should give stable outcomes when the true underlying distribution deviates slightly from the ideal distribution. In the context of EnKF, the ideal distribution refers to the Gaussian forecast distribution under which the EnKF gives unbiased analysis. In parameter estimation, let $g(\hat{F}_n)$ be the estimator of parameter from the empirical distribution \hat{F}_n , and $g(F)$ denotes its population counterpart, which is usually the large-sample limit of $g(\hat{F}_n)$. Suppose the true distribution is $(1 - r)F + rG$, a contaminated version of F , for some small $r > 0$. Then the estimator becomes $g((1 - r)F + rG)$. The robustness of g at F means that no matter what G looks like, $g((1 - r)F + rG)$ should be close to $g(F)$ as long as r is small. The quantification of this idea leads to the *Gâteaux derivative* and the *influence function*.

The Gâteaux derivative and the influence function

Following the above notation, the estimator can be viewed as a function of r , the amount of contamination. The *Gâteaux derivative* of g at F in the direction of G is defined by

$$\nu(G, F; g) = \lim_{r \rightarrow 0^+} \frac{g((1 - r)F + rG) - g(F)}{r}. \quad (6)$$

Intuitively, the Gâteaux derivative measures approximately how g is affected by an infinitesimal contamination of shape G on F .

If $G = \delta_t$ is a point mass at t , then one can define

$$\text{IF}(t; F, g) = \nu(\delta_t, F; g),$$

which is the *influence function* of g at F . There is a close analogy between the influence function and Green's function. In both cases, the general solution to a linear problem is a superposition of the solution to point mass problems. It can be shown that, under appropriate conditions, (see Bickel and Doksum, ch. 7.3),

$$\nu(G, F; g) = \int \text{IF}(t; F, g) dG(t). \quad (7)$$

As a result, the function $\text{IF}(\cdot; F, g)$ reflects the robustness of g at F . An important criterion in designing robust estimators is a bounded influence function:

$$\sup_t |\text{IF}(t; F, g)| < \infty.$$

Intuitively, this means that distorting any small proportion of the data can not have a big impact on the outcome.

c. Comparison from the robustness perspective: analytical results

In our study, the parameter, and hence the estimator, is a distribution. For any fixed \mathbf{x} , \mathbf{y} , the Gâteaux derivatives of the conditional densities at \mathbf{x} are⁵, under Model (5),

$$\nu(G, F; f_s(\mathbf{x}|\mathbf{y})) = \lim_{r \rightarrow 0^+} \frac{f_{s,r}(\mathbf{x}|\mathbf{y}) - f_{s,0}(\mathbf{x}|\mathbf{y})}{r} = \left. \frac{\partial}{\partial r} f_{s,r}(\mathbf{x}|\mathbf{y}) \right|_{r=0} \quad (8)$$

⁵In this paper we use $f(\cdot) = F'(\cdot)$ as the density function of $F(\cdot)$, whenever possible. E.g., $f_{s,r}(\cdot|\mathbf{y})$ is the density function of $F_{s,r}(\cdot|\mathbf{y})$. For succinctness, we will use $f_{s,r}$ instead of $f_{s,r}(\cdot|\mathbf{y})$ without confusion.

for the stochastic filter, and

$$\nu(G, F; f_d(\mathbf{x}|\mathbf{y})) = \lim_{r \rightarrow 0^+} \frac{f_{d,r}(\mathbf{x}|\mathbf{y}) - f_{d,0}(\mathbf{x}|\mathbf{y})}{r} = \left. \frac{\partial}{\partial r} f_{d,r}(\mathbf{x}|\mathbf{y}) \right|_{r=0} \quad (9)$$

for the deterministic filter. In our contaminated Gaussian model, the ideal distribution is $F = N(0, \mathbf{I})$ and $G = N(t, \mathbf{S})$ is the contamination distribution. Recall again that $f_{s,0} = f_{d,0} = f_{o,0}$, then equations (8) and (9) are comparing $f_{s,r}(\mathbf{x}|\mathbf{y})$ and $f_{d,r}(\mathbf{x}|\mathbf{y})$ with $f_{o,0}(\mathbf{x}|\mathbf{y})$ respectively.

However, the quantities in (8) and (9) involve not only \mathbf{x} but also \mathbf{y} , the random observation. In order to take all \mathbf{x} as well as the randomness of \mathbf{y} into account, we integrate the square of the Gâteaux derivatives and take expectation over \mathbf{y} under its marginal distribution when $r = 0$, which is $N(0, \mathbf{I} + \mathbf{R})$. Finally, the quantities indicating the robustness of the EnKFs are

$$E_{\mathbf{y}} \left(\int \nu^2(G, F; f_s(\mathbf{x}|\mathbf{y})) d\mathbf{x} \right) = E_{\mathbf{y}} \left[\int \left(\left. \frac{\partial}{\partial r} f_{s,r}(\mathbf{x}|\mathbf{y}) \right|_{r=0} \right)^2 d\mathbf{x} \right] \quad (10)$$

for the stochastic filter, and

$$E_{\mathbf{y}} \left(\int \nu^2(G, F; f_d(\mathbf{x}|\mathbf{y})) d\mathbf{x} \right) = E_{\mathbf{y}} \left[\int \left(\left. \frac{\partial}{\partial r} f_{d,r}(\mathbf{x}|\mathbf{y}) \right|_{r=0} \right)^2 d\mathbf{x} \right] \quad (11)$$

for the deterministic filter.

On the other hand, note that

$$\frac{\partial}{\partial r} \left[\int (f_{s,r}(\mathbf{x}|\mathbf{y}) - f_{s,0}(\mathbf{x}|\mathbf{y}))^2 d\mathbf{x} \right] = 2 \int (f_{s,r}(\mathbf{x}|\mathbf{y}) - f_{s,0}(\mathbf{x}|\mathbf{y})) \frac{\partial}{\partial r} f_{s,r}(\mathbf{x}|\mathbf{y}) d\mathbf{x},$$

and

$$\begin{aligned} & \frac{\partial^2}{\partial r^2} \left[\int (f_{s,r}(\mathbf{x}|\mathbf{y}) - f_{s,0}(\mathbf{x}|\mathbf{y}))^2 d\mathbf{x} \right] \\ &= 2 \int \left(\frac{\partial}{\partial r} f_{s,r}(\mathbf{x}|\mathbf{y}) \right)^2 d\mathbf{x} + 2 \int (f_{s,r}(\mathbf{x}|\mathbf{y}) - f_{s,0}(\mathbf{x}|\mathbf{y})) \frac{\partial^2}{\partial r^2} f_{s,r}(\mathbf{x}|\mathbf{y}) d\mathbf{x}. \end{aligned}$$

Evaluate the above derivatives at $r = 0$, we have

$$\left. \frac{\partial}{\partial r} \left[E_{\mathbf{y}} \int (f_{s,r}(\mathbf{x}|\mathbf{y}) - f_{s,0}(\mathbf{x}|\mathbf{y}))^2 d\mathbf{x} \right] \right|_{r=0} = 0.$$

and

$$\left. \frac{\partial^2}{\partial r^2} \left[\int (f_{s,r}(\mathbf{x}|\mathbf{y}) - f_{s,0}(\mathbf{x}|\mathbf{y}))^2 d\mathbf{x} \right] \right|_{r=0} = 2 \int \left(\left. \frac{\partial}{\partial r} f_{s,r}(\mathbf{x}|\mathbf{y}) \right|_{r=0} \right)^2 d\mathbf{x}.$$

Taking expectation over \mathbf{y} ,

$$E_{\mathbf{y}} \left[\int \left(\left. \frac{\partial}{\partial r} f_{s,r}(\mathbf{x}|\mathbf{y}) \right|_{r=0} \right)^2 d\mathbf{x} \right] = \frac{1}{2} \left. \frac{\partial^2}{\partial r^2} \left[E_{\mathbf{y}} \int (f_{s,r}(\mathbf{x}|\mathbf{y}) - f_{s,0}(\mathbf{x}|\mathbf{y}))^2 d\mathbf{x} \right] \right|_{r=0},$$

As a result, the quantity defined in (10) has a straightforward interpretation: It is the second derivative of the expected square of L_2 distance between $f_{s,r}$ and $f_{s,0}$. The same argument also holds for the deterministic filter. So a smaller value in (10) (or (11)) indicates a slower change in the updated distribution when r changes from zero to non-zero.

Our main theoretical results are summarized in the following theorems:

Theorem 2. *In model (5), we have*

(i) *For all \mathbf{R}, \mathbf{S}*

$$\lim_{\|\mathbf{t}\|_2 \rightarrow \infty} E_{\mathbf{y}} \int \nu^2(G, F; f_s(\mathbf{x}|\mathbf{y})) d\mathbf{x} = \infty, \quad \text{and} \quad 0 < \lim_{\|\mathbf{t}\|_2 \rightarrow \infty} \frac{E_{\mathbf{y}} \int \nu^2(G, F; f_s(\mathbf{x}|\mathbf{y})) d\mathbf{x}}{E_{\mathbf{y}} \int \nu^2(G, F; f_d(\mathbf{x}|\mathbf{y})) d\mathbf{x}} < 1; \quad (12)$$

(ii) *For all \mathbf{R}, t ,*

$$\lim_{\|\mathbf{S}\|_2 \rightarrow \infty} E_{\mathbf{y}} \int \nu^2(G, F; f_s(\mathbf{x}|\mathbf{y})) d\mathbf{x} = \infty, \quad \text{and} \quad 0 < \lim_{\|\mathbf{S}\|_2 \rightarrow \infty} \frac{E_{\mathbf{y}} \int \nu^2(G, F; f_s(\mathbf{x}|\mathbf{y})) d\mathbf{x}}{E_{\mathbf{y}} \int \nu^2(G, F; f_d(\mathbf{x}|\mathbf{y})) d\mathbf{x}} < 1; \quad (13)$$

(iii) For all t, \mathbf{S} ,

$$\lim_{\|\mathbf{R}\|_2 \rightarrow 0} E_{\mathbf{y}} \int \nu^2(G, F; f_s(\mathbf{x}|\mathbf{y})) d\mathbf{x} = \infty, \quad \text{and} \quad \lim_{\|\mathbf{R}\|_2 \rightarrow 0} \frac{E_{\mathbf{y}} \int \nu^2(G, F; f_s(\mathbf{x}|\mathbf{y})) d\mathbf{x}}{E_{\mathbf{y}} \int \nu^2(G, F; f_d(\mathbf{x}|\mathbf{y})) d\mathbf{x}} = 0. \quad (14)$$

Proof. The proof is included in Appendix B. □

Parts (i) and (ii) of Theorem 2 indicate that neither of the two filters has bounded Gâteaux derivative over all possible contaminations. However, when the contamination is wild, the stochastic filter is more stable than the deterministic filter. Loosely speaking, when there are outliers in the forecast ensemble, the Kalman filter will suffer from its non-robustness due to the use of the sample mean and sample covariance matrix. The deterministic filter is affected more because its rigid shifting and re-scaling (in order to make the exact covariance) leaves no chance to correct the outliers, while the stochastic filter uses a “softer” method to adjust the ensemble mean and covariance by using random perturbations. It is thus more resilient to outliers because there is some chance that the outliers are partially corrected by the random perturbations. This effect can also be seen in the top right plot of Figure 2. Moreover, it also implies that, in the multivariate case, when the contamination is wild, the deviation in the updated density is largely determined by the magnitude, not the orientation, of t and/or \mathbf{S} . As shown later in Section 4, the asymptotic result also holds even for moderately large choices of $\|t\|_2$ and $\|\mathbf{S}\|_2$.

Part (iii) indicates that stochastic filter is more stable when the observation is accurate. This result nicely supports the intuitive argument in Lawson and Hansen (2004): the convolution with a Gaussian random perturbation in the stochastic filter makes the updated ensemble closer to Gaussian while the deterministic might push the edge-members in the

ensemble to be far-outliers and have the major component in the mixture overly tight.

The case that $\|\mathbf{R}\|_2 \rightarrow \infty$ is particularly interesting. Intuitively, a very large $\|\mathbf{R}\|_2$ indicates a very non-informative observation. Thus the conditional distribution should be close to the forecast distribution. As a result, one should do little change on the forecast ensemble when $\|\mathbf{R}\|_2$ is large. This intuition suggests choosing the orthogonal matrix $\mathbf{U} = \mathbf{I}$ in the deterministic filter, the benefit of which can be seen through Theorem 3:

Theorem 3. *If in (3) we choose $\mathbf{U} = \mathbf{I}$, then for all t, \mathbf{S} ,*

$$0 < \lim_{\|\mathbf{R}\|_2 \rightarrow \infty} E_{\mathbf{y}} \int \nu^2(G, F; f_s(\mathbf{x}|\mathbf{y})) d\mathbf{x} < \infty, \quad \text{and} \quad \lim_{\|\mathbf{R}\|_2 \rightarrow \infty} \frac{E_{\mathbf{y}} \int \nu^2(G, F; f_s(\mathbf{x}|\mathbf{y})) d\mathbf{x}}{E_{\mathbf{y}} \int \nu^2(G, F; f_d(\mathbf{x}|\mathbf{y})) d\mathbf{x}} = 1. \quad (15)$$

Otherwise, we have

$$\lim_{\|t\|_2 \rightarrow \infty} \lim_{\|\mathbf{R}\|_2 \rightarrow \infty} \frac{E_{\mathbf{y}} \int \nu^2(G, F; f_s(\mathbf{x}|\mathbf{y})) d\mathbf{x}}{E_{\mathbf{y}} \int \nu^2(G, F; f_d(\mathbf{x}|\mathbf{y})) d\mathbf{x}} = 0, \quad (16)$$

and

$$\lim_{\|\mathbf{S}\|_2 \rightarrow \infty} \lim_{\|\mathbf{R}\|_2 \rightarrow \infty} \frac{E_{\mathbf{y}} \int \nu^2(G, F; f_s(\mathbf{x}|\mathbf{y})) d\mathbf{x}}{E_{\mathbf{y}} \int \nu^2(G, F; f_d(\mathbf{x}|\mathbf{y})) d\mathbf{x}} = 0. \quad (17)$$

Proof. See Appendix B. □

Theorem 3 is easy to understand. Intuitively, when \mathbf{R} is large, we have $\mu^a \approx \mu^f$ and $\mathbf{P}^a \approx \mathbf{P}^f$ in the Kalman filter. Here $\mathbf{U} = \mathbf{I}$ implies $\mathbf{A} \approx \mathbf{I}$, which means making little change on the forecast ensemble. In Section 4 we will see that the choice of $\mathbf{U} = \mathbf{I}$ does beat other choices even for moderately large \mathbf{R} , \mathbf{S} and t . The issue of choosing the orthogonal matrix in the square root filter has been discussed in Sakov and Oke (2007), which mainly focuses on the right-multiplication case. Theorem 3 suggests a stable choice of the left-multiplying

orthogonal matrix which means the corresponding right-multiplying orthogonal matrix is stable due to the correspondence between the left and right-multiplication in unbiased square root filters (Livings et al. 2008) if $p < n$.

Remark 4. Theorems 2 and 3 concern the effects caused by a large t , \mathbf{S} and \mathbf{R} separately, by means of sending one quantity to infinity while keeping others fixed. In fact, these quantities do interact in the optimal and EnKF updates, which will affect the comparison in a much more complicated manner. Although in this more interesting case analytical results seem hard to derive, we do think these theorems provide some qualitative view of the comparison as we will see in the numerical experiments.

d. Connection to bias comparison

The robustness tells us about the stability of the filters when the data distribution is nearly ideal. However, as mentioned earlier, a more direct comparison would be to just look at the bias, that is, the difference between the limiting distribution of the updated ensemble ($f_{s,r}$ and $f_{d,r}$), and the optimal conditional distribution ($f_{o,r}$). A first observation is that when r is small, then $f_{o,r} \approx f_{o,0}$, i.e., $f_{o,r}$ would mostly be as if there is no contamination at all, as long as \mathbf{y} is not too far from 0 or not too close to t , which is often the case when $\|t\|_2 \gg 0$ and \mathbf{y}° is randomly drawn from f_r . This can be seen from the fact that

$$F_{o,r} = (1 - \pi(r))N(\mu_{o,1}^a, \mathbf{P}_{o,1}^a) + \pi(r)N(\mu_{o,2}^a, \mathbf{P}_{o,2}^a), \quad (18)$$

where, letting $\phi(\mathbf{x}; \mu, \mathbf{P})$ be the density of $N(\mu, \mathbf{P})$ at \mathbf{x} ,

$$\pi(r) = \frac{r\phi(\mathbf{y}; t, \mathbf{S} + \mathbf{R})}{r\phi(\mathbf{y}; t, \mathbf{S} + \mathbf{R}) + (1 - r)\phi(\mathbf{y}; 0, \mathbf{I} + \mathbf{R})},$$

and, for $j = 1, 2$, with the convention that $\mu_1^f = 0$, $\mathbf{P}_1^f = \mathbf{P}^f$, $\mu_2^f = t$, and $\mathbf{P}_2^f = \mathbf{S}$,

$$\mathbf{K}_j = \mathbf{P}_j^f(\mathbf{P}_j^f + \mathbf{R})^{-1}, \quad \mu_j^a = (\mathbf{I} - \mathbf{K}_j)\mu_j^f + \mathbf{K}_j\mathbf{y}, \quad \mathbf{P}_j^a = (\mathbf{I} - \mathbf{K}_j)\mathbf{P}_j^f.$$

For the proof of (18), we refer the reader to Bengtsson et al. (2003) and references therein.

As a result, when $\|t\|_2 \gg 0$, and \mathbf{y} not far from 0, we have $\pi(r)/r \approx 0$.

As a result, although it might be difficult to compare $f_{s,r}$ ($f_{d,r}$) with $f_{o,r}$, comparing $f_{s,r}$ ($f_{d,r}$) with $f_{o,0}$ can give some rough idea for the hard comparison. Note further that $f_{s,0} = f_{d,0} = f_{o,0}$, which means that $f_{s,r} - f_{o,r} \approx f_{s,r} - f_{o,0} = f_{s,r} - f_{s,0}$. That is, robustness actually indicates small bias. In Section 4 we present simulation results to verify this idea.

A limitation of our analysis to this point is that the L_2 distance provides only partial information about the deviation of the analysis distribution from the optimal (Bayes) update. In fact, data assimilation is best evaluated by 1) the distance between the analysis center and the true posterior center and 2) the size of the analysis covariance which needs to be large enough to have the analysis ensemble cover a substantial proportion of the true posterior distribution including its center. These two criteria are labeled in Sacher and Bartello (2009) as “accuracy” and “stability” respectively (recall that in this paper the notion of “stability” is different). In the context of large ensemble behavior, the analysis center is almost the same for the stochastic filter and the deterministic filter. Therefore they should perform similarly in this aspect given they are starting from the same forecast ensemble. On the other hand, although both filters have the same second order statistics, the updated ensemble is distributed differently for a non-Gaussian prior. This difference will affect the future forecast ensemble and hence the filter performance in sequential applications, which needs to be explored further.

Another class of criteria are higher order moments since in a non-Gaussian distribution the higher moments contains much information about the error distribution. In the next subsection we consider the third moment as another measure of performance to support our previous results.

e. The third moment

The third moment is an indication of the skewness of the distribution. Therefore it seems a natural criterion beyond the first two moments to evaluate the updated ensemble. Lawson and Hansen (2004) also considered the ensemble skewness in their experiments. Here for presentation simplicity we consider the one dimensional model given by (5).

Assuming model (5), let $M_s(\mathbf{y}) = \int \mathbf{x}^3 f_s(\mathbf{x}|\mathbf{y})d\mathbf{x}$ be the third moment of the limiting updated distribution given by the stochastic filter and similarly define $M_d(\mathbf{y})$ for the deterministic filter. Then we have the following theorem:

Theorem 5. *Under model (5), if both $X_f \in \mathbb{R}^1$ and $Y \in \mathbb{R}^1$, then*

(i) *For all \mathbf{S} and \mathbf{y}*

$$\lim_{|t| \rightarrow \infty} |\nu(G, F; M_s(\mathbf{y}))| = \infty, \quad \text{and} \quad \lim_{|t| \rightarrow \infty} \frac{|\nu(G, F; M_s(\mathbf{y}))|}{|\nu(G, F; M_d(\mathbf{y}))|} < 1. \quad (19)$$

(ii) *For all t and \mathbf{y} ,*

$$\lim_{|\mathbf{S}| \rightarrow \infty} |\nu(G, F; M_s(\mathbf{y}))| = \infty, \quad \text{and} \quad \lim_{|\mathbf{S}| \rightarrow \infty} \frac{|\nu(G, F; M_s(\mathbf{y}))|}{|\nu(G, F; M_d(\mathbf{y}))|} < 1. \quad (20)$$

Proof. See Appendix B. □

These results are similar to those in the previous theorems, except that the third moment is a scalar which allows us to derive results for each value of \mathbf{y} . The intuition behind Theorem 5 can be seen from Figure 2, where the deterministic filter tends to produce two components which are less spread and further away from each other than in the stochastic filter. As a result, the deterministic filter puts a little more density in the region which are likely outliers (the bump near $=7$ on the bottom right plot). Despite maintaining the right mean and covariance, these outliers will have a substantial impact on the higher moments as shown in Theorem 5. The empirical comparison of the bias of the third moments is provided in Section 4.

4. Simulation results

In this section we present simulation results comparing the performance of the two versions of ensemble Kalman filters. As we will see later, the simulations do support the analytical results and intuitive discussion in Section 3c and 3d.

a. The 1-dimensional case

In the 1-dimensional case, n random samples are drawn from $F_r = (1 - r)F + rG$ as described in model (5), under different combinations of model parameters $(r, \mathbf{R}, t, \mathbf{S})$ as defined in Section 3a. Both versions of EnKF are applied to the same random sample and observation from which the optimal conditional distribution is calculated. We first check the expected square of L_2 distance as a measure of bias as a direct verification of Theorem 2

and 3, then we look at the third moment to further confirm our results.

The expected square of L_2 distance

Once all the parameters in Model (5) is specified, for any value of \mathbf{y} , the functions $f_{s,r}(\mathbf{x})$, $f_{d,r}(\mathbf{x})$ and $f_{o,r}(\mathbf{x})$ can be calculated analytically. The expected square of L_2 distances

$$E_{\mathbf{y}} \int (f_{s,r}(\mathbf{x}) - f_{o,r}(\mathbf{x}))^2 d\mathbf{x}, \quad \text{and} \quad E_{\mathbf{y}} \int (f_{d,r}(\mathbf{x}) - f_{o,r}(\mathbf{x}))^2 d\mathbf{x} \quad (21)$$

are calculated numerically. That is, \mathbf{y} is simulated many times, and for each simulated value of \mathbf{y} the above integrals are calculated numerically and averaged. In Table 1, we set $t = 8$, $\mathbf{S} = 1$, the same setup as in Figure 2. Actually the simulation is quantifying the difference between the density curves shown in Figure 2, except that it takes further expectation over all possible values of \mathbf{y} . Three different values of \mathbf{R} are chosen according to its relative size with $\mathbf{P}_r^f = \text{var}(\mathbf{x}|F_r)$. This result supports the analysis in Section 3c and the intuition in Section 3d: when r is small, $f_{s,r}$ is closer to $f_{o,r}$. Moreover, it seems that the asymptotic statement can be extended to much larger value of r , e.g., $r = 0.5$ as shown in Table 1. The expectation over \mathbf{y} is approximated by averaging over 1000 simulated values of \mathbf{y} (standard deviations are shown in the parentheses).

The third moment

The EnKF forces the updated ensemble to have the correct first and second moment, therefore the third moment becomes a natural criterion of comparison. The empirical third moments of the two updated ensembles are compared with the optimal third moment which

is calculated analytically.

Here, instead of taking expectation over \mathbf{y} , we investigate the impact of the value \mathbf{y} on the comparison. That is, we look at all $\mathbf{y} \in \left[-2(1 + \mathbf{R})^{\frac{1}{2}}, 2(1 + \mathbf{R})^{\frac{1}{2}}\right]$, which covers a majority of probability mass in F_r . In the experiment, $(\mathbf{R}, \mathbf{S}) \in \{1/4, 1, 4\}^2$, and $t \in \{1, 10, 30, 50, 100\}$. We choose $(n, r) = (500, 0.05)$. Several representative pictures are displayed in Figure 3. We see that for small t , both filters give very small bias for almost the whole range of \mathbf{y} , and when t gets bigger, the stochastic filter gives smaller bias for a wide range of \mathbf{y} , which covers the majority of probability mass of its distribution. Moreover, the difference is enhanced by larger values of \mathbf{R} and \mathbf{S} .

b. 2-dimensional case

In the 2-dimensional case, our theory claims that it is the magnitude of the matrices that determine the amount of deviation. However, in the finite sample simulation, it seems necessary to consider not only the magnitude, but also different orientations of the matrices. We consider two instances:

- *Orientation 1*: $\mathbf{P} = \mathbf{I}_2$, $\mathbf{R} = c_1 \mathbf{R}_0$ and $\mathbf{S} = c_2 \mathbf{S}_0$, where $(c_1, c_2) \in \{1/4, 1, 4\}^2$ tunes the magnitude of \mathbf{R} and \mathbf{S} , where $\mathbf{R}_0 = \text{diag}(1.5, 1)$, and \mathbf{S}_0 is a simulated 2 by 2 Wishart matrix:

$$\mathbf{S}_0 = \begin{pmatrix} 1.15 & 0.14 \\ 0.14 & 0.70 \end{pmatrix}.$$

- *Orientation 2*: In this case we consider a contamination distribution G with very different shape from F , i.e., \mathbf{S}_0 that has very different orientation from $\mathbf{P} = \text{cov}_F(\mathbf{x})$.

Here we choose \mathbf{P} to be, up to a scaling constant, the covariance matrix of the stationary distribution of the first two coordinates in the Lorenz 63 system, and \mathbf{S}_0 is obtained by switching the eigenvalues of P .

$$\mathbf{P} = \begin{pmatrix} 1.06 & 1.05 \\ 1.05 & 1.35 \end{pmatrix}, \quad \mathbf{S}_0 = \begin{pmatrix} 1.35 & -1.05 \\ -1.05 & 1.06 \end{pmatrix}.$$

Here \mathbf{S}_0 has the same eigenvectors as Σ , but with the eigenvalues switched. That is,

$$\mathbf{P} = \mathbf{Q} \begin{pmatrix} d_1 & 0 \\ 0 & d_2 \end{pmatrix} \mathbf{Q}^T, \quad \mathbf{S}_0 = \mathbf{Q} \begin{pmatrix} d_2 & 0 \\ 0 & d_1 \end{pmatrix} \mathbf{Q}^T,$$

where \mathbf{Q} is a orthogonal matrix and $d_1 = 0.15$, $d_2 = 2.27$ are eigenvalues of \mathbf{P} and \mathbf{S}_0 .

The other settings are the same as above expect that $\mathbf{R}_0 = \mathbf{I}_2$.

The contour of the two Gaussian densities are plotted in Figure 4.

In the deterministic algorithm we try two choices of \mathbf{U} in (3). The first is simply to choose $\mathbf{U} = \mathbf{I}$. The second choice is based on the “ensemble adjustment Kalman filter” (EAKF) proposed by Anderson (2001). Similar to the 1-dimensional case, the expectation over \mathbf{y} is approximated by averaging over 120 simulated \mathbf{y} . Standard deviations are shown in the parentheses. Some representative results are summarized in Table 2, where $r = 0.05$, $c_2 = 1$, and $t = (10, 10)$ (other values make no qualitative difference).

Recall that c_1 indicates the size of \mathbf{R} . We can see that for small c_1 , the stochastic filter is remarkably less biased, agreeing with the experiments in Lawson and Hansen (2004). Also note that in Model (5) both the forecast and the analysis distribution are a mixture of two Gaussian components, where the major component (i.e., the one with a weight close to 1) contains mostly “normal” ensemble members whereas the minor component (the one

whose weight is close to 0) contains mostly ensemble members that are likely outliers. When the observation is accurate, the optimal filter puts more weight on the major Gaussian component. On the other hand neither of the two EnKFs adjusts the component weights in the analysis. The two components in the analysis distribution given by the deterministic filter are less spread than those given by the stochastic filter. In order to have the same covariance, the less spread components have to be further away from each other. As a result, the outliers tends to be even more outlying in the deterministic update. An instance of this intuition can be seen in the right panel of Figure 2 where the deterministic filter always produces a small bump in the right tail, especially for small observation errors.

Another interesting observation is the comparison of the choices of the rotation matrix \mathbf{U} . For small observation noise, the difference is negligible. One can imagine that when the observation is accurate, the optimal analysis distribution tends to be closer to a Gaussian, whose distribution is determined by the first two moments, therefore the rotation does not make too much difference. While when c_1 gets bigger, the analysis ensemble becomes much less Gaussian and the choice $\mathbf{U} = \mathbf{I}$ shows significant advantage as compared with the EAKF, agreeing with Theorem 3. This basically says that when the observation is very uninformative, there is no need to change, and hence no need to rotate, the ensemble.

Moreover, the results shown in Table 2 also confirm the theory in that only the magnitude of the contamination matters since similar behavior is observed for two very different shapes of contamination distribution.

5. Conclusion

We have studied the large-ensemble performance of ensemble Kalman filters using the robustness approach. In the contaminated Gaussian model, the updated distribution is another mixture with two components, where the stochastic filter is more stable against small model violation due to the fact that its main component in the updated distribution is closer to that of the optimal filter. Our theoretical results are supported by intensive simulation over a wide range of the model parameters, agreeing with the empirical findings in Lawson and Hansen (2004), where the intuitive argument says that the deterministic shifting and re-scaling exaggerates the dispersion of some ensemble members.

Although our study focuses on the large-ensemble behavior under a classical model, our method can be extended in at least two directions. First, the influence function theory enables one to study other shapes of contamination, rather than Gaussian. Second, in geophysical studies the model deviation might come from the observation, instead of the state variable. In other words, the modeling error could come from the mis-specification of the distribution of the observation error. The approach developed in this paper is applicable to analysis of situations where the observation error is not exactly Gaussian.

The choice of the orthogonal matrix \mathbf{U} in the deterministic filter is an unsettled issue in data assimilation literature. Our L_2 -based stability criterion gives an answer to this question which is intuitively reasonable: you do almost nothing when the observation is uninformative.

In practice, there are many factors determining which filtering method to use, such as the computational constraints, the modeling error, the particular prediction task, and the specific shapes of the forecasting distribution and error distribution, etc. But this cannot

be done before we fully understand the properties of all the candidates. We hope our study contributes to that understanding.

Acknowledgments.

The authors would like to thank Dr. J. Hansen for helpful comments. Lei and Bickel are supported by NSF grant DMS 0605236.

APPENDIX A

Large-ensemble behavior of EnKFs

Following the discussion in Section c, we have:

Proposition 6. *As $n \rightarrow \infty$, we have*

$$\hat{F}_s \Rightarrow F_s, \quad \hat{F}_d \Rightarrow F_d,$$

where F_s and F_d are the distribution functions of $(\mathbf{I} - \mathbf{KH})\mathbf{x}^f + \mathbf{K}(\mathbf{y} + \epsilon)$ and $\mu^a + \mathbf{A}(\mathbf{x}^f - \mu^f)$, respectively.

Our theoretical result on comparing the stochastic and deterministic filters are based on F_s and F_d .

Proof. We show the weak convergence of \hat{F}_s . The proof for \hat{F}_d is similar.

Let J be a random index uniformly drawn from $\{1, \dots, n\}$. Let $\hat{Z}_n = (\mathbf{I} - \hat{\mathbf{K}}\mathbf{H})\mathbf{x}^{f(J)} + \hat{\mathbf{K}}(\mathbf{y} + \epsilon^{(J)})$ and $Z_n = (\mathbf{I} - \mathbf{KH})\mathbf{x}^{f(J)} + \mathbf{K}(\mathbf{y} + \epsilon^{(J)})$. Then $\hat{Z}_n \sim \hat{F}_s$, and $Z_n \sim F_s$, so it is enough to show that $\hat{Z}_n - Z_n \xrightarrow{P} 0$.

Consider the random variable $W = \mathbf{H}\mathbf{x}^f - \mathbf{y} - \epsilon$. For any $\xi > 0$, $\delta > 0$, one can find an M large enough such that $P(\|W\|_2 \geq M\xi) \leq \delta/2$. On the other hand, since $\hat{\mathbf{K}} - \mathbf{K} \xrightarrow{P} 0$, one can find $N_{\xi, \delta}$ such that $P(\|\hat{\mathbf{K}} - \mathbf{K}\|_2 \geq 1/M) \leq \delta/2$ whenever $n \geq N_{\xi, \delta}$. Then for all

$n \geq N_{\xi, \delta}$, we have

$$\begin{aligned}
P\left(\|\hat{Z}_n - Z_n\|_2 \geq \xi\right) &= P\left(\|(\hat{\mathbf{K}} - \mathbf{K})(\mathbf{H}\mathbf{x}^{f(J)} - \mathbf{y} - \epsilon^{(J)})\| \geq \xi\right) \\
&\leq P\left(\|\hat{\mathbf{K}} - \mathbf{K}\|_2 \geq 1/M\right) + P\left(\|\mathbf{H}\mathbf{x}^{f(J)} - \mathbf{y} - \epsilon^{(J)}\|_2 \geq M\xi\right) \\
&= P\left(\|\hat{\mathbf{K}} - \mathbf{K}\|_2 \geq 1/M\right) + P\left(\|\mathbf{H}\mathbf{x} - \mathbf{y} - \epsilon\|_2 \geq M\xi\right) \\
&\leq \delta/2 + \delta/2 \\
&= \delta.
\end{aligned}$$

□

Remark 7. In Proposition 6, there is nothing special about Gaussianity, so the result holds for any random variable \mathbf{x}^f such that $E\mathbf{x}^f = \mu^f$, $\text{Var}(\mathbf{x}^f) = \mathbf{P}^f$.

APPENDIX B

Proofs of the main theorems

Proof of Theorem 2

We give a sketchy proof for part (i), the argument applies similarly to other parts.

We first consider the simpler case: $t = \rho t_0$, where $\|t_0\|_2 = 1$.

Letting $\mathbf{K} = (\mathbf{I} + \mathbf{R})^{-1}$, $\mathbf{B} = \mathbf{I} - \mathbf{K}$, $\Gamma = tt^T + \mathbf{S} - \mathbf{I}$, $\mathbf{A} = \mathbf{A}(0) = \mathbf{B}^{\frac{1}{2}}\mathbf{U}$ for some orthogonal \mathbf{U} , and $V_s = \mathbf{B}\Gamma\mathbf{B}^T - \mathbf{A}\Gamma\mathbf{A}^T$, then, in the deterministic filter, we have

$$\begin{aligned} & \left. \frac{\partial}{\partial r} f_{d,r}(\mathbf{x}) \right|_{r=0} \\ &= \left[-\frac{1}{2} \text{tr}(\mathbf{B}^{-1}V_s) + (\Gamma\mathbf{K}\mathbf{y} + \mathbf{B}^{-1}(\mathbf{B} - \mathbf{A})t)^T(\mathbf{x} - \mathbf{K}\mathbf{y}) \right. \\ & \quad \left. + \frac{1}{2}(\mathbf{x} - \mathbf{K}\mathbf{y})^T \mathbf{B}^{-1}V_s \mathbf{B}^{-1}(\mathbf{x} - \mathbf{K}\mathbf{y}) - 1 \right] \phi(\mathbf{x}; \mathbf{K}\mathbf{y}, \mathbf{B}) + \phi(\mathbf{x}; \mathbf{K}\mathbf{y} + \mathbf{A}t, \mathbf{A}\mathbf{S}\mathbf{A}^T). \end{aligned}$$

Then it can be shown, via some algebra, that

$$E_{\mathbf{y}} \int \left(\left. \frac{\partial}{\partial r} f_{d,r}(\mathbf{x}) \right|_{r=0} \right)^2 d\mathbf{x} = C \cdot a_d(t_0)\rho^4 + P_d(\rho) + e^{-\kappa_d\rho^2} Q_d(\rho), \quad (\text{B1})$$

where $P_d(\rho)$ and $Q_d(\rho)$ are polynomials of degree 3; $C > 0$ is a constant depending only on

\mathbf{B} ; $\kappa_d > 0$ is a constant; and

$$a_d(t_0) = \frac{1}{2} \text{tr}(t_0 t_0^T \mathbf{K} t_0 t_0^T \mathbf{B}) + \frac{1}{16} E \left(z^T \left(\mathbf{B}^{\frac{1}{2}} t_0 t_0^T \mathbf{B}^{\frac{1}{2}} - \mathbf{U} t_0 t_0^T \mathbf{U}^T \right) z \right)^2. \quad (\text{B2})$$

On the other hand, in the stochastic filter,

$$\begin{aligned} & \left. \frac{\partial}{\partial r} f_{s,r}(\mathbf{x}) \right|_{r=0} \\ &= [(\Gamma \mathbf{K} \mathbf{y})^T (\mathbf{x} - \mathbf{K} \mathbf{y}) - 1] \phi(\mathbf{x}; \mathbf{K} \mathbf{y}, \mathbf{B}) + \phi(\mathbf{x}; \mathbf{K} \mathbf{y} + \mathbf{B} t, \mathbf{B} \mathbf{S} \mathbf{B}^T + \mathbf{K} \mathbf{R} \mathbf{K}^T). \end{aligned}$$

Similarly,

$$E_{\mathbf{y}} \int \left(\left. \frac{\partial}{\partial r} f_{s,r}(\mathbf{x}) \right|_{r=0} \right)^2 d\mathbf{x} = C \cdot a_s(t_0) \rho^4 + P_s(\rho) + e^{-\kappa_d \rho^2} Q_s(\rho), \quad (\text{B3})$$

where $P_s(\rho)$ and $Q_s(\rho)$ are polynomials of degree 3; C is the same constant as in (B1); $\kappa_s > 0$ is a constant; and

$$a_s(t_0) = \frac{1}{2} \text{tr}(t_0 t_0^T \mathbf{K} t_0 t_0^T \mathbf{B}). \quad (\text{B4})$$

Note that $\|\mathbf{B}^{\frac{1}{2}} t_0\|_2 < \|\mathbf{U} t_0\|_2$, for all $t_0 \neq 0$. Therefore,

$$\lim_{\rho \rightarrow \infty} E_{\mathbf{y}} \int \nu^2(G, F; f_s(\mathbf{x}|\mathbf{y})) d\mathbf{x} = \infty, \quad \text{and} \quad \lim_{\rho \rightarrow \infty} \frac{E_{\mathbf{y}} \int \nu^2(G, F; f_s(\mathbf{x}|\mathbf{y})) d\mathbf{x}}{E_{\mathbf{y}} \int \nu^2(G, F; f_d(\mathbf{x}|\mathbf{y})) d\mathbf{x}} = \frac{a_s(t_0)}{a_d(t_0)} < 1.$$

The statement of Theorem 2 (i) follows easily via a standard argument using the compactness of the set $\{t_0 \in \mathbb{R}^p : \|t_0\|_2 = 1\}$.

The proofs for part (ii) and (iii) are simply repeating the argument above on \mathbf{S} and \mathbf{R} , respectively.

Proof of Theorem 3

The argument is essentially the same as in the proof of Theorem 2. Starting from the easy facts:

$$\lim_{\|\mathbf{R}\|_2 \rightarrow \infty} \mathbf{K} = \mathbf{0}, \quad \lim_{\|\mathbf{R}\|_2 \rightarrow \infty} \mathbf{B} = \mathbf{I}, \quad \text{and} \quad \lim_{\|\mathbf{R}\|_2 \rightarrow \infty} \mathbf{A} = \mathbf{U},$$

then

$$\begin{aligned} & \lim_{\|\mathbf{R}\| \rightarrow \infty} \frac{\partial}{\partial r} f_{d,r}(\mathbf{x}) \Big|_{r=0} \\ &= \left[((\mathbf{I} - \mathbf{U})t)^\top \mathbf{x} + \frac{1}{2} \mathbf{x}^\top (\Gamma - \mathbf{U}\Gamma\mathbf{U}^\top) \mathbf{x} - 1 \right] \phi(\mathbf{x}; 0, \mathbf{I}) + \phi(\mathbf{x}; \mathbf{U}t, \mathbf{U}\mathbf{S}\mathbf{U}^\top), \end{aligned}$$

and

$$\lim_{\|\mathbf{R}\| \rightarrow \infty} \frac{\partial}{\partial r} f_{s,r}(\mathbf{x}) \Big|_{r=0} = -\phi(\mathbf{x}; 0, \mathbf{I}) + \phi(\mathbf{x}; t, \mathbf{S}).$$

The rest of the proof is simply repeating the argument for the proof of Theorem 2 (i) and (ii).

Proof of Theorem 5

The result is straight forward if one realizes that $F_{s,r}$ and $F_{d,r}$ are both Gaussian mixtures with two components. One can calculate analytically the parameters of each component. Then straight calculus gives:

$$\begin{aligned} & \nu(G, F; M_s(\mathbf{y})) \\ &= \beta^3 t^3 + (3\alpha^3 \beta + 6\alpha\beta^2) t^2 + (3\alpha^2 \beta + 3\beta^2 + 3\beta^3(\mathbf{S} - 1)) t + (\mathbf{S} - 1)(3\alpha^3 \beta + 6\alpha\beta^2), \end{aligned} \quad (\text{B5})$$

and

$$\begin{aligned} & \nu(G, F; M_d(\mathbf{y})) \\ &= \beta^{\frac{3}{2}} t^3 + (3\alpha^3 \beta + 6\alpha\beta^2) t^2 + \left(3\alpha^2 \beta + 3\beta^2 - 3\beta^{\frac{3}{2}} + 3\beta^{\frac{1}{2}} \mathbf{S} \right) t + (\mathbf{S} - 1)(3\alpha^3 \beta + 6\alpha\beta^2), \end{aligned} \quad (\text{B6})$$

where

$$\alpha = \frac{\mathbf{y}}{1 + \mathbf{R}}, \quad \beta = \frac{\mathbf{R}}{1 + \mathbf{R}}.$$

Then the results in Theorem 5 follows immediately because $0 < \beta < 1$ for all \mathbf{R} .

REFERENCES

- Anderson, J., 2001: An ensemble adjustment kalman filter for data assimilation. *Monthly Weather Review*, **129**, 2884–2903.
- Anderson, J. L., 2003: A local least squares framework for ensemble filtering. *Monthly Weather Review*, **131**, 634–642.
- Anderson, J. L., 2007: Exploring the need for localization in ensemble data assimilation using a hierarchical ensemble filter. *Physica D*, **230**, 99–111.
- Anderson, J. L. and S. L. Anderson, 1999: A monte carlo implementation of the nonlinear filtering problem to produce ensemble assimilations and forecasts. *Monthly Weather Review*, **127**, 2741–2758.
- Bengtsson, T., C. Snyder, and D. Nychka, 2003: Toward a nonlinear ensemble filter for high-dimensional systems. *Journal of Geophysical Research*, **108(D24)**, STS2.1–STS2.10.
- Bickel, P. J. and K. A. Doksum, to appear: Mathematical Statistics, basic ideas and selected topics. Volume II.
- Bishop, C. H., B. Etherton, and S. J. Majumdar, 2001: Adaptive sampling with the ensemble transformation kalman filter. part i: theoretical aspects. *Monthly Weather Review*, **129**, 420–436.

- Evensen, G., 1994: Sequential data assimilation with a non-linear quasi-geostrophic model using monte carlo methods to forecast error statistics. *J. Geophys. Res.*, **99(C5)**, 10 143–10 162.
- Evensen, G., 2003: The ensemble kalman filter: theoretical formulation and practical implementation. *Ocean Dynamics*, **53**, 343–367.
- Evensen, G., 2007: *Data assimilation: the ensemble Kalman filter*. Springer.
- Furrer, R. and T. Bengtsson, 2007: Estimation of high-dimensional prior and posterior covariance matrices in kalman filter variants. *Journal of Multivariate Analysis*, **98**, 227–255.
- Hampel, F., E. Ronchetti, P. Rousseeuw, and W. Stahel, 1986: *Robust Statistics: The Approach Based on Influence Functions*. John Wiley.
- Houtekamer, P. L. and H. L. Mitchell, 1998: Data assimilation using an ensemble kalman filter technique. *Monthly Weather Review*, **126**, 796–811.
- Lawson, G. W. and J. A. Hansen, 2004: Implications of stochastic and deterministic filters as ensemble-based data assimilation methods in varying regimes of error growth. *Monthly Weather Review*, **132**, 1966–1981.
- Livingston, D. M., S. L. Dance, and N. K. Nicols, 2008: Unbiased ensemble square root filters. *Physica D*, **237**, 1021–1028.
- Lorenz, E. N., 1963: Deterministic nonperiodic flow. *Journal of the Atmospheric Sciences*, **20**, 130–141.

- Ott, E., et al., 2004: A local ensemble kalman filter for atmospheric data assimilation. *Tellus*, **56A**, 415–428.
- Sacher, W. and P. Bartello, 2008: Sampling errors in ensemble Kalman filtering. part i: theory. *Monthly Weather Review*, **136**, 3035–3049.
- Sacher, W. and P. Bartello, 2009: Sampling errors in ensemble Kalman filtering. part ii: application to a barotropic model. *Monthly Weather Review*, **137**, 1640–1654.
- Sakov, P. and P. R. Oke, 2007: Implications of the form of the ensemble transformation in the ensemble square root filters. *Monthly Weather Review*, **136**, 1042–1053.
- Tippett, M. K., J. L. Anderson, C. H. Bishop, T. M. Hamill, and J. S. Whitaker, 2003: Ensemble square root filters. *Monthly Weather Review*, **131**, 1485–1490.
- Whitaker, J. S. and T. M. Hamill, 2002: Ensemble data assimilation without perturbed observations. *Monthly Weather Review*, **130**, 1913–1924.

List of Tables

- 1 Mean square L_2 distance to the true conditional distribution in 1-D, with
 $t = 8, \mathbf{S} = 1.$ 37
- 2 Mean square L_2 distance to the true conditional distribution in 2-D, with
 $t = (10, 10), r = 0.05, c_2 = 1.$ 38

TABLE 1. Mean square L_2 distance to the true conditional distribution in 1-D, with $t = 8$, $\mathbf{S} = 1$.

		$\mathbf{R} = 0.25\mathbf{P}_r^f$	$\mathbf{R} = \mathbf{P}_r^f$	$\mathbf{R} = 4\mathbf{P}_r^f$
r=0.05	Sto.	0.369(0.195)	0.435(0.105)	0.112(0.037)
	Det.	0.409(0.405)	0.586(0.137)	0.150(0.051)
r=0.1	Sto.	0.255(0.112)	0.286(0.094)	0.099(0.029)
	Det.	0.356(0.350)	0.464(0.161)	0.150(0.054)
r=0.5	Sto.	0.117(0.034)	0.124(0.006)	0.055(0.005)
	Det.	0.240(0.156)	0.199(0.064)	0.050(0.018)

TABLE 2. Mean square L_2 distance to the true conditional distribution in 2-D, with $t = (10, 10)$, $r = 0.05$, $c_2 = 1$.

		$c_1 = 1/4$	$c_1 = 1$	$c_1 = 4$	$c_1 = 16$
Orient. 1	Sto.	.035(.040)	.041(.039)	.043(.031)	.040(.027)
	Det. ($\mathbf{U} = \mathbf{I}$)	.462(.114)	.183(.093)	.100(.071)	.065(.055)
	Det. (EAKF)	.454(.111)	.183(.094)	.105(.075)	.086(.058)
Orient. 2	Sto.	.066(.142)	.047(.071)	.049(.056)	.050(.051)
	Det. ($\mathbf{U} = \mathbf{I}$)	.492(.224)	.204(.119)	.114(.098)	.077(.079)
	Det. (EAKF)	.500(.207)	.208(.118)	.128(.101)	.103(.085)

List of Figures

- 1 The scatter plots of the previous updated ensemble (left) and the forecast ensemble (right) in the Lorenz 63 system (simulated using fourth order Runge-Kutta method with step size 0.05, propagated 4 steps). 40
- 2 The density plots for $F_{o,r}$ (solid); $F_{s,r}$ (dotted) and $F_{d,r}$ (dash-dotted). Parameters: $t = 8$, $\mathbf{S} = 1$, $\mathbf{R} = k\mathbf{P}_r^f$. $k = 0.25$ (top row); $k = 1$ (middle row); $k = 4$ (bottom row). $r = 0.5$ (left column); $r = 0.05$ (right column). 41
- 3 The conditional third moments. Horizontal coordinate: the observation \mathbf{y} ; vertical coordinate: $E_{F_{o,r}}\mathbf{x}^3$ (solid), $E_{\hat{F}_{s,r}}\mathbf{x}^3$ (dotted) and $E_{\hat{F}_{d,r}}\mathbf{x}^3$ (dash-dotted). Parameters: $t = 1$ (top row), $t = 10$ (second row), $t = 50$ (third row), $t = 100$ (bottom row); $\mathbf{R} = \mathbf{S} = 1$ (left column), $\mathbf{R} = 1, \mathbf{S} = 4$ (middle column), $\mathbf{R} = \mathbf{S} = 4$ (right column). 42
- 4 The contour of the densities of the two components in Orientation 2 (up to shift). Left: $N(0, \mathbf{P})$; right: $N(0, \mathbf{S})$. The levels are (from inner to outer): 0.2, 0.15, 0.1, 0.05. 43

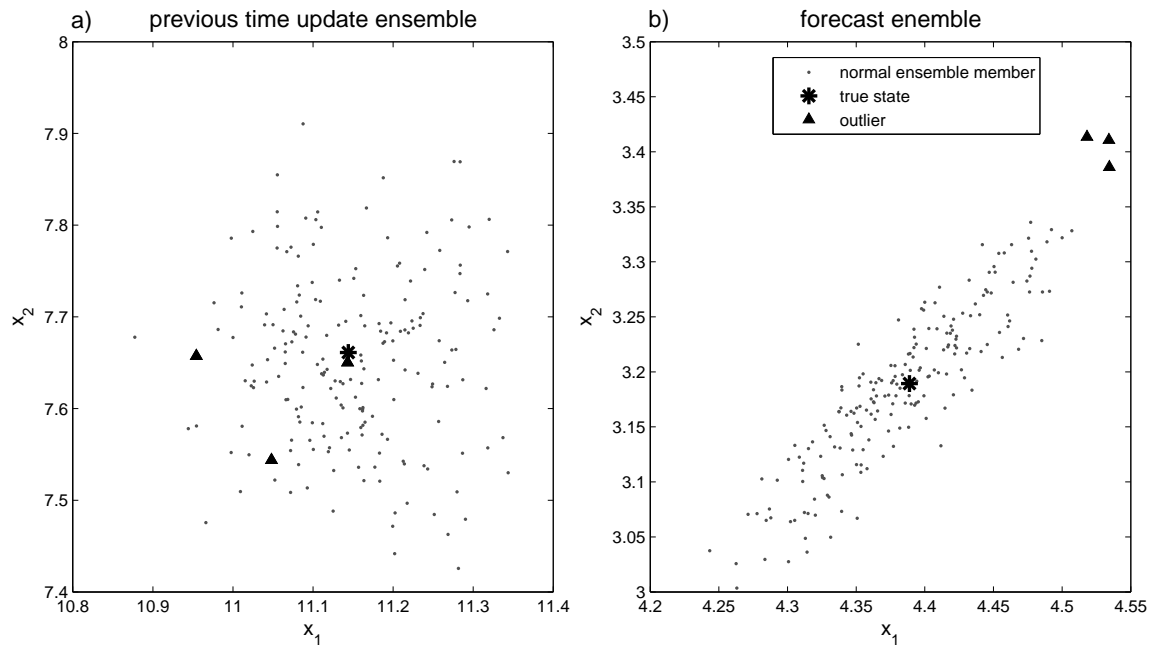


FIG. 1. The scatter plots of the previous updated ensemble (left) and the forecast ensemble (right) in the Lorenz 63 system (simulated using fourth order Runge-Kutta method with step size 0.05, propagated 4 steps).

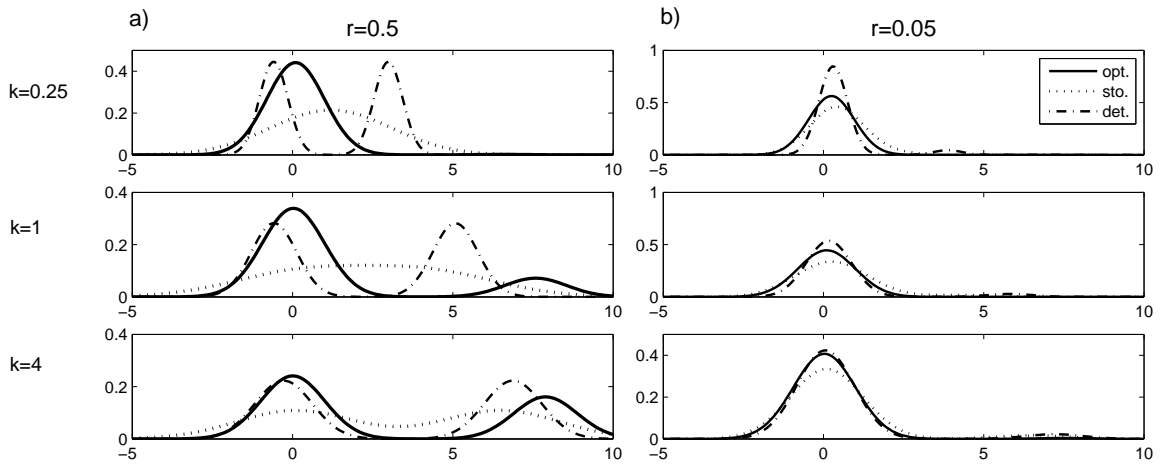


FIG. 2. The density plots for $F_{o,r}$ (solid); $F_{s,r}$ (dotted) and $F_{d,r}$ (dash-dotted). Parameters: $t = 8$, $\mathbf{S} = 1$, $\mathbf{R} = k\mathbf{P}_r^f$. $k = 0.25$ (top row); $k = 1$ (middle row); $k = 4$ (bottom row). $r = 0.5$ (left column); $r = 0.05$ (right column).

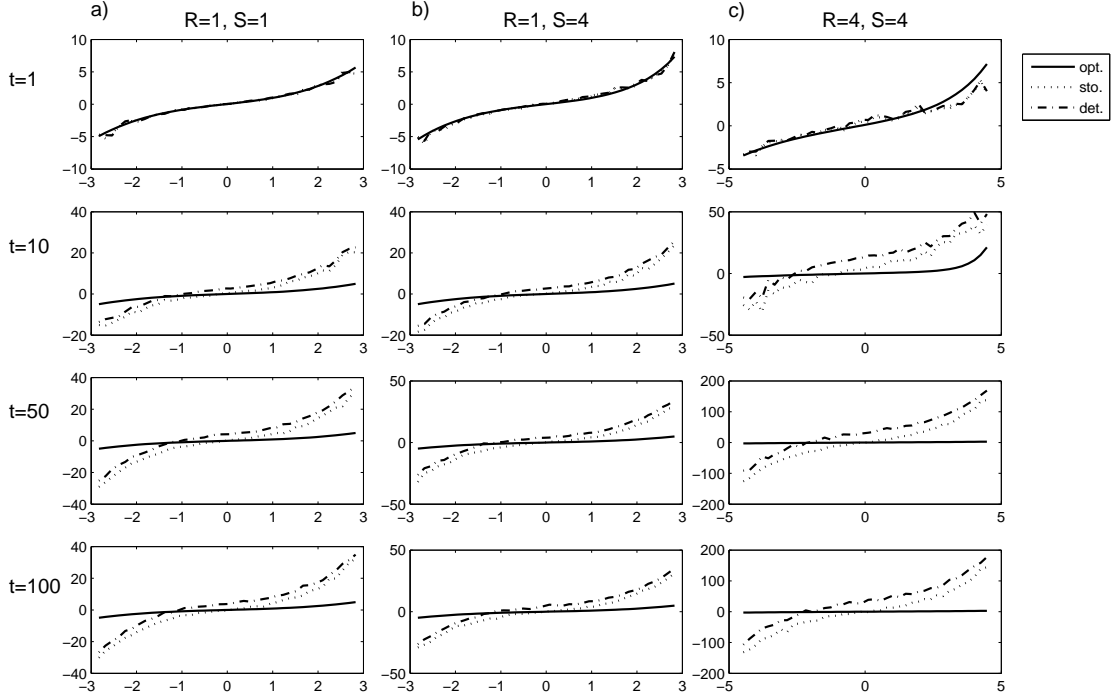


FIG. 3. The conditional third moments. Horizontal coordinate: the observation \mathbf{y} ; vertical coordinate: $E_{F_{0,r}} \mathbf{x}^3$ (solid), $E_{\hat{F}_{s,r}} \mathbf{x}^3$ (dotted) and $E_{\hat{F}_{d,r}} \mathbf{x}^3$ (dash-dotted). Parameters: $t = 1$ (top row), $t = 10$ (second row), $t = 50$ (third row), $t = 100$ (bottom row); $\mathbf{R} = \mathbf{S} = 1$ (left column), $\mathbf{R} = 1, \mathbf{S} = 4$ (middle column), $\mathbf{R} = \mathbf{S} = 4$ (right column).

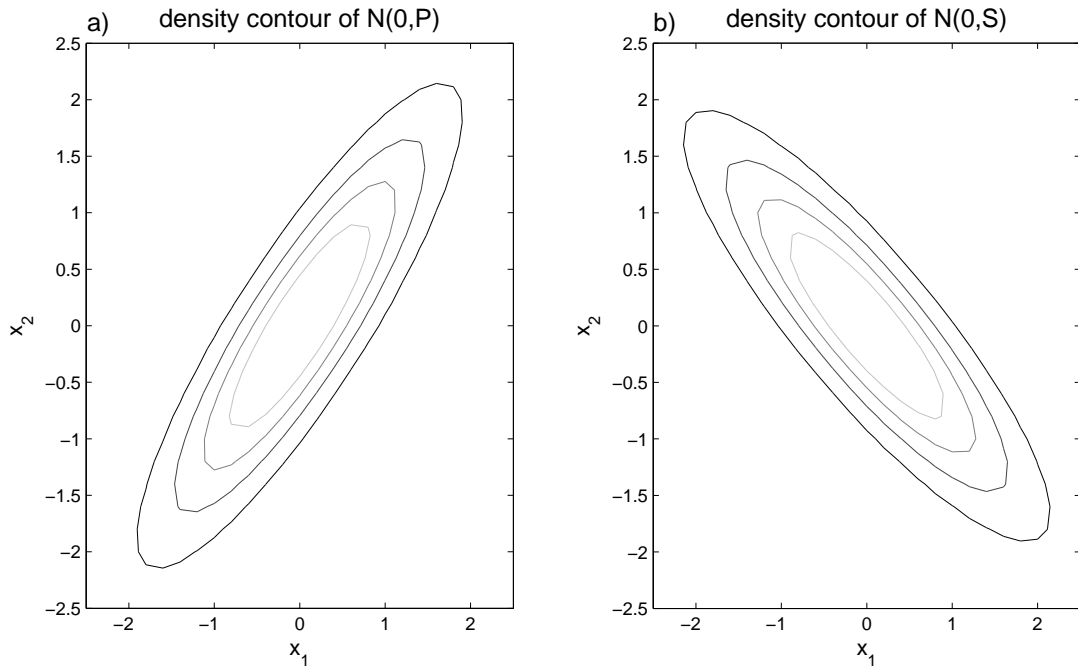


FIG. 4. The contour of the densities of the two components in Orientation 2 (up to shift). Left: $N(0, \mathbf{P})$; right: $N(0, \mathbf{S})$. The levels are (from inner to outer): 0.2, 0.15, 0.1, 0.05.