

Department of Mathematics and Statistics

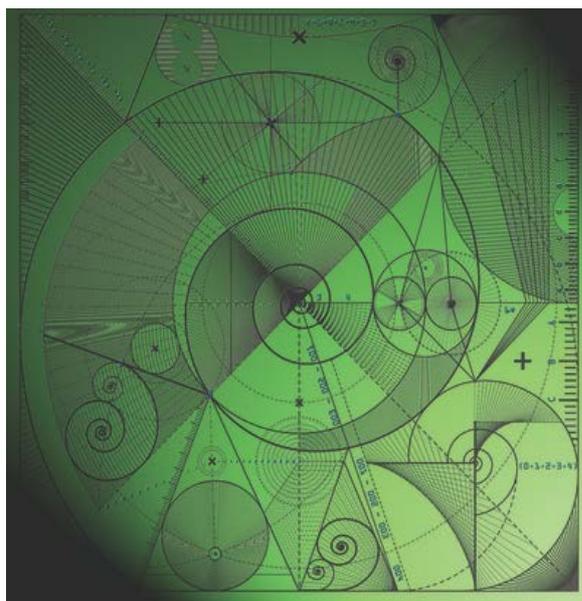
Preprint MPS-2014-08

3 March 2014

Conditioning and Preconditioning of the Optimal State Estimation Problem

by

S.A. Haben, A.S. Lawless and N.K. Nichols



Conditioning and Preconditioning of the Optimal State Estimation Problem

S. A. Haben, A. S. Lawless and N. K. Nichols

Department of Mathematics and Statistics

University of Reading

Abstract

Optimal state estimation is a method that requires minimising a weighted, nonlinear, least squares objective function in order to obtain the best estimate of the current state of a dynamical system. Often the minimisation is non-trivial due to the large scale of the problem, the relative sparsity of the observations and the nonlinearity of the objective function. To simplify the problem the solution is often found via a sequence of linearised objective functions. The condition number of the Hessian of the linearised problem is an important indicator of the convergence rate of the minimisation and the expected accuracy of the solution. In the standard formulation the convergence is slow, indicating an ill-conditioned objective function. A transformation to different variables is often used to ameliorate the conditioning of the Hessian by changing, or preconditioning, the Hessian. There is only sparse information in the literature for describing the causes of ill-conditioning of the optimal state estimation problem and explaining the effect of preconditioning on the condition number. This paper derives descriptive theoretical bounds on the condition number of both the unpreconditioned and preconditioned system in order to better understand the conditioning of the problem. We use these bounds to explain why the standard objective function is often ill-conditioned and why a standard preconditioning reduces the condition number. We also use the bounds on the preconditioned Hessian to understand the main factors that affect the conditioning of the system. We illustrate the results with simple numerical experiments.

Keywords Optimal state estimation, variational data assimilation, nonlinear least squares, condition number, preconditioning, correlation matrices, circulant matrices

1 Introduction

In dynamical systems, the aim of state estimation is to find the most likely current or future state of the system, given noisy, possibly indirect, observations. In many applications, such as numerical weather prediction (NWP), the number of observations is sparse relative to the dimension of the state space and so additional information, such as a prior estimate of the initial state of the system, is often required to guarantee a unique solution. The optimal state, called the ‘analysis’, minimises a weighted nonlinear least-squares objective function, measuring the distance between the state trajectory and the observations and between the initial state and the prior estimate, weighted by the covariance of the errors in the observations and the prior respectively. In the meteorology community this optimization problem is referred to as four-dimensional variational data assimilation or 4DVar [25]. The analysis is optimal in the sense

that, under certain assumptions, it provides the maximum a posteriori Bayesian estimate of the state of the system [22], [18]. Once the analysis is obtained the dynamical model is applied to predict future states.

The model and the observation operator, which maps model states to observations, are often nonlinear and therefore, to improve computational efficiency, the objective function is linearised around the current state estimate and a Gauss-Newton process is applied [10]. The linearised problem is solved by an inner iteration, called the *inner loop*, and then used to update the solution to the full nonlinear problem in an outer-loop. This inner-outer iteration process is repeated until the desired accuracy is obtained. The solution to the linearised problem is often found using gradient methods such as a conjugate gradient or quasi-Newton method [24]. However the accuracy of the solution and the rate of convergence of these iterative methods depends on the condition number, i.e. the ratio of the largest and smallest eigenvalues, and on the clustering of the eigenvalues, of the Hessian of the linearised objective function [9]. A large condition number often implies slow convergence and inaccurate results, whereas a small condition number means the gradient method converges quickly. Understanding the conditioning of the problem helps to identify sources of ill-conditioning and indicates how improvements in conditioning can be achieved.

The conditioning of the problem can be improved by transforming to new state variables. This is a form of *preconditioning*. A common type of preconditioning used for optimal state estimation in numerical weather prediction uses a square root of the prior error covariance matrix and has been shown to significantly reduce the condition number and time of convergence of the gradient methods [7], [12], [19]. This approach is equivalent to transforming to a set of variables with errors that are initially uncorrelated. Previous research has largely focused on the conditioning of the least-squares optimization (4DVar) problem in an experimental setting or in simplified systems [1], [26]. The conditioning of the state estimation problem without a prior estimate has also been considered in the literature on power systems and it has been shown that the observation errors and positions have an important effect on the conditioning of the problem [5] [23], [20]. Haben *et. al.* [12], [13], presented theoretical results for a one-dimensional periodic single-variable discretized system with observations made only at a one time point (a 3DVar problem). In these papers it was shown how dense accurate observations increase the condition number of the preconditioned Hessian. Additionally it was shown that the ill-conditioning of the covariance matrix of the prior errors was a major cause of the ill-conditioning of the unpreconditioned Hessian.

The focus of this paper is to prove theoretical bounds on the conditioning of the linearised objective function, with and without preconditioning, for a general state estimation or 4DVar data assimilation problem. In particular we establish the results of Haben *et. al.* [12], [13] as a special case of the results presented in here. We also present numerical results to illustrate the theory. We begin in Section 2 by providing background to the problem. In Section 3 we derive our main theoretical results for both the preconditioned and unpreconditioned case and then illustrate these results experimentally in Section 4 using simple numerical systems. Finally in Section 5 we summarise the results of this paper.

2 Optimal State Estimation

The aim of state estimation is to find the best estimate of the initial state of the system, $\mathbf{x}_0 \in \mathbb{R}^N$ (called *the analysis*), at time t_0 , given a prior estimate \mathbf{x}_0^b (called *the background*) and measurements $\mathbf{y}_i \in \mathbb{R}^{p_i}$ at time t_i ($i = 0, \dots, n$), taken within a time window $[t_0, t_n]$, and subject to the state space equations

$$\mathbf{x}_i = \mathcal{M}(t_i, t_0, \mathbf{x}_0), \quad (1)$$

$$\mathbf{y}_i = \mathcal{H}_i(\mathbf{x}_i) + \boldsymbol{\delta}_i, \quad (2)$$

for $i = 0, \dots, n$. The notation is as follows:

- the N model states at time t_i are denoted by the vector $\mathbf{x}_i \in \mathbb{R}^N$;
- the non-linear operator $\mathcal{M}(t_i, t_0, \cdot) : \mathbb{R}^N \rightarrow \mathbb{R}^N$, describes the evolution of the states from time t_0 to time t_i ;
- the non-linear operator $\mathcal{H}_i : \mathbb{R}^N \rightarrow \mathbb{R}^{p_i}$ relates the system states to the observations at time t_i and may include transformations and grid interpolations;
- the vector $\boldsymbol{\delta}_i \in \mathbb{R}^{p_i}$ describes the errors between the observed data and the model predictions of the data.

In many applications, such as numerical weather prediction, there are fewer total observations ($\sum_{i=0}^n p_i$) than the number of state variables (N). The prior estimate is therefore included with the aim of regularising the problem [22], [16]. We will assume the number of observations is less than the state dimension throughout this paper.

The errors ($\mathbf{x}_0 - \mathbf{x}_0^b$) in the background and $\boldsymbol{\delta}_i$ in the observations are assumed to be random with mean zero and symmetric positive-definite covariance matrices \mathbf{B} and \mathbf{R}_i , respectively. In addition, the observational errors are assumed to be temporally uncorrelated and uncorrelated with the errors in the background. The optimal estimate of the state of the system, \mathbf{x}_0^a , at time t_0 is found by minimising the following objective function with respect to \mathbf{x}_0

$$J(\mathbf{x}_0) = \frac{1}{2}(\mathbf{x}_0 - \mathbf{x}_0^b)^T \mathbf{B}^{-1}(\mathbf{x}_0 - \mathbf{x}_0^b) + \frac{1}{2} \sum_{i=0}^n (\mathcal{H}_i(\mathbf{x}_i) - \mathbf{y}_i)^T \mathbf{R}_i^{-1}(\mathcal{H}_i(\mathbf{x}_i) - \mathbf{y}_i), \quad (3)$$

subject to the model forecast equations (1)–(2). If the errors in the background and in the observations are assumed to have Gaussian probability distributions, then the solution to the optimization problem is equal to the *maximum a posteriori Bayesian estimate* of the system states at the initial time [18].

In this paper we concentrate on the case where the observation and model operators are linear, in which case we can write the objective function as

$$\tilde{J}(\mathbf{x}_0) = \frac{1}{2}(\mathbf{x}_0 - \mathbf{x}_0^b)^T \mathbf{B}^{-1}(\mathbf{x}_0 - \mathbf{x}_0^b) + \frac{1}{2}(\hat{\mathbf{H}}\mathbf{x}_0 - \hat{\mathbf{y}})^T \hat{\mathbf{R}}^{-1}(\hat{\mathbf{H}}\mathbf{x}_0 - \hat{\mathbf{y}}), \quad (4)$$

where we assume \mathbf{x}_i satisfies the linear model equations

$$\mathbf{x}_i = \mathbf{M}(t_i, t_0) \mathbf{x}_0 \equiv \mathbf{M}_i \mathbf{x}_0, \quad (5)$$

and

$$\begin{aligned}\hat{\mathbf{H}} &= [\mathbf{H}_0^T, (\mathbf{H}_1\mathbf{M}_1)^T, \dots, (\mathbf{H}_n\mathbf{M}_n)^T]^T, \\ \hat{\mathbf{y}} &= [\mathbf{y}_0^T, \mathbf{y}_1^T, \dots, \mathbf{y}_n^T]^T.\end{aligned}\tag{6}$$

The matrices \mathbf{M}_i and \mathbf{H}_i , $i = 0, \dots, n$, are linear evolution and observation operators respectively and $\hat{\mathbf{R}}$ is a block diagonal matrix with diagonal blocks equal to \mathbf{R}_i . The minimum of (4) can be found using iterative gradient methods such as the conjugate gradient method [9].

If the operators $\mathcal{M}(t_i, t_0, \cdot)$ and \mathcal{H}_i are nonlinear then the analysis is found by a double-loop algorithm, or inner-outer iteration procedure, where the problem is linearised around the current estimate of the model trajectory \mathbf{x}_i , $i = 0, \dots, n$, satisfying the nonlinear forecast model (1). An increment, $\delta\mathbf{x}_0$, to the current estimate of the analysis is then calculated by minimising the linearised objective function subject to the linearised model equations in an *inner-loop* [3]. The increment is then used to update the current estimate in an *outer-loop*. This double-loop process is equivalent to an approximate Gauss-Newton method [17], [10]. The inner-loop is often solved using a gradient method and is the main source of the computational cost of the minimisation. The linearised cost function has the same first-order Hessian as (3) and therefore the condition number analysis in this paper applies (to first order) whether the model and observation operators are linear or non-linear. Hence for the remainder of the paper we suppose that the model operators are linear.

2.1 Condition Number

A measure of the accuracy and efficiency with which the optimal state estimation problem can be solved is given by the *condition number* of the Hessian

$$\mathbf{S} = \mathbf{B}^{-1} + \hat{\mathbf{H}}^T \hat{\mathbf{R}}^{-1} \hat{\mathbf{H}},\tag{7}$$

of the linearized objective function (4) [9]. For any normal matrix, \mathbf{S} , the condition number in the ℓ_2 -norm is defined to be

$$\kappa(\mathbf{S}) = \|\mathbf{S}\|_2 \|\mathbf{S}^{-1}\|_2 \equiv \frac{|\lambda_{\max}(\mathbf{S})|}{|\lambda_{\min}(\mathbf{S})|},\tag{8}$$

where $\lambda_{\max}(\mathbf{S})$ and $\lambda_{\min}(\mathbf{S})$ denote the maximum and minimum (by moduli) eigenvalues of the matrix respectively. When \mathbf{S} is positive definite then all the eigenvalues are real and positive [9]. A Hessian with a large condition number is referred to as *ill-conditioned* and indicates that the solution to the linearized least-squares problem (4) is sensitive to relatively small perturbations in the data of the system. Additionally, ill-conditioning of the Hessian can have a detrimental effect on the convergence rates of the gradient methods used to solve the minimisation problem. For example, for the conjugate gradient method, the error in the computed solution after k iterations is bounded in proportion to $\left(\frac{(\sqrt{\kappa(\mathbf{S})} - 1)}{(\sqrt{\kappa(\mathbf{S})} + 1)}\right)^k$, which shows a potentially slow convergence for an ill-conditioned system [9]. Alternatively, a small condition number $\kappa(\mathbf{S}) \approx 1$ will lead to a rapid convergence of the conjugate gradient method.

2.2 Preconditioning

A common method for reducing the condition number of the objective function (4) is to use a linear transformation to change the variables [9]. The process of changing the condition number of the system is known as *preconditioning*. The condition number is minimised when the square root of the inverse of the Hessian is used as the change of variables transformation. However this is generally not practical due to the dimension of the problem and the complexity of the \mathbf{B} , $\hat{\mathbf{R}}$ and $\hat{\mathbf{H}}$ matrices. Instead, the symmetric square root of the covariance matrix of the errors in the prior estimates, $\mathbf{B}^{1/2}$, is often used [2], [19], [16]. The errors in the new variables $\mathbf{z}_0 = \mathbf{B}^{-1/2}\mathbf{x}_0$, are now uncorrelated, with unit variances, giving a prior error covariance matrix equal to the identity matrix.

In terms of the new variables, we aim to minimize the transformed objective function

$$\hat{J}(\mathbf{z}_0) = \frac{1}{2}(\mathbf{z}_0 - \mathbf{z}_0^b)^T(\mathbf{z}_0 - \mathbf{z}_0^b) + \frac{1}{2}(\hat{\mathbf{H}}\mathbf{B}^{1/2}\mathbf{z}_0 - \hat{\mathbf{y}})^T\hat{\mathbf{R}}^{-1}(\hat{\mathbf{H}}\mathbf{B}^{1/2}\mathbf{z}_0 - \hat{\mathbf{y}}), \quad (9)$$

with respect to \mathbf{z}_0 , where $\mathbf{z}_0^b = \mathbf{B}^{-1/2}\mathbf{x}_0^b$.

The effect of the variable transform is symmetrically to precondition the Hessian (7) with the square root of the error covariance matrix of the prior. The Hessian of the preconditioned objective function (9) is now given by

$$\hat{\mathbf{S}} = \mathbf{I}_N + \mathbf{B}^{1/2}\hat{\mathbf{H}}^T\hat{\mathbf{R}}^{-1}\hat{\mathbf{H}}\mathbf{B}^{1/2}, \quad (10)$$

where \mathbf{I}_m denotes the $m \times m$ identity matrix throughout the paper.

In general there are fewer observations than states of the system and therefore the matrix $\mathbf{B}^{1/2}\hat{\mathbf{H}}^T\hat{\mathbf{R}}^{-1}\hat{\mathbf{H}}\mathbf{B}^{1/2}$ is not of full rank, but is positive semi-definite. It follows that the smallest eigenvalue of (10) is unity and the condition number of the preconditioned Hessian is equal to its largest eigenvalue.

The aim of this paper is to prove theoretical bounds on the condition number of the unpreconditioned and preconditioned Hessians (7) and (10) respectively. The bounds enable the conditioning of the unpreconditioned and preconditioned Hessians to be compared and help to identify the main factors that affect the conditioning of the objective functions. This work extends the theoretical results presented in earlier work [12] [13], which examine the case of a discrete periodic single-variable system defined on a one dimensional grid with observations taken at only one time step. The proofs derived in this paper apply to more general cases where observations are taken over a time window and include, as special cases, proofs of the results summarized in previous papers. We illustrate the theory with numerical examples in a simplified system using common covariance structures and models.

3 Theory

3.1 Background Results

Before we can derive the algebraic bounds on the condition number of the Hessians, (7) and (10), of the unpreconditioned and preconditioned objective functions respectively, we require some basic results on circulant matrices.

For a periodic single-variable system discretized on a one-dimensional domain with equal spacing between grid points, many covariance and linear forecast models have a circulant structure. The eigenvalues for circulant matrices have a convenient form which makes them, and hence the condition number, simple to calculate. We exploit this useful property for producing our theoretical bounds. In more general cases, where the domain is not periodic, the autocovariance matrices will be Toeplitz instead. However when the dimension of the state space N is large these Toeplitz matrices and their properties can be approximated by circulant matrices [11], [15].

A circulant matrix has the form of a Toeplitz matrix where each row is a cyclic permutation of the previous row. Let $\mathbf{c} = [c_0, c_1, c_2, \dots, c_{N-1}]$ denote the top row of a $N \times N$ circulant matrix \mathbf{C} . Then the eigenvalues of \mathbf{C} are equal to the discrete Fourier transforms of the coefficients of the first row of the matrix [11] and can be written

$$\nu_m = \sum_{k=0}^{N-1} c_k e^{-2\pi i m k / N}. \quad (11)$$

The corresponding eigenvectors are given by the discrete exponential function,

$$\mathbf{v}_m = \frac{1}{\sqrt{N}} (1, e^{-2\pi i m / N}, \dots, e^{-2\pi i m (N-1) / N})^T. \quad (12)$$

Since circulant matrices are normal matrices we can explicitly calculate the condition number of \mathbf{C} from the definition (8) by taking the eigenvalues with the largest and smallest magnitude calculated using equation (11) [11].

A loose upper bound can be placed on the eigenvalues of a circulant matrix

$$|\nu_m| \leq \|\mathbf{C}\|_1 = \sum_{k=0}^{N-1} |c_k|. \quad (13)$$

If the circulant matrix has only positive coefficients then the largest eigenvalue is $\nu_0 = \sum_{k=0}^{N-1} c_k$. Assuming \mathbf{C} is a correlation matrix then $c_k \in [-1, 1]$, for $k = 0, \dots, N-1$, and the eigenvalue with the largest magnitude is at most N , the dimension of the state vector. However, unless all errors are strongly correlated ($|c_k| \approx 1$) this is likely to be a large overestimate.

3.2 Theory: Conditioning of the Hessian

In this section we consider the conditioning of the unpreconditioned Hessian (7). We make the following basic assumptions:

- A1. The covariance matrix for the errors in the prior estimate is of the form $\mathbf{B} = \sigma_b^2 \mathbf{C} \in \mathbb{R}^{N \times N}$, where \mathbf{C} is a symmetric, positive definite correlation matrix and $\sigma_b > 0$ is the standard deviation of the prior estimate errors.
- A2. The observation error covariance matrices at each time step are given by $\mathbf{R}_i = \sigma_o^2 \mathbf{I}_{p_i} \in \mathbb{R}^{p_i \times p_i}$, for $i = 0, \dots, n$, where $p_i \neq 0$, and $\sigma_o > 0$ is the standard deviation of the observation errors. Additionally $r \equiv \sum_{i=0}^n p_i < N$.

Assumption A1 implies that each component of the prior error has variance σ_b^2 . Assumption A2 assumes that observation errors are spatially and temporally uncorrelated with the same variances, σ_o^2 . These are reasonable assumptions where the states of the system represent values of the same variable at different spatial points and where the observations are all made with the same accuracy at every time. Where the observation errors satisfy assumption A2, $\hat{\mathbf{R}}$ will be block diagonal with blocks \mathbf{R}_i on the diagonal. In this case we can write $\hat{\mathbf{R}} = \text{diag}(\mathbf{R}_0, \mathbf{R}_1, \dots, \mathbf{R}_n)$. The final assumption ensures the total number of observations is less than the size of the state space we are estimating.

In the following theorem we derive the first theoretical bounds on the condition number of the unpreconditioned Hessian.

Theorem 1 *Let $\mathbf{B} \in \mathbb{R}^{N \times N}$ and $\hat{\mathbf{R}} = \text{diag}(\mathbf{R}_0, \mathbf{R}_1, \dots, \mathbf{R}_n) \in \mathbb{R}^{r \times r}$ be the prior error and observation error covariance matrices, respectively, satisfying assumptions A1 and A2. Additionally let $\hat{\mathbf{H}} \in \mathbb{R}^{r \times N}$ be the observation operator defined by (6). Then the following bounds hold on the condition number of the Hessian $\mathbf{S} = \mathbf{B}^{-1} + \hat{\mathbf{H}}^T \hat{\mathbf{R}}^{-1} \hat{\mathbf{H}}$:*

$$\frac{\kappa(\mathbf{C})}{\left(1 + \frac{\sigma_b^2}{\sigma_o^2} \lambda_{\max}(\mathbf{C}) \lambda_{\max}(\hat{\mathbf{H}}^T \hat{\mathbf{H}})\right)} \leq \kappa(\mathbf{S}) \leq \kappa(\mathbf{C}) \left(1 + \frac{\sigma_b^2}{\sigma_o^2} \lambda_{\min}(\mathbf{C}) \lambda_{\max}(\hat{\mathbf{H}}^T \hat{\mathbf{H}})\right), \quad (14)$$

where $\lambda_{\max}(\mathbf{A})$ and $\lambda_{\min}(\mathbf{A})$ represent the maximum and minimum (in modulus) eigenvalue, respectively, of the matrix \mathbf{A} .

Proof. To bound the condition number of the Hessian (7) we bound the maximum and minimum eigenvalues of \mathbf{S} . Suppose \mathbf{A}_1 and \mathbf{A}_2 are $N \times N$, symmetric, positive semi-definite matrices and that we label their eigenvalues such that $0 \leq \lambda_N(\mathbf{A}_i) \leq \lambda_{N-1}(\mathbf{A}_i) \leq \dots \leq \lambda_2(\mathbf{A}_i) \leq \lambda_1(\mathbf{A}_i)$ for $i = 1, 2$. Then from [9] the following bounds hold

$$\lambda_m(\mathbf{A}_1) + \lambda_N(\mathbf{A}_2) \leq \lambda_m(\mathbf{A}_1 + \mathbf{A}_2) \leq \lambda_m(\mathbf{A}_1) + \lambda_1(\mathbf{A}_2), \quad (15)$$

for $m = 1, \dots, N$. By assumption A1 and A2, $\mathbf{B} = \sigma_b^2 \mathbf{C}$ and $\hat{\mathbf{H}}^T \hat{\mathbf{R}}^{-1} \hat{\mathbf{H}} = \sigma_o^{-2} \hat{\mathbf{H}}^T \hat{\mathbf{H}}$ are both symmetric positive semi-definite matrices and therefore (15) applies with $\mathbf{A}_1 = \sigma_b^{-2} \mathbf{C}^{-1}$ and $\mathbf{A}_2 = \hat{\mathbf{H}}^T \hat{\mathbf{R}}^{-1} \hat{\mathbf{H}} = \sigma_o^{-2} \hat{\mathbf{H}}^T \hat{\mathbf{H}}$. In addition we have

$$\lambda_N(\hat{\mathbf{H}}^T \hat{\mathbf{R}}^{-1} \hat{\mathbf{H}}) = 0, \quad (16)$$

since $\hat{\mathbf{H}}^T \hat{\mathbf{R}}^{-1} \hat{\mathbf{H}}$ is not full rank by assumption A2. Hence the maximum and minimum eigenvalues of the Hessian satisfy

$$\sigma_b^{-2} \lambda_1(\mathbf{C}^{-1}) \leq \lambda_1(\mathbf{S}) \leq \sigma_b^{-2} \lambda_1(\mathbf{C}^{-1}) + \sigma_o^{-2} \lambda_1(\hat{\mathbf{H}}^T \hat{\mathbf{H}}), \quad (17)$$

and

$$\sigma_b^{-2} \lambda_N(\mathbf{C}^{-1}) \leq \lambda_N(\mathbf{S}) \leq \sigma_b^{-2} \lambda_N(\mathbf{C}^{-1}) + \sigma_o^{-2} \lambda_1(\hat{\mathbf{H}}^T \hat{\mathbf{H}}), \quad (18)$$

respectively. Combining (17) and (18) and the fact that $\lambda_N(\mathbf{A}) = \lambda_{\min}(\mathbf{A})$ and $\lambda_1(\mathbf{A}) = \lambda_{\max}(\mathbf{A})$ for any positive semi-definite symmetric matrix \mathbf{A} , we establish the following bounds on the condition number of the Hessian

$$\frac{\kappa(\mathbf{C})}{\left(1 + \frac{\sigma_b^2}{\sigma_o^2} \lambda_{\max}(\mathbf{C}) \lambda_{\max}(\hat{\mathbf{H}}^T \hat{\mathbf{H}})\right)} \leq \kappa(\mathbf{S}) \leq \kappa(\mathbf{C}) \left(1 + \frac{\sigma_b^2}{\sigma_o^2} \lambda_{\min}(\mathbf{C}) \lambda_{\max}(\hat{\mathbf{H}}^T \hat{\mathbf{H}})\right), \quad (19)$$

which completes the proof. \square

An alternative lower bound, which is easier to calculate explicitly, can be obtained using more restrictive assumptions

- A3. The observation operator is the same at each time step, that is, $\mathbf{H}_i = \mathbf{H} \in \mathbb{R}^{q \times N}$, where $p_i = q$, for $i = 0, \dots, n$, and all observations are direct observations of individual states.
- A4. The forecast model is assumed to be time invariant with $\mathbf{M}_i := \mathbf{M}^i$ for $i = 1, \dots, n$, for some circulant matrix $\mathbf{M} \in \mathbb{R}^{N \times N}$.
- A5. The symmetric positive-definite error covariance matrix $\mathbf{B} \in \mathbb{R}^{N \times N}$, and hence also its inverse, are circulant.

A consequence of assumption A3 is that $\mathbf{H}_i^T \mathbf{H}_i = \mathbf{H}^T \mathbf{H} \in \mathbb{R}^{N \times N}$ for $i = 0, \dots, n$ is a diagonal matrix with the k^{th} diagonal entry equal to one if the k^{th} position is observed or zero otherwise. Assumptions A4 and A5 mean that we can explicitly calculate the updated lower bound in the following theorem by using (11). We shall see in Section 4 that for a single periodic variable defined on a one-dimensional equally-spaced grid, the background covariance and forecast model matrices are circulant in many examples.

Under the additional assumptions (A3)-(A5) we can derive the following theoretical bounds on the condition number of the Hessian (7).

Theorem 2 *Let \mathbf{B} , $\hat{\mathbf{R}}$ and $\hat{\mathbf{H}}$ satisfy assumptions A1-A5 where $\hat{\mathbf{H}}$ is the generalised observation operator defined by (6). Then the following bounds hold on the condition number of the Hessian $\mathbf{S} = \mathbf{B}^{-1} + \hat{\mathbf{H}}^T \hat{\mathbf{R}}^{-1} \hat{\mathbf{H}}$:*

$$\left(\frac{1 + \frac{q}{N} \frac{\sigma_b^2}{\sigma_o^2} \lambda_{\min}(\mathbf{C}) \gamma_{\min}}{1 + \frac{q}{N} \frac{\sigma_b^2}{\sigma_o^2} \lambda_{\max}(\mathbf{C}) \gamma_{\max}} \right) \kappa(\mathbf{C}) \leq \kappa(\mathbf{S}) \leq \kappa(\mathbf{C}) \left(1 + \frac{\sigma_b^2}{\sigma_o^2} \lambda_{\min}(\mathbf{C}) \lambda_{\max}(\hat{\mathbf{H}}^T \hat{\mathbf{H}}) \right), \quad (20)$$

where $\lambda_{\max}(\mathbf{A})$ and $\lambda_{\min}(\mathbf{A})$ represent the maximum and minimum (in modulus) eigenvalues of a matrix $\mathbf{A} \in \mathbb{R}^{N \times N}$ respectively and we have $\gamma_{\max} := \sum_{j=0}^n |\lambda_{\max}(\mathbf{M})|^{2j}$ and $\gamma_{\min} := \sum_{j=0}^n |\lambda_{\min}(\mathbf{M})|^{2j}$.

Proof. From Theorem 1 the upper bound on the condition number is automatically obtained. The lower bound is achieved using the *Rayleigh quotient*, which for a Hermitian matrix $\mathbf{A} \in \mathbb{C}^{N \times N}$ and non-zero vector $\mathbf{v} \in \mathbb{C}^N$ is defined by

$$R_{\mathbf{A}}(\mathbf{v}) = \frac{\mathbf{v}^H \mathbf{A} \mathbf{v}}{\mathbf{v}^H \mathbf{v}}, \quad (21)$$

where the superscript H denotes complex conjugate transpose. The Rayleigh quotient is a bounded function satisfying [9]

$$\lambda_N(\mathbf{A}) \leq R_{\mathbf{A}}(\mathbf{v}) \leq \lambda_1(\mathbf{A}), \quad (22)$$

where \mathbf{A} has eigenvalues $\lambda_N(\mathbf{A}) \leq \lambda_{N-1}(\mathbf{A}) \leq \dots \leq \lambda_2(\mathbf{A}) \leq \lambda_1(\mathbf{A})$. Consider the Rayleigh quotient of $\hat{\mathbf{H}}^T \hat{\mathbf{H}}$ at an eigenvector \mathbf{v} of the symmetric, positive definite matrix \mathbf{B}^{-1} (which by assumption A5 is of the form (12)). By assumption A3 and A4, $\mathbf{H}_j^T \mathbf{H}_j = \mathbf{H}^T \mathbf{H}$ and $\mathbf{M}_j = \mathbf{M}^j$ for $j = 0, 1, \dots, n$ and therefore

$$\mathbf{v}^H (\hat{\mathbf{H}}^T \hat{\mathbf{H}}) \mathbf{v} = \mathbf{v}^H \left(\sum_{j=0}^n (\mathbf{M}^j)^T \mathbf{H}^T \mathbf{H} \mathbf{M}^j \right) \mathbf{v} = \sum_{j=0}^n (\lambda^H(\mathbf{M}))^j (\lambda(\mathbf{M}))^j \mathbf{v}^H \mathbf{H}^T \mathbf{H} \mathbf{v}. \quad (23)$$

The last equality holds because \mathbf{M} is circulant by assumption A4 and hence has the same eigenvectors as \mathbf{B}^{-1} , giving

$$\mathbf{M}^j \mathbf{v} = (\lambda(\mathbf{M}))^j \mathbf{v}, \quad (24)$$

where $\lambda(\mathbf{M})$ is the eigenvalue of \mathbf{M} corresponding to the eigenvector \mathbf{v} . Since the eigenvectors are of the form (12) and the observations are only of individual states, then we have

$$\mathbf{v}^H \mathbf{H}^T \mathbf{H} \mathbf{v} = \frac{q}{N}, \quad (25)$$

using assumption A3. Combining (25) and (23) we obtain

$$\mathbf{v}^H (\hat{\mathbf{H}}^T \hat{\mathbf{H}}) \mathbf{v} = \frac{q}{N} \sum_{j=0}^n |\lambda(\mathbf{M})|^{2j}. \quad (26)$$

We can then use (26) to put new bounds on the Hessian, \mathbf{S} , given by (7). Let \mathbf{v}_{\max} denote the eigenvector associated with $\lambda_{\max}(\mathbf{B}^{-1}) = \lambda_1(\mathbf{B}^{-1})$ and some eigenvalue $\lambda_\alpha(\mathbf{M})$ of \mathbf{M} . Applying the Rayleigh quotient to \mathbf{S} we obtain

$$\lambda_{\max}(\mathbf{S}) \geq \mathbf{v}_{\max}^H \mathbf{S} \mathbf{v}_{\max} = \mathbf{v}_{\max}^H \mathbf{B}^{-1} \mathbf{v}_{\max} + \sigma_o^{-2} \mathbf{v}_{\max}^H (\hat{\mathbf{H}}^T \hat{\mathbf{H}}) \mathbf{v}_{\max} \quad (27)$$

$$= \sigma_b^{-2} \lambda_{\max}(\mathbf{C}^{-1}) + \frac{q}{N} \sigma_o^{-2} \sum_{j=0}^n |\lambda_\alpha(\mathbf{M})|^{2j} \quad (28)$$

$$\geq \sigma_b^{-2} \lambda_{\max}(\mathbf{C}^{-1}) + \frac{q}{N} \sigma_o^{-2} \sum_{j=0}^n |\lambda_{\min}(\mathbf{M})|^{2j}, \quad (29)$$

where $\lambda_{\min}(\mathbf{M})$ is the minimum (in modulus) eigenvalue of \mathbf{M} . Similarly, let \mathbf{v}_{\min} be the eigenvector corresponding to $\lambda_{\min}(\mathbf{B}^{-1}) = \lambda_N(\mathbf{B}^{-1})$ and some eigenvalue $\lambda_\beta(\mathbf{M})$ of \mathbf{M} . Then the Rayleigh quotient can be used to show

$$\lambda_{\min}(\mathbf{S}) \leq \sigma_b^{-2} \lambda_{\min}(\mathbf{C}^{-1}) + \frac{q}{N} \sigma_o^{-2} \sum_{j=0}^n |\lambda_{\max}(\mathbf{M})|^{2j}, \quad (30)$$

where $\lambda_{\max}(\mathbf{M})$ is the maximum (in modulus) eigenvalue of \mathbf{M} . Now define

$$\gamma_{\max} = \sum_{j=0}^n |\lambda_{\max}(\mathbf{M})|^{2j}, \quad (31)$$

$$\gamma_{\min} = \sum_{j=0}^n |\lambda_{\min}(\mathbf{M})|^{2j}, \quad (32)$$

and combine the bounds (29) and (30) to give

$$\kappa(\mathbf{S}) \geq \left(\frac{1 + \frac{q}{N} \frac{\sigma_b^2}{\sigma_o^2} \lambda_{\min}(\mathbf{C}) \gamma_{\min}}{1 + \frac{q}{N} \frac{\sigma_b^2}{\sigma_o^2} \lambda_{\max}(\mathbf{C}) \gamma_{\max}} \right) \kappa(\mathbf{C}), \quad (33)$$

which completes the proof. \square

The theoretical bounds for the special case where observations are made at only one time step, presented in [12], follow automatically from Theorem 2, as shown in the following corollary.

Corollary 3 *Let $n = 0$, and let \mathbf{B} , $\hat{\mathbf{R}} \equiv \mathbf{R} = \sigma_o^2 \mathbf{I}_q$ and $\hat{\mathbf{H}} \equiv \mathbf{H}$ satisfy assumptions A1-A5. Then the following bounds hold on the condition number of $\mathbf{S} = \mathbf{B}^{-1} + \mathbf{H}^T \mathbf{R}^{-1} \mathbf{H}$:*

$$\left(\frac{1 + \frac{q}{N} \frac{\sigma_b^2}{\sigma_o^2} \lambda_{\min}(\mathbf{C})}{1 + \frac{q}{N} \frac{\sigma_b^2}{\sigma_o^2} \lambda_{\max}(\mathbf{C})} \right) \kappa(\mathbf{C}) \leq \kappa(\mathbf{S}) \leq \left(1 + \left(\frac{\sigma_b^2}{\sigma_o^2} \right) \lambda_{\min}(\mathbf{C}) \right) \kappa(\mathbf{C}), \quad (34)$$

where $\lambda_{\max}(\mathbf{C})$ and $\lambda_{\min}(\mathbf{C})$ are the maximum and minimum (in modulus) eigenvalues respectively of the matrix \mathbf{C} .

Proof. The proof follows directly from Theorem 2, since $\lambda_{\max}(\mathbf{H}^T \mathbf{H}) = 1$ and $\gamma_{\max} = 1 = \gamma_{\min}$ in the case $n = 0$. \square

The bounds presented in Theorem 1, Theorem 2 and Corollary 3 demonstrate the influence of the conditioning of the prior error covariances on the conditioning of the Hessian (7). In particular, the bounds provide us the following further details about the condition number of the Hessian:

1. If we vary the ratio σ_b^2/σ_o^2 in the bounds (14) and (20), while fixing all other variables, then as $\sigma_b^2/\sigma_o^2 \rightarrow 0$, both the lower bounds and upper bounds on the condition number converge to $\kappa(\mathbf{C})$. In this case the observations become much less accurate than the prior estimate or, conversely, the accuracy of the prior estimate becomes much greater than the accuracy of the observations. In either case the solution is dependent primarily on the prior estimate because the observations provide little or no constraint on the the solution to the problem. Hence the conditioning depends essentially on the conditioning of the prior error correlation matrix.
2. If the ratio $\sigma_b^2/\sigma_o^2 \rightarrow \infty$, while all other variables are fixed, then the upper bound in (14) and (20) grows linearly with σ_b^2/σ_o^2 . In this case the observations become much more accurate than the prior estimate or, conversely, the accuracy of the prior estimate becomes much worse than the accuracy of the observations. In the limit we are then trying to fit model trajectories to perfectly accurate observations, and the prior places no constraint on the problem. As the regularization provided by the prior reduces, we expect the state estimation (data assimilation) problem to become more ill-posed and harder to solve and, in general, we expect the condition number to increase.
3. If the prior errors become strongly positively correlated, the matrix \mathbf{C} becomes singular, since all its components converge to unity. In the limit $\lambda_{\max}(\mathbf{C})$ is bounded away from zero and $\lambda_{\min}(\mathbf{C}) \rightarrow 0$, and the upper and lower bounds in (14) and (20) converge to $\kappa(\mathbf{C}) \rightarrow \infty$.

4. In the special case where the prior errors are all uncorrelated and $\mathbf{C} = \mathbf{I}_N$, with \mathbf{I}_N being the $N \times N$ identity matrix, then all eigenvalues of \mathbf{C} are unity and the exact condition number $\kappa(\mathbf{S}) = 1 + \frac{\sigma_b^2}{\sigma_o^2} \lambda_{\max}(\hat{\mathbf{H}}^T \hat{\mathbf{H}})$, which is equal to the upper bounds in (14) and (20). In this case the upper bound on the conditioning of the Hessian is strict.

In conclusion, the conditioning of the state estimation problem (4) is strongly dependent on the conditioning of the prior error covariance matrix. With commonly arising prior error covariance matrices, it was shown in [14] that for large correlation length-scales, these matrices are very ill-conditioned and lead to a poorly conditioned Hessian of (4). This is consistent with previous results on variational data assimilation that suggest that the error covariances of the prior estimates are the cause of slow convergence in the minimization of the objective function [19]. In Section 4 we further illustrate the effect of an ill-conditioned prior error covariance matrix on the conditioning of the optimal state estimation problem using simplified numerical experiments.

3.3 Conditioning of the Preconditioned System

In this section we consider the effect of preconditioning the Hessian with the square root of the error covariance matrix of the prior estimate. The following theorem derives new theoretical bounds on the condition number of the preconditioned Hessian (10).

Theorem 4 *Let $\mathbf{B} = \sigma_b^2 \mathbf{C} \in \mathbb{R}^{N \times N}$ and $\hat{\mathbf{R}} = \text{diag}(\mathbf{R}_0, \mathbf{R}_1, \dots, \mathbf{R}_n) \in \mathbb{R}^{r \times r}$ be the prior and observation error covariance matrices, respectively, satisfying assumptions A1 and A2. Additionally let $\hat{\mathbf{H}} \in \mathbb{R}^{r \times N}$ be the observation operator defined by (6). Then the following bounds hold on the condition number of the preconditioned Hessian $\hat{\mathbf{S}} = \mathbf{I}_N + \mathbf{B}^{1/2} \hat{\mathbf{H}}^T \hat{\mathbf{R}}^{-1} \hat{\mathbf{H}} \mathbf{B}^{1/2}$:*

$$1 + \frac{1}{r} \frac{\sigma_b^2}{\sigma_o^2} \sum_{k,l=1}^r \{\hat{\mathbf{H}} \mathbf{C} \hat{\mathbf{H}}^T\}_{k,l} \leq \kappa(\hat{\mathbf{S}}) \leq 1 + \frac{\sigma_b^2}{\sigma_o^2} \|\hat{\mathbf{H}} \mathbf{C} \hat{\mathbf{H}}^T\|_{\infty}, \quad (35)$$

where $\{\mathbf{A}\}_{k,l}$ represents the $(k,l)^{\text{th}}$ entry of the matrix \mathbf{A} .

Proof. Since there are fewer observations than variables in the state space ($r < N$), the Hessian $\hat{\mathbf{S}}$ is just a low rank update of the identity matrix and its smallest eigenvalue is unity. The condition number of the Hessian is then equal to the largest eigenvalue of $\hat{\mathbf{S}}$. Let $\mathbf{E} = \hat{\mathbf{R}}^{-1/2} \hat{\mathbf{H}} \mathbf{B}^{1/2}$. The matrices $\mathbf{E}^T \mathbf{E} = \mathbf{B}^{1/2} \hat{\mathbf{H}}^T \hat{\mathbf{R}}^{-1} \hat{\mathbf{H}} \mathbf{B}^{1/2}$ and $\mathbf{E} \mathbf{E}^T = \hat{\mathbf{R}}^{-1/2} \hat{\mathbf{H}} \mathbf{B} \hat{\mathbf{H}}^T \hat{\mathbf{R}}^{-1/2}$ have the same non-zero eigenvalues and therefore the Hessian $\hat{\mathbf{S}}$ has the same non-unit eigenvalues as the symmetric, positive definite matrix

$$\mathbf{G} = \mathbf{I}_r + \hat{\mathbf{R}}^{-1/2} \hat{\mathbf{H}} \mathbf{B} \hat{\mathbf{H}}^T \hat{\mathbf{R}}^{-1/2} = \mathbf{I}_r + \frac{\sigma_b^2}{\sigma_o^2} \hat{\mathbf{H}} \mathbf{C} \hat{\mathbf{H}}^T. \quad (36)$$

For any ℓ_p -norm $\|\cdot\|$, $|\lambda_{\max}(\mathbf{A})| \leq \|\mathbf{A}\|$ where $\lambda_{\max}(\mathbf{A})$ is the maximum (in modulus) eigenvalue of matrix $\mathbf{A} \in \mathbb{R}^{r \times r}$. Therefore letting $\mathbf{A} = \hat{\mathbf{H}} \mathbf{C} \hat{\mathbf{H}}^T$ we obtain

$$\kappa(\hat{\mathbf{S}}) = \lambda_{\max}(\mathbf{G}) \leq 1 + \frac{\sigma_b^2}{\sigma_o^2} \|\hat{\mathbf{H}} \mathbf{C} \hat{\mathbf{H}}^T\|_{\infty}, \quad (37)$$

which establishes the upper bound.

The lower bound is established by applying the Rayleigh quotient of \mathbf{G} with the unit vector $\mathbf{y} = \frac{1}{\sqrt{r}}(1, 1, \dots, 1)^T \in \mathbb{R}^r$,

$$R_{\mathbf{G}}(\mathbf{y}) = \mathbf{y}^T \mathbf{G} \mathbf{y} = 1 + \frac{1}{r} \frac{\sigma_b^2}{\sigma_o^2} \sum_{k,l=1}^r \{\hat{\mathbf{H}}\mathbf{C}\hat{\mathbf{H}}^T\}_{k,l}. \quad (38)$$

Since $\lambda_{\max}(\mathbf{G}) \geq R_{\mathbf{G}}(\mathbf{y})$ for any $\mathbf{y} \in \mathbb{R}^{r \times r}$, this completes the proof. \square

The bounds on the condition number of the preconditioned Hessian for the special case of observations at only one time step, derived in [13], can be found in the following corollary to Theorem 4.

Corollary 5 *Let $n = 0$, and let \mathbf{B} , $\hat{\mathbf{R}} \equiv \mathbf{R} = \sigma_o^2 \mathbf{I}_q$ and $\hat{\mathbf{H}} \equiv \mathbf{H}$ satisfy assumptions A1-A5. Then the following bounds on the condition number of $\hat{\mathbf{S}} = \mathbf{I}_N + \mathbf{B}^{1/2} \mathbf{H}^T \mathbf{R}^{-1} \mathbf{H} \mathbf{B}^{1/2}$ hold*

$$1 + \frac{1}{q} \frac{\sigma_b^2}{\sigma_o^2} \sum_{i,j \in K} \{\mathbf{C}\}_{i,j} \leq \kappa(\hat{\mathbf{S}}) \leq 1 + \frac{\sigma_b^2}{\sigma_o^2} \|\mathbf{H}\mathbf{C}\mathbf{H}^T\|_{\infty}, \quad (39)$$

where K are indices of the state variables that are observed.

Proof. From Theorem 4 with $n = 0$ we obtain the bounds

$$1 + \frac{1}{q} \frac{\sigma_b^2}{\sigma_o^2} \sum_{k,l=1}^q \{\mathbf{H}\mathbf{C}\mathbf{H}^T\}_{k,l} \leq \kappa(\hat{\mathbf{S}}) \leq 1 + \frac{\sigma_b^2}{\sigma_o^2} \|\mathbf{H}\mathbf{C}\mathbf{H}^T\|_{\infty}. \quad (40)$$

Since $\mathbf{H}\mathbf{C}\mathbf{H}^T$ is simply the matrix \mathbf{C} with rows and columns removed at the unobserved positions it follows that

$$\sum_{k,l=1}^q \{\mathbf{H}\mathbf{C}\mathbf{H}^T\}_{k,l} = \sum_{i,j \in K} \{\mathbf{C}\}_{i,j}, \quad (41)$$

where K are indices of the state variables that are observed. \square

Before discussing the implications of the bounds we first note that the matrix $\hat{\mathbf{H}}\mathbf{C}\hat{\mathbf{H}}^T$ which appears in the upper and lower bounds (35) can also be written in the form $\hat{\mathbf{H}}\tilde{\mathbf{C}}\hat{\mathbf{H}}^T = \sigma_b^{-2} \tilde{\mathbf{H}}\tilde{\mathbf{B}}\tilde{\mathbf{H}}^T$, where $\tilde{\mathbf{H}}$ is the block diagonal matrix consisting of $n + 1$ blocks equal to \mathbf{H}_i , $i = 0, \dots, n$, and $\tilde{\mathbf{B}} = \sigma_b^2 \tilde{\mathbf{C}}$ is the four-dimensional error covariance matrix associated with the background state vector $(\mathbf{x}_0^{bT}, \mathbf{x}_1^{bT}, \dots, \mathbf{x}_n^{bT})^T$. Here \mathbf{x}_i^b denotes the state vector at time t_i , $i = 1, \dots, n$ evolved from the prior state estimate, \mathbf{x}_0^b , using the dynamical model (5) [14]. Since $\hat{\mathbf{H}}\tilde{\mathbf{C}}\hat{\mathbf{H}}^T$ is simply $\tilde{\mathbf{C}}$ with rows and columns deleted at positions that are unobserved, we refer to this as the *reduced* error covariance matrix.

The reduced error covariance matrix $\tilde{\mathbf{H}}\tilde{\mathbf{C}}\tilde{\mathbf{H}}^T$ plays a key role in the condition number of the Hessian (10). In particular, the lower bound in (35) is linearly related to the *average* row sum $\frac{1}{r} \sum_{k,l=1}^r \{\tilde{\mathbf{H}}\tilde{\mathbf{C}}\tilde{\mathbf{H}}^T\}_{k,l}$ whereas the upper bound is related to the absolute maximum row sum $\|\tilde{\mathbf{H}}\tilde{\mathbf{C}}\tilde{\mathbf{H}}^T\|_{\infty}$. In fact the lower and upper bounds are identical if all entries of $\tilde{\mathbf{H}}\tilde{\mathbf{C}}\tilde{\mathbf{H}}^T$ are positive and its row sums are identical. The dependence on the reduced error covariance matrix implies further details about the condition number of the preconditioned Hessian.

1. The number and positions of the observations are important to the conditioning of the preconditioned problem. In particular, if we assume that the correlations in the prior error covariance matrix decrease with increased distance between grid points and also that the linear model \mathbf{M} acts to ensure that the coefficients of the correlation matrix \mathbf{C} remain positive and decrease monotonically with distance, then increasing the distance between observations will imply smaller entries in the reduced error covariance matrix and thus smaller sums in the upper and lower bounds in (35) and, potentially, a smaller condition number. The assumptions apply, for instance, in the case where the model is an advection equation and the prior error covariance has a Gaussian or SOAR structure [14].
2. Additionally, under the same assumptions, if we have fewer observations at fewer time steps, then there will be fewer entries in the reduced error covariance matrix, implying smaller sums in the bounds and hence a smaller condition number of the Hessian (10).
3. Finally, it follows from the dependence of the bounds (35) on the ratio σ_b^2/σ_o^2 that the accuracy of the observations is also important to the conditioning of the problem. In particular, increasing the accuracy of the observations, where $\sigma_o^2 \rightarrow 0$ while the other variables remain fixed, implies an increase in the bounds and a potential increase in the conditioning of the Hessian. In the limit, the model trajectories must fit exactly to the data, as in the unpreconditioned case, and the problem becomes much harder to solve and hence more ill-posed.

The bounds (35) and (39) are quite general and do not require the more restrictive assumptions A3-A5 used in the unpreconditioned case. Additionally the bounds do not depend on the condition number of the background error covariance matrix but simply on a summation of the coefficients of a four-dimensional background error covariance matrix. In [12] it was shown that, in the case where observations are only made at a single time step, preconditioning brought a dramatic reduction in the condition number of the Hessian compared to the unpreconditioned case. Contrary to intuition, however, the bounds also show that in the preconditioned case, as well as in the unpreconditioned case, increasing the accuracy and density of observations is likely to make the conditioning of the problem increase and the estimation problem harder to solve accurately.

4 Numerical Experiments

In this section we illustrate the effect of varying different parameters and properties of the state estimation problem on the condition number of the unpreconditioned (7) and preconditioned (10) Hessians. We apply the theoretical bounds derived in Section 3.2 and 3.3, respectively, to explain these effects. Throughout this Section we consider a dynamical system where the state vector consists of a single periodic variable discretized at equally spaced grid points on a one-dimensional domain. As shown in Section 3.2 the prior error covariance plays a influential role in the conditioning of the preconditioned and unpreconditioned Hessian and so we first introduce and describe some of the properties of a common prior error covariance matrix.

4.1 Condition Number of Error Covariance Matrices

In this section we assume that the prior error covariance matrix is of the form $\mathbf{B} = \sigma_b^2 \mathbf{C} \in \mathbb{R}^{N \times N}$ where \mathbf{C} denotes the error correlation matrix and σ_b^2 is the error variance. By definition (8), the condition number $\kappa(\mathbf{B}) = \kappa(\mathbf{C})$. We use the second-order auto-regressive correlation (SOAR) function [4], defined by

$$\rho_S(r) = \left(1 + \frac{|d|}{L}\right) \exp\left(-\frac{|d|}{L}\right), \quad (42)$$

to model the correlation structure, where $L > 0$ is the correlation length-scale and $0 \leq d \in \mathbb{R}$ is the distance between two points on the real line. The SOAR function is commonly used to define correlations in meteorological applications [4]. For a periodic variable we identify the values of the variable at two points $-D$ and D . However, the function (42), which defines a valid correlation function on the real line, may no longer define valid correlation models on the finite interval, since the corresponding Fourier transforms are not necessarily positive [29], [8], [28]. We transform to a valid correlation model on the circle by replacing the distance along the great circle by the chordal distance

$$d = 2a \sin(\theta/2), \quad (43)$$

where θ is the angle between two points on the circle and a is the radius. This guarantees that the corresponding correlation matrix is positive definite [30, Sec. 22.5]. Applying the transform (43) to the SOAR correlation function and sampling at evenly spaced points on the circle s_i , $i = 1, \dots, N$, produces the SOAR correlation matrix \mathbf{C}_S on the circle with elements given by

$$(\mathbf{C}_S)_{i,j} = \left(1 + \frac{|2a \sin(\theta_{i,j}/2)|}{L}\right) \exp\left(-\frac{|2a \sin(\theta_{i,j}/2)|}{L}\right) \quad (44)$$

where $i, j = 1, \dots, N$ and $\theta_{i,j}$ is the angle between the points s_i and s_j on the circle. We note that the resultant correlation matrix is circulant and therefore has eigenvalues given by (11).

Length-scale	0.05	0.1	0.15	0.2	0.25	0.3	0.35
Condition Number	5.96	58.1	265	807	1963	3978	7328

Table 1: The condition number of the SOAR correlation matrix as a function of different correlation length-scales.

Table 1 shows the condition number of $\mathbf{C} = \mathbf{C}_S$, for different length-scales, L , where the correlation function is sampled at $N = 500$ equally spaced grid points on the interval $[-25, 25]$. The table shows that the condition number of the correlation matrix increases as a function of the correlation length-scale. As shown there is a large increase in the condition number as the length-scale increases. An increase in the length-scale from 0.1 to 0.2 causes an increase in the condition number by about 750 whereas the same increase in length-scale from 0.2 to 0.3 causes an increase by over 3000 in the condition number. Similar results also hold for other common auto-correlation matrices, see [14, Chap. 5]. Since all coefficients of the SOAR correlation

matrix are positive, we find from (11) that the largest eigenvalue satisfies

$$\lambda_{\max}(\mathbf{C}) = \|\mathbf{C}\|_{\infty} = \sum_{k=0}^{N-1} c_k, \quad (45)$$

and therefore increases slowly as a function of L and is bounded by $N = 500$ since $|c_k| \leq 1$. It is the decrease in the smallest eigenvalue that causes the increase in condition number of the correlation matrix [14, Chap. 5]. For the remainder of this paper we define the background correlation matrix using the SOAR correlation matrix (44).

4.2 Numerical Example: Advection Model

To compare the conditioning of the unpreconditioned and preconditioned Hessians (7) and (10), respectively, we assume a simple linear advection equation for our dynamical forecast model (5) throughout the rest of this section. We discretise using the upwind scheme described at the k^{th} grid point at time t_{j+1} by [21, Chap. 4]

$$U_k^{j+1} = U_k^j - c \frac{\Delta t}{\Delta x} (U_k^j - U_{k-1}^j) = U_k^j - \nu (U_k^j - U_{k-1}^j), \quad (46)$$

where $c = 0.3$ is the speed of advection and $U_0 = U_N$. We assume there are $N = 500$ grid points with a fixed spatial spacing of $\Delta x = 0.1$ and that the time step is $\Delta t = 0.1$. (Here $\nu = c \frac{\Delta t}{\Delta x} \in (0, 1)$ so the finite difference equation satisfies the CFL condition and is therefore convergent [21]). In matrix form (46) can be written $\mathbf{U}^{j+1} = \mathbf{M}\mathbf{U}^j$ where $\mathbf{U}^j = (U_1^j, \dots, U_N^j)^T$. The linear forecast model matrix \mathbf{M} is a circulant matrix with top row $(1 - \nu, -\nu, 0, \dots, 0, -\nu)$ and satisfies $\mathbf{M}_j = \mathbf{M}^j$. We observe at the same randomly spatially distributed grid points at three different time steps $t_0 = 0, t_1 = 3\Delta t$ and $t_2 = 6\Delta t$, giving 60 observations in total. Finally, we fix the background error variance to be $\sigma_b^2 = 1$ and the observation error variances to be $\sigma_o^2 = 1$. With these criteria the corresponding Hessian satisfies the assumptions A1-A5 and therefore the bounds (20) and (35) derived in Theorems 2 and 4, respectively, hold.

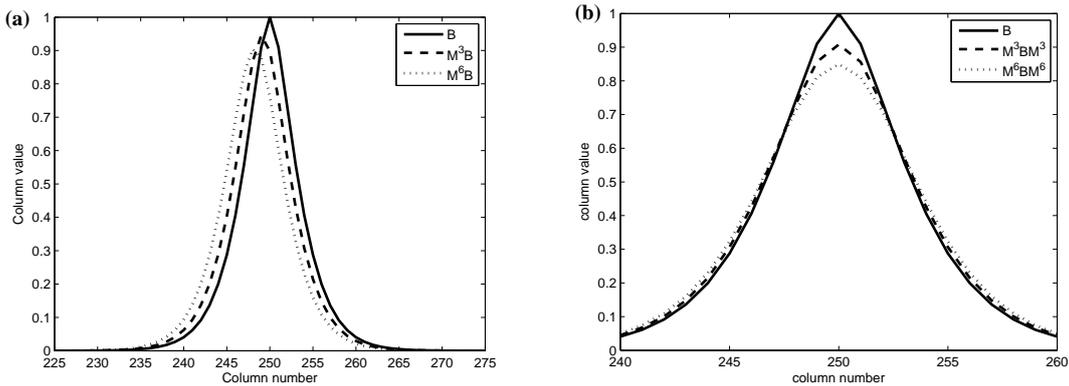


Figure 1: The coefficients of the 250th row of the off-diagonal (a) and diagonal blocks (b) of the 4D-prior error covariance matrix using the SOAR correlation matrix.

Since \mathbf{C} and \mathbf{M} are circulant, then each block of the 4D-prior error covariance matrix $\tilde{\mathbf{C}}$ is circulant. By plotting rows of each block of $\tilde{\mathbf{C}}$ we notice all coefficients are less than or equal to one. For example, Figure 1 shows the 250th row of the matrix diagonal and off-diagonal blocks $\mathbf{M}^j\mathbf{C}(\mathbf{M}^k)^T$ and $\mathbf{M}^j\mathbf{C}$, respectively, for $j, k = 0, 3, 6$ in the case where the length-scale of the correlations is given by $L = 0.2$. From Theorem 4 the maximum size of the condition number of the preconditioned Hessian is then given by

$$\kappa(\hat{\mathbf{S}}) \leq 1 + \frac{\sigma_b^2}{\sigma_o^2} \|\hat{\mathbf{H}}\hat{\mathbf{C}}\hat{\mathbf{H}}^T\|_\infty \leq 61. \quad (47)$$

If the conditioning of the unpreconditioned system is driven by the background error covariance matrices, as indicated by the bounds derived in Theorem 1 and 2, then we expect the preconditioning to significantly reduce the condition number of the problem in the cases where the length-scales of the prior error correlations are sufficiently large. Figure 2 compares the actual condition number and the theoretical bounds on the conditioning of the (a) unpreconditioned and (b) preconditioned Hessians as a function of length-scale using the SOAR correlation matrix. For this experiment the conditioning of the unpreconditioned system follows the upper bound, which demonstrates that the bound is strict. From these results it can be seen that, even for short length-scales, the preconditioning improves the conditioning of the Hessian by orders of magnitude. Similar results hold for Hessians constructed using other common correlation matrices to define the prior error covariance matrix (see [14, Chap. 7] for more details.).

A comparison of the magnitude of the condition number of the unpreconditioned Hessian and the condition number of the corresponding prior error covariance matrix with the same length-scale, given in Table 1, reveals that the conditioning of the Hessian is closely coupled to the conditioning of the prior error covariance matrix, as predicted by the theoretical bounds found in Theorem 2. For example, at length-scale $L = 0.25 = 2.5\Delta x$ the condition number of the Hessian is approximately 1900 whereas the condition number of the SOAR correlation matrix for the same length-scale as shown in Table 1 is 1963.

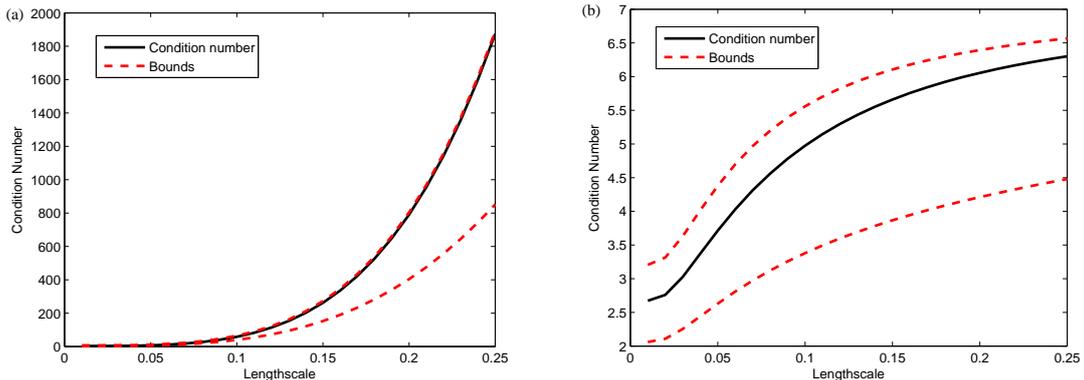


Figure 2: Condition number of the Hessian (solid line) together with the bounds (dashed) as a function of length-scale for SOAR correlation matrix in the (a) unpreconditioned and (b) preconditioned case.

For the preconditioned Hessian, the condition number is much smaller than the absolute upper bound predicted by (47) and is much better conditioned than the unpreconditioned Hessian. For instance, at length-scale $L = 0.25$ the condition number of the unpreconditioned Hessian is approximately 1900, whereas for the preconditioned Hessian it is around 6. The conditioning of the preconditioned system increases as the length-scale increases, which can be explained by the increase in the bounds (35). The larger length-scale increases the coefficients of the matrix $\tilde{\mathbf{C}}$ and therefore the size of the row sums of the coefficients of $\tilde{\mathbf{H}}\tilde{\mathbf{C}}\tilde{\mathbf{H}}^T = \hat{\mathbf{H}}\hat{\mathbf{C}}\hat{\mathbf{H}}^T$ in the upper and lower bounds (35).

4.3 The Effect of Observations on the Conditioning of the Preconditioned Hessian

We now consider the conditioning of the preconditioned Hessian for the numerical advection forecast model in more detail. The bounds for the preconditioned Hessian (35) and (39) identify the accuracy and positioning of observations as important to the conditioning of the preconditioned objective function.

Assuming the same data as for the experiment shown in Figure 2, we consider the effect of changing the observation accuracy on the condition number of the Hessian. We use the SOAR correlation matrix and fix the correlation length-scale to $L = 0.2$, but vary the observation variance. Table 2 shows the effect of changing the observation accuracy on the condition number of the preconditioned Hessian. As demonstrated in section 3.3, the bounds (35) are linearly related to the inverse of the observation variance and hence we expect the condition number of the Hessian to increase as the observation variance decreases and the accuracy of the observations increases. This is confirmed by the results of the numerical experiment, as seen in Table 2. For instance, a doubling in the accuracy of the observations from a variance of 0.1 to 0.05 roughly doubles the condition number of the Hessian from 51.55 to 102.11. Similar results also hold where other common prior error covariance matrices and observation locations are used (see [14]).

Obs Variance	0.01	0.05	0.10	0.50	1.00	2.00	5.00	10.00
Condition Number	506.53	102.11	51.55	11.11	6.06	3.53	2.01	1.51

Table 2: The condition number of the preconditioned Hessian as a function of the observation error variance using SOAR correlation matrices.

We now consider the condition number of the preconditioned Hessian as a function of the separation of the observations. From the definition of the correlation matrix (44) the coefficients in each block of $\{\tilde{\mathbf{C}}\}_{i,j}$ monotonically decrease as the distance $|i - j|$ increases, as shown in Figure 1. The upper and lower bounds on the Hessian (10) depend on sums of the elements of the matrix $\tilde{\mathbf{H}}\tilde{\mathbf{C}}\tilde{\mathbf{H}}^T$, which is viewed as a ‘reduced’ covariance matrix. The reduced matrix is simply the 4D covariance matrix $\tilde{\mathbf{C}}$ with all non-observed rows and columns deleted. As the separation of the observations increases, the elements of the reduced matrix become smaller in magnitude due to the decrease in the coefficients (or covariance) with distance. We therefore

expect the conditioning of the problem to decrease as the separation of the observations increases or the density decreases. We illustrate this with our numerical model.

We fix the observation error variances to $\sigma_o^2 = 1$ and assume that $q = 20$ observations are made at grid points at each of the time steps $t_0 = 0, t_1 = 3\Delta t$ and $t_2 = 6\Delta t$ with uniform spacing between adjacent observations. We consider the condition number as the uniform spacing is increased. Table 3 shows the results of the experiment. As expected from the theoretical bounds (14), increasing the spacing between the observations reduces the size of the condition number of the Hessian. Since the coefficients of the covariance matrix $\tilde{\mathbf{C}}$, given by (44), decrease with an increase in the distance between sampling points, the condition number of the preconditioned Hessian becomes smaller with larger distances and decreased density of observations. Additionally, as predicted, the condition number is larger for larger length-scales corresponding to the increase in the size of the coefficients of $\tilde{\mathbf{C}}$. Similar results hold for the preconditioned Hessians using other common prior error covariance matrices (See [14, Chap. 7]).

Spacing	1	2	3	4	5	6	7	8	9	10
Condition Number ($L = 0.2$)	22.0	12.5	8.9	6.9	5.8	5.1	4.6	4.3	4.0	3.9
Condition Number ($L = 0.3$)	29.6	17.6	12.5	9.8	8.1	7.0	6.2	5.6	5.1	4.8
Condition Number ($L = 0.5$)	39.8	26.3	19.3	15.2	12.6	10.8	9.4	8.4	7.6	7.0

Table 3: The condition number of the preconditioned Hessian as a function of the number of spaces between observations for different correlation length-scales L .

The results of this section indicate that less accurate and less dense observations reduce the conditioning of the preconditioned Hessian and hence may increase the rate of convergence of the iterative solver used to find the optimal state estimate. These results appear to be counter-intuitive as it would indicate that better observational data leads to a more inaccurate (numerical) solution. This may be explained by the fact that highly accurate, dense observations put tighter restrictions on the optimisation problem and so the problem becomes more difficult to solve whilst accurately satisfying the constraints. In practice there must a balance between satisfactorily solving the physical problem (by collecting many accurate data) and the numerical problem (as typified by the condition number).

5 Discussion

In state estimation the conditioning of the objective function plays an important role in determining the accuracy of the numerical solution and the speed of convergence of the iterative methods used to solve the problem. If the Hessian of the objective function has a large condition number, then we say the problem is ill-conditioned and the iterative method may be slow to converge. The problem can be reformulated with a variable transform which *preconditions* the problem to one with a smaller condition number.

In this paper we have examined the conditioning of an optimal state estimation (data assimilation) problem and shown how preconditioning with a standard change of variables affects this conditioning. The main results presented in this work are new theoretical bounds on the

condition number of the Hessian of the objective function in both the unpreconditioned and preconditioned forms. The bounds derived identify the main sources of ill-conditioning in both systems and explain how preconditioning can improve the conditioning of the problem. In particular, we found that the condition number of the unpreconditioned Hessian is proportional to the condition number of the prior error covariance matrix. Hence an ill-conditioned prior error covariance matrix can produce an ill-conditioned Hessian. The bounds on the preconditioned system showed that preconditioning using the prior error covariance matrix can produce a significant reduction in the condition number of the Hessian. Additionally, the distribution, quantity and accuracy of the observations play key roles in the conditioning of the preconditioned Hessian, with more accurate and dense observations creating a more ill-conditioned problem.

We presented results from numerical experiments in order to demonstrate the effect of the various factors on the condition number of the Hessians, as indicated by the bounds. We presented the SOAR covariance matrix, which is commonly used in variational data assimilation, and showed that the conditioning of this matrix becomes very ill-conditioned for only relatively small increases in correlation length-scale. We then demonstrated that this prior error covariance matrix resulted in the ill-conditioning of the unpreconditioned Hessian and that preconditioning dramatically reduced the conditioning, as predicted by the theoretical bounds. We also illustrated the reduction in the conditioning of the preconditioned system as we increased the separation between observations and reduced the accuracy of the observations, as expected from our theoretical results. We remark that the conclusions derived from the theory presented here have also been found to hold for experimental data from the high-dimensional, multi-variable Met Office Numerical Weather Prediction data assimilation system [13].

A simple, natural extension to this problem would be to consider more general observation operators which incorporate interpolation and to introduce correlations into the observation errors. Very recently, extra preconditioning, in addition to the variable transform via the matrix \mathbf{B} , has been considered [27], [6] for use in optimal state estimation. Further exploration and analysis of these, and other, preconditioning techniques, following the theoretical approach presented here, may be valuable in order to produce further improvements in the conditioning of the problem.

Acknowledgements This research has been supported in part by the National Centre for Earth Observation, the UK Engineering and Physical Sciences Research Council and the Met Office.

References

- [1] E. Andersson, M. Fisher, R. Munro and A. McNally, *Diagnosis of background errors for radiances and other observable quantities in a variational data assimilation scheme, and the explanation of a case of poor convergence*, Q. J. R. Met Soc., **126**, 1455–1472, 2000.
- [2] R. Bannister, *A Review of forecast error covariance statistics in atmospheric variational data assimilation. II: Modelling the forecast error covariance statistics*, Q. J. R. Met. Soc., **134**, 1971–1996, 2008.
- [3] P. Courtier, J.-N. Thépaut and A. Hollingsworth, *A strategy for operational implementation of 4D-Var, using an incremental approach*, Q. J. R. Met. Soc., **120**, 1367–1387, 1994.

- [4] R. Daley, *Atmospheric Data Analysis*, Cambridge University Press, 1993.
- [5] R. Ebrahimiyan and R. Baldrick, *State Estimator Condition Number Analysis*, IEEE Transactions on Power Systems, **16**, pp. 273-279, 2001.
- [6] M. Fisher, J. Nocedal, Y. Tremolet and S.J. Wright *Data assimilation in weather forecasting: a case study in PDE-constrained optimization*, Optim. Eng., **10**, 409–426, 2009.
- [7] P. Gauthier, C. Charette, L. Fillion, P. Koclas and S. Laroche *Implementation of a 3D Variational Data Assimilation System at the Canadian Meteorological Centre. Part 1: The Global Analysis*, Atmosphere-Ocean, **37**, 103–156, 1999.
- [8] T. Gneiting, *Simple tests for the validity of correlation function models on the circle*, Statistics and Probability Letters, **39**, 119–122, 1998.
- [9] G. H. Golub and C. F. Van Loan, *Matrix Computations, third edition*, Johns Hopkins University Press, 1996.
- [10] S. Gratton, A.S. Lawless and N.K. Nichols, *Approximate Gauss-Newton methods for non-linear least-squares problems*, SIAM J. Optim., **18**, 106–132, 2007.
- [11] R.M. Gray, *Toeplitz and Circulant matrices: A Review*, Foundations and Trends® on Communications and Information Theory, **2**, pp.155–239, 2006.
- [12] S.A. Haben, A.S. Lawless and N.K. Nichols, *Conditioning and preconditioning of the variational data assimilation problem*, Computers and Fluids, **46**, 252–256, 2011.
- [13] S.A. Haben, A.S. Lawless and N.K. Nichols, *Conditioning of incremental variational data assimilation, with application to the Met Office system*, **64**, Tellus A, 782-792, 2011.
- [14] S.A. Haben, *Conditioning and Preconditioning of the Minimisation Problem in Variational Data Assimilation*, PhD Thesis, Department of Mathematics and Statistics, University of Reading, 2011. <http://www.reading.ac.uk/web/FILES/maths/HabenThesis.pdf>
- [15] S. B. Healy and A. A. White, *Use of discrete Fourier transforms in the 1D-Var retrieval problem*, Q. J. R. Met Soc., **131**, 63–72, 2005.
- [16] C. Johnson, B.J. Hoskins, N.K. Nichols, *A singular vector perspective of 4D-Var: Filtering and interpolation*, Q. J. R. Met Soc., **131**, 1–19, 2005.
- [17] A.S. Lawless, S. Gratton and N.K. Nichols, *An investigation of incremental 4D-Var using non-tangent linear models*, Q. J. R. Met. Soc., **131**, 459–476, 2005.
- [18] A.C. Lorenc, *Optimal nonlinear objective analysis*, Q. J. R. Met. Soc., **114**, 205–240, 1988.
- [19] A.C. Lorenc, *Development of an Operational Variational Assimilation Scheme*, J. Met. Soc. Japan, **75**, 339–346, 1997.
- [20] J. Meng and C. L. DeMarco, *Application of Optimal Multiplier Method in Weighted Least-Squares State Estimation Part I: Theory*, University of Wisconsin-Madison, 1999.

- [21] K. W. Morton and D. F. Mayers, *Numerical Solution of Partial Differential Equations*, Cambridge University Press, 1994.
- [22] N.K. Nichols, *Data Assimilation: Aims and Basic Concepts* In: R. Swinbank, V. Shutyaev and W.A. Lahoz, (ed.) *Data Assimilation for the Earth System*. Kluwer Academic, pp. 9-20, 2003.
- [23] S. Pajic, *Power System State Estimation and Contingency Constrained Optimal Power Flow - A Numerically Robust Implementation*, PhD Thesis, Worcester Polytechnic Institute, 2007.
- [24] F. Rawlins, S.P. Ballard, K.J. Bovis, A.M. Clayton, D. Li, G.W. Inverarity, A.C. Lorenc and T.J. Payne, *The Met Office global four-dimensional variational data assimilation scheme*, Q. J. R. Met. Soc., **133**, 347–362, 2007.
- [25] O. Talagrand and P. Courtier, *Variational assimilation of meteorological observations with the adjoint vorticity equation. I: Theory*, Q. J. R. Met. Soc., **113**, 1311–1328, 1987.
- [26] Y. Tremolet, *Incremental 4D-Var convergence study*, Tellus A, **59**, 706–718, 2007.
- [27] J. Tshimanga, S. Gratton, A. T. Weaver and A. Sartenaer *Limited-memory preconditioners, with application to incremental four-dimensional variational data assimilation*, Q. J. R. Met. Soc., **134**, 751-769, 2008.
- [28] Rudolf O. Weber and Peter Talkner, *Some Remarks on Spatial Correlation Function Models*, Monthly Weather Review, **121**, 2611–2617, 1993.
- [29] Andrew T. A. Wood, *When is a truncated covariance function on the line a covariance function on the circle?*, Statistics and Probability Letters, **24**, 157–164, 1995.
- [30] A. M. Yaglom, *Correlation Theory of Stationary and Related Random Functions I. Basic Results*, Springer-Verlag, 1986.