

---

---

**A study of dynamical instability and  
filter stability using ensemble  
Kalman filter**

---

---

*A thesis*

*submitted to the*

Tata Institute of Fundamental Research, Mumbai

*for the degree of*

DOCTOR OF PHILOSOPHY

*in*

PHYSICS

*by*

**SHASHANK KUMAR ROY**

INTERNATIONAL CENTRE FOR THEORETICAL SCIENCES

TATA INSTITUTE OF FUNDAMENTAL RESEARCH

Bengaluru, India.

Submitted in January 2024

[Final version submitted in January 2025]



# Declaration

This thesis is a presentation of my original research work. Wherever contributions of others are involved, every effort is made to indicate this clearly, with due reference to the literature, and acknowledgement of collaborative research and discussions. The work was done under the guidance of Professor Samriddhi Sankar Ray and Professor Amit Apte at the International Centre for Theoretical Sciences- Tata Institute of Fundamental Research, Bengaluru.



Shashank Kumar Roy

January 31, 2025

In my capacity as supervisor of the candidate's thesis, I certify that the above statements are true to the best of my knowledge.



Prof. Samriddhi Sankar Ray



# Acknowledgements

To begin, this thesis would not have been possible without my *Guru ji*, Professor *Amit Apte* whose mentorship was fundamental in my scientific endeavors. It was his indispensable guidance, the degree of freedom, and support that allowed me to refine and advance my scientific temperament through this journey. In times of doubt and confusion, he was always there to patiently understand and discuss them without any prejudice. I really admire this quality, something I want to deeply assimilate in my own personality.

I owe a debt of gratitude to Professor *Samriddhi Sankar Ray*, my Supervisor, for supporting me through the different phases of my academic journey at ICTS. Special thanks to Professor *Vishal Vasan*, from whom I learned a lot during the last year of my PhD, engaging in invigorating discussions and receiving honest feedback on both scientific and general issues of a graduate student.

Throughout my coursework, I cherished my time bonding with an amazing group of batchmates: *Soumyadip, Prashant, Arnab, Pinak, Pronobesh, Sujan, Srikant, Vijay, Monica* and *Ashwin*. I really admire them for their dedication and discipline.

I thank *Pinak*, a collaborator and friend, who has been my companion for evening tea and lovely long strolls through GKVK and NCBS. From him, I learned how to be hyper-focused on getting difficult things done while being humble about it. I thank *Manisha* for being very supportive and lending a helping hand within and beyond academic assistance, playing the role of a friend and sometimes, as an elder sister. I made quite some friends which made my life inside and outside of ICTS, full of fun happenings and gossip. Faces that lit me up with their bright and welcoming smile were *Anup, Uddepta, Srashti, Divya, Basudeb, Sparsh, Aditya Thorat, Ankush, Alan, Akash, Harshit, Mukesh, Souvik, Shivam, Shalabh*, and the list goes on. I extend my heartfelt thanks and warm regards to all of you.

I now move on to thank people in my life where no amount of words can succeed in portraying my exact feelings. Yet, I will attempt to bring the sublime into a few words about my roots – my family. Without them, none of the things would have been possible for me the way they are today. I am deeply grateful to my mother and father for their unwavering support in my decisions, even when I was uncertain about what was best for me. My parents are the fundamentals of my belief system. They never let the feeling of discouragement come near me no matter how many times I failed. I learned a lot of

life lessons from my father, who worked in the coal mines with extreme hard work and resilience, day and night all his life with the sole purpose of getting his children educated in the best possible way. My mother taught me the difference between easy and right choices. It's her love and determination that speaks through me. Due to their understanding and remarkable foresight, they made some of the most challenging decisions that I now consider to be ahead of their time. I marvel at their courage and wisdom that never conformed to the imposed structures of the society in which we lived. In the long journey of mine, it's them who made all the sacrifices. Their endless love and encouragement make me feel like the most fortunate person in the world.

I am very grateful to my sister *Deepshikha*, who has always been the constant forward push of my life. I watched her work hard towards her goals which instilled in me the desire to study science and pursue a career in it. She taught me the value of relentless pursuit of one's goals, fighting distractions while balancing responsibility. To have a loving and protective sister like her is a blessing. I also show my endless love to my nephew *Rivaansh* who brings me tremendous joy and happiness.

I consider myself extremely lucky to have found a very loving, caring, and supportive partner in *Kirti*. Meeting her was one of the most beautiful things that happened to me during my initial days in Bangalore. Whenever I had setbacks, she was always there to motivate me, cheer me up, and never let my spirits drop low. I thank you from the core of my heart as your support over the years has been crucial in pursuing this difficult journey. Together with *Aastha*, we have lived so many beautiful moments and will continue doing so.

The story of all the people and their contributions to my life cannot be summarized so quickly. However, the thesis must be submitted on a specific date.

Lastly, I remain forever indebted to ICTS. I am also grateful to NCBS and its library, a beautiful and productive place which saw me roaming around over-caffeinated from early mornings to late nights, running codes and writing this thesis. I remember coming to ICTS as a naive and ignorant person, concerned about my career and thinking about science. Now, ICTS—its values, its people, and its community— are forever an integral part of my life's trajectory. It is here that I met beautiful people, full of love, friendship, compassion, and support.

*With profound love and gratitude,  
I dedicate this thesis to my dearest mother and father,*

*KANCHANMALA DEVI*

*&*

*SRIKANT ROY*





# Abstract

Estimating and predicting the state of a chaotic dynamical system which evolves over time from partial and indirect observations is a challenging problem. Also known as nonlinear filtering or data assimilation in earth sciences, this problem is solved using sequential algorithms which approximate the state and the associated uncertainty. Bayesian formalism allows one to define the filtering distribution, the conditional distributions of the state based on all available information. Since the state is unknown, the initial distribution for any filtering algorithm starts with a choice for the initial distribution of the state. It is crucial that the influence of the wrong initial condition on the filter estimated distributions diminish over time, eventually converging to same distribution - a property known as filter stability. Developing a method to understand and directly investigate this property across various numerical filtering algorithms is of great utility.

The first problem of this thesis addresses the problem of nonlinear filter stability numerically. Using ensemble Kalman filter, a widely utilized bayesian filtering algorithm, we directly investigate nonlinear filter stability by performing data assimilation on chaotic dynamical systems such as Lorenz-63 and Lorenz-96 as our choice for the testbed models. With the help of Sinkhorn algorithm, we compute the distances between two different filtering distributions starting from different initial conditions. The numerical approach developed here for studying non-linear filter stability using a distance on the space of probability distributions has been a novel contribution out of this thesis.

The second problem studied in the thesis is related to dynamical instability and their computation for a dynamical system. We focus on the Lyapunov vectors such as backward and covariant Lyapunov vectors which are strongly related to the growth of error and uncertainty in state estimation problems. However, computing these vectors requires the knowledge of the the dynamics and a true underlying trajectory which is unknown in the setting of data assimilation. We propose to first estimate the underlying trajectory and then use Genelli's algorithm to recover the Lyapunov vectors themselves. We also investigate the sensitivity of these vectors themselves by systematically replacing the actual trajectory with approximate trajectories. This allows us to understand the limitations of the numerical algorithms for computing Lyapunov vectors from estimated trajectory.



# Contents

<b>1</b>	<b>Introduction</b>	<b>17</b>
1.1	A general introduction to this thesis . . . . .	19
1.2	Understanding nonlinear filter stability numerically using Wasserstein distance	19
1.3	Computing Lyapunov vectors from filter approximated trajectory and their sensitivity . . . . .	22
1.4	Structure of the thesis . . . . .	25
<b>2</b>	<b>Bayesian filtering for dynamical systems using ensemble Kalman filter</b>	<b>27</b>
2.1	Nonlinear filtering for dynamical systems . . . . .	27
2.2	The state estimation problem for a discrete-time dynamical systems . . . .	31
2.3	Background: Bayesian filtering theory . . . . .	32
2.4	Background: The Kalman Filter – a recursive solution to filtering of linear dynamical systems . . . . .	34
2.5	Ensemble Kalman filter: extending Kalman filters to nonlinear dynamical systems with monte-carlo approach . . . . .	38
2.5.1	Curse of dimensionality in EnKF . . . . .	40
2.5.2	Inflation . . . . .	41
2.5.3	Localization . . . . .	41
2.6	Chaotic dynamical systems used in this thesis . . . . .	42
2.6.1	Lorenz 63: A 3-dimensional chaotic ODE . . . . .	43
2.6.2	Lorenz 96: A 40-dimensional chaotic ODE . . . . .	44
2.7	Some performance metric for numerical filtering algorithms . . . . .	45
2.7.1	Accuracy . . . . .	45
2.7.2	Reliability . . . . .	46
2.7.3	Stability . . . . .	47
2.8	Twin experiments and results . . . . .	47
2.8.1	A brief overview of the experimental setup . . . . .	48
2.8.2	Ensemble size versus Time of divergence . . . . .	48
2.8.3	The effect of localization and inflation on filter divergence . . . . .	50
2.8.4	Localization scale versus RMSE for different localization functions .	51

2.8.5	Effect of model error in data assimilation . . . . .	52
2.9	Summary . . . . .	56
<b>3</b>	<b>Numerical filter stability of EnKF using Sinkhorn divergence</b>	<b>59</b>
3.1	An introduction to nonlinear filter stability . . . . .	60
3.2	Mathematical definition of filter stability . . . . .	61
3.3	Background: Wasserstein distance on the space of probability distributions	64
3.3.1	The optimal transport problem and Wasserstein distance . . . . .	65
3.3.2	Entropy-regularized optimal transport and Sinkhorn divergence . . . . .	67
3.3.3	Sinkhorn-Knopp algorithm . . . . .	68
3.3.4	Important metric for filter stability . . . . .	70
3.4	Results and discussion . . . . .	71
3.4.1	Understanding numerical properties of Sinkhorn divergence . . . . .	71
3.4.2	Results for Lorenz-96 system . . . . .	73
3.4.3	Dependence of the filter stability w.r.t observation gap . . . . .	76
3.4.4	Dependence of the filter stability w.r.t observation noise . . . . .	79
3.5	Summary . . . . .	82
<b>4</b>	<b>Computing Lyapunov instabilities of a dynamical system using data assimilation</b>	<b>83</b>
4.1	Introduction . . . . .	84
4.2	Theory and computation of Lyapunov vectors . . . . .	85
4.2.1	Definition and importance of covariant Lyapunov vectors (CLV) . . . . .	86
4.2.2	Computation of Lyapunov vectors . . . . .	88
4.2.3	Data-based algorithm to calculate the Lyapunov vectors . . . . .	91
4.2.4	Lyapunov Vectors from perturbed trajectory . . . . .	92
4.2.5	Models . . . . .	93
4.2.6	Details of computing Lyapunov Vectors from filtered and perturbed trajectories . . . . .	94
4.2.7	Comparison metrics . . . . .	95
4.3	Results and discussion . . . . .	96
4.3.1	BLVs and CLVs computed from assimilated trajectories . . . . .	97
4.3.2	Dependence on perturbation strength for Lorenz-63 . . . . .	100
4.3.3	Dependence on dimension and perturbation strength for Lorenz-96 . . . . .	102
4.3.4	Oseledets' subspaces spanned by the LVs for different perturbation strengths . . . . .	103
4.4	Summary . . . . .	106
	<b>Bibliography</b>	<b>109</b>





# List of Publications

- [1] Pinak Mandal, Shashank Kumar Roy, and Amit Apte. Stability of nonlinear filters- numerical explorations of particle and ensemble Kalman filters. In 2021 Seventh Indian Control Conference (ICC), pages 307–312. IEEE, 2021 <https://doi.org/10.1109/ICC54714.2021.9703185>
- [2] Pinak Mandal, Shashank Kumar Roy, and Amit Apte. Probing robustness of nonlinear filter stability numerically using sinkhorn divergence. *Physica D: Nonlinear Phenomena*, 451:133765, 2023 <https://doi.org/10.1016/j.physd.2023.133765>
- [3] Shashank Kumar Roy and Amit Apte. Computation of covariant Lyapunov vectors using data assimilation. Submitted to *Nonlinear Processes in Geophysics*, <https://doi.org/10.5194/egusphere-2023-2168>





# Chapter 1

## Introduction

Weather prediction is one of the most important problems in the world. It is also a very difficult task due to the sheer complexity of the governing set of equations itself combined with the coupled dynamics of different components of the Earth system, such as the land, atmosphere, the oceans [11, 34]. Based on fundamental physical theories and observed empirical behavior, large and complex numerical models of the ocean and atmosphere are built, which are then used to perform the task of predicting the future of the system in numerical weather prediction or NWP. At any given time, the state of the ocean and the atmosphere is typically described and governed by a set of state variables such as temperature, pressure, humidity, and wind speed at different points in space and time. Some of the aforementioned variables can be observed via sensors and measurement devices contaminated by measurement noise. Other unknown external factors also contribute to measurement noise distribution, making uncertainty their inherent attributes.

Real observations provide a direct link to the hidden state of the underlying system by sampling it at different points in space and time. However, it is impossible to produce estimates of all the state variables only from the observations that are limited in time and space, i.e., they are sparse. Numerical models, on the other hand, are based on theory, formulated in order to understand the spatiotemporal evolution of the underlying system [53]. They are developed to numerically approximate the real dynamics of the system, trading off complexity with accuracy. These models come to life with certain simplifying assumptions, approximations, and discretization and are solved in time making a specific choice of the numerical scheme. They also contain unavoidable modeling errors that are attributed to a finite-scale representation and unknown and unresolved dynamics that are present in the real system but absent in the numerical models[26, 41].

With a numerical model as the best representative of the system at large, the objective of operational weather, ocean, and climate prediction centers around the world is to estimate the current state and use their most accurate estimates to forecast future weather. The typical models employed in numerical weather prediction are complex, high-dimensional,

nonlinear set of equations where the system's state at any time is represented by a large state vector of the order of  $O(10^9)$  [36]. The problem of estimating the state is further complicated by the highly nonlinear and chaotic nature of these systems, where small errors in the initial conditions grow exponentially fast in time. Hence, even if we assume that the models are perfect, ignoring the error due to numerical representation, the forward integration of these models starting from the initial condition soon departs from reality.

The shortcomings of each of the two sources of information must be considered in order to balance the quality of information coming from these two distinct sources. The problem of inverting observations to estimate the underlying state is a poorly posed problem [7]. Data assimilation plays a key role in consistently combining information from different sources of observations from various sources, such as satellites, radars, and ground-based sensors, into numerical models to improve the accuracy of weather predictions. The goal of data assimilation is to provide an estimate of the true state of the system [45] based on a set of sparse and noisy observations. Data assimilation can often be referred to as an interpolation method for dynamical systems, where numerical models play a crucial role in the interpretation and interpolation of observations to estimate the full state of the system [53, 36].

Data Assimilation research addresses two problems of inverse modeling, filtering, which is to estimate the true state of the system, and prediction which aims to predict the future of the system with uncertainty estimates. The two sources of information with uncertainties at hand are dynamical models and observation data. The former captures the underlying physics of the system and faithfully represents our best objective knowledge of the system. On the other hand, the noisy and partial observations from the system are incomplete to describe the full state of the system. Combined together, they produce more accurate estimates of the state and also predict future observations, which is better statistically than either just the model or the observations. The information coming from real observations improves the model outcomes on which the next forecast cycle is then based. This complete cycle is what constitutes a general assimilation cycle.

The history of data assimilation originates from the control theory perspective [42, 54], which led to the development of variational data assimilation methods. The goal here is to obtain a trajectory which minimizes a cost function collectively representing the misfit with the observations while being consistent with the dynamical equations. Another recent perspective stemming from the Bayesian approach is now well recognized and adopted in data assimilation in geosciences [13, 4]. A few examples of sequential Bayesian filtering algorithms, such as the ensemble Kalman filter and particle filter which attempt to represent the posterior distribution via an ensemble approach and implement Bayesian posterior approximations. The state-of-the-art algorithms for the variational data assimilation include the 3DVAR and 4DVAR algorithms. In this thesis, we focus on the Bayesian

filtering algorithms, specifically the ensemble Kalman filters which we apply to nonlinear dynamical systems to study the research problems addressed in this thesis.

## 1.1 A general introduction to this thesis

In the context of data assimilation for deterministic dynamical systems, the focus of this thesis is twofold: (1) to develop numerical methods to assess the stability of filtering algorithms and (2) to develop data-based numerical methods to compute Lyapunov vectors using filtering algorithms. Data assimilation or filtering is the method of combining observations with dynamics in order to sequentially estimate the state of a dynamical system [26, 13]. The associated uncertainties depend on the unstable directions along which errors grow exponentially and are important for the goal of state estimation. Lyapunov exponents and vectors are the fundamental tools used in the study of chaotic systems, which describe the asymptotic behavior of infinitesimal errors [70]. For chaotic and dissipative systems in high dimensions [80], these subspaces are much smaller compared to the size of the system itself. Their information is crucial for analysis and predictability. We study the sensitivity of the Lyapunov vectors in response to perturbations introduced into the underlying trajectory. This allows us to understand the limitations on the accuracy of such vectors computed from estimated trajectories obtained using a filtering algorithm.

In the same spirit that sensitivity to the initial conditions is important for chaotic dynamical systems, we study the problem of filter stability [74, 75], which deals with the sensitivity to the initial condition of a filtering algorithm. In practice, we often do not know the initial distribution and use a different distribution to start the filtering algorithm. A measure of the robustness of a filtering algorithm is reflected in understanding how the posterior conditional distribution of the state becomes independent of the choice of the initial distribution over time. In this thesis, we develop a method based on the Sinkhorn algorithm to numerically study the stability of a nonlinear filtering algorithm. We illustrate the proposed method using the ensemble Kalman filter, where we approximate the distance between Monte Carlo samples representing different filtering distributions.

## 1.2 Understanding nonlinear filter stability numerically using Wasserstein distance

**Bayesian filtering theory** The Bayesian filtering problem is defined as the sequential estimation of the conditional distribution of the state of a dynamical system given the history of observations up to that time [26]. A filtering problem consists of a dynamical equation governing the state and observations, which are the functions of the state. We

assume that the system dynamics is Markovian and is observed indirectly through noisy and incomplete observations over time. We state the filtering problem in a discrete-time setting for state  $x_k \in \mathbb{R}^d$  and the observation  $y_k \in \mathbb{R}^m$ , given by,

$$x_k = f_{k-1}(x_{k-1}); \quad y_k = h_k(x_k) + \epsilon_k \quad (1.1)$$

where,  $f_{k-1}$  is the propagator from time  $t_{k-1}$  to  $t_k$ ,  $h_k$  is the observation operator and  $\epsilon_k$  are *i.i.d* Gaussian errors in the observations with distribution  $\mathcal{N}(0_m, \sigma^2 I_m)$ . Since the true initial state  $x_0$  is unknown, the filtering starts with  $\pi(x_0)$  as the initial distribution at time  $t_0$ . As observations  $y_k$  arrive sequentially in time, the filter distribution  $\pi(x_k|y_{0:k})$  at time  $t_k$  is computed using the likelihood of observation and the prior distribution  $\pi(x_k|y_{0:k-1})$  capturing the flow-dependent uncertainty via Bayes theorem [13]. Different numerical filtering algorithms are used to obtain an estimate for the filtering distribution, denoted by  $\hat{\pi}_k(\mu)$ , using different approximations involved in Bayesian posterior computation.

Kalman filter provides an iterative closed-form solution [77] to the Bayesian filtering problem in the scenario when the dynamic and measurement operator is linear with the observation noise being Gaussian. A simple approximation for the case of nonlinear dynamics was introduced by G. Evensen [34] where an ensemble or collection of particles are used to represent the distribution and operations of the standard Kalman filtering are performed by using ensemble estimates of corresponding quantities. This ensemble representation is powerful for dimension reduction, leading to computational feasibility for systems in large dimensions. Additional ad-hoc procedures such as inflation and localization [16], have made EnKF more applicable and popular in high-dimensional operational data assimilation problems.

**Numerical filter stability** In this work, we introduce a practical way to assess filter stability, an important property of filtering algorithms. Filter stability aims to understand the effect of incorrect initialization of numerical filtering algorithms. In data assimilation research, many studies have relied on these models as test models to understand and demonstrate filtering algorithms and their intercomparison [7]. A numerical filtering algorithm requires the initialization of the probability distribution  $\pi(x_0)$ , an arbitrary choice that may be different from the truth distribution. The stability of the filter is a property ascertaining that starting with different initial distributions, the filter converges to the same posterior distribution asymptotically [75] making the filtering robust to the wrong choice of  $\pi(x_0)$ .

In our definition [64], a numerical filtering algorithm is stable if a numerical filter, starting with two different initial distributions  $\mu$  and  $\nu$ , yields  $\pi_n(\mu)$  and  $\pi_n(\nu)$  as the respective filtering distributions at time  $t_n$ , then asymptotically, we have

$$\lim_{n \rightarrow \infty} \mathbb{E}[D(\hat{\pi}_n(\mu), \hat{\pi}_n(\nu))] = 0, \quad (1.2)$$

where  $D$  is a distance on  $\mathcal{P}(\mathbb{R}^d)$ , the space of probability measures on  $\mathbb{R}^d$ . By studying the equation 1.2 for a numerical filter using an appropriate distance metric  $D$ , we can directly study the problem of filter stability numerically by choosing different initial distributions for the filter and how the distance between the corresponding filtering distributions varies over time.

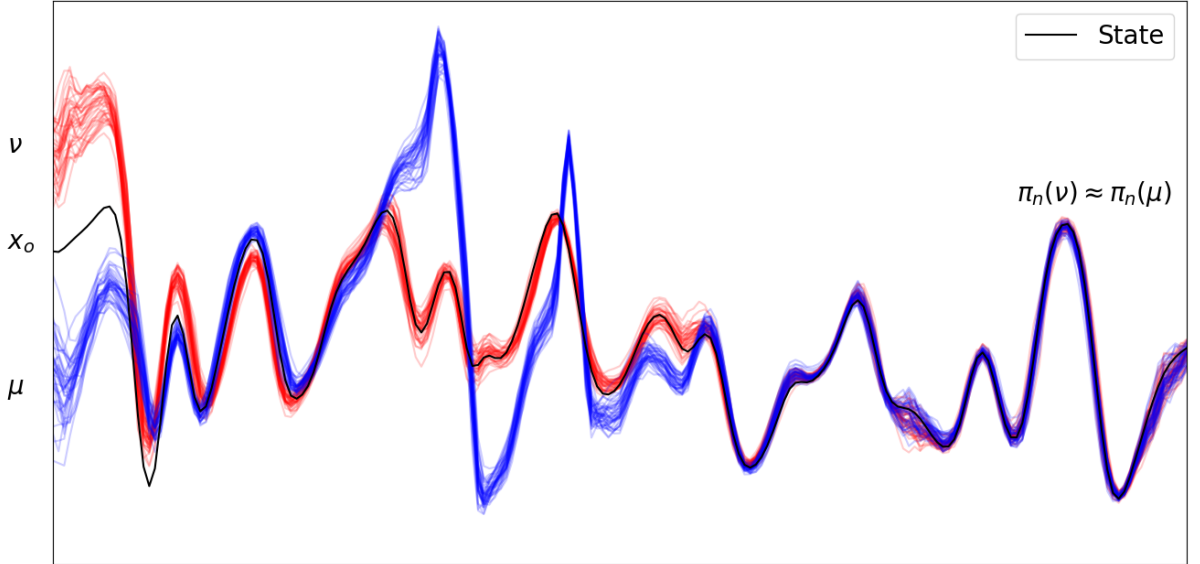


Figure 1.1: A schematic illustrating filter stability. The red and the blue trajectories represent the two different filtering ensemble over time with initial distribution  $\nu$  and  $\mu$ .

**Wasserstein distance and Sinkhorn divergence** In this work, we have used an approximation of the Wasserstein-p distance [6, 38], denoted as  $W_p$ , in the space of probability distributions to numerically study filter stability. Sinkhorn divergence [37], denoted by  $S_\varepsilon$ , is computationally cheaper due to entropy regularization [35] and can be solved by various iterative methods, making it useful to approximate  $W_p$  between two probability distributions. We use Sinkhorn-divergence to approximate  $W_2$  in the limit of  $\varepsilon \rightarrow 0$  [35]

$$\lim_{\varepsilon \rightarrow 0} \sqrt{S_\varepsilon} = W_2(\mu, \nu). \quad (1.3)$$

We study the distance between filtering distributions, defined in equation (1.2), where the two different filtering distributions are obtained using ensemble Kalman filter with different initial distributions over time using  $S_\varepsilon$ . In the following discussion, we refer to  $D_\varepsilon$  as the distance, which is the  $\sqrt{S_\varepsilon}$  defined in equation (1.3).

**Results with ensemble Kalman filter for Lorenz-96** We investigate the filter stability using Lorenz-96 [60] for  $d = 10$ , forcing constant  $f=10$  with observation gap  $g = 0.1$  units of time. We use three different initial conditions and study  $\mathbb{E}[D_\varepsilon(\hat{\pi}_n(\mu_i), \hat{\pi}_n(\mu_j))]$ ,  $i \neq j$  as a function of time  $n$  for initial conditions  $\mu_i$ , with the expectation taken by averaging over 10 observation realizations. We also fit an exponential curve of the following form

between Sinkhorn divergence and time [83],

$$\mathbb{E}[D_\varepsilon(\pi_n(\mu_0), \pi_n(\mu_b))] = a \exp(-\lambda t) + c \quad (1.4)$$

where  $t = \text{assimilation step} \times \text{observation gap} = ng$  and compute the rates of convergence in different cases. With the recent development of the Sinkhorn algorithm, we propose a method to study the stability of filtering algorithms. We demonstrate this method by applying filtering to a partially observed dynamical system using EnKF. This allows us to directly assess the stability by computing the expected value, averaged over multiple observation realizations, of the distance between filtering distributions as a function of time. We also studied extensively the dependence of stability properties on two main parameters, the time between observations and the observational error covariance. Furthermore, we find that the distance between two filters initialized with distinct initial conditions has an empirically linear relationship with the  $l_2$ -error of the filter mean.

### 1.3 Computing Lyapunov vectors from filter approximated trajectory and their sensitivity

Characterizing the directions along which the errors grow rapidly is an important property of a chaotic dynamical system. The associated uncertainties depend on the unstable directions along which the errors grow exponentially and are important for the goal of state estimation. Lyapunov exponents and vectors describe the asymptotic behavior of infinitesimal errors [33]. The spectrum Lyapunov exponents summarize the effect of all possible perturbations on the initial condition for the system, and the values of the exponents describe the average growth rate of the respective perturbations over time. The Lyapunov vectors, or LVs, are local objects spanning the tangent space for a specific point in the phase space. The computation of LV requires the knowledge of the dynamics and the underlying trajectory along which these vectors are obtained. We investigate the computation of Lyapunov vectors for deterministic dynamical systems using partial and noisy observations and data assimilation.

**Lyapunov vectors and subspaces** We consider an autonomous continuous-time dynamical system in  $R^n$  together with the equation governing the evolution of infinitesimal perturbations in tangent space, given by,

$$\dot{x}_t = f(x_t), \quad \dot{Q}_t = J(x_t)Q_t, \quad J_{ij}(x_t) = \frac{\partial f_i(x_t)}{\partial x_{t,j}} \quad (1.5)$$

where  $Q \in R^{n \times n}$  satisfies the equation for the Fundamental matrix which contains the set of perturbation vectors along the columns.  $J_{ij}(x_t)$  is the jacobian matrix evaluated at  $x_t$ ,

and  $f_i$  and  $x_{t,j}$  denote the  $i^{\text{th}}$  and  $j^{\text{th}}$  component of  $f(x_t)$  and  $x_t$  respectively.

If  $z_k$  and  $z_l$  denotes the infinitesimal perturbations at time  $t_k$  and  $t_l$ , the tangent linear propagator from time  $t_k$  to  $t_l$  can be written as

$$\mathcal{M}_{k,l} = Q_l Q_k^{-1} \quad \text{with the property that} \quad z_l = \mathcal{M}_{k,l} z_k. \quad (1.6)$$

Under suitable conditions, the Oseledec's theorem implies the existence of the following limits:

$$\lambda(z, t_k) := \lim_{l \rightarrow \infty} \frac{1}{|t_l - t_k|} \log \frac{\|\mathcal{M}_{k,l} z\|}{\|z\|}, \quad \text{and} \quad \lambda(z, t_l) := \lim_{k \rightarrow -\infty} \frac{1}{|t_l - t_k|} \log \frac{\|\mathcal{M}_{k,l} z\|}{\|z\|}, \quad (1.7)$$

A non-increasing tuple  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n$  of Lyapunov exponents summarizes the global, asymptotic rate of change of linear perturbations around a trajectory. The forward and backward LVs are the eigenvectors of the forward and backward Oseledec matrices, denoted by  $\Phi_l^+$  and  $\Phi_l^-$  respectively [88], which are defined as

$$\Phi_l^- = \lim_{k \rightarrow -\infty} \frac{1}{|t_l - t_k|} \log [\mathcal{M}_{k,l}^T \mathcal{M}_{k,l}], \quad \Phi_l^+ = \lim_{k \rightarrow \infty} \frac{1}{|t_k - t_l|} \log [\mathcal{M}_{k,l}^T \mathcal{M}_{k,l}]. \quad (1.8)$$

Using the intersection between the subspaces spanned by the FLVs and the BLVs, one can define a norm-independent set of LVs called covariant Lyapunov vectors or CLVs [52]. These vectors have different features from their orthogonal counterparts and have been computed with recent algorithms [39, 88] as can be seen in figure 1.2 where there is additional information contained in their mutual angles.

**Lyapunov vectors from partial and noisy observations** We focus on computing Lyapunov vectors for deterministic dynamical systems using partial and noisy observations and data assimilation. When the initial condition for a given dynamical system is unknown and the system is accessible only by partial and noisy observations, identification of the underlying trajectory is a challenge. Filtering techniques can be used to compute an approximate trajectory using the observations of the system and the numerical model. We propose to obtain the vectors and their subspaces from the estimated trajectory obtained from a data assimilation algorithm since it is close to the true trajectory. We take the approach of using data assimilation to compute the LVs using Ginelli's algorithm [39] from a filter-estimated trajectory. We propose to use the state estimated from a numerical filter as the best estimate of the unknown state  $x_{t_k}$  on the underlying trajectory and carry out the computation of Lyapunov vectors and Oseledec's subspaces. We use the ensemble Kalman filter or EnKF [34, 7] as the filtering algorithm, which is described in section 2.5 in detail.

A trajectory estimated with DA using any filtering or smoothing algorithm does not meet the requirements of being a dynamical trajectory and at any time  $t_k$ , even the most

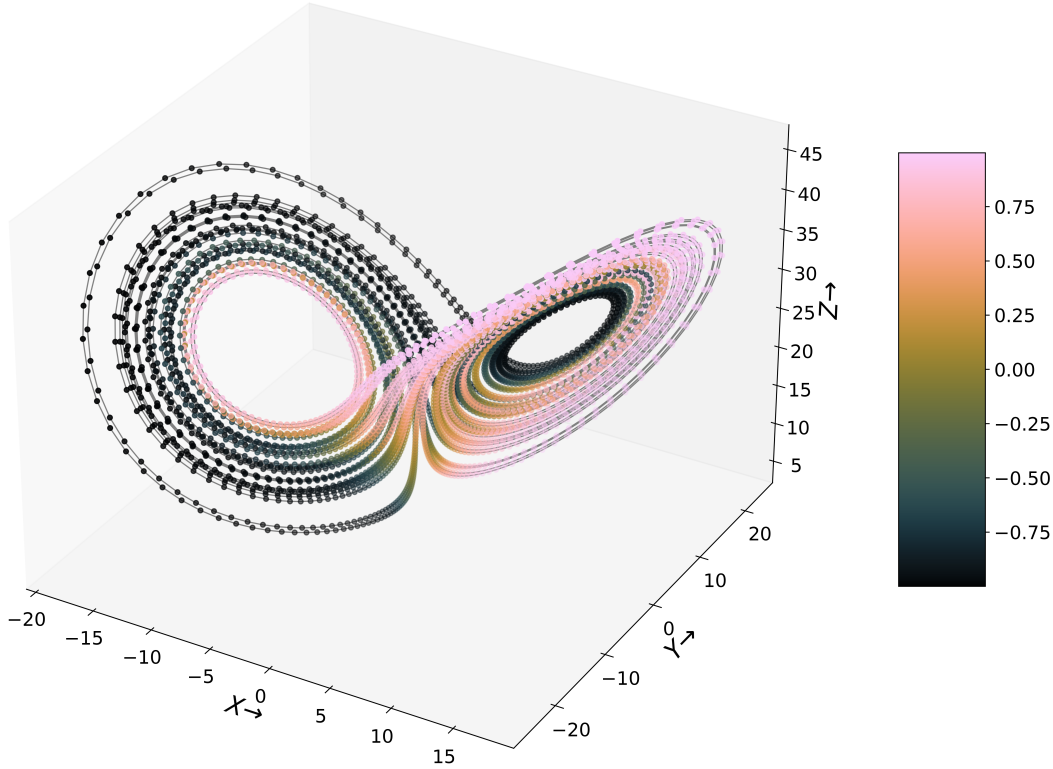


Figure 1.2: The attractor of Lorenz-63 system where the color indicates the cosine of the angle between the 1<sup>st</sup> and the 2<sup>nd</sup> CLV.

accurate estimate, denoted by  $x_k^a$ , has minor discrepancies  $e_k^a$  from the actual state  $x_k$ , which is expressed as  $e_k^a = x_k^a - x_k$ . The amount of inaccuracy in the calculation of LVs depends on the magnitude of the error  $e_k^a$ . This naturally leads to the question of the sensitivity of the Lyapunov vectors in response to perturbations introduced into the underlying trajectory. We explore this problem using the notion of perturbed trajectory where the perturbations follow some unknown error statistics, we break down our problem into two sub-parts.

- i. Firstly, we present a data-based algorithm for computing the LVs using the state estimates obtained from a data assimilation method, namely the ensemble Kalman filter (though any assimilation method may be used in place of EnKF). This algorithm can be used to produce the full spectrum or a subset of LVs, either backward or covariant. This data-based method does not, by itself, give any bounds on how close these vectors may be to the true vectors associated with the trajectory that is being observed. To understand this aspect, we are naturally led to the second aim of this paper.
- ii. We present the other main contribution, which is an extensive numerical exploration



of the sensitivity of the LVs to perturbations of the underlying trajectory. We do this by using the same algorithm mentioned above but with a noisy trajectory and then comparing the true LVs with the approximate one obtained from the perturbed trajectory. We used the principal subspace angles between the Oseledets spaces in order to quantify this discrepancy.

We numerically computed LVs using Ginelli’s algorithm but with noise added to a trajectory, and then compared the true LVs with the approximate one obtained from the perturbed trajectory. We investigate the stability of the LVs from a more general perspective by adding perturbations to a trajectory and recovering the vectors and the subspaces from such perturbed trajectories. Using principal angles [43], we quantify the discrepancy for the Oseledets subspaces. We find that this sensitivity is consistent with and helps explain the errors in the approximate Lyapunov vectors from the estimated trajectory of the filter. These results allow us to understand the limitations on the accuracy of such vectors computed from estimated trajectories obtained using a filtering algorithm. We have developed numerical methods for computing the Lyapunov vectors from assimilated trajectories obtained using a filtering algorithm. It is important to study the limitation on the accuracy of different Lyapunov vectors and their subspaces recovered using the assimilated trajectory. We understand this by adding perturbations to a dynamical trajectory and systematically studying the sensitivity of LVs and their subspaces. We find that the errors in the approximated LVs implicitly depend on the errors in the estimated trajectory. Our systematic study reveals that individual vectors can be quite sensitive for dynamical systems with high dimensions as opposed to low dimensions. We also find that the Oseledet subspaces defined by the LVs computed from the approximate trajectory are less sensitive than the individual vectors. Our results provide an understanding and limitations on the accuracy of such vectors computed from estimated trajectories obtained using a filtering algorithm.

## 1.4 Structure of the thesis

The thesis is structured as follows: we first present an introduction to nonlinear filtering for dynamical systems under the Bayesian formalism. We then state the problem of sequential state estimation using partial and noisy observations for a linear dynamical system. We present a derivation of the Kalman filtering equation, the optimal solution to the problem of filtering for linear dynamical systems. We introduce EnKF, the Monte Carlo approximation of the standard Kalman filtering and the ad-hoc approach which make EnKF applicable in high-dimensional settings. In section 2.6, we introduce Lorenz-63 and Lorenz-96, the two chaotic dynamical systems that we use for performing the numerical experiments in the context of the problems addressed in the following chapters.

We conclude chapter 2 by presenting some twin experiments in data assimilation using standard EnKF on Lorenz-96 illustrating different ideas such as inflation and localization together with well-known metrics of filtering algorithms.

In chapter 3 begins with an introduction to the problem of nonlinear filter stability in data assimilation. In section 3.2, we present the mathematical definition of filter stability using a distance on the space of probability distributions. We introduce the Wasserstein distance, an optimal transport based distance, and present the algorithm to approximate it via the Sinkhorn divergence. Using EnKF, we then demonstrate our numerical approach for numerical filter stability by applying it to Lorenz-96. We present the results for two important cases- where we study the filter stability for the case of fixed observation gap and fixed observation noise. The first part of this work aims to study whether filtering algorithms are affected by the wrong choice of initial distributions at the beginning.

In chapter 4, we discuss the problem of estimating Lyapunov vectors from an estimated trajectory using data assimilation. To approximate the Lyapunov vectors using the estimate of the underlying trajectory obtained from the filter mean. In section 4.2, we describe the theory and the method of computing the vectors using Ginelli's algorithm. In section 4.3, we first present the results about the errors in the approximated LV in subsection 4.3.1 when computed from the filter estimated trajectory for Lorenz-63 and Lorenz-96 systems. Following this, in subsection 4.3.2 and 4.3.3, we present an extensive study of the sensitivity of these approximate LVs and the corresponding Oseledets' subspaces to added perturbations of different strengths. Finally, in section 4.3.4 principal angles, we demonstrate that the Oseledets' subspaces defined by the LVs computed from the approximate trajectory are less sensitive than the individual vectors.

# Chapter 2

## Bayesian filtering for dynamical systems using ensemble Kalman filter

We begin this chapter by introducing filtering theory in section 2.1 from a historical perspective on the estimation of signals from noisy measurements in time. In section 2.2, we set up the filtering problem for discrete-time dynamical systems. Under this setup, we introduce the mathematical foundations of Bayesian filtering theory. In section 2.2, we discuss the sequential filtering problem, followed by section 2.4, where we introduce Kalman filtering for the linear dynamical system, which provides an optimal solution to the Bayesian filtering problem under certain assumptions. We provide a derivation of the Kalman filtering equations and gain matrices following the approach of. In the following section 2.5, we present the ensemble Kalman filter (EnKF) which is motivated with simplifications of the standard Kalman filter and explain its application to nonlinear dynamical systems. We then discuss two practical approaches in the following subsections 2.5.2 and 2.5.3 that make it possible for EnKF to work in high-dimensional systems. In section 2.6, we introduce two dynamical systems, namely Lorenz-63 and Lorenz-96, which we use to perform twin experiments and simulations in the rest of the thesis. We present the results of some twin experiments in section 2.8, where we use EnKF to perform data assimilation with partial and noisy observations performed on the Lorenz systems. We conclude this chapter with an additional discussion of some of the introductory numerical results obtained from twin experiments performed using EnKF.

The presentation of different topics in this chapter is based on the following books [14, 54, 7, 36], all of which are comprehensive in their treatment of this topic.

### 2.1 Nonlinear filtering for dynamical systems

Inspired by different fields of physics where understanding the non-linear behavior of a system is required, the field of dynamical systems evolved as a mathematical formalism to

study the instability properties of a system. Astronomy, astrophysics, systems biology, fluid dynamics, and weather prediction are all research areas where the scientific quest involves modeling systems and analyzing data originating from real-time experiments.

Dynamical systems are ubiquitous in engineering and natural sciences, where mathematical models are used to define how the system evolves over time. Modeling such problems involves borrowing ideas from various physical theories and empirical experimentation using mathematical tools such as differential equations and PDEs. The vast area of numerical weather prediction and geosciences abound in creating and analyzing such complex models. The underlying system evolves in time as the state evolves in phase space. In the presence of uncertainty related to important parameters and variables in the models, it is natural to introduce a probabilistic approach. Dynamical systems, together with the probabilistic state-space approach, also lead to a dynamical system in the space of probability distributions [54]. Similarly to how the dynamical equations describe how the system evolves in the phase space, the corresponding probability distributions change in time, and their evolution is captured by the corresponding Fokker-Planck equation, which incorporates the knowledge of the original dynamical system. Depending on the location in the phase space, the system has different dynamical correlations and hence is not a stationary process. Hence, in these problems, methods that rely on the assumptions of stationarity are out of scope for modeling the distribution in time.

The complexity of these models ranges from simple ones such as discrete maps, and deterministic, continuous ordinary differential equation models, to highly complicated non-linear and stochastic partial-differential continuous models, often in very high dimensions with a large number of unknown parameters. To calibrate the parameters of the model itself or to make certain predictions about the behavior of the system in the future, one must produce estimates that are accurate and are able to capture the knowledge of the model and the information coming from the observations at the same time. In general, both the state and the parameters can be estimated together using filtering techniques. Using joint distributions of the state and the parameter, we solve the new filtering problem by simply extending the state vectors to include unknown parameters. However, in this thesis, we only focus on the state estimation problem and hence assume the correct parameter values.

The measurements obtained by observing these systems is low-dimensional compared to the state of the system and typically contains errors of different kinds. Hence as the data is both sparse and noisy, it is difficult to simply determine the full state of the system from merely using the relationship between the two. To address this, filtering theory provides a systematic way to remove the effects of measurement noise and extract useful information about the state of the system. The goal is to estimate the true state of the system using observations made on the system by running an algorithm called a filter.

Filtering theory is a mathematical framework for estimating the state of a dynamical system from a set of observations supplied with a set of governing equations [26, 36, 45]. The approach is to use the flow-dependent uncertainty provided by the governing dynamical equations of the system along with a model of observations which tells us how the observations are related to the true state. This problem is also known as the state estimation or filtering theory [26].

In 1949, Norbert Wiener, in his classical work on interpolation, extrapolation, and smoothing on time-series data [87], addressed the problem of estimating and therefore predicting the underlying signal where the observation data are corrupted by noise. He developed the solution to the problem based on spectral decomposition based on the assumption of second-order stationarity of the signal and additive noise. The continuous-time formulation leads to what is now called the Wiener-Hopf integral equation, while the discrete-time formulation is widely known as the Wiener filter [44]. In the following discussion, we focus on the Wiener problem in discrete time, which is relevant to the discrete-time filtering problem in the sections later in this chapter. Although we may have the underlying process to be continuous, we can assume that some continuous-time system is being sampled at a fixed interval of time  $\delta t$  for simplicity.

Wiener's filtering problem in discrete time is as follows: Suppose that we have received a sequence of observations in time  $\{y_0, y_1, y_2, \dots, y_n\}$ , where  $y_i = x_i + \eta_i$  are the measurements that corrupted the signal  $x_i$  by additive noise  $\eta_i$ . The goal is to determine the sequence  $\{x_0, x_1, x_2, \dots, x_n\}$  of the underlying states of the given the observations accompanied by certain prior assumptions. The phrase *filtering the signal* means separating the underlying process signal from the noisy observations by modeling the dynamics of the process. A simple linear model in which the future  $x_{k+1}$  is the desired result and the input consists of the last  $p$  previous time steps  $x_k$  is

$$x_{k+1} = \sum_0^p a_i x_i \quad (2.1)$$

where, the set of coefficients  $\{a_i\}$  need to be determined. Obtaining a solution to this problem using the least squares method seeks to minimize the mismatch between the filter output and the desired signal, which determines the coefficients of the Wiener filter. However, this linear filter assumes that both the signal and the noise are stationary and their statistics are known. Its performance is also sensitive to the parameters for these statistics. These assumptions are often not satisfied for real observations, which makes them unsuitable for many applications.

Overcoming the drawbacks of Wiener's problem, Kalman, in his seminal paper [44], described his optimal solution to Wiener's problem in a discrete-time linear setting, extending it to include the non-stationary process and multivariate signals. He modeled

the underlying signal using a higher-dimensional linear dynamical system of dimension  $d$ , with  $d > p$ . Using ideas from the dynamical systems and conditional expectations of the distribution of the state, he approached the problem of modeling the underlying dynamical system via a state transition matrix and an observation operator. Since the algorithms were trying to filter out the true signal from the measurements, the algorithms are now known as filtering algorithms.

Any general filtering algorithm, often referred to simply as a *filter*, describes a way of combining the observed data with a mathematical model of the system to produce an estimate of the true state of the system at each time step. These algorithms account for the error in the measurements and also rely on the statistical properties of the system to provide the best possible estimate of the true state of the system. In this chapter, our focus is on Bayesian filtering theory and the filters which are based on approximating the Bayesian posterior computations using some observation likelihood for dynamical systems.

Before we dive into our discussion of sequential state estimation and Bayesian filtering, we briefly describe the variational data assimilation formalism; the other approach to DA. Variational data assimilation algorithms are routinely used in operational centers for ocean and weather prediction. This is based on the least-squares formulation of DA that relies on an objective or loss function, a quantity that represents the mismatch between observations over time and a given background state of the system. Sasaki introduced the variational approach in one of his seminar papers and the different formalisms for applying tools from numerical variational analysis to the problem of numerical weather prediction [78, 27].

The 4D-Var approach includes a time component in its objective function. The solution at the initial time determines a best-fit trajectory over an observation window, where we have a sequence of observations. The methods to solve 4D-Var data assimilation draw heavily on optimization theory and use adjoint-based gradient optimization methods [36]. The optimal solution is an initial condition which, when propagated forward in time, generates a model trajectory. The two types of formalism included the weak and the strong constraint 4D-Var, with the former being formulated for a perfect model scenario and the latter relaxed this assumption thus accounting for model errors. Some examples are NCEP, ECMWF, and UKMet Office which use some version of 4D-Var for global numerical weather prediction [3, 36].

In the following section, we introduce a discrete-time dynamical system model and then formally define the filtering problem with all the mathematical details. We then proceed to describe the state estimation problem of determining the probability distributions on the sequence of states and given the observation sequence up to a certain time  $k$ .

## 2.2 The state estimation problem for a discrete-time dynamical systems

We now introduce the filtering problem for the most common scenario where we have sequential observations arriving in time. A general discrete-time dynamical system can be used to model both processes that occur in discrete-time steps and continuous-time processes where measurements are sampled at certain fixed intervals.

Assume that we have a linear dynamical system where the evolution of the hidden state of the system  $\mathbf{x}_k$  over time. What we observe indirectly through measurements is a sequence of observations  $\{\mathbf{y}_0, \mathbf{y}_1, \dots, \mathbf{y}_k\}$  at time  $t_0, t_1, \dots, t_k$ . The state-space formulation of this problem can be written in terms of a discrete-time propagator, which is a matrix that maps the state from time  $t_{k-1}$  to  $t_k$  with a linear observation operator that specifies the relationship between the measurement  $\mathbf{y}_k$  with the state  $\mathbf{x}_k$  at any time  $t_k$ .

$$\mathbf{x}_{k+1} = \mathbf{F}_k(\mathbf{x}_k) + \eta_k \quad (2.2)$$

$$\mathbf{y}_k = \mathbf{H}(\mathbf{x}_k) + \epsilon_k \quad (2.3)$$

where,

$\mathbf{x}_k \in R^d$ : state of the system.

$\mathbf{F}_k \in R^{d \times d}$ : discrete propagator from  $t_{k-1}$  to  $t_k$ .

$\eta_k \in R^d$ : model error with distribution  $\mathcal{N}(0, \mathbf{Q}_k)$ .

$\mathbf{H} \in R^{p \times d}$ : the observation operator of the model.

$\mathbf{y}_k \in R^p$ : the observation obtained at time  $t_k$  with  $p < d$ .

$\epsilon_k \in R^p$ : observation error with distribution  $\mathcal{N}(0, \mathbf{R}_k)$ .

The general state estimation problem for the setup described by equations (2.2)-(2.3) is as follows: determine the sequence of states  $\{\mathbf{x}_0, \mathbf{x}_1, \dots, \mathbf{x}_k\}$  given  $\{\mathbf{y}_0, \mathbf{y}_1, \dots, \mathbf{y}_k\}$ . Inverse problems of this nature where the state lies in a high-dimensional space, while observations are low-dimensional and noisy, are poorly posed in the Hadamard sense [7, Ch. 2, p 34] as one cannot simply invert the observations to estimate the state. But from a probabilistic view, multiple states are probable, and hence it is possible to quantify a distribution on the state space and quantify all the possible states under certain conditions.

From a probabilistic perspective,  $\{\mathbf{x}_k, \mathbf{y}_k\}$  are random variables with a joint distribution. The given set of observed values  $\mathbf{y}_k$  are realizations of a random variable; it is possible to determine the probability of the simultaneous occurrence of various values  $\mathbf{x}_k$  if there

is some statistical or dynamic relationship. We can determine the probability of the occurrence of a particular observed value  $\mathbf{y}_k$  given some fixed value of the underlying state  $\mathbf{x}_k$  and the information of the noise distribution.

Our goal is to estimate the probability distribution of the sequence of states  $\mathbf{x}_{0:k}$  given the sequence of observations  $\{\mathbf{y}_{0:k}\}$ . This estimation problem can be restructured and categorized into three different but related problems: (i) filtering, (ii) smoothing and (iii) prediction. Filtering corresponds to the use of past records up to the present time  $t_k$  to estimate the distribution  $\pi(\mathbf{x}_k|\mathbf{y}_{0:k})$ . This is a relevant problem in practical scenarios where only information history is available to us and we want to estimate/predict the present in the best possible way. The problem of smoothing takes into account future observations in order to determine the state estimate of the past i.e.  $\pi(\mathbf{x}_j|\mathbf{y}_{0:k})$  for  $j < k$ . Smoothing results in more accurate estimates, as they contain information from future observations as well. The prediction step aims to find the distribution of the future state based on the history of observations  $\pi(\mathbf{x}_j|\mathbf{y}_{0:k})$  with  $j > k$ .

## 2.3 Background: Bayesian filtering theory

In this section, we now discuss Bayesian filtering theory, which closely follows the discussion in [77, Ch.4, p.54]. We note that we use  $\pi(\mathbf{x})$  and  $\rho(\mathbf{x})$  to represent the probability distribution and its associated probability density function, respectively.

If we place a prior distribution on the states  $\pi(\mathbf{x}_{0:k})$  with the corresponding density given by  $\rho(\mathbf{x}_{0:k})$ , then using the observation likelihood  $\rho(\mathbf{y}_{0:k}|\mathbf{x}_{0:k})$  obtained under a given observation model, then using Bayes theorem, the joint posterior density of the states  $\mathbf{x}_{0:k}$ , given by

$$\rho(\mathbf{x}_{0:k}|\mathbf{y}_{0:k}) = \frac{1}{\mathcal{Z}} \rho(\mathbf{y}_{0:k}|\mathbf{x}_{0:k}) \rho(\mathbf{x}_{0:k}) \quad (2.4)$$

where  $\mathcal{Z}$  is the normalization factor given by

$$\mathcal{Z} = \int \rho(\mathbf{y}_{1:k}|\mathbf{x}_{0:k}) \rho(\mathbf{x}_{0:k}) d\mathbf{x}_{0:k} \quad (2.5)$$

When the set of observations  $\{\mathbf{y}_k\}$  arrives serially in time, then in long time, as  $k \rightarrow \infty$ , the computation of the posterior in (2.4) distributions becomes intractable. A simple trade-off is that when the joint posterior of the whole state sequence is not required, one can find specific conditional distributions of interest as follows. Under the assumption that the hidden states  $\mathbf{x}_k$  follow Markovian dynamics, i.e., given  $\mathbf{x}_{k-1}$ ,  $\pi(\mathbf{x}_k|\mathbf{x}_{0:k-1}) = \pi(\mathbf{x}_k|\mathbf{x}_{k-1})$  together with conditional independence of the observation  $\mathbf{y}_k$  given the state  $\mathbf{x}_k$  i.e. , the distribution of the observations is independent of all the past and future states  $\pi(\mathbf{y}_k|\mathbf{x}_{0:k}, \mathbf{y}_{1,k}) = \pi(\mathbf{y}_k|\mathbf{x}_k)$ , the joint estimation problem can be broken into a sequential estimation of conditional state distribution, also known as the Bayesian



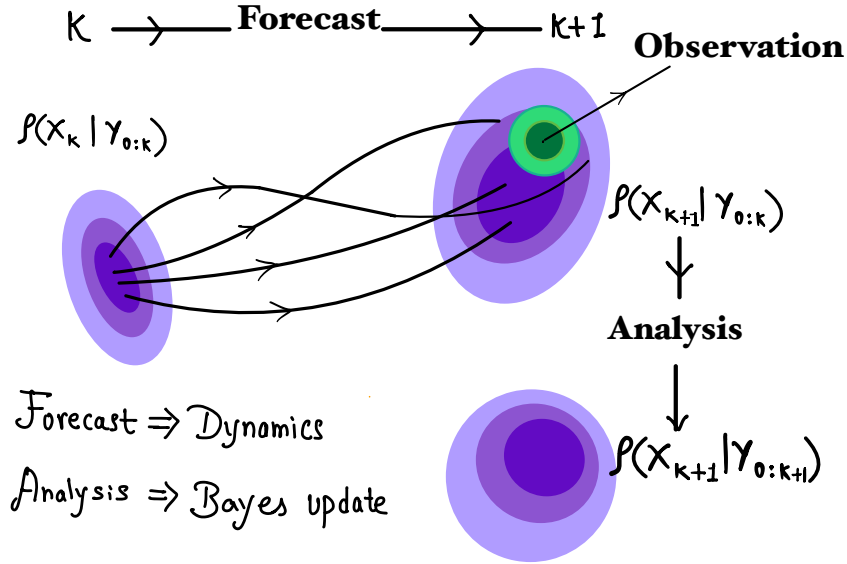


Figure 2.1: The schematic diagram representing the two steps of Bayesian filtering algorithm.

filtering problem.

Under an assumed model of the dynamics and likelihood of the observations, Bayesian filtering deals with the sequential estimation of the distribution of the state of the system conditioned on the sequence of observations  $y_{0:k}$  up to certain time  $t_k$ , known as the filtering or the analysis distribution [26]. This is followed by prediction, where the goal is to predict the future distribution, known as the forecast distribution of the state based on the system dynamics by accounting for possible sources of error at the future time  $t_{k+1}$ .

Bayesian filtering starts with an initial distribution  $\pi(\mathbf{x}_0)$  with  $\rho(x_0)$  as the density function for the state at  $t_0$ . At time  $t_k$ , we represent the filtering or analysis distribution by  $\pi(\mathbf{x}_k | \mathbf{y}_{0:k})$ , the distribution of the state conditioned on the history of observations up to time  $t_k$ . This is then propagated in time by solving the fokker-planck equation [42] to obtain the prior or forecast distribution at time  $t_{k+1}$ , which is the probability distribution of the state conditioned on observations only upto time  $t_k$ . At time  $t_{k+1}$ , the the forecast or the predictive distribution and is given by  $\pi(\mathbf{x}_{k+1} | \mathbf{y}_{0:k})$ , Using the Chapman-Kolmogorov equation [26], the corresponding forecast distribution density at time  $t_{k+1}$  is given by

$$\rho(\mathbf{x}_{k+1} | \mathbf{y}_{0:k+1}) = \int \rho(\mathbf{x}_{k+1} | \mathbf{x}_k) \rho(\mathbf{x}_k | \mathbf{y}_{0:k}) d\mathbf{x}_k \quad (2.6)$$

When the observation arrives at time  $t_{k+1}$ , one computes the posterior distribution  $\pi(\mathbf{x}_{k+1} | \mathbf{y}_{0:k+1})$  and it's density represented by  $\rho(\mathbf{x}_{k+1} | \mathbf{y}_{0:k+1})$  using Bayes theorem.

$$\rho(\mathbf{x}_{k+1} | \mathbf{y}_{0:k+1}) = \frac{\rho(\mathbf{y}_{k+1} | \mathbf{x}_{k+1}, \mathbf{y}_{0:k}) \rho(\mathbf{x}_{k+1} | \mathbf{y}_{0:k})}{\rho(\mathbf{y}_{k+1} | \mathbf{y}_{0:k})} \quad (2.7)$$

Assuming that given  $\mathbf{x}_{k+1}$ , the observation  $\mathbf{y}_{k+1}$  is conditionally independent of all previous observations and states, we have  $\pi(\mathbf{y}_k|\mathbf{x}_{0:k}, \mathbf{y}_{1:k}) = \pi(\mathbf{y}_k|\mathbf{x}_k)$ , which results in a simplified expression for the observation likelihood. Using the observation model and the measurement noise distribution, one can obtain the observation likelihood at time  $t_{k+1}$  given by  $\pi(\mathbf{y}_{k+1}|\mathbf{x}_{k+1})$  and equation (2.7) simplifies into

$$\rho(\mathbf{x}_{k+1}|\mathbf{y}_{0:k+1}) = \frac{\rho(\mathbf{y}_{k+1}|\mathbf{x}_{k+1})\rho(\mathbf{x}_{k+1}|\mathbf{y}_{0:k})}{\rho(\mathbf{y}_{k+1}|\mathbf{y}_{0:k})} \quad (2.8)$$

Figure 2.1 represents the schematic picture of the two steps involved in Bayesian filtering depicting the density updates in time. In practice, any Bayesian filtering algorithm tries to numerically approximate equation (2.6) and (2.8) under different assumptions and approximations in order to solve a sequential estimation problem.

## 2.4 Background: The Kalman Filter – a recursive solution to filtering of linear dynamical systems

We now present the optimal solution of the Bayesian filtering problem for the case of a linear dynamical system which was introduced earlier in section 2.2. These set of equations were developed by Kalman in his seminal work [44], and are now widely known as Kalman filter equations. They represent a set of update equations used to obtain the filtering and the predictive distributions in equation (2.8) and (2.6) respectively, as a closed-form recursive solution. We later present a complete derivation of Kalman filter equations using the best linear unbiased estimator, also known as BLUE, the statistical approach to obtaining the Kalman filter equations. Another approach using the Bayesian formulation to arrive at the Kalman filter equations involves products of Gaussian distributions and integration of the normalization terms, for which we refer to [77, Ch. 4 ,p. 56].

For the case of a linear dynamical system with linear observation operator and Gaussian measurement noise as introduced in section 2.2, when we start with initial distribution as Gaussian, all the resulting filtering and predictive distributions at any time  $t_k$  are Gaussian. Since we can completely specify a Gaussian distribution by its mean and covariance, we find that we can obtain equations (2.8) and (2.6) in terms of updating the mean and covariance sequentially over time. The updated mean and covariance at time  $t_{k+1}$  are related to the mean and covariance of the predictive distribution at the same time. These updates are computed recursively in two steps, starting with a prediction step from time  $t_k$ , where we use  $\pi(\mathbf{x}_k|\mathbf{y}_{1:k})$  and the hidden dynamics to obtain the predictive distribution  $\pi(\mathbf{x}_{k+1}|\mathbf{y}_{1:k})$ . When the observation  $\mathbf{y}_{k+1}$  at time  $t_{k+1}$  is available, we update the mean and covariance again to obtain the filtering distribution  $\pi(\mathbf{x}_{k+1}|\mathbf{y}_{1:k+1})$  using Bayesian principles.

We now define a notation for our convenience in the case of the Gaussian probability density function. For a  $d$ -dimensional Gaussian distribution  $\mathcal{N}(\mathbf{m}, \mathbf{P})$  in which is parameterised by  $\mathbf{m}$  and  $\mathbf{P}$ , let the associated probability density be represented by  $\rho_{\mathcal{N}}(\mathbf{x}; \mathbf{m}, \mathbf{P})$ , the expression for which is given by

$$\rho_{\mathcal{N}}(\mathbf{x}; \mathbf{m}, \mathbf{P}) = \frac{1}{\sqrt[2]{2\pi \det(\mathbf{P})}} \exp\left(-\frac{1}{2}(\mathbf{x} - \mathbf{m})^T \mathbf{P}^{-1}(\mathbf{x} - \mathbf{m})\right). \quad (2.9)$$

We assume that the prior distribution of the state  $\mathbf{x}_0$  at time  $t_0$  is a Gaussian distribution, that is,  $\mathbf{x}_0 \sim \mathcal{N}(\mathbf{m}_0, \mathbf{P}_0)$  where  $\mathbf{m}_0$  and  $\mathbf{P}_0$  denote the mean and covariance of the Gaussian distribution. Starting at time  $t_k$ , the dynamics of the systems completely specify the future distribution at time  $t_{k+1}$ . Using the Markovian assumption, which follows from the linear dynamics, the conditional distribution of the state  $\mathbf{x}_k$  given the state  $\mathbf{x}_{k-1}$  is given by

$$\pi(\mathbf{x}_{k+1}|\mathbf{x}_k) = \mathcal{N}(\mathbf{F}_k \mathbf{x}_k, \mathbf{Q}_k), \quad (2.10)$$

where  $\mathbf{Q}_k$  is the covariance of the process noise or model error. When the process noise or model error is zero, i.e. they have deterministic dynamics since the previous state  $x_k$  completely determines the next state  $\mathbf{x}_{k+1}$ . In this case we obtain a dirac-delta distribution function  $\delta(\mathbf{x}_{k+1} - \mathbf{F}_k \mathbf{x}_k)$  which is centered on  $\mathbf{F}_k \mathbf{x}_k$ .

At time  $t_{k+1}$ , we assume that the forecast or the prior distribution  $\pi(\mathbf{x}_{k+1}|\mathbf{y}_{0:k})$  is represented by  $\mathcal{N}(\mathbf{m}_{k+1}^f, \mathbf{P}_{k+1}^f)$ . To find the mean and covariance of this prior distribution at time  $t_{k+1}$ , a simple approach is to find the mean and covariance using  $E[\mathbf{x}_{k+1}|\mathbf{y}_{0:k}] = \mathbf{F}_k \mathbf{m}_k^a$  and  $E[(\mathbf{x}_{k+1} - \mathbf{m}_{k+1}^f)(\mathbf{x}_{k+1} - \mathbf{m}_{k+1}^f)^T] = \mathbf{F}_k \mathbf{P}_k^a \mathbf{F}_k^T$ .

Thus, the forecast distribution  $\pi(\mathbf{x}_{k+1}|\mathbf{y}_{0:k}) = \mathcal{N}(\mathbf{m}_{k+1}^f, \mathbf{P}_{k+1}^f)$  where,

$$\begin{aligned} \mathbf{m}_{k+1}^f &= \mathbf{F}_k \mathbf{m}_k^a \\ \mathbf{P}_{k+1}^f &= \mathbf{F}_k \mathbf{P}_k^a \mathbf{F}_k^T \end{aligned} \quad (2.11)$$

We also assume that the measurement error  $\epsilon_{k+1}$  in the observation is assumed to be distributed normally, i.e.  $\epsilon_{k+1} \sim \mathcal{N}(0, \mathbf{R}_{k+1})$  with the fact that  $\mathbf{y}_{k+1} = \mathbf{H}\mathbf{x}_{k+1} + \epsilon_{k+1}$ , the likelihood of observing  $\mathbf{y}_k$  given the state  $\mathbf{x}_k$  is given by,

$$\pi(\mathbf{y}_{k+1}|\mathbf{x}_{k+1}) = \mathcal{N}(\mathbf{H}\mathbf{x}_{k+1}, \mathbf{R}_{k+1}).$$

The distribution  $\pi(\mathbf{y}_k | \mathbf{y}_{0:k-1}) = \mathcal{N}(\mathbf{H}\mathbf{m}_k^f, \mathbf{H}\mathbf{P}_k^f \mathbf{H}^T + \mathbf{R}_k)$  is the normalization distribution which appears in the denominator in equation.

We now have all the essential ingredients in order to solve the sequential filtering problem. Let us denote the filtering distribution at time  $t_{k+1}$  by  $\pi(\mathbf{x}_{k+1}|\mathbf{y}_{0:k+1}) = \mathcal{N}(\mathbf{m}_{k+1}^a, \mathbf{P}_{k+1}^a)$ .

Using the bayes threom in equation (2.8) and substituting all the Gaussian density functions defined in equation (2.9), we obtain

$$\rho_{\mathcal{N}}(\mathbf{x}_{k+1}; \mathbf{m}_{k+1}^a, \mathbf{P}_{k+1}^a) = \frac{\rho_{\mathcal{N}}(y_{k+1}; \mathbf{H}\mathbf{x}_{k+1}, \mathbf{R}_{k+1}) \rho_{\mathcal{N}}(x_{k+1}; \mathbf{m}_{k+1}^f, \mathbf{P}_{k+1}^f)}{\rho_{\mathcal{N}}(y_{k+1}; \mathbf{H}\mathbf{m}_{k+1}^f, \mathbf{H}\mathbf{P}_{k+1}^f\mathbf{H}^T + \mathbf{R}_{k+1})}. \quad (2.12)$$

where  $\rho_{\mathcal{N}}(\mathbf{x}_{k+1}; \mathbf{m}_{k+1}^a, \mathbf{P}_{k+1}^a)$  denotes the resulting density for the filtering distribution  $\mathcal{N}(\mathbf{m}_{k+1}^a, \mathbf{P}_{k+1}^a)$ .

The mean and covariance  $\mathbf{m}_{k+1}^a$  and  $\mathbf{P}_{k+1}^a$  of the filtering distribution at time  $t_{k+1}$  can be expressed in terms of the mean and covariance of the predictive distribution  $\mathbf{m}_{k+1}^f$  and  $\mathbf{P}_{k+1}^f$  and are given by,

$$\mathbf{m}_{k+1}^a = \mathbf{m}_{k+1}^f + \mathbf{K}_{k+1} \left[ \mathbf{y}_{k+1} - \mathbf{H}\mathbf{m}_{k+1}^f \right] \quad (2.13)$$

$$\mathbf{P}_{k+1}^a = (\mathbf{I} - \mathbf{K}_{k+1}\mathbf{H}) \mathbf{P}_{k+1}^f \quad (2.14)$$

where, the matrix  $\mathbf{K}_{k+1}$  is called the Kalman gain matrix given by

$$\mathbf{K}_{k+1} = \mathbf{P}_{k+1}^f \mathbf{H}^T \left[ \mathbf{H}\mathbf{P}_{k+1}^f \mathbf{H}^T + \mathbf{R}_{k+1} \right]^{-1} \quad (2.15)$$

This recursive solution is an online solution, i.e., with more observation over time, it computes the best state estimate and its distributions in closed-form only using the update equations (2.11) and (2.13) for mean and covariance.

We now present a derivation of the Kalman gain matrix in (2.15) which closely follows the notation in [14]. Let us define the state estimate at time  $k$  before observation, which is denoted by  $\hat{x}_k^f$  and the observation as  $\hat{x}_k$ , where we use the "hat" to denote the estimate. Note that these state estimates are random vectors and are characterized by their relative distributions.

We start at time  $t_{k+1}$ , where the distribution of our state is  $\mathcal{N}(\mathbf{m}_{k+1}^f, \mathbf{P}_{k+1}^f)$ . Let us denote the error in our estimate by  $\varepsilon_k$ , given by

$$\varepsilon_{k+1}^f = \mathbf{x}_{k+1} - \hat{\mathbf{x}}_{k+1}^f \quad (2.16)$$

$$\mathbf{P}_{k+1}^f = \mathbf{E} \left[ \varepsilon_{k+1}^f (\varepsilon_{k+1}^f)^T \right] \quad (2.17)$$

Our goal is to improve our state estimate using the information from the observation  $y_k$ . The difference between true and predicted observations at any time  $t_{k+1}$  is called innovation, which is denoted by  $I_{k+1} = \mathbf{y}_{k+1} - \mathbf{H}\mathbf{x}_{k+1}$ . This is the correction, which is

weighted and added to the the initial state estimate at time  $t_{k+1}$  before the observation  $\mathbf{y}_{k+1}$  was assimilated. Our goal is to determine a matrix  $\mathbf{K}_{k+1}$  such that, it produces a correction in the previous state estimate  $\mathbf{x}_{k+1}^f$  of the following form,

$$\mathbf{x}_{k+1}^a = \mathbf{x}_{k+1}^f + \mathbf{K}_{k+1} \left( \mathbf{y}_{k+1} - \mathbf{H}\mathbf{x}_{k+1}^f \right). \quad (2.18)$$

We now establish the criteria to determine the optimal combination of  $\mathbf{x}_{k+1}^f$  and  $\mathbf{y}_{k+1} - \mathbf{H}\mathbf{x}_{k+1}^f$  which is determined by  $\mathbf{K}_{k+1}$ . We first compute the expected error covariance matrix corresponding to our new estimate  $\mathbf{x}_{k+1}^a$  which is given by,

$$\begin{aligned} \varepsilon_{k+1}^a &= \mathbf{x}_{k+1} - \mathbf{x}_{k+1}^a \\ &= \mathbf{x}_{k+1} - \left( \mathbf{x}_{k+1}^f + \mathbf{K}_{k+1} (\mathbf{y}_{k+1} - \mathbf{H}\mathbf{x}_{k+1}^f) \right) \\ &= \mathbf{x}_{k+1} - \mathbf{x}_{k+1}^f - \mathbf{K}_{k+1} (\mathbf{H}\mathbf{x}_{k+1}^f + \epsilon_{k+1} - \mathbf{H}\mathbf{x}_{k+1}^f) \\ &= \mathbf{x}_{k+1} - \mathbf{x}_{k+1}^f - \mathbf{K}_{k+1} \mathbf{H} \left( \mathbf{x}_{k+1} - \mathbf{x}_{k+1}^f \right) - \mathbf{K}_{k+1} \epsilon_{k+1} \\ &= \varepsilon_{k+1}^f - \mathbf{K}_{k+1} \mathbf{H} \varepsilon_{k+1}^f - \mathbf{K}_{k+1} \epsilon_{k+1} \\ &= (\mathbf{I} - \mathbf{K}_{k+1} \mathbf{H}) \varepsilon_{k+1}^f - \mathbf{K}_{k+1} \epsilon_{k+1} \end{aligned} \quad (2.19)$$

Since error in the state estimate  $\varepsilon_{k+1}^f$  and the measurement error  $\epsilon_{k+1}$  are uncorrelated,  $\mathbf{E} \left[ \varepsilon_{k+1}^f (\epsilon_{k+1})^T \right] = 0$ , the posterior covariance is given by,

$$\begin{aligned} \mathbf{P}_{k+1}^a &= \mathbf{E} \left[ \varepsilon_{k+1}^a (\varepsilon_{k+1}^a)^T \right] \\ &= (\mathbf{I} - \mathbf{K}_{k+1} \mathbf{H}) \mathbf{E} \left[ \varepsilon_{k+1}^f (\varepsilon_{k+1}^f)^T \right] (\mathbf{I} - \mathbf{K}_{k+1} \mathbf{H})^T + \mathbf{K}_{k+1} \mathbf{E} \left[ \varepsilon_{k+1} \varepsilon_{k+1}^T \right] \mathbf{K}_{k+1} \\ &= (\mathbf{I} - \mathbf{K}_{k+1} \mathbf{H}) \mathbf{P}_{k+1}^f (\mathbf{I} - \mathbf{K}_{k+1} \mathbf{H})^T + \mathbf{K}_{k+1} \mathbf{R}_{k+1} \mathbf{K}_{k+1}^T. \end{aligned} \quad (2.20)$$

The optimal criterion to determine the gain matrix  $\mathbf{K}_{k+1}$  in equation (2.18) is such that the expected error covariance is minimized. Since the trace of a covariance matrix represents the total error variance summed over the individual components, we set the trace of  $\mathbf{P}_{k+1}^a$ ,

$$\frac{d \operatorname{tr}(\mathbf{P}_{k+1}^a)}{d \mathbf{K}_{k+1}} = 0 \quad (2.21)$$

which involves derivative of a scalar with respect to a matrix.

We first expand the expression of analysis covariance, given by,

$$\mathbf{P}_{k+1}^a = \mathbf{P}_{k+1}^f - \mathbf{K}_{k+1} \mathbf{H} \mathbf{P}_{k+1}^f - \mathbf{P}_{k+1}^f \mathbf{H}^T \mathbf{K}_{k+1}^T + \mathbf{K}_{k+1} \left( \mathbf{H} \mathbf{P}_{k+1}^f \mathbf{H}^T + \mathbf{R}_{k+1} \right) \mathbf{K}_{k+1}^T. \quad (2.22)$$

Using matrix differentiation rules, the derivative of the equation (2.22) results in the following,

$$\frac{d\text{Trace}(\mathbf{P}_{k+1}^a)}{d\mathbf{K}_{k+1}} = -2 \left( \mathbf{H}\mathbf{P}_{k+1}^f \right)^T + 2\mathbf{K}_{k+1} \left( \mathbf{H}\mathbf{P}_{k+1}^f \mathbf{H}^T + \mathbf{R}_{k+1} \right) \quad (2.23)$$

From equation (2.21) and (2.23), we get

$$\mathbf{K}_{k+1} = \mathbf{P}_{k+1}^f \mathbf{H}^T \left( \mathbf{H}\mathbf{P}_{k+1}^f \mathbf{H}^T + \mathbf{R}_{k+1} \right)^{-1}. \quad (2.24)$$

Not only does the optimal combination generated by the Kalman gain minimize the mean-square estimation error, but it also simplifies the expression in equation (2.22) in an even simpler form, given by,

$$\mathbf{P}_{k+1}^a = (\mathbf{I} - \mathbf{K}_{k+1}\mathbf{H}) \mathbf{P}_{k+1}^f \quad (2.25)$$

Let us now understand the role of the Kalman gain matrix and its interpretation in two different cases as follows. When observations tend to be perfect,  $\mathbf{R}_k \rightarrow \mathbf{0}$ , and  $\lim_{\mathbf{R}_k \rightarrow 0} \mathbf{K}_k = \mathbf{H}^\dagger$ , the observations are heavily trusted. In the limit  $\lim_{\mathbf{P}_k^f \rightarrow 0} \mathbf{K}_k = 0$ , the filter relies heavily on the model estimates, ignoring the measurements.

## 2.5 Ensemble Kalman filter: extending Kalman filters to nonlinear dynamical systems with monte-carlo approach

From the above discussion, we know that Kalman filters provide closed-form solutions to the Bayesian filtering equations in the case where both the dynamics and measurement models are linear with Gaussian measurement error distribution. In this section, we discuss the ensemble Kalman filter or EnKF, a Monte Carlo approximation of the original Kalman filter, which was introduced by G. Evensen in [34]. Since its introduction, there has been a lot of research and development on improving the standard version of EnKF. But all the different formulations of EnKF are built on these basic sets of ideas; for a more comprehensive discussion of EnKF and its variants, we refer to [7, Ch.6, p.153]. However, the above Kalman filter equations are true only for linear dynamical systems where the Gaussian distribution is preserved by linearity. Thus, it is not applicable in general to nonlinear systems and is difficult to handle for very high-dimensional systems. For nonlinear problems, linearisation can be done around the best present filter estimate, which can be used to propagate the covariance matrices, which leads to the formulation of an extended Kalman filter [77, ch.5, p.69]. But in the presence of nonlinearity in the

dynamics, there is no restriction on why the evolution of the distribution of states should be Gaussian.

EnKF uses the Monte Carlo approach to implement the Kalman filter equations (2.11) and (2.13). Since a probability distribution can also be represented by an ensemble of states, EnKF represents the probability distributions via their ensemble, each of size  $N$ , which is sampled at the beginning at  $t_0$  from the initial distribution. In the case of Gaussian distributions, the ensemble can be easily generated using the mean and covariance matrix. In between the observations, each member of the ensemble is integrated forward in time using the dynamics, hence propagating the probability distribution from one time step to the other until the next observation arrives. At any time, the spread of the ensemble is a measure of uncertainty in the mean of the ensemble. To determine the ensemble representing the filtering distribution, the mean and covariance required in equations (2.11) and (2.13) are then estimated from the ensemble. Each of the individual members of the ensemble is then updated via a rule inspired by the original Kalman filter equations. In order to account for the observation error statistics, we generate  $N$  realizations of the observation by adding randomly generated noise using the knowledge of the noise distribution. This scheme is also known as perturbed observation EnKF [34].

We generally choose the size of the ensemble so that the estimated covariance matrix of the ensemble is close to the true covariance. Definitely, as we increase the number of members in the group, the more accurately it represents the pdf of the initial and forecast and will represent the exact distribution in the limit  $N \rightarrow \infty$ . However, as we increase the ensemble size, we often see a saturation in RMSE beyond a particular number  $N$ . Using this information, we can choose a cutoff on the size of the ensemble since a further increase in ensemble size is insignificant to the errors.

Despite its advantages, the EnKF without further improvements suffers from important limitations. EnKF assumes that the errors in the observations and the model are independent and identically distributed (i.i.d.), which may not be the case in practice. Additionally, the EnKF suffers from sampling errors, especially in high-dimensional systems. An artifact of using a small ensemble size in covariance estimation for EnKF leads to long-range correlations and an underestimation of the covariance. Localization is the procedure of eliminating such spurious correlations. To deal with the underestimation of covariance matrices, we artificially inflate the covariance matrix before updating the ensemble with the latest observation and this is referred to as Inflation. Together, they play an important role in making the EnKF highly applicable and wildly popular for high-dimensional data assimilation problems in spatially extended systems. Now we discuss these two diagnostic procedures that help overcome the shortcomings of EnKF and make it a practical data assimilation algorithm. For a detailed account of EnKF, we refer to (Chapter Ensemble Kalman Filter: Current Status and Potential) [53]. For a nice historical account of the

---

**Algorithm 1:** EnKF with covariance localization in state-space.  $\circ$  denotes Hadamard product.

---

Initialize  $N$  particles  $\{x_0^i\}_{i=1}^N$  according to the initial distribution and set  $x_0^{i,a} = x_0^i$   
Set  $\rho$  as the Gaspari-Cohn localization matrix [16].

```

for  $k = 1, \dots, n$  do
  for  $i = 1, \dots, N$  do
     $x_k^{i,f} \leftarrow f_g(x_{k-1}^{i,a})$ 
   $m_k^f \leftarrow \frac{1}{N} \sum_i x_k^{i,f}$ 
   $P_k^f \leftarrow \rho \circ \frac{\sum_i (x_k^{i,f} - m_k^f)(x_k^{i,f} - m_k^f)^\top}{N-1}$ 
   $K \leftarrow P_k^f H^\top [H P_k^f H^\top + R_k]^{-1}$ 
  for  $i = 1, \dots, N$  do
    Sample  $\eta_k^i \sim \mathcal{N}(0_q, \sigma^2 I_q)$ 
     $y_k^i \leftarrow y_k + \eta_k^i$ 
     $x_k^{i,a} \leftarrow x_k^{i,f} + K [y_k^i - H x_k^{i,f}]$ 
   $\hat{\pi}_k \leftarrow \frac{1}{N} \sum_{i=1}^N \delta_{x_k^{i,a}}$ 

```

---

ensemble forecasting in numerical weather prediction, see [23, 72].

### 2.5.1 Curse of dimensionality in EnKF

Finite ensemble representation of probability distributions in higher dimensions suffer from what is known as the ‘Curse of dimensionality’ [8] that makes many computations numerically intractable in higher dimensions. This arises simply due to the fact that in high dimensions, as the volume expands exponentially, we need exponential increase in the sample size to represent the distributions. Simple EnKF fails without additional procedures and assumptions. The failure of EnKF can be understood considering the fact that the distance between two random sample points drawn from a distribution scales with the dimension. Distant observables are weakly correlated for short time scales. The low number of samples introduces long-distance spurious correlations in covariance. As a result, the ensemble updates degrade and result in a failure to estimate the trajectory over time. Obtaining accurate background covariance requires a large number of ensemble members. Increasing the ensemble size increases the demand for computational resources.

We illustrate this issue numerically by drawing three samples of size  $N = 15, 70,$  and  $300.$  from a standard multivariate normal distribution. In figure 2.2, we plot the empirical covariance computed from the samples, with presence of non-zero off-diagonal entries.



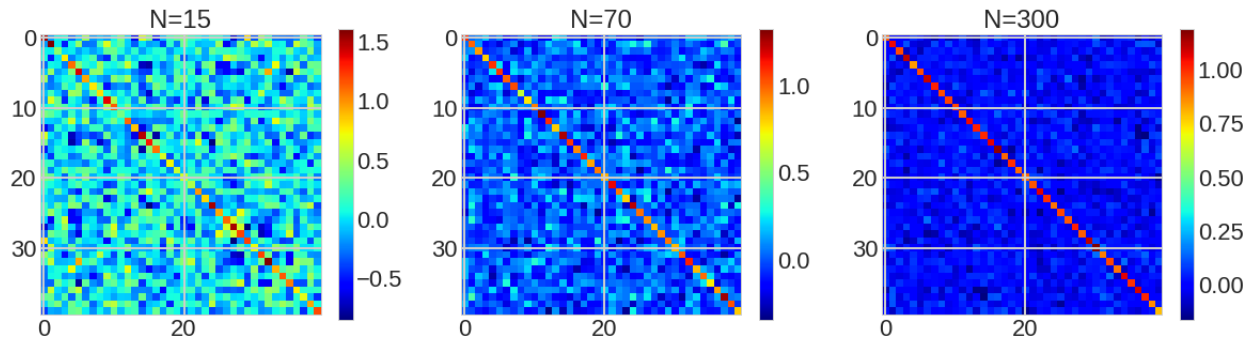


Figure 2.2: Long-distance spurious correlations due to finite sample-size  $N$ . The original covariance is  $I_d$ , the identity matrix.

## 2.5.2 Inflation

Ensemble Kalman filter with a small number of ensemble members for a high-dimensional system encounters the issue of underestimating the forecast error covariance. This leads to overconfident forecasts and, as a result, the observations are given lower weights in the analysis or the update step. An ad hoc way to resolve this issue is to pre-multiply the forecast ensembles before updating by a factor  $\alpha$  such that the entire covariance matrix is multiplied by  $\alpha^2$ .

$$\mathbf{P}_k^f \rightarrow \alpha^2 \mathbf{P}_k^f \quad (2.26)$$

This small modification has been shown to improve forecasts and analysis. However, choosing the “right” value of  $\alpha$  is often ad hoc. The right value of  $\alpha$  is chosen based on empirical methods using trial and error and is tuned according to the performance of the filter.

## 2.5.3 Localization

As discussed earlier, estimating the covariance using a small ensemble size leads to spurious correlations among points which are far in physical space. A simple approach to address this issue is an ad hoc method called localization. The idea behind this heuristic approach is the locality hypothesis: Two points that are far apart are independent or uncorrelated at small time scales. Hence, the correlations between any two points must decay with distance eventually decaying to zero. Since the empirical ensemble covariance is only an estimator of the true covariance, we can pre-process the covariance so that it does not have spurious correlations beyond a certain length scale. Physically, this means that only observations within a certain distance from a grid point actually contribute to the modification of the value of the ensemble members at the respective point. For a more comprehensive discussion about localization, see [16].

A simple way to achieve this is to modify the Kalman gain matrix by pre-multiplying the forecast covariance matrix by another positive semi-definite matrix  $\rho$  called a localization matrix. This localization matrix can be constructed using any function that is symmetric and zero outside the localization length scale. The type of localization function and its respective length scale, denoted by  $\ell$ , decide the number of points in the grid beyond which we assume the correlations to be spurious in nature. We use the modified forecast covariance matrix in the Kalman filtering equations where the modified forecast covariance is obtained by the Schur of the forecast covariance  $\mathbf{P}_k^f$  with a localization matrix  $\rho$

$$\mathbf{P}_k^f \rightarrow \rho \circ \mathbf{P}_k^f, \quad (2.27)$$

where  $\circ$  denotes the pointwise multiplication of the two matrices. In principle, the localization functions used on the basis of how correlations of nearby variables depend on each other depend on the dynamics and are local in phase space. We will revisit some localization functions and their matrices in subsection 2.8.4, where we conduct twin experiments to analyze their impact on filtering performance.

## 2.6 Chaotic dynamical systems used in this thesis

Modern computers have opened a new paradigm in which we can study and understand the behavior of complex systems by simulating them according to our needs. The set of equations which govern or approximate the dynamics is then used to produce solutions using numerical methods [62]. Different equations and models can be defined that attempt to emulate the physical system, giving new ways to understand them. Surprisingly, chaotic dynamical systems pose a more difficult problem in this respect, where starting from two similar states, the behavior of the system leads to totally distinct states in the future. This defining feature of deterministic chaos is the sensitivity of the flow of the system to errors in the initial condition, leading to exponential departure of the initially nearby trajectories.

Simple models help to understand complex phenomena, where we only retain some of the features while still being able to reproduce desired behavior. Dynamical systems that replicate rich dynamics and behaviors have been used periodically to study and investigate predictability, a fundamental question [58]. Toy models serve as a simple and safe playground for testing and developing ideas for a larger goal, such as numerical weather prediction. Since their introduction to the data assimilation community, Lorenz models have been one of the well-established and widely used toy models in the research community to perform numerical experiments using different assimilation algorithms [62, 53, 16]. We now discuss the two Lorenz models, both of which are chaotic, dissipative, and have global

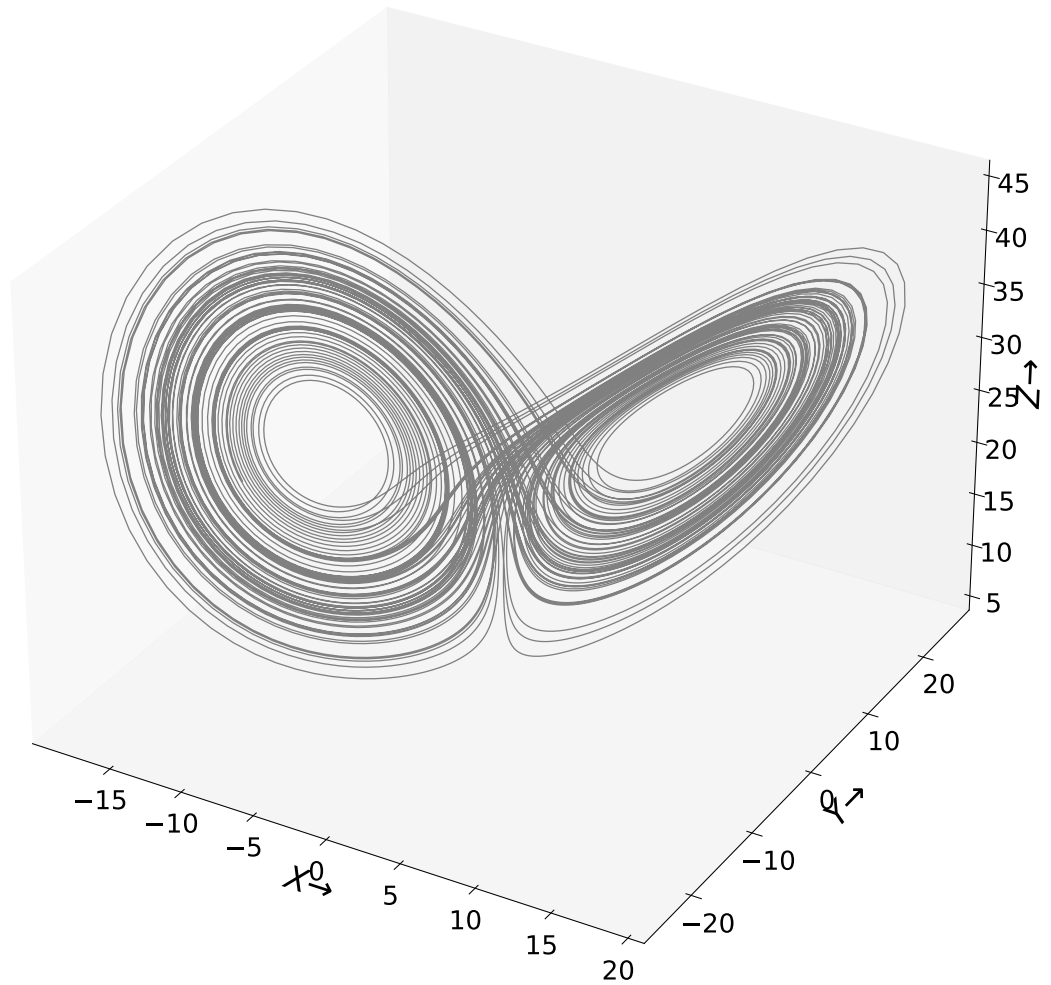


Figure 2.3: The butterfly shaped attractor of Lorenz-63 system for the specific parameter values  $(\sigma, \rho, \beta) = (10, 28, 8/3)$ .

attractor sets. We will refer to these dynamical systems for experiments in chapter 3 for results on nonlinear filter stability and for the computation of Lyapunov vectors in chapter 4.

### 2.6.1 Lorenz 63: A 3-dimensional chaotic ODE

Proposed by Edward Lorenz in the year 1963 [59] as a simplistic model for understanding atmospheric convection; the Lorenz-63 model is a three-dimensional continuous-time dynamical system. As a simple model, it has contributed to various insights into the theory and computation of non-linear and chaotic systems. The coupled nonlinear ODE is

given by,

$$\begin{aligned}\dot{x} &= \sigma(y - x) \\ \dot{y} &= \rho x - y - xz \\ \dot{z} &= -\beta z + xy\end{aligned}$$

with the parameters  $(\sigma, \rho, \beta) = (10, 28, 8/3)$  for which the system exhibits rich chaotic behavior and has the well-known butterfly-shaped attractor which has a fractal dimension less than 3. In figure 2.3, we visualize this attractor by plotting a trajectory obtained after long integration in time.

### 2.6.2 Lorenz 96: A 40-dimensional chaotic ODE

In an attempt to devise a small set of  $d$ -dimensional dissipative and chaotic differential equations, Lorenz proposed another model in 1995, now popularly known as the Lorenz-96 ODE [61]. It is a phenomenological model that is deterministic and chaotic in nature and has been extensively studied and used to test and model chaotic systems and their predictability. It is a system of coupled ordinary differential equations that are continuous in time and discrete in a space-periodic lattice with  $d$  grid points, indexed by  $\mathbf{X}_1, \dots, \mathbf{X}_d$ , is given by the following set of equations;

$$\dot{\mathbf{X}}_i = -\mathbf{X}_i - \mathbf{X}_{i-1}(\mathbf{X}_{i-2} - \mathbf{X}_{i+1}) + \mathbf{F} \quad (2.28)$$

where,  $\mathbf{X}_i$  denotes the variable at the lattice point  $i$  with periodic boundary conditions  $\mathbf{X}_{k+d} = \mathbf{X}_k$  and  $\mathbf{F}$  denotes a constant external forcing. For increasing values of  $\mathbf{F}$ , the behavior of the system changes from stable to weakly and strongly chaotic. For a detailed exposition on the behavior of this model under different regimes of forcing and dimensions, we refer to [62, 49, 84]. For the specific value of forcing  $\mathbf{F} = 8$  in dimension  $d = 40$ , it is a hyperchaotic system with 13 positive Lyapunov exponents and its Kaplan-Yorke dimension is close to 28.4. Due to its rich dynamical properties, it has served as a computationally tractable model for evaluation and analysis of different data assimilation algorithms to benchmark their performance (see, e.g., [17]) before being applied to very large-scale atmospheric models. Figure 2.4 shows a contour plot of a typical trajectory of the system.

Using non-linear, dissipative, and external forcing terms, this autonomous equation is said to mimic the circulation of the earth's atmosphere in an over-simplified manner [62]. Despite their simplicity, these models have had a major effect on the advancement of the theory of dynamical systems, particularly because of their chaotic nature in any number of dimensions.

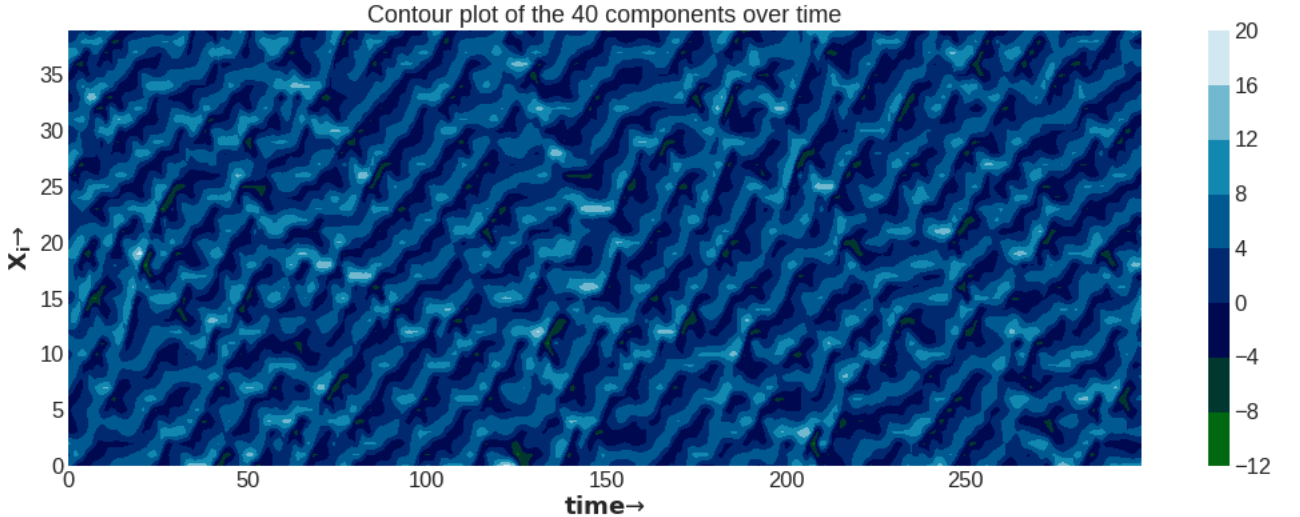


Figure 2.4: The state trajectory as a contour plot for Lorenz-96 system for  $d = 40$  obtained using Runge-Kutta scheme with time step= 0.1 and forcing  $\mathbf{F} = 8$

## 2.7 Some performance metric for numerical filtering algorithms

Assessing the quality of any general filtering algorithm asks the following question: How good is the given filter at extracting information about the underlying signal? This is necessary because evaluating the performance of filtering algorithms is essential, as it helps drive the selection of the best one for specific tasks [54].

We now turn to discuss some of the important metrics in DA which are used to evaluate different filtering algorithms. We describe three well-known metrics methods that are used for this purpose.

### 2.7.1 Accuracy

To quantify the quality of the state reconstruction from noisy observations, the natural question that arises is as follows. What is the error between the reconstructed state and the true underlying state? We compare the trajectory estimated by the filtering algorithm with the actual trajectory to answer this question. The quality of this reconstruction is calculated by the quantity called room mean-square error or simply RMSE, which is given by,

$$\text{RMSE} = \sqrt{\frac{1}{n} \sum_{k=0}^n (\mathbf{x}_k^{\text{true}} - \hat{\mathbf{x}}_k)^2} \quad (2.29)$$

where  $\mathbf{x}_k^{\text{true}}$  is the best state estimate and  $n$  is the number of assimilation steps.

## 2.7.2 Reliability

Rank histograms are a method to analyze the quality of any ensemble forecasting algorithm. The reliability of an ensemble forecast implies that the probability distribution represented by the ensemble adequately captures the uncertainty about the mean of the forecast. Real observations should be seen as a statistically independent draw from the forecast distribution at any time. This implies that the actual observation obtained at time  $t_k$  should be a random sample drawn from the observation distribution corresponding to this forecast ensemble [2].

The idea behind estimating the reliability of the forecast ensembles produced by a model uses the notion of ranking the ensemble members. In general, we can use rank histograms to study specific components or any scalar quantity which is a function of the state which can be used to determine the ranks. We represent a forecast ensemble at some time by  $\{\mathbf{x}^{i,f}\}_{i=1}^N$ , and the true state by  $\mathbf{x}^{true}$ , where we omit the time index for simplifying the notation. The ordered ensemble is represented by  $\{\mathbf{x}^{(i),f}\}_{i=1}^N$  where the ordered ensemble can be ranked to generate a bin based on the ranks after sorting in increasing order of their respective values. At any time  $t_k$ , we expect the observed state  $\mathbf{x}_k$  is equally likely to be in any of the bins generated by the rank partitioning of the  $N$  members of the forecast ensemble at the same time  $t_k$ . Mathematically, we can write,

$$P(\mathbf{x}^{(i),f} \leq \mathbf{x}^{true} < \mathbf{x}^{(i+1),f}) = \frac{1}{N+1}. \quad (2.30)$$

Because the observations are irregularly sampled and sparse in space and time for most of the real physical systems of interest, what is observed is the true state mapped to the observation space. Therefore, we map the forecast ensembles from the state space to the observation space with added observation noise given by  $\mathbf{y}^{i,f} = \mathbf{H}\mathbf{x}^{i,f} + \epsilon^i$ , where  $\epsilon^i$  are drawn using the observation noise distribution. We represent the corresponding ordered observation ensemble as  $\{\mathbf{y}^{(i),f}\}_{i=1}^N$  where the members of the forecast observation ensemble are sorted in an increasing order of their values. Since we have  $N$  members of the ensemble, we have  $N$  points that partition the real line into  $N+1$  bins. Now, if the recorded observation is denoted by  $\mathbf{y}^o$ , the observation is equally likely to lie in any of these bins,

$$P(\mathbf{y}^{(i),f} \leq \mathbf{y}^o < \mathbf{y}^{(i+1),f}) = \frac{1}{N+1} \quad (2.31)$$

To study the reliability of the ensemble forecast, we plot a histogram of different components of the state vectors using the counts per bin of the observation over a large number of forecast steps. For a reliable ensemble forecast, the histogram obtained corresponds to a uniform distribution. We will use this idea to discuss some of the results in the context of twin experiments using EnKF in section 2.8.5.

### 2.7.3 Stability

In previous sections, we introduced the notion of Bayesian filtering where we start with a guess distribution  $\pi(\mathbf{x}_0)$  at  $t_0$ . In reality, the true distribution used in the beginning by any filtering algorithm is based on some ideas of climatology and is far from the actual distribution. Hence it is important that whatever initial distribution we initialize our filter with, the filtering distribution  $\pi(\mathbf{x}_n|\mathbf{y}_{0:k})$  must eventually be independent of the choice of  $\pi(\mathbf{x}_0)$ . Thus, a filtering algorithm is stable if the filtering distributions converge to each other over time. We discuss filter stability and how to numerically study stability of a filtering algorithm in general, and EnKF in particular, which is the main objective of chapter 3. Using EnKF algorithm and the chaotic models introduced in this chapter, we will address the filter stability problem numerically, which is the main content of the next chapter.

## 2.8 Twin experiments and results

In the previous section, we introduced some of the important ways to measure the performance of a filtering algorithm. In reality, assessing the quality of the reconstruction of the underlying signal is difficult; in weather prediction problems, it is impossible to know the infinite-dimensional true underlying state of a real system such as the ocean or atmosphere, and what is generally given to us is their manifestation through the real observations, which are themselves noisy and sparse. Twin experiments allow us to build confidence in different algorithms and benchmark them in terms of their reconstruction quality before finally deploying the algorithms in practice; see [16] for an overview in the context of geosciences. They serve as practical tools for benchmarking various data assimilation algorithms by assimilating artificially generated observations from an artificial truth obtained by integrating the model itself. To use our desired filtering algorithms on a numerical model in twin experiments, we first generate a model trajectory that serves as the true underlying trajectory of the system. Synthetic observations are then generated from this trajectory using the observation operator, mimicking the sequential arrival of the observations. Using filtering techniques on these observations, we reconstruct the state of the system over time and obtain a filter-estimated trajectory. The observations are assimilated to the model to obtain forecast and analysis states which can then be used to compare with the artificial true state using an appropriate metric. For a comprehensive overview of twin experiments in the context of EnKF, where all components of the state are observed, we refer to [7, Ch. 6, p. 172].

### 2.8.1 A brief overview of the experimental setup

In this section, we give a brief overview of the twin experiments which were performed using EnKF and partial observations from the full state. We describe in detail the different steps involved in the implementation of the experiments and study the performance of EnKF from the perspective of two important metrics such as RMSE and the rank histograms introduced earlier in section 2.7.

We start the assimilation experiment with the generation of an ensemble of states which are sampled from a Gaussian distribution. The first member is chosen to represent the true initial state and serves as the initial condition for the model to generate the true trajectory. On this true trajectory, we generate 10 observation realizations using the observation operator and the observation error covariance. The assimilation experiment consists of a cycle of two steps repeated over time: the forecast step and the update step. The forecast at some time  $t_k$  is performed by forward integrating the initial ensemble forward in time from time  $t_{k-1}$ . This gives us the ensemble representing the forecast distribution or the prior distribution at time  $t_k$  in the Bayesian sense. For the L96-40 dimensional system, the Lyapunov time scale is 0.5. Observations are assimilated at every 0.1 interval of time. The update step follows the perturbed observation EnKF algorithm, which generates an ensemble of observations from the real observation  $y_k$  using the noise statistics. As discussed in sub-section 2.5.3 and 2.5.2, we use EnKF with localization and inflation procedures, and their parameters are specified at the start of the assimilation experiments. A multiplicative inflation factor  $\alpha$  is used to inflate the covariance matrix by pre-multiplying the forecast covariance matrix by  $\alpha^2$  before assimilation of the observations. When inflation is not implemented,  $\alpha$  is set to 1. We run the experiments for 800 assimilation steps, with a time gap of 0.1 at which the observations are used. We now discuss the numerical results of applying EnKF to perform twin experiments using the Lorenz-96 ODE that was introduced earlier in section 2.6.2.

### 2.8.2 Ensemble size versus Time of divergence

We start our discussion of twin experiments focusing on the issue of filter divergence. A filter is said to diverge from the true trajectory if the information coming from the measurements is ignored by the filter over time. Filtering distributions become overconfident and ignore observation corrections over time. This leads to gradual departure of the filter trajectory from the true underlying trajectory, eventually to a completely different trajectory.

In figure 2.5 above, what is observed is that the filter starts to depart from the true-state trajectory after some time, gradually ignoring all the observations, and diverging from the actual trajectory. In general, the time at which the trajectory diverges from the truth may be different for the different components of the state. In our discussion here, we denote



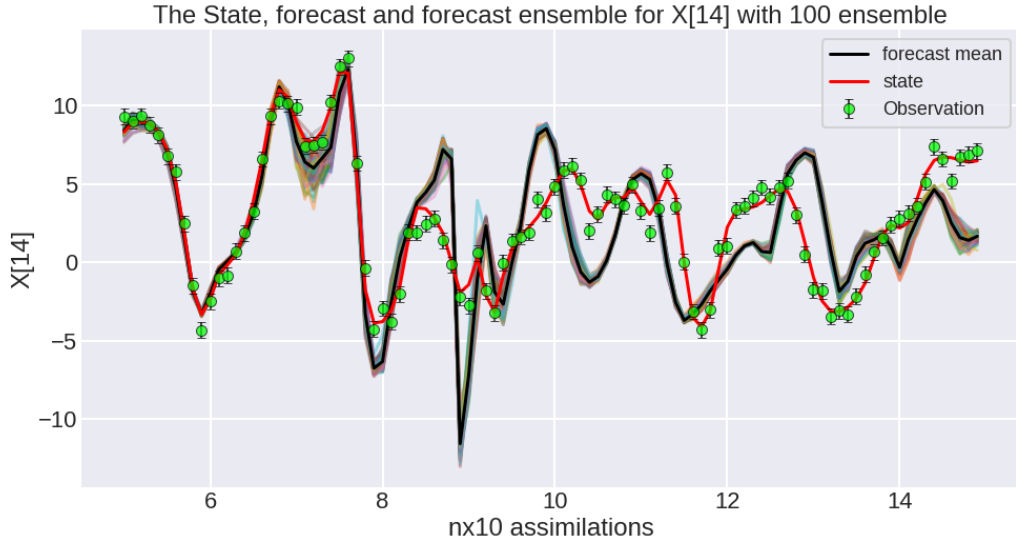


Figure 2.5: A particular component diverging from the true trajectory after finite assimilation steps using EnKF on Lorenz-96.

the average time over the components as a filter divergence time. For the same component plotted in figure 2.5, we plot the absolute error over time below in figure 2.6.

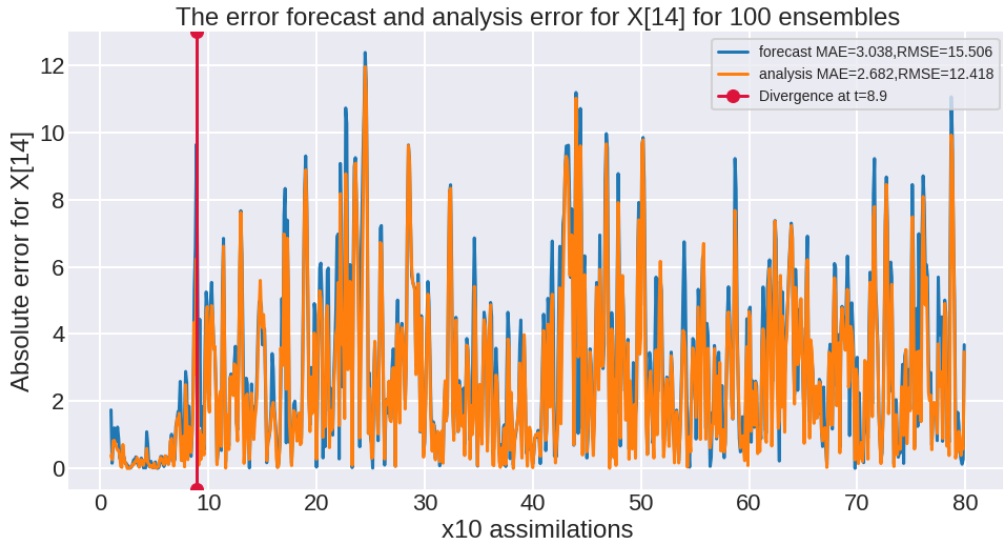


Figure 2.6: Absolute error over time for the same component when the filter divergence occurs. The  $X$ -axis shows the number of assimilation steps.

Filter divergence is a common issue for EnKF when the ensemble size is quite small compared to the dimension of the state space [7]. When the ensemble collapses, the diagonal elements of the forecast ensemble covariance  $\mathbf{P}_k^f$  approaches 0. This leads to the Kalman gain matrix  $\mathbf{K}_k$  adding no correction to the ensemble, and hence the filter

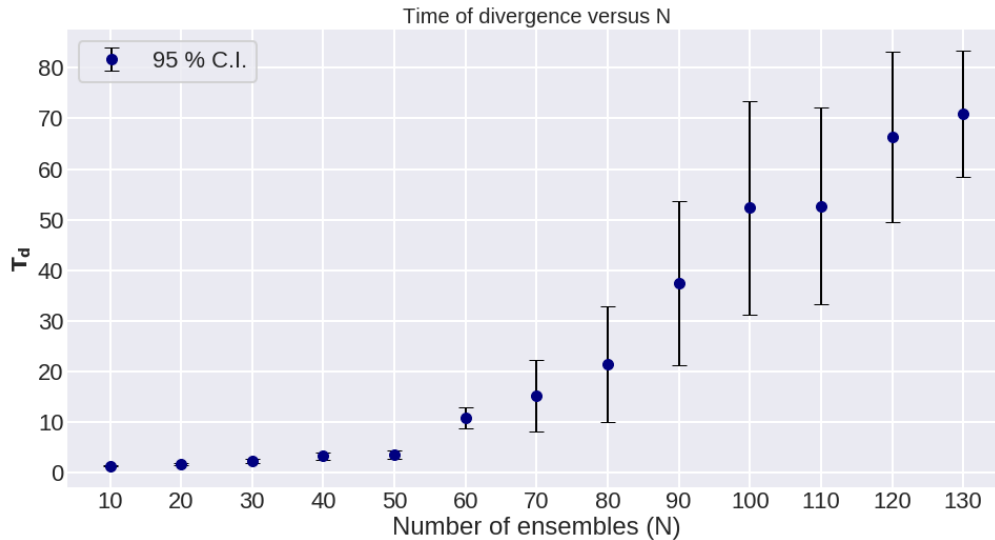


Figure 2.7: Time of divergence versus ensemble size. The error bars represent the confidence interval corresponding to the 10 observation realizations, which illustrates the variability.

estimates depart from the true state over time.

Using this set of experiments, we demonstrate the issue of filter divergence, and our objective is to highlight the behavior of EnKF before any additional procedures, such as localization and inflation. The detailed procedure of this experiment and the choice of parameters are as follows. We plot the average departure time over the different components which is here termed the time of divergence  $T_d$ . To see how this time of divergence depends on the number of ensembles, we plot this time versus the ensemble number. We perform assimilations with one and the same trajectory but with 10 different observation realizations. We gradually increase the size of the ensemble in steps of 10. Since there is no localization or inflation, we can capture the significant effect of increasing the size of the ensemble  $N$  on the divergence of the filters.

### 2.8.3 The effect of localization and inflation on filter divergence

We now discuss an experiment similar to subsection 2.8.2 where we want to show how filter divergence is affected by the introduction of localization and inflation. Figure 2.8. This time, we set our localization function to be concave with the radius set to 3 and plot the filter divergence time  $T_d$  for different ensemble sizes  $N$ .

To investigate the effect of inflation on the time of filter divergence, we now add inflation to our localization experiments. In this experiment, we only choose small ensemble sizes where the role of inflation is evident. For this experiment, we fix the localization type to concave and the radius to 3. We gradually increase the inflation factor  $\alpha$  from 1.00 to 1.30 in steps of 0.02 and calculate the time of filter divergence averaged over 10 observation

realizations, with the confidence interval 95%, In figure 2.9, we plot the filter divergence time versus the inflation factor  $\alpha$  for ensemble size 10, 15 and 20, respectively.

## 2.8.4 Localization scale versus RMSE for different localization functions

We now implement localization in our assimilation experiments. Both the localization function used and the corresponding scale are important. Here, we use the following three types of localization functions, which can be used to define the elements of the localization matrix  $\rho$  for a chosen localization function as follows:

$$[\rho]_{ij} = \begin{cases} 0, & \text{for } |i - j| \geq \ell \\ \frac{2\exp(-\frac{|i-j|}{\ell})}{\exp(|i-j|)+1}, & \text{for convex} \\ \frac{\ell^2 - |i-j|^2}{\ell^2}, & \text{for concave} \\ \frac{\ell - |i-j|}{\ell}, & \text{for linear} \end{cases} \quad (2.32)$$

To study the best localization function for the system and the corresponding length scale, we can use RMSE. Since using large numbers of ensembles is computationally expensive, we vary the ensemble numbers from 10 to 20 in steps 5. For each localization function, we vary the localization scale in steps of 1. We then repeat it for the convex, concave, and linear localization functions, respectively. Averaging the RMSE over the observation realizations, we plot the figure 2.11 below.

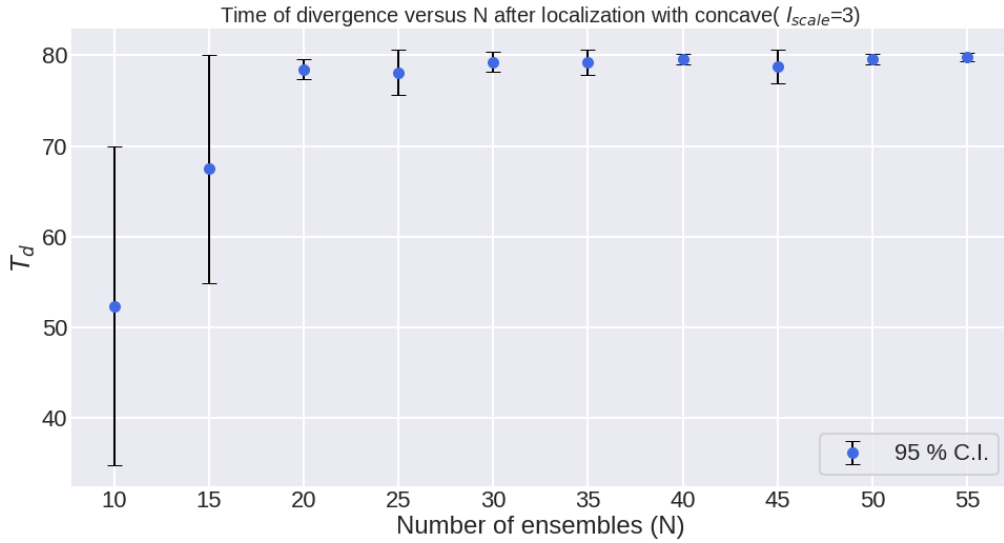


Figure 2.8: Plot of time of divergence versus ensemble size

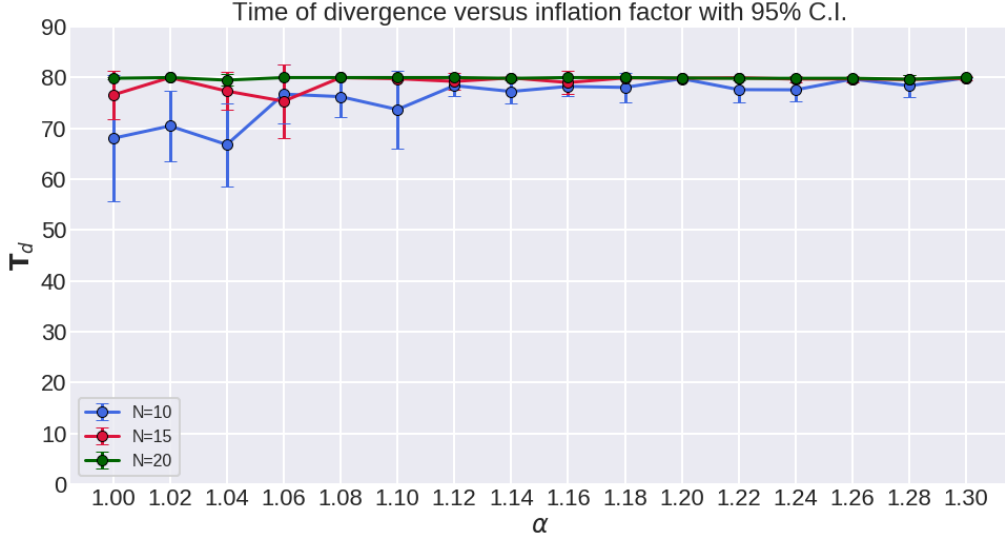


Figure 2.9: Averaged time of divergence over 10 observation realizations for different values of inflation factor. The localization is fixed at radius 3 for the concave type function.

### 2.8.5 Effect of model error in data assimilation

In this section, we focus on scenarios where we have a model error, which is closer to real applications. Different numerical approximations and parameterizations for unresolved processes that are unknown or are computationally very expensive to solve, are neglected in order to make data assimilation computationally feasible. Thus, real data assimilation always works with an approximate model and accounts for the model error for the additional uncertainty in the model forecasts in an appropriate way. Collectively, it can be understood as the presence of model error in the data assimilation arising under different conditions, and the model still has skills in estimating the state, but with additional uncertainties. The experiments and results presented in this section focus on understanding and trying to mitigate the effect of unaccounted model error present in the assimilation.

We assume that the true model for the physical system is the Lorenz II level system, which has two different variables, a grid of fast-scale variables  $\mathbf{Y}$  and a grid of slow-scale variables  $\mathbf{X}$ . The model that we use for assimilation is the Lorenz-I level system, introduced earlier in section 2.6.2, and that we only observe the slow-scale variables  $\mathbf{X}$  from the system. The situation here corresponds to the scenarios where the unresolved scale dynamics of  $\mathbf{Y}$  affect the dynamics of the resolved scale dynamics of  $\mathbf{X}$  via some coupling, but they are not part of the dynamical model in our data assimilation.

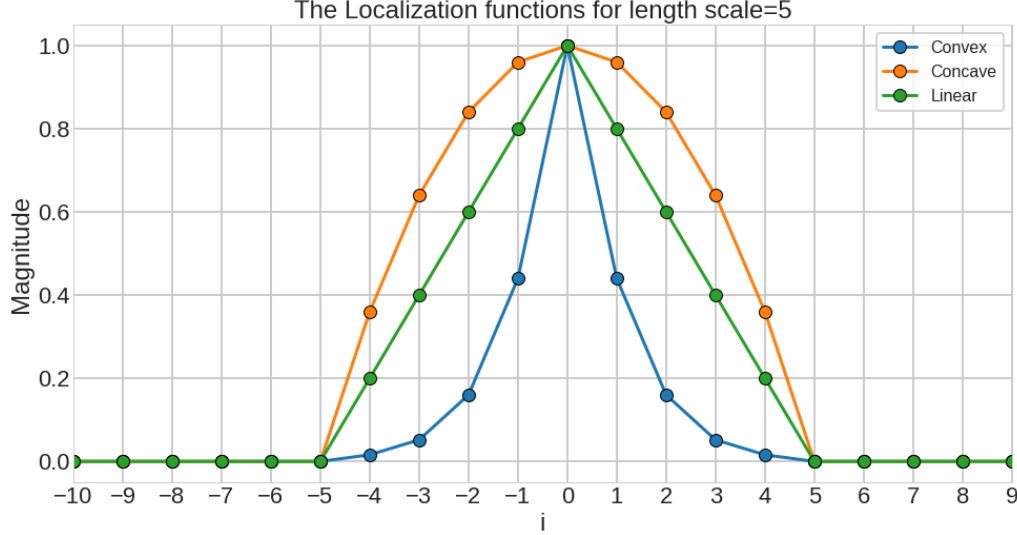


Figure 2.10: We plot the three different types of localization function introduced in equation (2.32). The x-axis indicates  $i$ , which is the index for the distance on a discrete grid assuming the origin is fixed at 0.

The Lorenz II level system is given by the following coupled ODE,

$$\frac{d\mathbf{X}_k}{dt} = \mathbf{X}_{k-1}(\mathbf{X}_{k-2} - \mathbf{X}_{k+1}) - \mathbf{X}_k + \mathbf{F} - \frac{hc}{b} \sum_j \mathbf{Y}_{k,j} \quad (2.33)$$

$$\frac{d\mathbf{Y}_{k,j}}{dt} = cb\mathbf{Y}_{k,j+1}(\mathbf{Y}_{k,j+2} - \mathbf{Y}_{k,j-1}) - c\mathbf{Y}_{k,j} + \frac{hc}{b}\mathbf{X}_k \quad (2.34)$$

where,  $\mathbf{X}_k$  are slow-scale variables,  $\mathbf{Y}_{k,j}$  are fast-scale variables interacting with  $\mathbf{X}_k$ . The parameters  $c$  and  $b$  are the temporal scale and the spatial scale ratio of the variables of the fast and slow scales with  $h$  being a coupling constant.  $\mathbf{F}$  is a constant forcing as is present in the original Lorenz-96 model.

To set up the experiment, we fix the parameters as follows:  $\dim(\mathbf{X}) = 20$ ,  $\dim(\mathbf{Y}) = 10$ ,  $c = 10$ ,  $b = 10$ ,  $h = 1$ ,  $\mathbf{F} = 15$ . We first generate the trajectory from the true system, i.e. the Lorenz II level system. We then generate observations using only the  $X$  slow-scale variables. To perform data assimilation with the observations, we consider the Lorenz I level system, as a truncated model to the full model that has only the slow-scale dynamics and observe 10 alternate components. We use 20 ensemble members with concave localization with  $\ell = 3$  without any inflation.

Since observations come from slow-scale variables  $\mathbf{X}_k$ , which are part of a coupled to a larger system, our goal is to see the effect of using an observation error covariance different from that used to generate them. We consider the uncertainty of the observations to be enhanced because of the presence of the model error. We plot the RMSE of both the observed and unobserved components versus  $\mu_u$ , where the assumed observation error

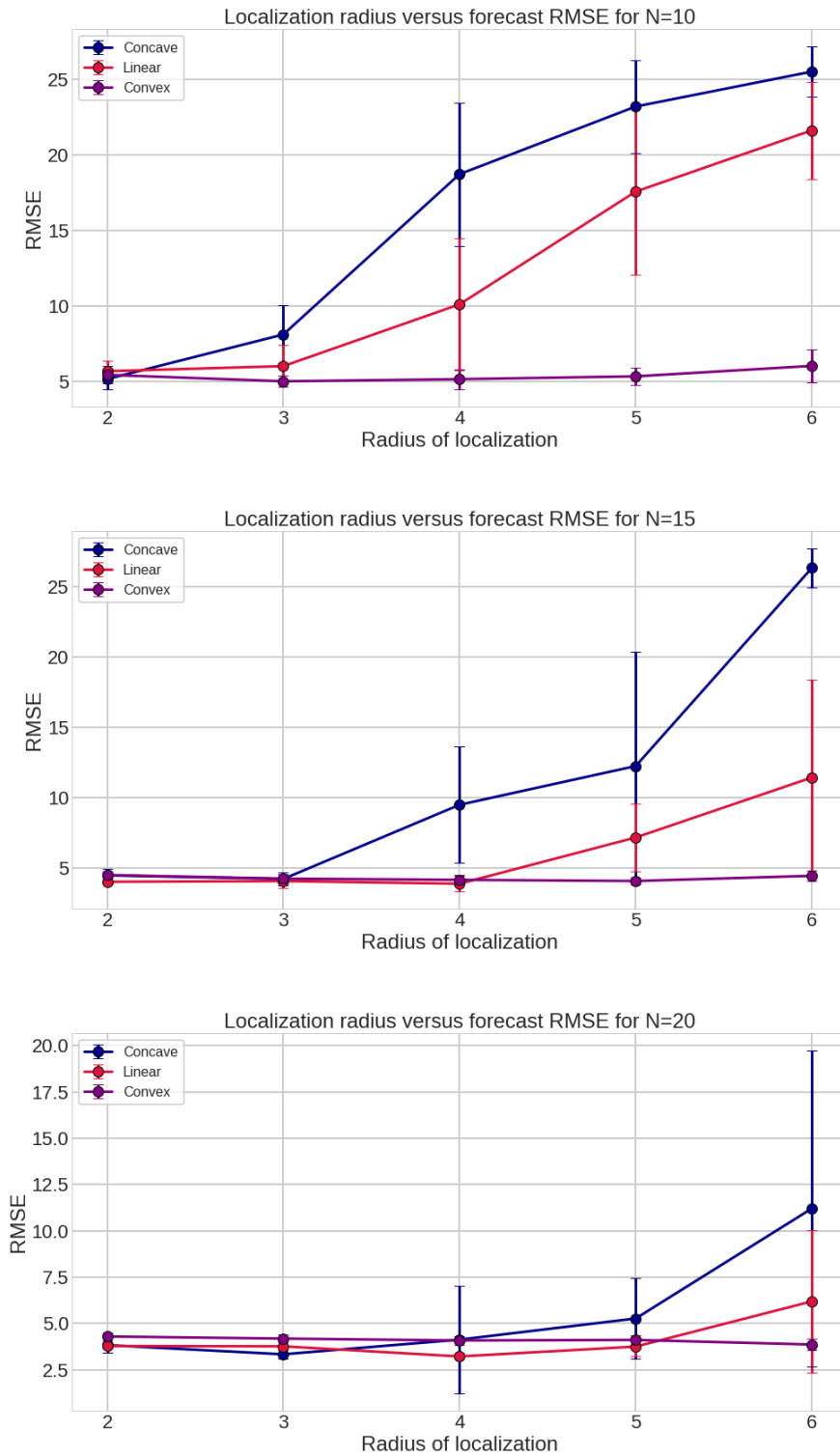


Figure 2.11: Effect of localization radius on RMSE for different ensemble size N. The three different colors represent different types of localization functions introduced in subsection 2.5.3.

covariance is of the form  $\mu_u I_d$ . Averaging the different observation realizations, the RMSE versus  $\mu_u$  plot is shown in figure 2.12. The asymptotic RMSE of the observed part is

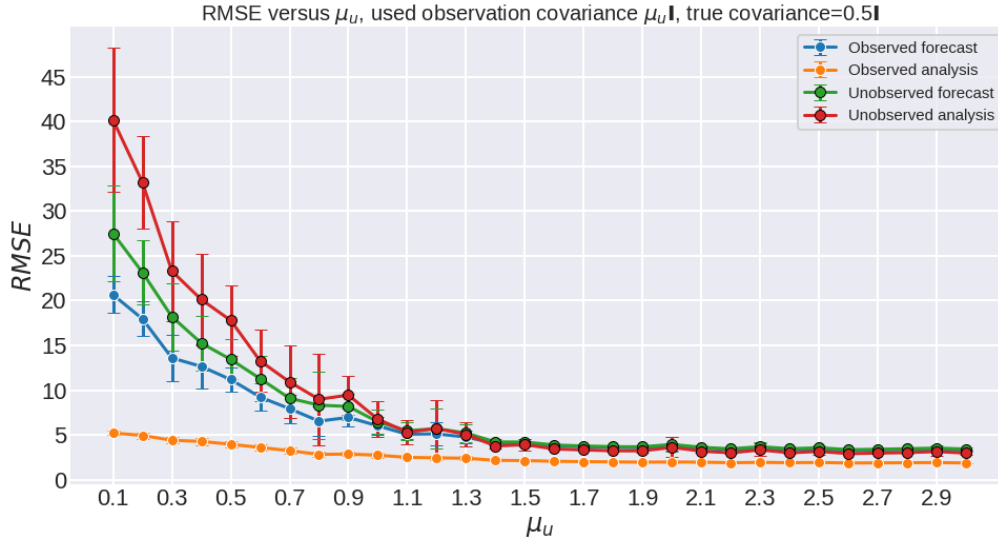


Figure 2.12: RMSE versus  $\mu_u$ , the assumed observation covariance used in assimilating the observations. The true observation covariance is  $0.5\mathbf{I}_{10}$

clearly less than the unobserved part of the state vector and their values are close to the observation error.

We also plot the bias-variance trade-off for different inflation factors. We quantify this by computing the ratio of the absolute error to the spread of the ensemble, averaged over the components. We calculate the absolute error in each component averaged over time and the 10 observation realizations. Figure 2.13 shows the plots for these quantities. Similarly, it is calculated for the ensemble spread for each component. The rank histogram, already introduced in subsection 2.7.2, is another important tool for understanding the reliability of the forecasts[Hamill, 2001]. Figure 2.14 shows the rank histogram of both an observed and an unobserved component, respectively, for different  $\mu_u$ . The expected shape of the histogram is flat when any member of the ensemble is equally likely to be a possible true state. For an unreliable ensemble, a U-shaped curve appears for  $\mu_u = 0.5$  (overconfident ensemble), while an inverted U-shaped curve appears for  $\mu_u = 2.0$  (underconfident ensemble).

For a reliable ensemble, the distance between the mean of the ensemble and the observation should be close to the ensemble spread. When this is not true, it is possible that the spread of the ensemble overestimates the uncertainty or underestimates the uncertainty when it is systematically smaller in the mean predictions. If the observations consistently lie outside the spread defined by the ensemble, the ensemble forecast is said to be under-dispersive. Similarly, When the observations are consistently within the spread defined by the ensemble, the ensemble forecast is said to be over-dispersive. In the presence of model errors, we must account for their overall effect on the forecasts. If ignored, the

forecasts tend to become overconfident, making them unreliable.

## 2.9 Summary

In this chapter, we introduced the Bayesian approach to nonlinear filtering and data assimilation for dynamical systems. We presented the toy dynamical systems used in this thesis and the filtering algorithm- ensemble Kalman filter (EnKF), the approximation of the original Kalman filter for nonlinear systems. The Bayesian filtering theory provides a recursive solution to estimate the posterior distribution of the state given the history of observations. The chapter then presents the Kalman filter, which is the optimal solution for the filtering problem when the dynamic and measurement models are linear and the observation noise is Gaussian. We also derive the Kalman filter equations using the best linear unbiased estimator (BLUE) approach and explain the role of the Kalman gain matrix in minimizing the trace of the estimated posterior covariance matrix. We then discuss the EnKF, which extends the Kalman filter to nonlinear dynamical systems by using an ensemble of particles to represent the distribution and perform the operations of the standard Kalman filter. We discuss the practical methods such as inflation, and localization that make EnKF applicable in high-dimensional operational data assimilation scenarios. This chapter also introduces the two chaotic dynamical systems, the Lorenz-63 and the Lorenz-96 models, that are used as standard test-bed chaotic dynamical systems for the numerical experiments related to data assimilation. Performance metrics for filtering algorithms, such as accuracy, reliability, and stability, explains how to measure them using RMSE and rank histograms.

We then discuss the utility of twin experiments in data assimilation research in order to systematically test and compare different data assimilation algorithms using synthetic observations from a known true trajectory. The chapter concludes with some twin experiments and results using EnKF and partial observations from the Lorenz-96 model. We presented results for different important parameters chapter illustrates the effects of ensemble size, localization function radius, inflation factor, observation error covariance, and model error on the filter divergence time using RMSE and rank histograms. The chapter also demonstrates the effectiveness of localization and inflation for improving the quality and reliability of the EnKF forecasts.



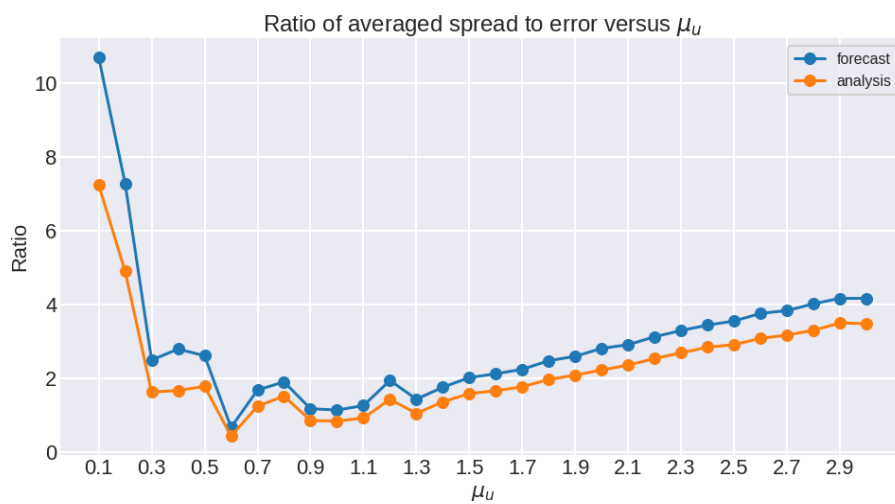
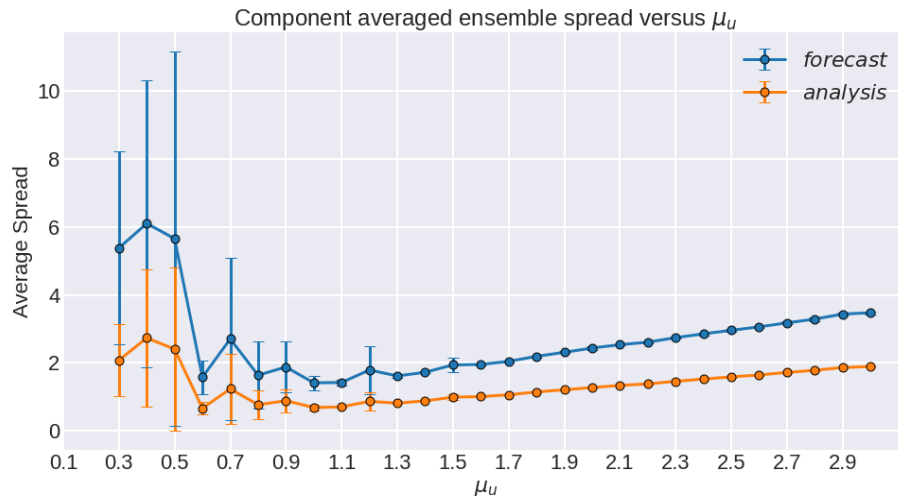
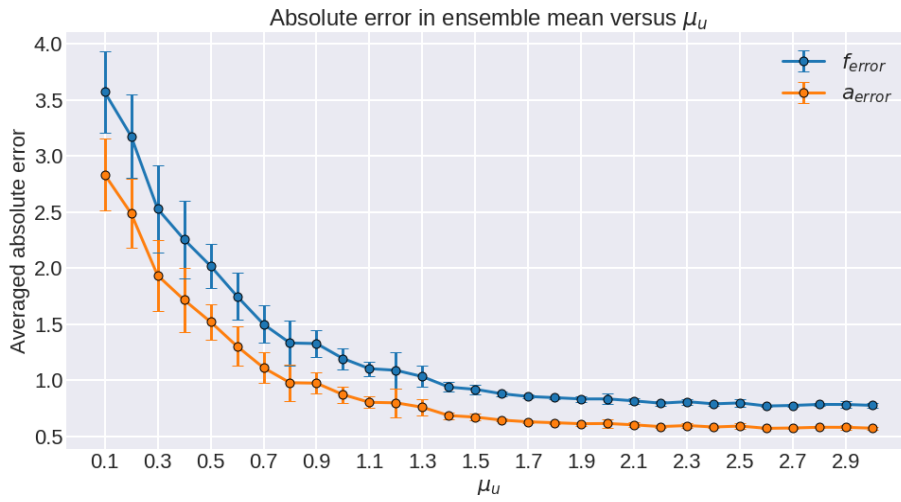


Figure 2.13: Average absolute error, average spread and their ratio versus assumed covariance  $\mu_u$  fo

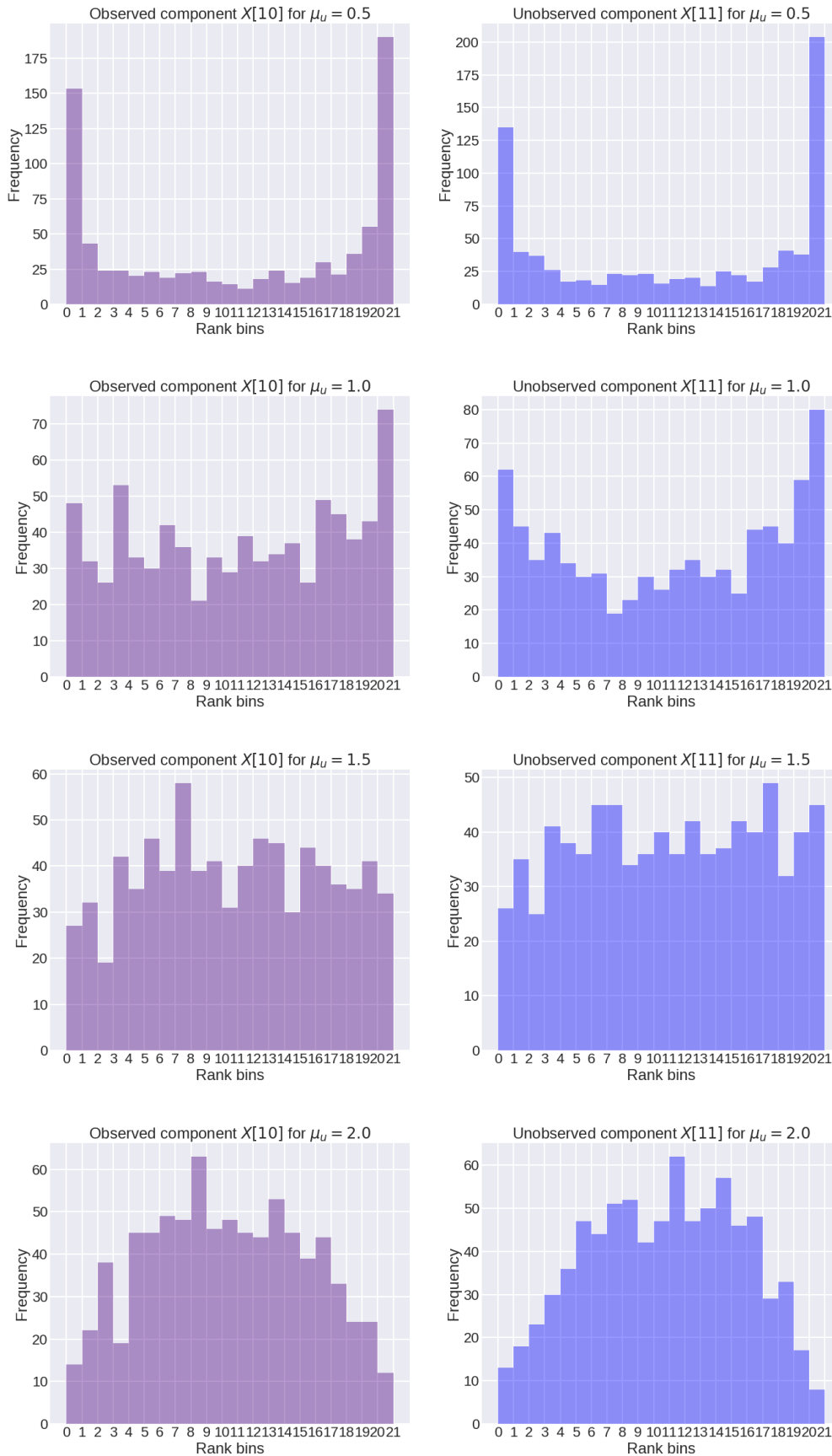


Figure 2.14: Rank histograms for an observed and an unobserved component for different values of  $\mu_u$  assumed in EnKF data assimilation. The U-shape appears for under-dispersive ensemble (1<sup>st</sup> row) and the inverted U-shape (4<sup>th</sup> row) for over-dispersive ensemble.

# Chapter 3

## Numerical filter stability of EnKF using Sinkhorn divergence

In chapter 2, we studied the Bayesian filtering problem where the goal is to estimate the conditional distribution of the state based on the history of the observations. The focus of this chapter is to study the problem of stability of such filtering algorithms. Filter stability is the property that the conditional posterior distribution computed sequentially over a long time is robust to the choice of the distribution made when initializing the filter. Specifically, we study nonlinear filter stability for a specific sequential filtering algorithm—the ensemble Kalman filters, which was introduced in section 2.5. Using EnKF for filtering the Lorenz-96 system, we study the stability numerically by directly using the notion of distances on the space of probability distributions. These numerical experiments are performed for deterministic dynamical systems using the recently developed Sinkhorn algorithm [29] to approximate the distance between probability distributions represented by Monte Carlo samples. Through twin experiments performed using Lorenz systems, we study the numerical convergence of the distance between the corresponding filtering distributions starting from different initial distributions. We show that the ensemble Kalman filter is stable for two different parameters that directly affect the numerical filters, namely, the observation gap and observation covariance. We also find that the Wasserstein distance between filters with two different initial conditions is proportional to the filter’s bias or RMSE, an empirical relationship between filter stability and filter convergence.

In section 3.1 we give a brief introduction to the problem of nonlinear filter stability from the perspective of data assimilation. We then define the problem of filter stability mathematically and introduce our definition in section 3.2. In section 3.3, we start by discussing the ideas of optimal transport-based Wasserstein distance and its approximation by Sinkhorn divergence. In subsection 3.3.3, we describe the Sinkhorn-knopp algorithm and develop our approach to study filter stability numerically based on the Sinkhorn algorithm. In section 3.4, we describe the computational aspects of the Sinkhorn algorithm relevant

to our method and discuss the results obtained using EnKF on the Lorenz-96 system. By approximating the distance between Monte Carlo samples representing different filtering distributions, our study allows us to address the problem of filter stability directly using numerical experiments. We are successful in demonstrating evidence for the exponential filter stability of our commonly used nonlinear filtering algorithms, namely the ensemble Kalman filters and the particle filter. In this thesis, we restrict our discussion to EnKF and refer to results in the context of particle filters which are present in our joint work [64, 65] when necessary.

### 3.1 An introduction to nonlinear filter stability

Data assimilation in Earth science has benefited from various numerical studies of filter algorithms and their performance. These studies focus on understanding algorithms' drawbacks and devising methods to improve them. These are important from the point of view of practical implementations and calibration before being used in real applications. To understand and evaluate different filtering algorithms, the data assimilation literature comprises of numerous such numerical studies. Twin experiments are regularly used to generate synthetic observations, and then data assimilation is performed by incorporating these observations into the numerical models to understand their efficacy under different conditions. However, most numerical studies have studied other measures of filter performance, such as the accuracy of filter mean and filter reliability tests using rank histograms [2].

In practice, a filtering algorithm is implemented numerically where they compute the conditional distribution of the state in the phase space using their own assumptions and approximation. Filtering algorithms are initialized with a certain choice of the distribution at the beginning, and the question of what is a good choice is difficult as the true state is unknown. A plausible choice used in operational forecasts is the climatological distribution of the variable, which has been generated using observation statistics for some components.

A Bayesian filtering algorithm aims to approximately represent the conditional distribution of the state based on the observation history. As the observations arrive serially, more information is available over time which results in improvement of the state estimates. Consequently the quality of the distribution of the state estimated by a filtering distribution must improve, overcoming the limitations due to the wrong choice at the time of initialization.

Previous works in the literature have studied the problem of filter stability in terms of RMSE where the effect of incorrect initial conditions on filtering algorithms has been studied to understand the behavior of errors over time. In reality, since the true state is not available, they have been studied using twin experiments for various filtering

algorithms. However, RMSE is not representative of the distance between probability distributions. Hence, convergence of quantities other than the distributions themselves cannot be used otherwise to define and illustrate filter stability. Other major numerical studies have focused on important quantities that can be studied numerically, such as the rank histogram, the statistics of innovation vectors for Bayesian filtering algorithms [2].

The problem of sensitivities to the initial distribution of a nonlinear filter has been studied in [25, 24, 28] as a question of the stability of the filter. A filter is said to be stable if, for any choice of initial distribution, the resulting filtering distribution in time is asymptotically the same. This is indeed a desirable property for any filter, since the initial distribution of the state is unknown in the beginning. Satisfying the property mentioned above, a reliable filter will soon “forget” about the arbitrary choice of the initial distribution [83]. Thus, filter stability is an extensively studied topic, but has mainly received attention in the context of stochastic dynamical systems [74, 68, 75, 21].

There are a few theoretical results on the stability of EnKF, such as in [30]. The other important results in the context of twin experiments are on *bias* and if the errors in the filter mean remain bounded in time [54, 55] for different systems. However, these studies do not provide any direct evidence of stability in the context of a filter. Therefore, our goal is to assess filter stability directly by using the filtering distributions in time used with the definition (3.2) provided in section 3.2, which we can test directly for any numerical filtering algorithm. We study the numerical stability of the filter by applying EnKF, which has been one of the foundations for different ensemble-based filtering techniques [13, 7, 45]. For EnKF which was introduced in section 2.5, we can use the ensemble representations of the filtering distributions directly in order to approximate the distance between filtering distributions obtained using different initialisations.

Another class of filters that solve the Bayesian filtering problem without imposing any assumption on the underlying distribution are particle filters. The Bayesian posterior calculations for particle filters are performed by weighted sampling, where the weights of the particles are modified based on the likelihood of observation, see [31] for an overview of particle filtering algorithms.

## 3.2 Mathematical definition of filter stability

In chapter 2, we defined the notation for different conditional distributions which arise in nonlinear filtering. We now specifically focus on the filtering distribution  $\pi(\mathbf{x}_n|\mathbf{y}_{0:n})$  which is the conditional distribution of the state  $\mathbf{x}_n$  given all observations upto time  $t_n$ . Since any numerical filtering algorithm starts with the probability distribution  $\pi(\mathbf{x}_0)$  of the state at  $t_0$ . The chosen distribution is often arbitrary and wrong since there are both unobserved components in the state whose climatology is unknown. Thus, the true

filtering distribution is unknown and unavailable for use in filter initialization, and the choice of  $\pi(\mathbf{x}_0)$  is made arbitrary. The actual distribution may be very different from  $\pi(\mathbf{x}_0)$ , a drawback that all numerical filtering algorithms share that must be overcome to allow the assimilation system to operate. *In the same spirit that sensitivity to the initial conditions is important for chaotic dynamical systems, we study the sensitivity to initial condition of a filtering algorithm.* It is critical that the conditional distribution of the state  $\pi(\mathbf{x}_n|\mathbf{y}_{0:n})$  at time  $t_n$  is independent of the original choice used to initialize the filter so that quantities estimated using the posterior are finally independent of our initial distribution  $\pi(\mathbf{x}_0)$ . To simplify notation, we represent the filtering distribution obtained at time  $t_n$  as  $\pi_n(\mu)$ , where  $\mu$  represents the initial distribution choice. A good filtering algorithm will forget the choice of initial condition over time, and how fast a given filter achieves stability is determined by a quantity called the filter stability rate [83].

In the setting of the filtering problem for a dynamical system introduced in the earlier chapter 2 in section 2.2, having discussed the notion of filter stability, we now describe the filter stability definition in [74], which explores the “true” filter stability for deterministic dynamics. We adapt their definition to numerical filter stability by denoting the numerical approximation of the true filter  $\pi$  as  $\hat{\pi}$ .

**Definition**[Stability-RA [74]] A numerical filter is stable if, for any measure  $\nu$  w.r.t which the true initial measure  $\mu$  is absolutely continuous i.e.,  $\mu \ll \nu$ , we have,

$$\lim_{n \rightarrow \infty} \mathbb{E} \left[ \int_{\mathbb{R}^d} h(x) \hat{\pi}_n(\mu, dx) - \int_{\mathbb{R}^d} h(x) \hat{\pi}_n(\nu, dx) \right] = 0, \quad (3.1)$$

for any bounded and continuous function  $h$ , where the above expectation is taken over the observation noise.

Although (3.1) accurately conveys the concept of filter stability, it can be improved from a computational point of view in the following ways. First, we note that computing the expectation for all bounded and continuous functions is not possible. From a practical point of view, it may not be possible to obtain  $\mu$ , so a definition where the expectation does not depend on  $\mu$  is preferable. We need a definition that can be computed and verified numerically in order to study the stability of the numerical filter, and hence, we propose the following definition in [65] which does not involve the true measure  $\mu$ .

A numerical filtering algorithm is stable if the filter, starting with two different initial distributions  $\mu$  and  $\nu$  has  $\pi_n(\mu)$  and  $\pi_n(\nu)$  as the corresponding filtering distribution at time  $n$ , then the following holds

$$\lim_{n \rightarrow \infty} \mathbb{E}[D(\hat{\pi}_n(\nu_1), \hat{\pi}_n(\nu_2))] = 0, \quad (3.2)$$

where  $D$  is the distance on  $\mathcal{P}(\mathbb{R}^d)$ , the space of probability measures on  $\mathbb{R}^d$  and the expectation is taken with respect to observation noise. By studying the equation 3.2 using

an appropriate distance metric  $D$ , we can directly study the numerical stability of the filter by choosing different initial distributions for the filter and how the distance between the corresponding filtering distributions vary over time. Definition (3.2) is quite general and is a stronger version of definition 3.1, see the Appendix in [65] for details.

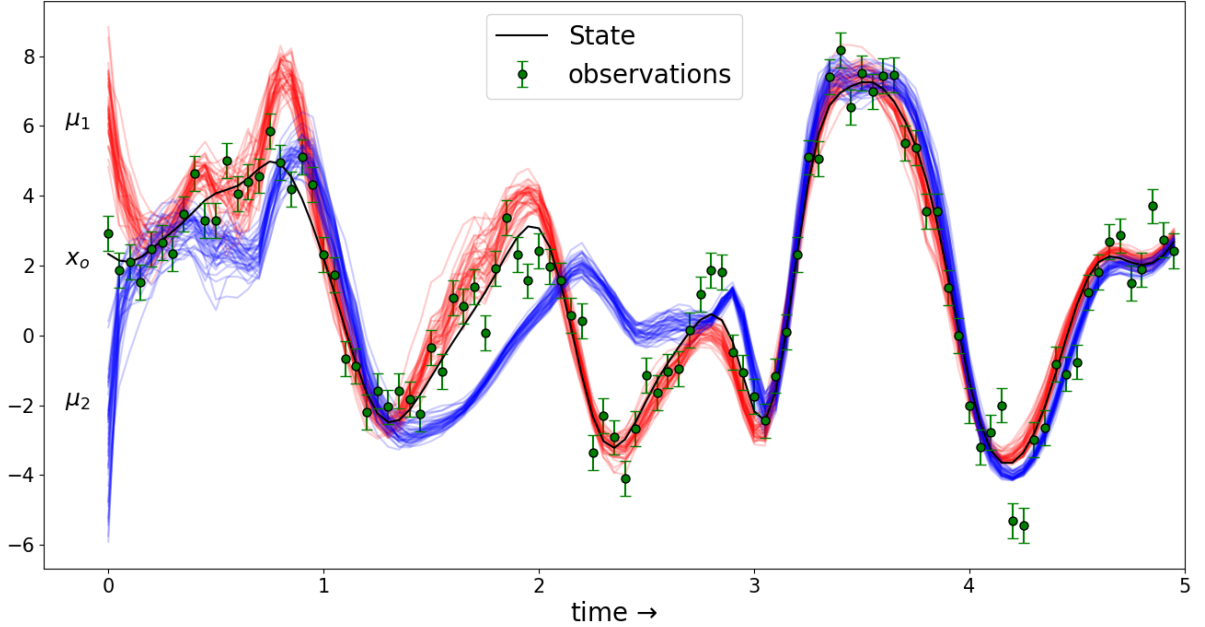


Figure 3.1: The filtering distributions for a single component are represented by an ensemble of trajectories over time.  $\mu_1$  and  $\mu_2$  are two biased distribution used to start the algorithm. The green dots are the observations and the state is the actual trajectory of one out of 40 components. The ensemble of trajectories are obtained using EnKF.

Figure 3.1 shows an experiment using EnKF on Lorenz-96 system assimilation experiment, how two different initial distributions evolve over time for a given set of observations. From a numerical perspective, definition (3.2) offer some difficulty in the sense that for practical purposes, we can only study the stability of a filter numerically using finite pairs of different initial distributions. However, this limitation does not prevent us from exploring the filter stability numerically, as the initial distributions we decide for the filter are generic. The randomness in  $x_0 \sim \mu$  is also present; however, we can address this by fixing the initial state  $x_0(\omega)$  as a realization, and then, the only randomness which is present in  $\hat{\pi}$  is attributed to the observations error distribution.

We now discuss the details of the distance and the related computational aspects. Using definition (3.2), we now turn to the choice of the distance  $D$  which we can use to compute the distance between filtering distributions over time.

### 3.3 Background: Wasserstein distance on the space of probability distributions

In general, we may be tempted to use any distance metric on  $\mathcal{P}(\mathbb{R}^d)$ , to investigate stability numerically such as Total variation distance or other popular pseudo-distance metric such as KL-divergence. However, in our study, we have used the Wasserstein metric, denoted by  $W_p$ , as our choice for the distance  $D$ . Wasserstein distance has some suitable properties, which makes it relevant for our approach. The Wasserstein distance formulation allows using well-known distances defined in the Euclidean space where the samples live, which are also used in defining the cost function.

A cost function is said to lift the distance from the feature space  $\mathcal{X}$  to the general family of measures  $M^+(\mathcal{X})$  if for all dirac masses  $\delta_x$  and  $\delta_y$ , we have

$$D(\delta_x, \delta_y) = \|x - y\|, \quad (3.3)$$

the cost function between two points  $x$  and  $y$  in the ambient space is their geometric distance. Wasserstein distances lift the standard metrics on  $\mathbb{R}^d$  to the space of probability distributions  $\mathcal{P}(\mathbb{R}^d)$  unlike KL-divergence or total variation distance. This gives them a geometrical interpretation that is not present in other distances on  $\mathcal{P}(\mathbb{R}^d)$ . For our purpose of computing distances, we use  $p = 2$  for no reason other than the familiarity of the 2-norm on Euclidean spaces. For a comparison of these distances the interested reader may see [6] where example 1 (learning parallel lines) depicts how the output of  $W_p$  can often be intuitive.

Another property of the Wasserstein distance is that it metrizes convergence in the space of probability distribution. A distance metric  $D$  on the space of probability distribution is said to metrize the convergence if for a sequence of distributions  $\{\alpha_n\}$  converges to  $\alpha$ , then following holds,

$$\alpha_n \rightarrow \alpha \iff D(\alpha_n, \alpha) \rightarrow 0. \quad (3.4)$$

The sequence  $\{\alpha_n\}$  can represent the discrete measure obtained by sampling  $\alpha$  which is a continuous distribution. This is important from the point of view of numerical filtering algorithms such as in EnKF where we represent the filtering distributions using an ensemble representation. In the following section, we discuss the theory of optimal transport using which one can define a distance on the space of probability distributions such as the Wasserstein distance.



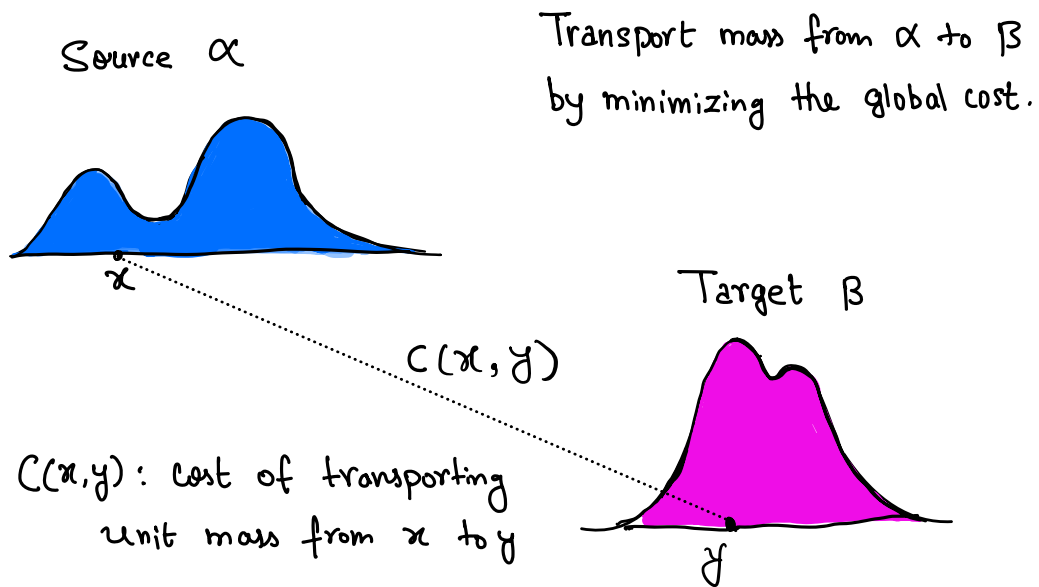


Figure 3.2: Schematic showing two distributions  $\alpha$  and  $\beta$ .  $c(x, y)$  is the cost of moving a unit mass from point  $x$  to  $y$ .

### 3.3.1 The optimal transport problem and Wasserstein distance

Optimal transport is defined as the minimum cost of morphing one probability distribution into another where the transport of mass comes with an associated cost defined on the ground metric. When the cost function is a distance on the ground metric, the optimal transport plan which minimizes the cost, this minimum cost defines a distance between two probability distributions.

Problems which involve learning probability distribution have focused on such distances and led to efficient algorithms for their computation. However, the origin of optimal transport dates back to 1781 when Gaspard Monge [86] was interested in finding the solution to the logistic problem of the most efficient way of transporting material from one region A to another region B. Monge's sought a deterministic transport map as the solution for this assignment which would determine the transport plan to move material from a set of points  $x_i$  in the source region to a point  $y_i$  in the target region B to minimize the total transportation cost [73]. This formulation however didn't allow for mass-splitting and existence and uniqueness of the solution didn't hold in general.

Leonid Kantorovich, in 1942, revisited the assignment problem via optimization for resource allocation problems in economics. His relaxation to the original formulation of Monge's transport problem, now known as the Kantorovich formulation [37], which was formulated and generalized for probability distributions. Kantorovich solved the optimal transport transport by introducing duality theory and linear programming [46, 48, 47].

We now turn to Kantorovich's formulation of optimal transport where the problem is

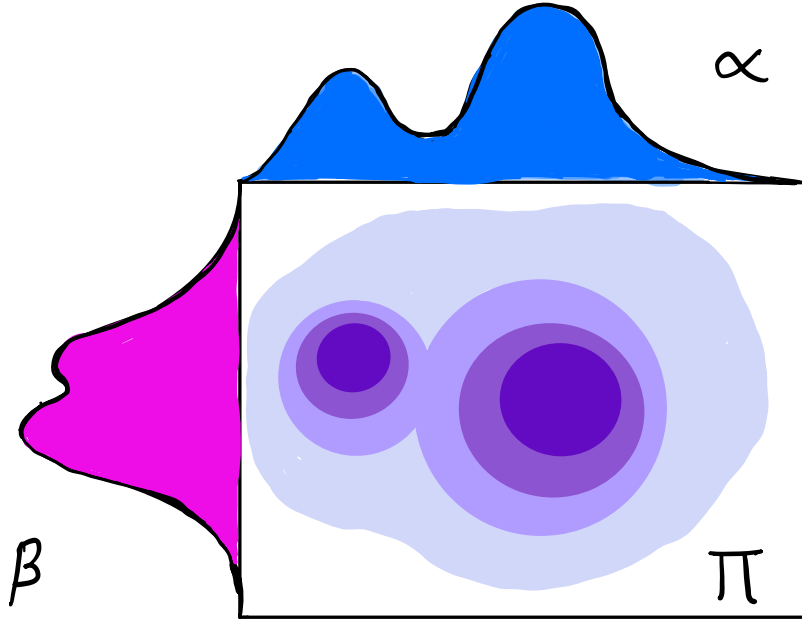


Figure 3.3: The joint distribution with the marginal distributions being  $\alpha$  and  $\beta$  respectively.

posed as finding the optimal transport plan which minimizes the cost of mass transportation [37] between two probability distributions under a well-defined transportation cost for each such movement between two points.

Consider two probability distributions  $\alpha, \beta \in \mathcal{M}_{\mathcal{X}}$ , where  $\mathcal{M}_{\mathcal{X}}$  represents the space of probability distribution on  $\mathcal{X}$ . Let  $c(x, y)$  be a function that defines the cost of transporting a unit mass from point  $x$  to point  $y$ , taking into account the geometry of the sample space. It may also be chosen as a function of the distance defined in the sample space. The cost of minimizing the transportation plan between two different distributions is equivalent to morphing from  $\alpha$  to  $\beta$ . This is obtained by solving the following optimization problem

$$\mathbf{OT}_c(\alpha, \beta) \stackrel{\text{def}}{=} \min_{\pi \in \mathbf{\Pi}} \int_{\mathcal{X} \times \mathcal{X}} c(x, y) d\pi(x, y), \quad (3.5)$$

where,  $\pi$  is a joint distribution which minimizes equation (3.5) and lies in the space of joint distributions over the product space  $(\mathcal{X}, \mathcal{X})$ . This space is denoted by  $\mathbf{\Pi}$ , where,

$$\mathbf{\Pi} = \left\{ \pi \mid \pi \in \mathcal{M}_{\mathcal{X} \times \mathcal{X}}^+, \int_{\mathcal{X}} d\pi(x, \cdot) = \alpha, \int_{\mathcal{X}} d\pi(\cdot, y) = \beta \right\},$$

the space of all joint probability distributions over  $\mathcal{X} \times \mathcal{X}$  with two marginal distribution being  $\alpha$  and  $\beta$  over  $\mathcal{X}$ .

Kantorovich also introduced a distance, known as the Kantorovich-Rubinstein [86] distance between probability measures according to which the distance between two

measures should be the optimal transport cost from one to the other, if the cost is chosen as the distance function.

When the cost of transporting a unit mass from point  $x$  to point  $y$  is of the form  $c(x, y) = d(x, y)^p$  where  $d(x, y)$  is a distance between two points  $x$  and  $y$ , equation 3.5 defines a distance on the space of probability distribution called the Wasserstein-p distance, given by,

$$W_p \stackrel{\text{def}}{=} [\text{OT}_c]^{1/p} = \left[ \min_{\pi \in \Pi} \int_{\mathcal{X} \times \mathcal{X}} d(x, y)^p d\pi(x, y) \right]^{1/p} \quad (3.6)$$

Wasserstein-p defines a distance only when the cost function is chosen in the above manner, since the conditions for satisfying triangle inequality are not satisfied otherwise [86].

Recently, optimal transport has received a tremendous amount of attention both from theoretical and computational perspectives in the field of machine learning, economics, and finance. Chief examples of success in deep learning endeavor are adversarial generative networks that have seen tremendous growth in the last few years [51, 76, 57]. Owing to the vast literature on this topic, we refer to [86, 73] for their comprehensive treatment of the theoretical and computational perspectives on optimal transport. Solving the Wasserstein distance for two discrete measures leads to a computationally expensive linear program that as must be solved in order to obtain the coupling between the two distributions  $\alpha$  and  $\beta$  in order to find their distance. When applied to two sampling distributions with both having sample size  $k$ , computing  $W_1$  is equivalent to solving a constrained linear programming problem(LPP) in  $n = k^2$  variables. Since LPPs take  $O(n^3)$  time to solve a problem with  $n$  variables, computing  $W_1$  takes  $O(k^6)$  time which is prohibitively expensive.

### 3.3.2 Entropy-regularized optimal transport and Sinkhorn divergence

The solution to second Wasserstein distance  $W_2$ , is obtained by computing the minimizer of equation (3.6) with  $c(x, y) = \|x - y\|_2^2$  for  $p=2$ . As mentioned above, solving the above optimization problem becomes computationally expensive in high dimension owing to the curse of dimensionality [8, 12]. Cuturi in [29] introduced an entropy regularization to equation (3.5) which facilitates a smoother solution and easier computational handling. The resulting dual optimization problem turns into a concave optimization problem, which can be solved by various iterative methods. Using parameters  $\varepsilon$ , which is the regularization coefficient, we now define the regularized optimal entropy transport [38] with the relative entropy as the regularization term in the following equation:

$$\text{OT}_\varepsilon(\alpha, \beta) \stackrel{\text{def}}{=} \min_{\pi \in \Pi} \left[ \int_{\mathcal{X} \times \mathcal{X}} \|x - y\|_2^2 d\pi(x, y) + \varepsilon \text{KL}(\pi | \alpha \otimes \beta) \right]. \quad (3.7)$$

Here,  $\mathbf{KL}(\pi|\alpha \otimes \beta)$  denotes the KL-divergence, computing the relative-entropy of the coupling or joint distribution  $\pi$  with respect to the product measure  $\alpha \otimes \beta$  on  $\mathcal{X} \times \mathcal{X}$ , given by,

$$\mathbf{KL}(\pi|\alpha \otimes \beta) = \int_{\mathcal{X} \times \mathcal{X}} \log \left( \frac{d\pi(x, y)}{d\alpha(x)d\beta(y)} - 1 \right) d\pi(x, y) + 1. \quad (3.8)$$

Entropy-regularized optimal transport can be used to approximate the Wasserstein distance between two probability measures. For a recent account of analytical results using entropy-regularized optimal transport for Gaussian distributions, see [79].

However, the minimizer of the equation (3.7) does not satisfy the properties of a distance and is also not zero for  $\alpha = \beta$  i.e.  $\mathbf{OT}_\varepsilon(\alpha, \alpha) \neq 0$ . A simple way around this problem was to redefine a quantity which overcomes this issue. by defining Sinkhorn-divergence, which is denoted as  $S_\varepsilon$ . For a fixed  $\varepsilon$ , the Sinkhorn divergence between two distributions  $\alpha$  and  $\beta$  is defined as

$$\mathbf{S}_\varepsilon(\alpha, \beta) \stackrel{\text{def}}{=} \mathbf{OT}_\varepsilon(\alpha, \beta) - \frac{1}{2} (\mathbf{OT}_\varepsilon(\alpha, \alpha) + \mathbf{OT}_\varepsilon(\beta, \beta)) \quad (3.9)$$

which satisfies  $\mathbf{S}_\varepsilon(\alpha, \alpha) = 0$ . Sinkhorn divergence satisfies the following properties:

- $\mathbf{S}_\varepsilon(\beta, \alpha) = \mathbf{S}_\varepsilon(\alpha, \beta) \geq \mathbf{S}_\varepsilon(\alpha, \alpha) = 0$
- $\mathbf{S}_\varepsilon(\alpha, \beta) = 0 \iff \alpha = \beta$
- For a sequence  $\{\alpha_n\} \rightarrow \alpha \iff \mathbf{S}_\varepsilon(\alpha_n, \alpha) \rightarrow 0$ .

Additionally, in the limit that  $\varepsilon \rightarrow 0$ , we have [35]

$$\lim_{\varepsilon \rightarrow 0} \sqrt{\mathbf{S}_\varepsilon(\alpha, \beta)} = W_2(\alpha, \beta) \quad (3.10)$$

which allows us to compute  $\mathbf{S}_\varepsilon$  for small enough  $\varepsilon$  in order to approximate  $W_2$  between two distributions.

### 3.3.3 Sinkhorn-Knopp algorithm

The dual optimization problem in equation (3.7) can be obtained via Lagrangian dual approach [37], and the resulting dual problem is given by

$$\mathbf{OT}_\varepsilon(\alpha, \beta) = \max_{f, g \in \mathcal{C}(\mathcal{X})} \left[ \int_{\mathcal{X}} f d\alpha + \int_{\mathcal{X}} g d\beta - \varepsilon \int_{\mathcal{X}^2} \left( \exp \left( \frac{f \oplus g - c}{\varepsilon} \right) - 1 \right) d(\alpha \otimes \beta) \right] \quad (3.11)$$

where  $f$  and  $g$  represent the dual potential function corresponding to the distribution  $\alpha$  and  $\beta$  respectively and  $c(x, y) = \|x - y\|_2^2$ , the  $l_2$ -distance function on  $\mathcal{X}$ . In real applications

and from the point of view of numerical filters, we know the distribution via their ensemble representations. We consider the above formulation of regularized optimal transport in equation (3.7) in a discrete setting where we work with probability measures instead of distributions. Since the filtering distributions can be represented using empirical measures, this helps us define the algorithm for any numerical filter and specifically to EnKF whose filter stability is studied in section 3.4.

For two empirical measures  $\hat{\alpha}$  and  $\hat{\beta}$  given by  $\hat{\alpha} = \sum_{i=1}^N \alpha_i \delta_{x_i}$ ,  $\hat{\beta} = \sum_{j=1}^M \beta_j \delta_{y_j}$  represented by using two i.i.d. samples  $(x_1, \dots, x_N)$  and  $(y_1, \dots, y_M)$  of size  $N$  and  $M$  which are drawn from their corresponding distribution  $\alpha$  and  $\beta$  respectively, equation (3.11) becomes

$$\mathbf{OT}_\varepsilon(\hat{\alpha}, \hat{\beta}) = \max_{f, g \in \mathcal{C}(\mathcal{X})} \sum_{i=1}^N \alpha_i f_i + \sum_{j=1}^M \beta_j g_j - \varepsilon \sum_{i,j} \alpha_i \beta_j \exp\left(\frac{f_i + g_j - c_{ij}}{\varepsilon}\right) + \varepsilon, \quad (3.12)$$

where  $f_i = f(x_i)$ ,  $g_j = g(y_j)$  and  $c_{ij} = \|x_i - y_j\|_2^2$  is the cost matrix. To show that this problem is a concave optimization problem, we first write down the first-order and second-order derivatives of  $\mathbf{OT}_\varepsilon$ . For simplifying the notations, we represent  $\mathbf{OT}_\varepsilon(\hat{\alpha}, \hat{\beta})$  by  $D$  and differentiate it w.r.t  $f_i$  and  $g_j$  to get

$$\frac{\partial D}{\partial f_i} = \alpha_i - \alpha_i \frac{1}{\varepsilon} \sum_j \beta_j \exp\left(\frac{f_i + g_j - c_{ij}}{\varepsilon}\right) \quad (3.13)$$

$$\frac{\partial D}{\partial g_j} = \beta_j - \beta_j \frac{1}{\varepsilon} \sum_i \alpha_i \exp\left(\frac{f_i + g_j - c_{ij}}{\varepsilon}\right) \quad (3.14)$$

$$\frac{\partial^2 D}{\partial f_i \partial f_j} = -\frac{1}{\varepsilon} \sum_j \beta_j \exp\left(\frac{f_i + g_j - c_{ij}}{\varepsilon}\right) \quad (3.15)$$

$$\frac{\partial^2 D}{\partial g_i \partial g_j} = -\frac{1}{\varepsilon} \sum_i \alpha_i \exp\left(\frac{f_i + g_j - c_{ij}}{\varepsilon}\right) \quad (3.16)$$

From the first order optimality conditions, we know that the first derivative in (3.13) and (3.14) become zero  $\forall i, j$ , and we get

$$f_i = -\varepsilon \log \left[ \sum_{j=1}^M \beta_j e^{\frac{g_j - c_{ij}}{\varepsilon}} \right] \quad \text{and} \quad g_j = -\varepsilon \log \left[ \sum_{i=1}^N \alpha_i e^{\frac{f_i - c_{ij}}{\varepsilon}} \right]. \quad (3.17)$$

From the second order optimality conditions, we clearly see that this is a concave optimization problem since both equation (3.15) and (3.16) are always negative. Thus, the solution to this optimization problem can be found by initializing  $f_i, g_j = 0 \forall i, j$  and iterating equation (3.17) until convergence. This allows us to compute the regularized optimal transport problem and the Sinkhorn divergence in equation (3.9). We describe the complete algorithm for computing  $\mathbf{S}_\varepsilon$  in details in algorithm 2.

---

**Algorithm 2:** Computation of  $S_\varepsilon$ 

---

**Input:**  $\{\alpha_i\}_{i=1}^N, \{x_i\}_{i=1}^N, \{\beta_j\}_{j=1}^M, \{y_j\}_{j=1}^M$ **Output:**  $S_\varepsilon \left( \sum_{i=1}^N \alpha_i \delta_{x_i}, \sum_{j=1}^M \beta_j \delta_{y_j} \right)$ Initialize  $f_i \leftarrow 0 \forall i = 1, \dots, N$  and  $g_j \leftarrow 0, \forall j = 1, \dots, M$ .iteration  $\leftarrow 0$ **while**  $\min\{L_1 \text{ relative errors in } f \text{ and } g\} > 0.1\%$  **do**    **for**  $i = 1, \dots, N$  **do**         $f_i \leftarrow -\varepsilon \log \left( \sum_{k=1}^M \beta_k \exp \left( \frac{1}{\varepsilon} g_k - \frac{1}{\varepsilon} \|x_i - y_k\|_2^2 \right) \right)$     **for**  $j = 1, \dots, M$  **do**         $g_j \leftarrow -\varepsilon \log \left( \sum_{k=1}^N \alpha_k \exp \left( \frac{1}{\varepsilon} f_k - \frac{1}{\varepsilon} \|x_k - y_j\|_2^2 \right) \right)$     iteration  $\leftarrow$  iteration + 1OT $_{\alpha,\beta} \leftarrow \sum_{i=1}^N \alpha_i f_i + \sum_{j=1}^M \beta_j g_j$ Initialize  $f_i \leftarrow 0 \forall i = 1, \dots, N$  and  $g_j \leftarrow 0, \forall j = 1, \dots, M$ .**while**  $L_1 \text{ relative error in } f > 0.1\%$  **do**    **for**  $i = 1, \dots, N$  **do**         $f_i \leftarrow \frac{1}{2} \left[ f_i - \varepsilon \log \left( \sum_{k=1}^N \alpha_k \exp \left( \frac{1}{\varepsilon} f_k - \frac{1}{\varepsilon} \|x_i - x_k\|_2^2 \right) \right) \right]$ **while**  $L_1 \text{ relative error in } g > 0.1\%$  **do**    **for**  $j = 1, \dots, M$  **do**         $g_j \leftarrow \frac{1}{2} \left[ g_j - \varepsilon \log \left( \sum_{k=1}^M \beta_k \exp \left( \frac{1}{\varepsilon} g_k - \frac{1}{\varepsilon} \|y_j - y_k\|_2^2 \right) \right) \right]$  $S_\varepsilon \leftarrow$  OT $_{\alpha,\beta} - \sum_{i=1}^N \alpha_i f_i - \sum_{j=1}^M \beta_j g_j$ 

---

### 3.3.4 Important metric for filter stability

We focus mainly on studying  $\mathbb{E}[D_\varepsilon(\pi_n(\mu_0), \pi_n(\mu_b))]$  and how it changes over time as more and more information arrives with the observations sequentially. For a nonlinear stochastic filtering problem under additional assumptions [83], it is proved that the filter is exponentially stable. We do not expect those assumptions to hold for the deterministic case, but this leads us to fit a curve that has exponential behaviour in time of the following form,

$$\mathbb{E}[D_\varepsilon(\pi_n(\mu_0), \pi_n(\mu_b))] = a \exp(-\lambda t) + c, \quad (3.18)$$

where  $t$ , is the assimilation time multiplied by the observation gap and  $c$  is a constant to which the distance saturates asymptotically, what one can expect from using finite sample size, as motivated in the above subsection 3.4.1.

We then study the convergence of the filter mean to the true trajectory over time. We

define the following quantity called scaled  $l_2$  error,

$$\begin{aligned} e_n(\nu) &\stackrel{\text{def}}{=} \frac{1}{\sqrt{d}} \left\| \mathbb{E}_{\hat{\pi}_n(\nu)}[x_n] - x_n^{\text{true}} \right\|_2, \\ &= \left[ \frac{1}{d} \sum_{i=1}^d \left( \frac{1}{N} \sum_{\alpha=1}^N x_n^{\alpha,i} - x_n^{\text{true},i} \right)^2 \right]^{1/2}, \end{aligned} \quad (3.19)$$

where  $d$  is the dimension of the system,  $N$  is the ensemble size, superscript  $\alpha$  and  $i$  denote the index for the ensemble member and the component respectively, and  $n$  is the assimilation step. We expect that for a reliable filter,  $\mathbb{E}[e_n^2(\nu)] \sim \sigma^2$  [55].

Another quantity that we plot is  $s_n(\nu)$  which captures the uncertainty in the best estimate of the state obtained by the filter, defined as

$$\begin{aligned} s_n(\nu) &\stackrel{\text{def}}{=} \left[ \frac{1}{d} \text{tr} \left[ \mathbb{E}_{\hat{\pi}_n(\nu)}[(x_n - \mathbb{E}_{\hat{\pi}_n(\nu)}[x_n]) (x_n - \mathbb{E}_{\hat{\pi}_n(\nu)}[x_n])^t] \right]^{1/2} \right], \\ &= \left[ \frac{1}{d} \sum_{i=1}^d \frac{1}{N-1} \sum_{\alpha=1}^N \left( x_n^{\alpha,i} - \sum_{\beta=1}^N x_n^{\beta,i} \right)^2 \right]^{1/2} \end{aligned} \quad (3.20)$$

It is the square root of the trace of the sample covariance of EnKF filtered ensemble at the corresponding assimilation step  $n$ . This quantity is often understood as a measure of reliability by comparing it with the  $L_2$  error. We note that we are unaware of any theoretical results on the asymptotic nature of this quantity.

## 3.4 Results and discussion

### 3.4.1 Understanding numerical properties of Sinkhorn divergence

In this section, we explore the effect of finite sample size in the numerical computation of Sinkhorn divergence. Our goal is to understand the effect of using empirical distributions for computing and interpreting the numerical results obtained. We first ask how the Sinkhorn distance  $S_\epsilon$  obtained using empirical distributions obtained by sampling one and the same distribution approaches zero with increasing sample size  $N$ . For this purpose, we use the Gaussian distribution as the true distribution from which we generate the sample since ensemble Kalman filters use the empirical mean and covariance while performing the posterior approximations of the filter.

We draw two samples  $\alpha_m^d = \frac{1}{m} \sum_{i=1}^m \delta_{x_i^{m,d}}$  and  $\beta_m^d = \frac{1}{m} \sum_{i=1}^m \delta_{y_i^{m,d}}$ , where  $\{x_i^{m,d}\}$  and  $\{y_i^{m,d}\}$  are two i.i.d. samples drawn from  $\mathcal{N}_d^\lambda = \mathcal{N}(0_d, \lambda I_d)$ . We choose  $\lambda = 0.1, 1.0$ . We then perform the computation of Sinkhorn divergence using the two samples from

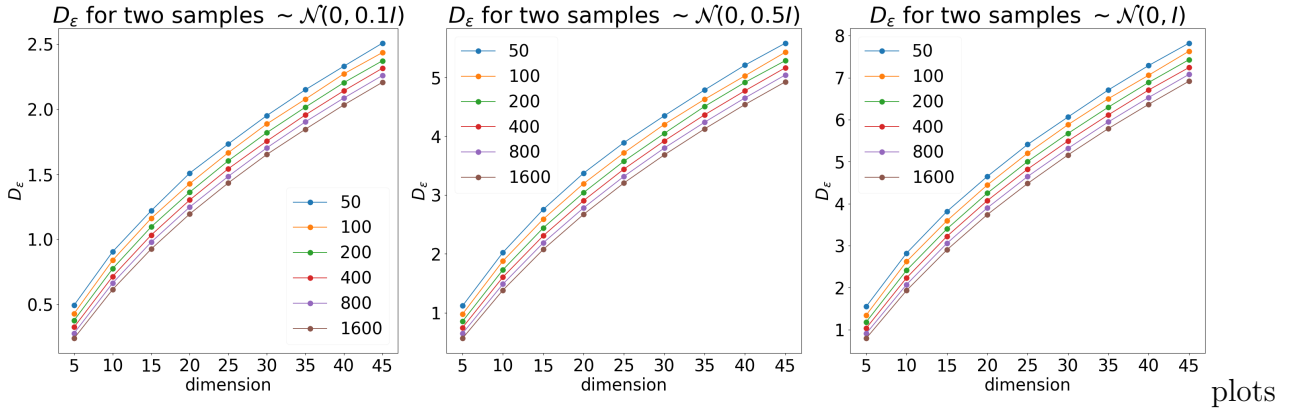


Figure 3.4: Average  $D_\varepsilon(\alpha_m^d, \beta_m^d)$  (over 20 realizations) where  $\alpha_m^d, \beta_m^d$  are two different sampling distributions with the same sample size  $m$  for the same underlying  $d$ -dimensional Gaussian  $\mathcal{N}(0_d, \lambda I_d)$

the above for Gaussian distribution, for which the exact Wasserstein distance  $W_2$  has a well-known analytical form, given by,

$$W_2(\mu_1, \mu_2)^2 = \|\mu_1 - \mu_2\|_2^2 + \text{trace}(C_1 + C_2 - 2(C_2^{1/2}C_1C_2^{1/2})^{1/2}), \quad (3.21)$$

where  $\mu_1 \sim \mathcal{N}(m_1, C_1)$  and  $\mu_2 \sim \mathcal{N}(m_2, C_2)$  are the two multivariate Gaussian distributions.

To understand the convergence of  $W_2$  between two filtering distributions to zero, for our definition of filter stability in equation (3.2), we will first numerically explore how close  $D_\varepsilon$  can get to zero using samples. In figure 3.4 we plot the average  $D_\varepsilon(\alpha_m^d, \beta_m^d)$  where  $\alpha_m^d = \frac{1}{m} \sum_{i=1}^m \delta_{x_i^{m,d}}$  and  $\beta_m^d = \frac{1}{m} \sum_{i=1}^m \delta_{y_i^{m,d}}$ , where  $\{x_i^{m,d}\}$  and  $\{y_i^{m,d}\}$  are both samples of the same underlying  $d$ -dimensional Gaussian distribution  $\mathcal{N}_d^\lambda := \mathcal{N}(0_d, \lambda I_d)$ . For small  $\lambda$ , we expect  $D_\varepsilon$  to behave in a similar fashion to that of  $\mathcal{N}_d^\lambda$  supported on a compact set.

- *Drop with increase in sample size* Theorem VI.3 in [64] explains the monotonic drop in average  $D_\varepsilon$  for a fixed dimension while increasing the sample size which results in the distance between empirical distributions approach the distance between the true distributions.
- *Rise with increase in dimension* We note that  $D_\varepsilon(\alpha_m^d, \beta_m^d)$  increases with increasing dimension  $d$ . This can be understood from the fact that as we go higher in dimension, larger sample sizes are required to accurately represent  $\mathcal{N}_d^\lambda$ . This results in  $\alpha_m^d, \beta_m^d$  becoming poorer estimators of the underlying distribution  $\mathcal{N}_d^\lambda$ , resulting in the increase in  $d$ .
- *Drop with decrease in covariance* We note that decreasing the covariance parameter  $\lambda$  has the opposite effect since, for fixed dimension  $d$  and sample size  $m$ , a distribution with smaller covariance is better estimated by the sample. Thus, the resulting



empirical distributions  $\alpha_m^d$  and  $\beta_m^d$  become better estimators of  $\mathcal{N}_d^\lambda$  as  $\lambda$  decreases, leading to a decrease in distance between them.

- *Support of our distributions* Since Lorenz systems are known to contain bounded attractors, the trajectories for both systems were obtained by integrating after a long transient period. Hence, we assume that true filtering distributions in this case are supported on a compact set. Consequently, in the filtering experiments shown later, the zero of the Sinkhorn algorithm shows qualitatively similar behavior (e.g., in figure 3.5) with respect to dimension as seen in figure 3.4.

### 3.4.2 Results for Lorenz-96 system

We now study the problem of numerical filter stability discussed above in the context of the ensemble Kalman filter. We now present the numerical explorations of the filter stability for EnKF. To obtain the initial condition for the attractor, we start from a random initial condition and integrate the ode for a long time consisting of  $10^5$  iterations. The terminal point is then used as  $x_0$ , the initial condition for generating the true underlying trajectory. To perform twin experiments in data assimilation, one generates synthetic observations based on this trajectory by adding random noise drawn from a given distribution of the measurement noise, which we assume to be the multivariate Gaussian distribution for simplicity. We generate 10 observation realizations for the same trajectory by adding noise at each time step.

First, we present the results related to Lorenz-96 in 10 and 40 dimensions in order to illustrate the stability of EnKF as a filtering algorithm. In both the cases, we observe half of the total components, i.e. for 40 dimensional case, we observe 20 alternate even indexed components and assimilate observations every 0.1 units of time. We use observation covariance  $R = 0.1I_5$  and  $1.0I_{20}$  respectively. In L96-10, we use ensemble sizes  $N = 50$  with localization and  $N=200$  without localization. This was done to see if the localization has an effect on filter stability. For L96-40, we perform assimilation with localization for  $N = 50$ .

We use three different initial distributions for this case, denoted by  $\mu_1, \mu_2, \mu_3$ , where

$$\begin{aligned}\mu_1 &= \mathcal{N}(x_0^{true}, 0.1 \times I_d), \\ \mu_2 &= \mathcal{N}(x_0^{true} + 2 \times 1_d, 0.5 \times I_d), \\ \mu_3 &= \mathcal{N}(x_0^{true} + 4 \times 1_d, 1.0 \times I_d).\end{aligned}\tag{3.22}$$

We perform 500 assimilation steps using EnKF with 10 different observation realizations. For each observation realization, we use different initial distributions in (3.22) for initializing

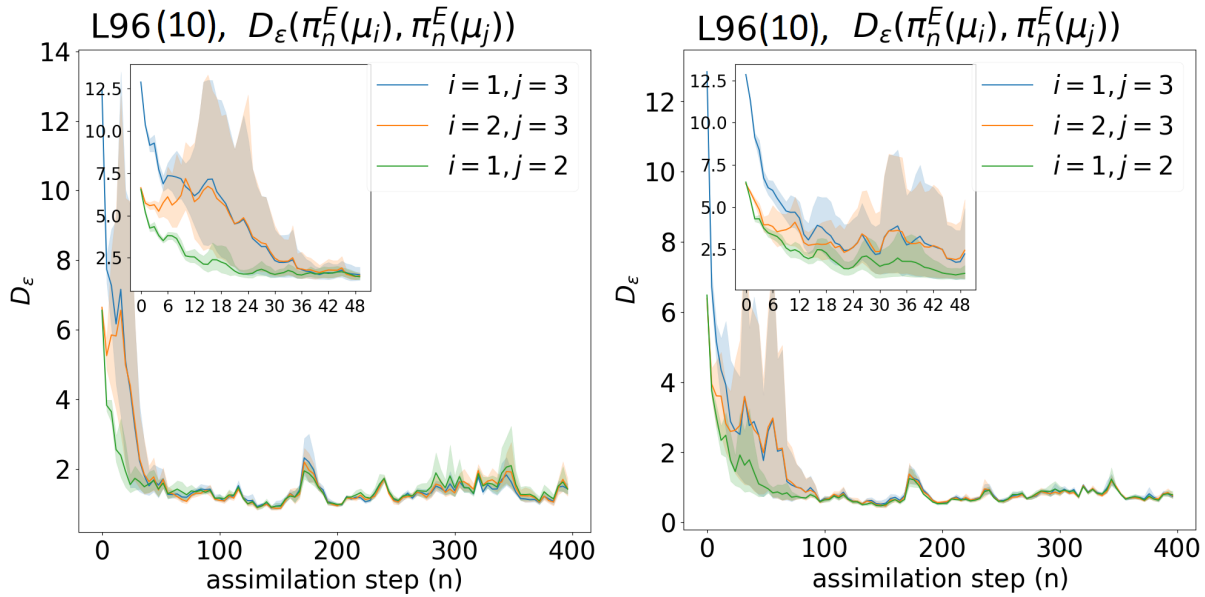


Figure 3.5:  $D_\varepsilon$  (averaged over 10 observation realizations, with one standard deviation confidence band) for EnKF for 10-dimensional L96 with  $N = 50$  with localization (left),  $N = 200$  without localization for observation covariance  $\sigma^2 = 0.1$  for pairs of initial distributions in equation 3.22. The inset shows the drop in  $D_\varepsilon$  for the first 50 assimilation steps.

EnKF. Since we have three different initial conditions, we have three corresponding pairs of distances. For each of the pairs  $(\mu_1, \mu_2)$ ,  $(\mu_2, \mu_3)$  and  $(\mu_1, \mu_3)$  and for a given observation realization, we compute the Sinkhorn divergence between the resulting filtering distribution as a function of time  $n$ . Figure 3.5 and 3.6 illustrate the stability for the three different pairs of initial distributions.

Here we use the notation  $\pi_n^E$  for  $\hat{\pi}_n$  obtained by the EnKF which was described in algorithm 1.

- **Drop in  $D_\varepsilon$  over time** From figures 3.5 and 3.6, we see that for every pair  $(\mu_i, \mu_j)$  of initial distributions  $D_\varepsilon(\pi_n^E(\mu_i), \pi_n^E(\mu_j))$  decreases with time rapidly within the first 50 assimilation steps and beyond 100 assimilation steps, the observation average of  $D_\varepsilon$  for filters with different pairs initial distributions are similar and have very little variance.
- **Variation with respect to observation realization** We see that the variation of  $D_\varepsilon$  for different observation realizations (shown by the shaded bands in figures 3.5 and 3.6) is significant for initial times (e.g.  $n < 100$  for the 10-dimensional L96 model). On the other hand, for larger times (approx.  $n > 100$ ), this variation is comparatively smaller.
- **Effect of localization** EnKF with a small ensemble size needs localization which, however, is an ad-hoc procedure to prevent filter divergence and may not approximate

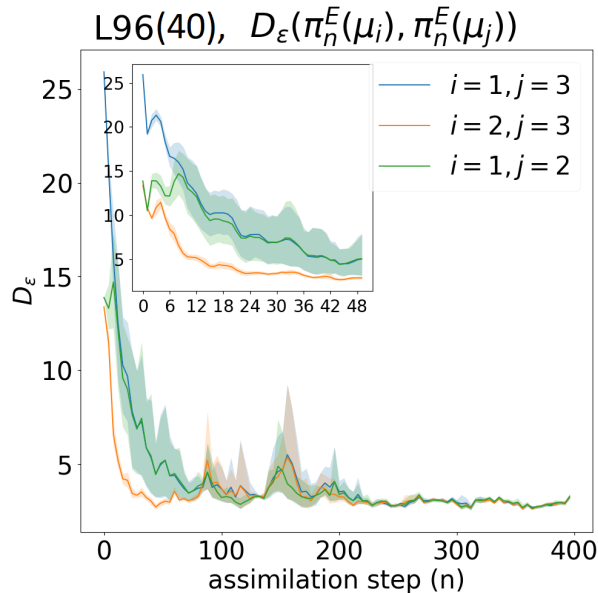


Figure 3.6:  $D_\varepsilon$  (averaged over 10 observation realizations, with one standard deviation confidence band) for EnKF with 40-dimensional L96 with  $N = 50$  with localization (right) with observation covariance  $\sigma^2 = 1.0$  for pairs of initial distributions in 3.22. The inset shows the drop in  $D_\varepsilon$  for the first 50 assimilation steps.

the true filter. Figure 3.5 for 10-dimensional L96 (left panel) and figure 3.6 for 40-dimensional L96 (right panel) shows that for  $N = 50$  with localization length 4, the EnKF is stable, whereas the right panel in figure 3.5 (with the same configuration as the left panel) for 10-dimensional L96 without localization, but with larger ensemble size  $N = 200$ . This indicates that localization does not affect EnKF’s stability properties.

We now systematically study the dependence of filter stability on two important parameters, namely the observation gap and the observation noise covariance. These two parameters are crucial for a numerical filtering algorithm as the former determines how frequently the new information coming from the observation is incorporated into the filter and the other determines the accuracy of the observations themselves, a direct measure of the observation quality. We perform experiments with Lorenz-96 in 10 dimensions by varying observation gaps and observation covariances as a parameter. We plot the average of the  $D_\varepsilon$  over 10 different realizations and also plot the points corresponding to different individual realizations to show their variability.

To study the filter stability rates, we use two Gaussian initial distributions as initial conditions to perform data assimilation; however, other choices are also possible. Since we know the true initial condition  $x_0$  in synthetic experiments, we use this to generate two different initial distributions, one denoted by  $\mu_0$  which is unbiased and precise, and the other denoted by  $\mu_b$  which is biased and imprecise. The exact distributions with their

$g$		<b>0.01</b>	<b>0.03</b>	<b>0.07</b>	<b>0.09</b>
a	<b>PF</b>	$5.367 \pm 0.080$	$7.077 \pm 0.055$	$8.672 \pm 0.084$	$9.54 \pm 0.40$
	<b>EnKF</b>	$9.28 \pm 0.13$	$10.08 \pm 0.27$	$11.11 \pm 0.32$	$10.76 \pm 0.37$
$\lambda$	<b>PF</b>	$10.73 \pm 0.76$	$4.423 \pm 0.058$	$2.203 \pm 0.021$	$1.392 \pm 0.052$
	<b>EnKF</b>	$3.904 \pm 0.085$	$3.56 \pm 0.15$	$4.24 \pm 0.21$	$2.95 \pm 0.17$
c	<b>PF</b>	$4.03425 \pm 0.00083$	$2.7524 \pm 0.0018$	$2.1362 \pm 0.0093$	$1.69 \pm 0.14$
	<b>EnKF</b>	$0.258 \pm 0.016$	$0.459 \pm 0.036$	$0.711 \pm 0.046$	$0.827 \pm 0.064$

Table 3.1: Parameters of the best-fit for the mean  $D_\varepsilon$  versus time as in (3.18) with associated confidence intervals for fixed observation covariance  $\sigma^2 = 0.4$  and different observation gap  $g$  shown in the top row.

parameters are as follows:

$$\begin{aligned}\mu_0 &= \mathcal{N}(x_0^{true}, 0.1I_d) \\ \mu_b &= \mathcal{N}(x_0^{true} + 4.01d, I_d)\end{aligned}\tag{3.23}$$

where,  $1_d$  and  $I_d$  represent a vector with all entries 1 and the identity matrix of order  $d \times d$  respectively.

We plot the Sinkhorn distance between the two distributions over time, with different initial measures used to initialize the numerical filter at  $n = 0$ . We now discuss the numerical stability experiments performed using EnKF for the Lorenz-96 model which we introduced in subsection 2.6.2 with their parameters in a detailed manner and describe the obtained results. We also have results for the particle filtering algorithm in [65] which was used in our joint work. where the exact implementation of the PF algorithm is present in Algorithm 2.

### 3.4.3 Dependence of the filter stability w.r.t observation gap

We first discuss the results of our numerical experiments with EnKF for a fixed observation covariance  $\sigma^2 = 0.4I_{10}$  with varying observation gap. As shown in figure 3.7 and 3.8, the mean  $D_\varepsilon$  falls exponentially over time until reaching a stationary value for both the filters. Table 3.1 shows the values of the coefficients of the best-fit of mean  $D_\varepsilon$  versus time, according to equation (3.18), for different observation gaps  $g$ .

For PF, we see that the rate  $\lambda$  decreases with increasing observation gap  $g$  but for EnKF, the rates are not significantly affected by the change in  $g$ . The highest Lyapunov exponent for the model (with the chosen parameter value) is approximately  $\lambda_{\max} = 1.7$  whereas the exponential rate  $\lambda$  for EnKF is seen to be in the range of (3.0, 4.2), close to  $2\lambda_{\max}$ , indicating a possible close relation between the dynamics and the EnKF that could be explored further. The exponential rate for the PF does not seem to show such a relation.

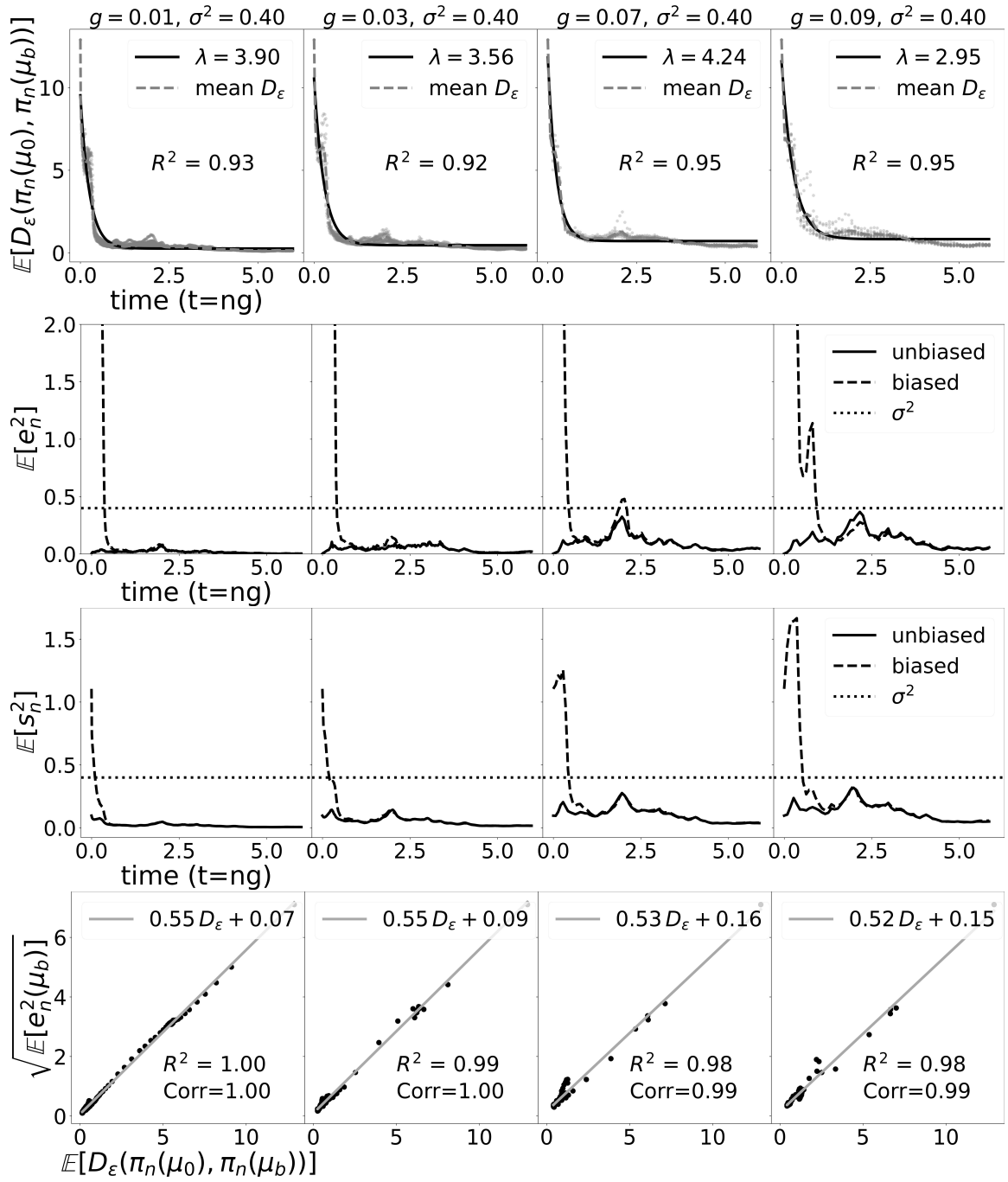


Figure 3.7: Results for EnKF with fixed observational error variance  $\sigma^2 = 0.4$ , and each column contains the results for different time between observations  $g = 0.01, 0.03, 0.07, 0.09$ . Row 1: Mean  $D_\epsilon$  versus time. The dots represent 10 different realisations. The solid line is the exponential best-fit line for the mean  $D_\epsilon$  as in (3.18). Row 2: Mean scaled  $l_2$  error from (3.19) versus time for the two initial distributions. Row 3: Mean uncertainty from (3.20) versus time for the two initial distributions. The constant dotted line in rows 2 and 3 shows the observational error variance  $\sigma^2$  for reference. Row 4: RMSE versus mean  $D_\epsilon$ . Pearson correlation coefficient between these two quantities is depicted alongside the goodness of fit for the best-fit line.

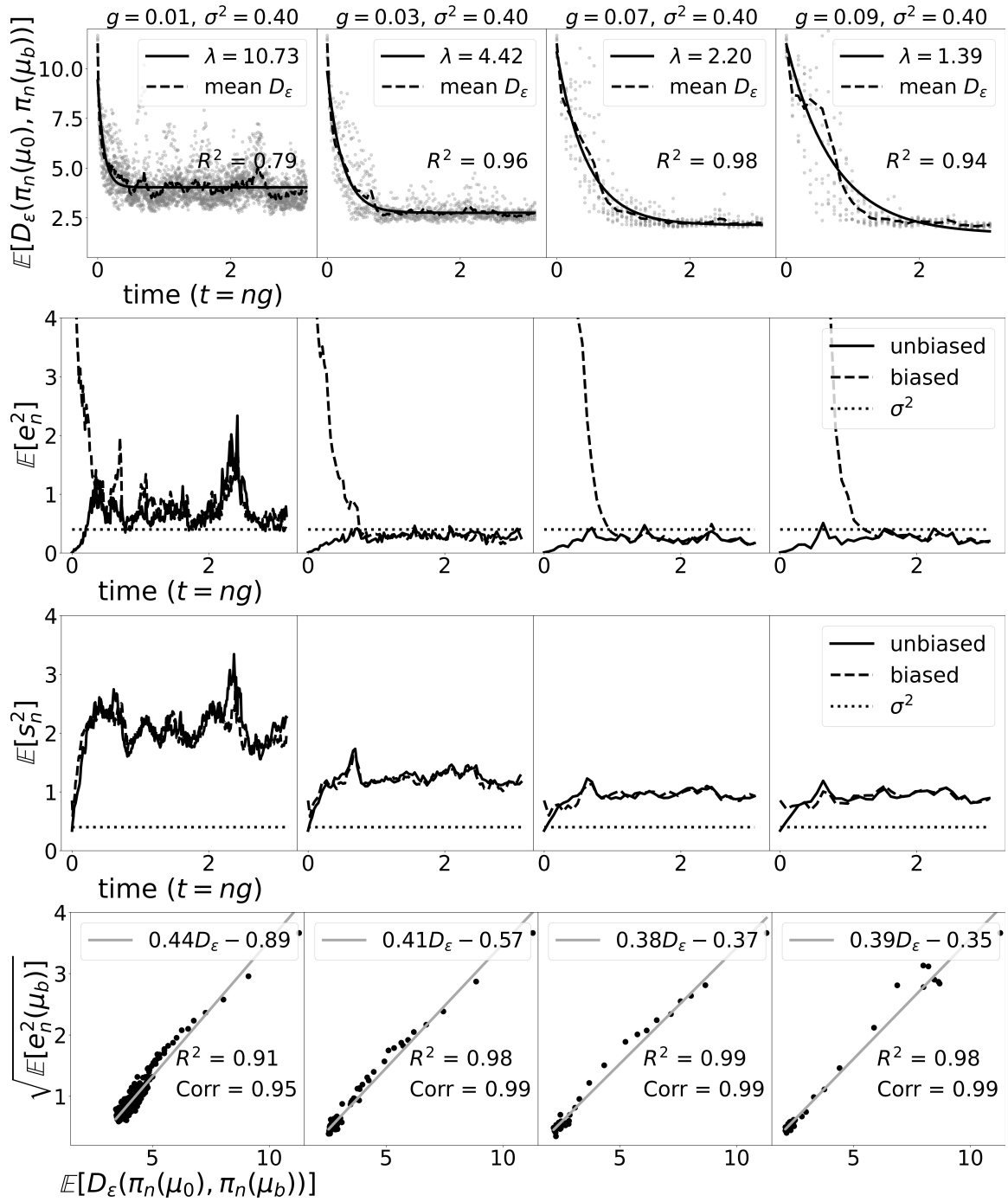


Figure 3.8: Results for PF with fixed observational error variance  $\sigma^2 = 0.4$ , and each column contains the results for different time between observations  $g = 0.01, 0.03, 0.07, 0.09$ . The different rows corresponds to are arranged in the same order as in figure 3.7. Row 1: Mean  $D_\epsilon$  versus time. Row 2: Mean scaled  $l_2$  error versus time. Row 3: Mean uncertainty versus time for the two initial distributions. Row 4: RMSE versus mean  $D_\epsilon$ .

Another difference between PF and EnKF may also be noted: with increasing observation gap, the stationary value  $c$  for the PF decreases whereas it increases for the EnKF. Also, the stationary values  $c$  of the  $D_\epsilon$  for EnKF are significantly lower compared to corresponding values for PF. These asymptotic values of  $D_\epsilon$  over time for both filters

can be explained by their mean posterior covariance using the following argument.

A characteristic of the numerical distance  $D_\varepsilon$  is that for two different *i.i.d.* samples drawn from the same probability distribution,  $D_\varepsilon$  has a nonzero positive value. Statistically,  $D_\varepsilon$  between two empirical measures approaches this value at which they are essentially representing the same distribution and cannot be distinguished. For a fixed dimension  $d$  and sample size  $N$ , this value increases with increasing covariance of the distribution [64, Figure 1 and discussion therein].

The mean posterior covariance trace is directly proportional to the  $s_n^2$ . With increasing observation gap  $g$ , the mean uncertainty decreases for PF while it increases for EnKF. Hence the asymptotic value of  $D_\varepsilon$  decreases with increasing observation gap for PF, while for the latter, it increases. We note that the previous paragraph explains the asymptotic value of  $D_\varepsilon$ , but the difference in the behaviour of the filter uncertainty  $s_n$  for PF and EnKF as a function of observation gap needs to be explored further.

The scaled  $l_2$  errors also reach an asymptotically constant value around the same time the corresponding filters stabilize in  $D_\varepsilon$ . The scatter plots for the RMSE against the  $D_\varepsilon$  shows strong correlation between them. This suggests that we can use the RMSE over time as a good indicator for the time when the filter stabilizes. Note that the methods in this paper give us a direct way to check whether a numerical filter is stable for a given dynamical and observational model, and the relation between the filter stability and the  $l_2$  error or bias  $e_n$  implies that a stable filter may be expected to be an accurate one.

We note that in the plots in the bottom row, the cluster at the bottom left corresponds to the time after which both RMSE and the  $D_\varepsilon$  have reached their stationary values. Even for two different biased initial distributions for the filter, there is a finite transient growth after which  $D_\varepsilon$  falls exponentially [64]. Although not shown here, the linear regime is still present in the scatter plot of the RMSE of either one of them versus the  $D_\varepsilon$  in those cases.

### 3.4.4 Dependence of the filter stability w.r.t observation noise

We now discuss the results of numerical experiments for fixed observation gap  $g = 0.05$ , with varying observation error variances  $\sigma^2 = 0.2, 0.4, 0.8$  and  $1.6$ . In figure 3.9 and 3.10, we again note the exponential decrease of the distance  $D_\varepsilon$  over time until it reaches a stationary value  $c$ . The parameter values obtained for the best fit for different observation covariances are shown in table 3.2.

In contrast with the case of varying observational gap, the exponential rates for the PF stability are not affected by the change in observational uncertainty. While the rates for EnKF are again close to twice the Lyapunov exponent, the rates for PF are smaller.

The scaled  $l_2$  error and the  $D_\varepsilon$  achieve their stationary value around the same time as in the former case of fixed observation. As expected, this asymptotic value  $c$  as well as the asymptotic values of the uncertainty  $s_n$  and the bias  $e_n$  all increase with increasing  $\sigma^2$  for

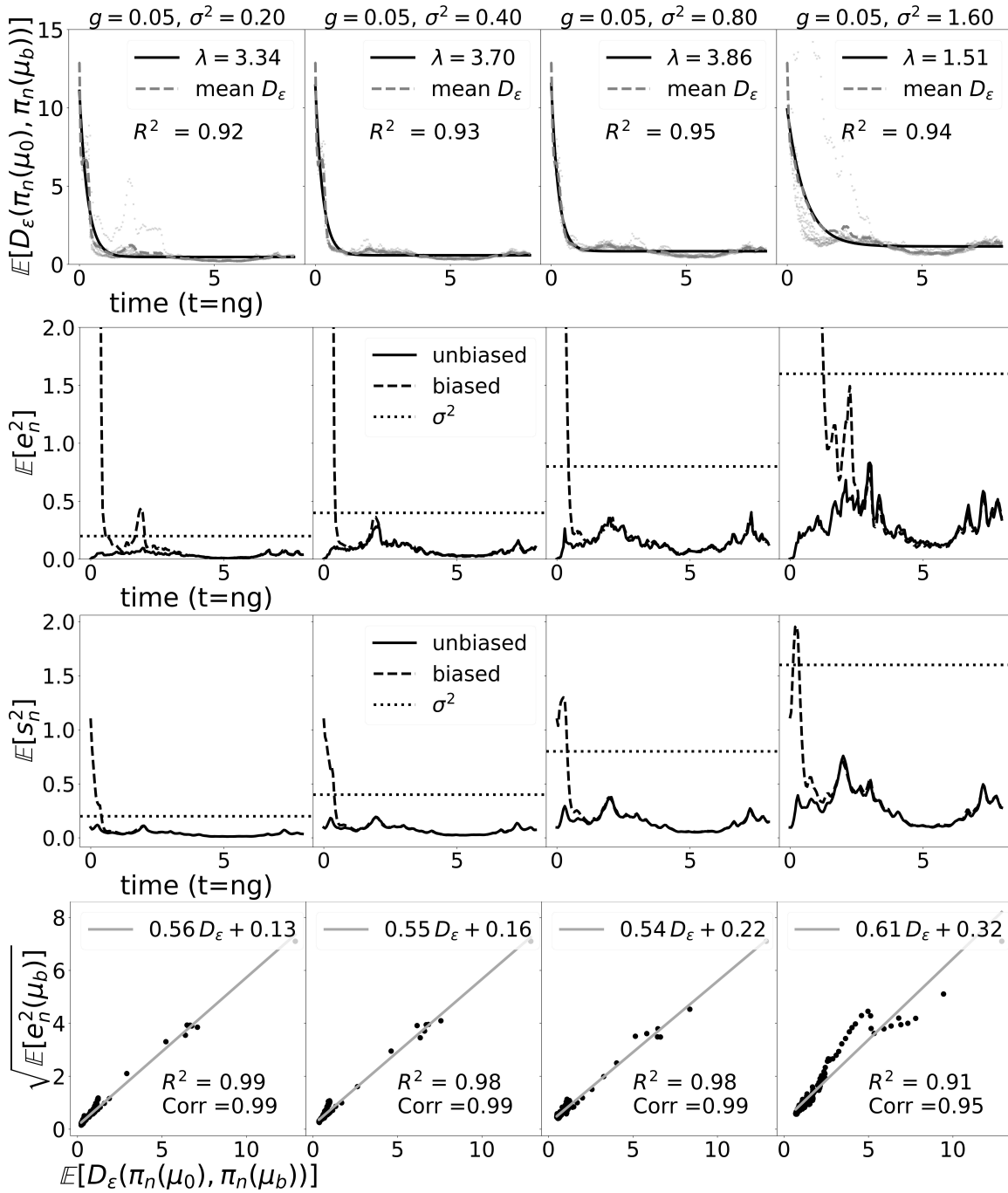


Figure 3.9: Showing the results for EnKF with fixed time between observations  $g = 0.05$ , and each column containing the results for different observational error variances  $\sigma^2 = 0.2, 0.4, 0.8, 1.6$ .

both PF and EnKF.

We also see near perfect correlation in the scatter plots for the RMSE versus mean  $D_\epsilon$  as for both PF and EnKF, we get a Pearson correlation coefficient very close to 1. We remark that in our numerical experiments, either with varying observational time interval or with varying observational covariance, we did not notice any relation between stability and posterior uncertainty or precision, i.e., there did not seem to be any relation between



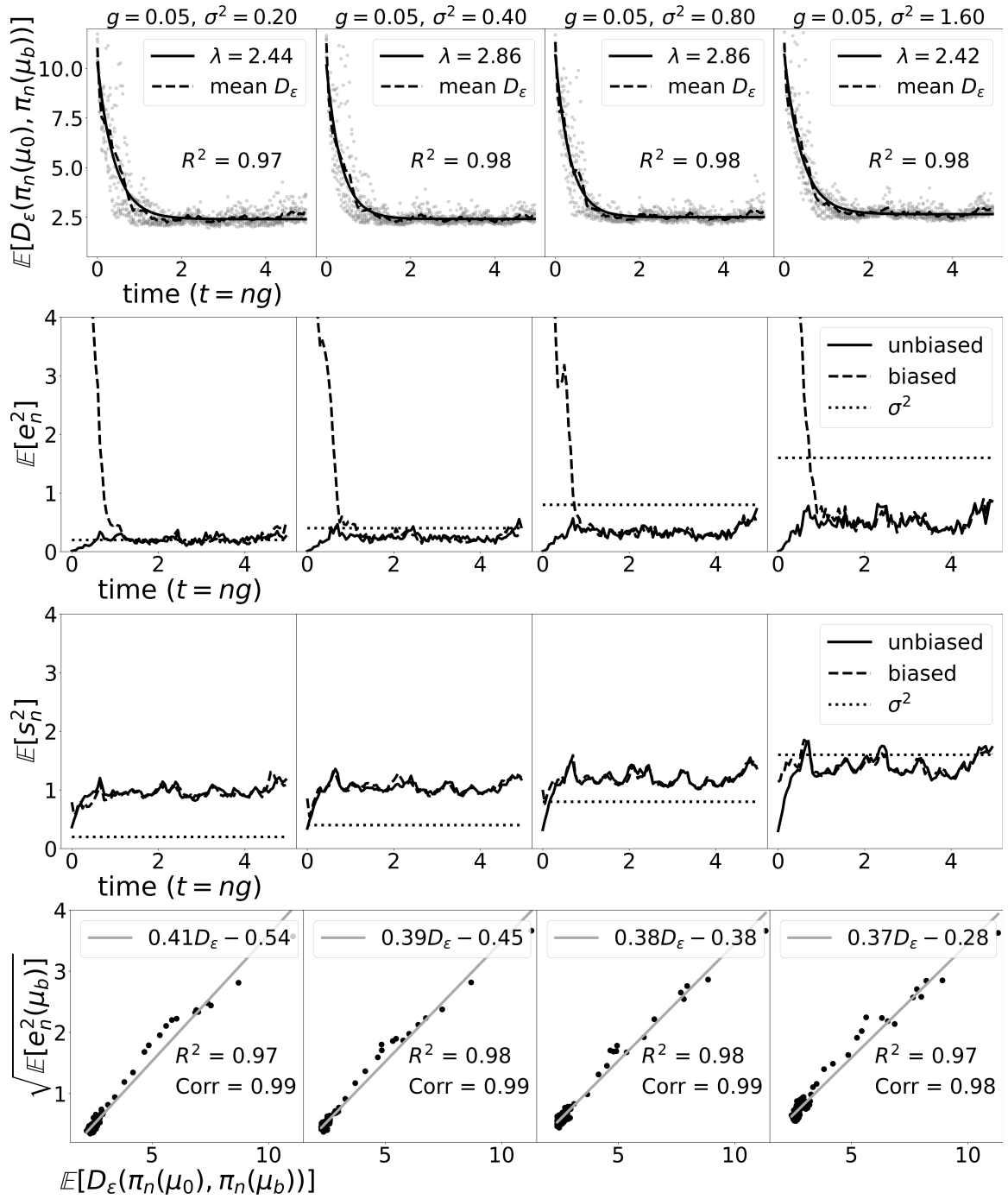


Figure 3.10: Same as in figure 3.9 showing the results for PF for fixed time between observations  $g = 0.05$ , and each column containing the results for different observational error variances of  $\sigma^2 = 0.2, 0.4, 0.8, 1.6$ .

$D_\epsilon$  and  $s_n$ .

$\sigma^2$		<b>0.2</b>	<b>0.4</b>	<b>0.8</b>	<b>1.6</b>
a	<b>PF</b>	$7.842 \pm 0.058$	$7.730 \pm 0.046$	$8.153 \pm 0.044$	$8.038 \pm 0.048$
	<b>EnKF</b>	$10.61 \pm 0.32$	$10.84 \pm 0.30$	$10.69 \pm 0.23$	$8.75 \pm 0.22$
$\lambda$	<b>PF</b>	$2.442 \pm 0.016$	$2.858 \pm 0.017$	$2.859 \pm 0.015$	$2.416 \pm 0.012$
	<b>EnKF</b>	$3.34 \pm 0.16$	$3.70 \pm 0.16$	$3.86 \pm 0.14$	$1.507 \pm 0.062$
c	<b>PF</b>	$2.4050 \pm 0.0022$	$2.4144 \pm 0.0015$	$2.5051 \pm 0.0014$	$2.6554 \pm 0.0019$
	<b>EnKF</b>	$0.470 \pm 0.039$	$0.579 \pm 0.035$	$0.838 \pm 0.027$	$1.148 \pm 0.041$

Table 3.2: Parameters of the best-fit for the mean  $D_\epsilon$  versus time as in (3.18) with associated confidence intervals for fixed observation gap  $g = 0.05$  and different observational error covariance  $\sigma^2$  shown in the top row.

### 3.5 Summary

The central focus of this chapter revolves around the comprehensive examination of the stability of nonlinear filters with a focus on the ensemble Kalman filter, utilizing a numerical approach. We employ the recently devised Sinkhorn algorithm to compute an approximation of the Wasserstein distance between Monte Carlo samples, each representing a filtering distribution. This method facilitates a direct evaluation of the stability of a filter by computing the expected value, averaged over multiple observation realizations, of the distance between filtering distributions as a function of time. Furthermore, we undertake an extensive exploration of the dependence of the stability characteristics of the filter with respect to two crucial parameters: the time between observations and the observational error covariance. Our results provide substantial numerical evidence of exponential stability for EnKF in the context of deterministic chaotic dynamical systems.

As the gap between observations increases, the rate of decay of the distance  $D_\epsilon$  between two filters decreases for the particle filter, while it stays approximately the same and close to twice the Lyapunov exponent for the ensemble Kalman filter. Investigating the connection between the chaotic characteristics of the system and exponential stability is an interesting direction for future work. Proving the stability of these numerical algorithms mathematically is another intriguing open field of study.

# Chapter 4

## Computing Lyapunov instabilities of a dynamical system using data assimilation

In the previous chapter, we studied the problem of nonlinear filter stability numerically using EnKF and the dynamical systems introduced in this thesis. The notion was to numerically understand the convergence of the filtering distribution over time. We now arrive at the second problem of computing the stability( instability) properties of the underlying dynamical system- the Lyapunov vectors and the exponents. The asymptotic instability directions of a chaotic dynamical system are defined by the Lyapunov vectors and their associated growth rates called the Lyapunov exponents. The objective of this chapter is to understand how good our reconstructed instability directions are when computed from approximations of a true underlying trajectory. Our focus is specifically on the two types of these vectors, namely backward Lyapunov vectors, the BLVs, and covariant Lyapunov vectors or the CLVs.

The organization of this chapter is as follows. In section 4.1, we start with the importance of instability directions and their utility in data assimilation. In section 4.2, we provide an overview of the theory and mathematical definitions of Lyapunov vectors and their subspaces in section 4.2.1. We then explain in detail Ginelli's algorithm which we use for the computation of covariant Lyapunov vectors in section 4.2.2. Our proposed data-based algorithm using data assimilation for computation of LVs is presented in section 4.2.3. We also explain our approach of computing the approximate LVs of perturbed trajectories in section 4.2.4. In section 4.2.7, we describe the metrics employed for comparison of the exact and perturbed or approximate CLV and the related Oseledet's subspaces. We present the details of the numerical implementation are presented in section 4.2.6. The numerical results are discussed in section 4.3 followed by a summary of conclusions and directions of further studies in section 4.4.

## 4.1 Introduction

The uncertainties in the state estimates, obtained using DA, depend quite crucially on the directions of the error growth, or in other words, the dynamical instabilities of the system. Lyapunov exponents and vectors are the fundamental tools used in the study of nonlinear and chaotic systems, in order to characterize the stability properties of a dynamical system with respect to perturbations along different directions. The paradigm called assimilation in unstable subspace, see, e.g., [82, 20, 18, 19, 80] uses the subspace spanned by the unstable and neutral Lyapunov vectors for producing the analysis or the update step of assimilation and shows promising improvements over the traditional algorithms. The importance of Lyapunov vectors in general and in the context of assimilation is also discussed extensively in recent works of [71, 85, 13, 40, 15] and others.

For a dynamical systems, a non-increasing tuple  $\{\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_s\}$  of Lyapunov exponents summarizes the global, asymptotic rate of change of linear perturbations around a trajectory. Existence of at least one positive Lyapunov exponent indicates exponential divergence of perturbations and is indicative of instabilities and chaotic dynamics. The associated directions in the tangent space are called Lyapunov vectors, which span the tangent space at a specific point in the phase space and contain information about the past and future evolution of local perturbations [33, 56]. Various methods have been proposed for computation of Lyapunov exponent, with or without the use of the dynamical equations and / or their linearization, or using a time series of observations of the system [9, 10, 1].

Lyapunov vectors that capture the asymptotic growth rate as  $t \rightarrow \pm\infty$  are called, respectively, forward or backward Lyapunov vectors (FLV / BLV). These FLV and BLV provide a orthonormal basis for a filtration of the tangent space, but they are not covariant with respect to the dynamics: the time evolution, under the linear dynamics, of an FLV (resp. BLV) at one time does not lead to an FLV (resp. BLV) at other time but leads to a linear combination of FLVs (resp. BLVs). Further, even though the asymptotic growth rate of the FLV in forward time is the Lyapunov exponent, it is not so in backward time (and similarly for BLV). Vectors that have both these properties, namely covariance with respect to the time evolution and asymptotic growth rates in forward and backward time, are called covariant Lyapunov vectors (CLV), which are the central focus of this work. In recent years, various methods for computing these Lyapunov vectors have been developed [88, 52, 39, 67].

The main aims of this chapter are twofold: one is to present a data-based algorithm for computation of Lyapunov vectors and the other is to study their sensitivity to noise. The first part - a data-based algorithm for calculating CLV presented in 4.2.3 along with the results presented in 4.3.1 - is largely motivated by the work of [81, 82] and subsequent developments of the AUS (assimilation in unstable subspace) methodology. Their work focuses mainly on BLVs whereas we consider a natural extension to compute the CLV as

well. Recently, another data-based algorithm has been proposed in [66], using the data to reconstruct the system dynamics and then use this reconstructed dynamics to compute CLVs. Our approach differs in a fundamental way, since our basic assumption - which is also common with all the data assimilation (DA) methods - is that a dynamical model is available but we do not know the trajectory that is being observed. This naturally leads to our proposed algorithm that uses a DA algorithm to compute an optimal estimate, which is then used to compute the CLV. We note that any algorithm, including our algorithm, for computation of CLV necessarily needs to be “offline” since it requires the backward evolution from far future. But one by-product of our algorithm that is indeed “online” (without the need to future observations) is the data-based computation of BLV using only the past observations.

The second part studies the sensitivity of the LVs to noisy perturbations of a trajectory. We study how well the LVs and the associated Oseledets’ subspaces spanned by these vectors are approximated when calculated by using a noisy trajectory instead of a true trajectory. We note that this is quite distinct from the question of continuity of the LVs with respect to the phase space, as has been studied in [50, sec. 19.02]; [5, 32, 63]. We also note that the noisy trajectories we use are neither shadowing trajectories nor are they solutions of a stochastic version of the deterministic dynamics, since the noise is added only at discrete observation times. This choice is motivated by the parallel with the data assimilated state estimates which also provide trajectories that are neither shadowing nor solutions of stochastic dynamics.

## 4.2 Theory and computation of Lyapunov vectors

In this section, we present a brief summary of the mathematical framework for defining the Lyapunov vectors, followed by a description of the method we used for computing them [39, 52]. We then discuss our proposed method for computing these vectors either (i) using state estimates obtained from the ensemble Kalman filter or (ii) using the perturbed trajectories of a dynamical system. We also discuss the metrics that we use, namely the subspace angles between Oseledets subspaces, to assess the accuracy of the approximate Lyapunov vectors obtained using these two methods - this is accuracy with respect to the Lyapunov vectors obtained from the exact trajectory (or rather, its numerical approximation by the RK4 algorithm that we use).

### 4.2.1 Definition and importance of covariant Lyapunov vectors (CLV)

We consider an autonomous continuous-time dynamical system represented by ODE of the form

$$\dot{x}_t = f(x_t), \quad \text{where, } x_t \in S \subseteq R^n \text{ and } f : S \rightarrow R^n \text{ is the vector field.} \quad (4.1)$$

Associated to such an ODE, the evolution for the infinitesimal perturbations in the tangent space is obtained by linearizing along a trajectory, thus leading to the following ODE for  $z_t \in R^n$ :

$$\dot{z}_t = J(x_t)z_t, \quad \text{where } J(x) \text{ is the jacobian matrix given by } J_{ij}(x_t) = \frac{\partial f_i(x_t)}{\partial x_{t,j}}, \quad (4.2)$$

where  $f_i$  and  $x_{t,j}$  denote the  $i^{\text{th}}$  and  $j^{\text{th}}$  component of  $f(x_t)$  and  $x_t$ . A fundamental matrix solution of this linear non-autonomous equation solves  $\dot{Q}_t = J(x_t)Q_t$  with any non-singular matrix  $Q_0 \in R^{n \times n}$  as an initial condition. In the discussion below, we consider the evolution of the trajectory and the corresponding perturbations at times  $\dots, t_k, t_{k+1}, \dots$ , and use the notation  $x_k \equiv x_{t_k}$ ,  $z_k \equiv z_{t_k}$ , etc. With this notation, using the fundamental matrix solution, the tangent linear propagator from time  $t_k$  to  $t_l$  can be written as

$$\mathcal{M}_{k,l} = Q_l Q_k^{-1} \quad \text{with the property that } z_l = \mathcal{M}_{k,l} z_k. \quad (4.3)$$

The eigenvectors and eigenvalues of  $\mathcal{M}_{k,l}$  for large  $l \rightarrow \infty$  and  $k \rightarrow -\infty$  capture the asymptotic stability properties of the perturbations around a trajectory and will be the primary focus in this paper. The existence of these limits is the main content of the Oseledets' multiplicative ergodic theorem [69, 70]. Specifically, under suitable conditions, the theorem implies the existence of the following limits:

$$\lambda^+(z, t_k) := \lim_{t_l \rightarrow \infty} \frac{1}{t_l - t_k} \log \frac{\|\mathcal{M}_{k,l} z\|}{\|z\|}, \quad \text{and} \quad \lambda^-(z, t_l) := \lim_{t_k \rightarrow -\infty} \frac{1}{t_k - t_l} \log \frac{\|\mathcal{M}_{k,l} z\|}{\|z\|}, \quad (4.4)$$

where  $\lambda^\pm(z, t_k)$  are called the Lyapunov characteristic exponents. Under appropriate regularity conditions, the two limits in (4.4) give the same set of Lyapunov exponents but with opposite sign and we drop the superscript  $\pm$  for the exponents. There are at most  $n$  distinct such exponents, which are usually ordered as  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_s$  with  $s \leq n$ . Oseledets' theorem also proves the existence of the Oseledets subspaces

$$\{0\} = S_{s+1}^+ \subsetneq S_s^+ \subsetneq \dots \subsetneq S_1^+ = R^n \quad (4.5)$$

with the property that  $\lambda(z, t_k) = \lambda_i$  when  $z \in S_i^+ \setminus S_{i+1}^+$ . Such vectors  $z$  satisfying this latter property are called the *forward Lyapunov vectors*. The Oseledets subspaces corresponding to the  $t_k \rightarrow -\infty$  limit in (4.4) are given by

$$\{0\} = S_0^- \subsetneq S_1^- \subsetneq \cdots \subsetneq S_s^- = R^n \quad (4.6)$$

and analogously define the *backward Lyapunov vectors*. Note that these subspaces depend on time, i.e. on  $t_k$  (resp.  $t_l$ ) for forward (resp. backward) Lyapunov vectors, so more precisely,  $S_j^+ = S_j^+(t_k, x_k)$  etc. For autonomous dynamical systems, this time dependence is only through the trajectory, i.e., on the phase space point  $x_k = x(t_k)$ . Thus, more precisely,  $S_j^+ = S_j^+(x_k)$ . Though this dependence was dropped above for simplicity of notation, exploring the dependence of the Oseledets subspaces and the Lyapunov vectors on the phase space trajectories is one of the main aims of this paper, as we discuss below.

The above discussion naturally raises the question of whether nearby points in phase space have Oseledets spaces that are “close” to each other in an appropriate metric. In other words, this is a question of continuity of these Oseledets spaces with respect to phase space and has been investigated theoretically in [50, 5, 32, 63], proving Hölder continuity of these spaces, and numerical methods for computing derivatives of CLVs has been developed recently in [22]. But we are not aware of any numerical study of the continuity of the Oseledets spaces with respect to perturbations to the full trajectory. The main focus of this paper is precisely to address this lacunae by numerically studying the sensitivity of the Lyapunov vector to perturbations in phase space. One of the key difficulties in these numerical investigations is explained in detail in section 4.2.2.

The forward and backward Lyapunov vectors are not mapped to each other under the action of the tangent linear operator, i.e., they are not covariant. Further they are also not invariant with respect to time reversal. Instead they satisfy the following property: if the forward (resp. backward) Lyapunov vectors at time  $t_k$  are arranged in columns of  $\Phi^+(t_k)$  [resp.  $\Phi^-(t_k)$ ], then

$$\mathcal{M}_{k,l} \Phi^+(t_k) = \Phi^+(t_l) L_{k,l} \quad \text{and} \quad \mathcal{M}_{k,l}^{-1} \Phi^-(t_l) = \Phi^-(t_k) R_{k,l}, \quad (4.7)$$

where  $L_{k,l}$  and  $R_{k,l}$  are, respectively, lower and upper triangular matrices. The diagonal elements of these matrices give the local stretching or contraction of the Lyapunov vectors. These properties motivate the numerical algorithms to calculate the BLV and FLV, e.g., see the review [52].

To summarize, the BLV and FLV are not covariant but form orthonormal bases of the Oseledets subspaces. On the other hand, these subspaces are covariant under the linear dynamics as indicated by (4.7). By looking for bases that may not be orthonormal, it is possible to find a set of basis vectors of these Oseledets’ spaces that are covariant with

respect to the dynamics and invariant with respect to time reversal. Such basis vectors are called *covariant Lyapunov vectors* (CLV). In particular, they have the property that the  $i$ -th covariant Lyapunov vector  $q_i(t_k)$  satisfies the dynamics given by,

$$q_i(t_l) = \mathcal{M}_{k,l} q_i(t_k) \quad \text{and} \quad \|\mathcal{M}_{k,k+l} q_i(t_k)\| \sim e^{\lambda_i t_l} \quad \text{for} \quad t_l \rightarrow \pm\infty. \quad (4.8)$$

The FLV satisfy the above only in the limit of  $t_l \rightarrow +\infty$  but not in the other limit of  $t_l \rightarrow -\infty$ , and similarly the BLV satisfy only one of the two limits above.

The existence of such covariant Lyapunov vectors is guaranteed by the following fact: the dimensions of  $i$ -th forward and backward Oseledets subspace  $S_i^+$  and  $S_i^-$  are, respectively,  $d_s + \dots + d_i$  and  $d_1 + \dots + d_i$ . Since the sum of these dimensions is  $n + d_i$ , their intersection has minimum dimension of  $d_i$ . We can see that the vectors that belong to this intersection  $S_i^+ \cap S_i^-$  satisfy the properties (4.8) and are the covariant Lyapunov vectors that we seek. Indeed the numerical algorithm presented below directly make use of the fact that they belong to this intersection.

## 4.2.2 Computation of Lyapunov vectors

We now discuss the dynamic algorithm introduced in the work of Ginelli [39] for the computation of BLVs and CLVs about a reference trajectory. We assume that we have the database consisting of states  $x_j$  sampled at time  $t_j = j\Delta t$  with  $\Delta t$  as the interval between two consecutive states from a fixed trajectory starting at  $j = 0$ . Since we have pre-computed trajectories, we first carry out time integration of the perturbation vectors in tangent space, using the state information from the trajectory whenever required in the evaluation of the jacobian in equation (4.2) to obtain the BLVs. More specifically, BLVs are required as an intermediate step since they provide a basis in the tangent space which is then used to represent the CLVs at any time  $t_j$ . This necessitates that the trajectory should be long enough to contain three distinct time intervals: (i) an initial ‘‘forward transient’’ time interval denoted as  $[0, I]$  to account for the forward transient required to converge to the BLV basis, (ii) the subsequent interval of interest denoted by  $[I, F]$  over which we want to obtain the BLVs (and later on the CLVs), and (iii) an additional ‘‘backward transient’’ interval denoted by  $[F, E]$  to account for the backward iterations to converge to give the CLVs. This is shown schematically in the left panel of figure 4.1.

We perform gram-schmidt re-orthonormalization via QR decomposition after every  $l\Delta t$  interval and use or store the vectors for integrating over the next interval. The number  $N_{0I}$  of times the QR-decomposition is performed in the forward transient time interval is given by the relation  $N_{0I}l\Delta t = I$ . Similarly for the other two intervals,  $N_{IF}$  and  $N_{FE}$  denote the number of times the QR-decomposition is performed, i.e.,  $N_{IF}l\Delta t = F - I$  and  $N_{FE}l\Delta t = E - F$ .



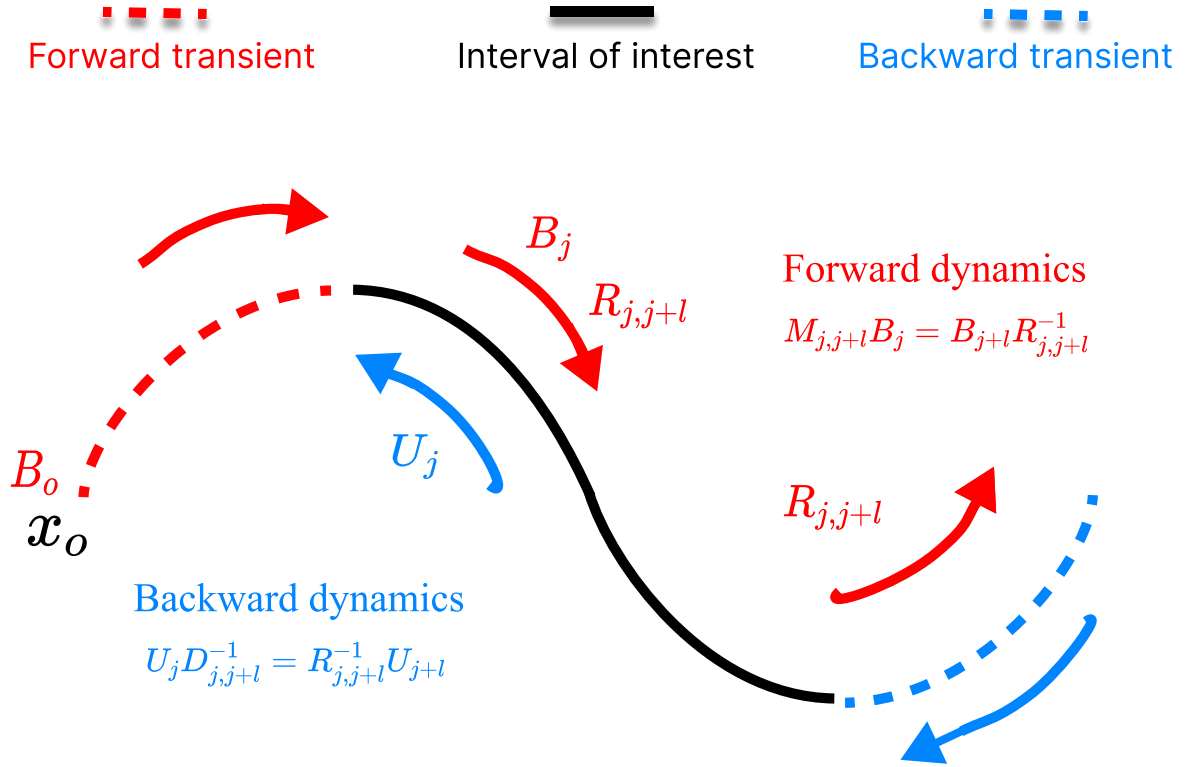


Figure 4.1: Schematic diagram of Ginelli's dynamic algorithm for computing covariant Lyapunov vectors.

Let  $B_j \in R^{n \times m}$  with  $m \leq n$  denote the matrix containing a set of orthogonal perturbation vectors along the columns at time  $t_j$ . For the forward transient interval, we initialize  $B_0$  as the initial condition for equation (4.2) and integrate over small time intervals  $[t_j, t_{j+l}]$  of length  $l\Delta t$ , at the end of which we perform the gram-schmidt re-orthonormalization of the columns of  $B_j$ . Recall that  $\mathcal{M}_{j,j+l}$  is the linear tangent propagator from time  $t_j$  to  $t_{j+l}$ . So the evolution of the perturbation vectors followed by re-orthonormalization satisfies the following relations,

$$\tilde{B}_{j+l} = \mathcal{M}_{j,j+l}B_j \quad (\text{evolution}) \quad \text{and} \quad \tilde{B}_{j+l} = B_{j+l}R_{j,j+l} \quad (\text{re-orthonormalization}), \quad (4.9)$$

for  $j = 0, l, \dots, (N_{0I} - 1)l$  where,  $R_{j,j+l}$  is the matrix containing the growth rates of the vectors over the interval  $[t_j, t_{j+l}]$  obtained by QR decomposition of  $\tilde{B}_{j+l}$ . The choice of  $l$  must be small enough so as to prevent the rank collapse of the columns of  $\tilde{B}_j$  at any time  $t_j$  where the matrices are stored. The length of the forward transient interval  $[0, I]$  is chosen to be long enough such that at time  $I$ , the columns of the matrix  $B_j$  provide a good approximation of the columns of  $\Phi_j^-$ , i.e., the BLVs, as defined in equation (4.2).

From this time onwards, we store both the matrices  $B_j$  and  $R_{j-l,j}$ , as computed in (4.9) for  $j = (N_{0I})l, \dots, (N_{0I} + N_{IF} - 1)l$ , thus covering the time interval  $[I, F]$  during which we want to compute the CLVs. Note that over this time interval, we already get the numerically approximate BLVs as the columns of the matrices  $B_j$ . The main idea is to note the following fact: if the CLVs at time  $t_j$  are arranged as columns of matrix  $C_j$ , then their relation to the BLV basis can be expressed as

$$C_j = B_j U_j, \quad (4.10)$$

for some  $U_j$  which is an upper triangular coefficient matrix. This is due to the fact that the  $i^{\text{th}}$  CLV lies in the span of the first  $i$  BLVs. Since the CLV are covariant, they satisfy the following relation:

$$C_{j+l} D_{j,j+l} = \mathcal{M}_{j,j+l} C_j, \quad (4.11)$$

where  $D_{j,j+l}$  is a diagonal matrix with the diagonal elements being the norms of the columns of the product on the right-hand side. Combining the above three relations (4.9)-(4.11), the backward evolution of the perturbation vectors in matrices  $U_j$  followed by re-normalization satisfies the following relations,

$$\tilde{U}_j = R_{j,j+l}^{-1} U_{j+l} \quad (\text{backward evolution}) \quad \text{and} \quad \tilde{U}_j = U_j D_{j,j+l}^{-1} \quad (\text{re-normalization}). \quad (4.12)$$

Hence the rest of the algorithm is the backward evolution aimed at calculating these matrices  $U_j$ .

At the end of the interval  $[I, F]$ , we further integrate forward the perturbation vectors using (4.9) for  $j = (N_{0I} + N_{IF})l, \dots, (N_{0I} + N_{IF} + N_{FE} - 1)l$ , thus covering the time interval  $[F, E]$ . During this time, we only store the matrices  $R_{j-l,j}$ . At the end of this backward transient interval, we set  $U_j$  for  $j = (N_{0I} + N_{IF} + N_{FE})l$  to be a generic upper triangular full-rank matrix  $U_F$ . This is to ensure that the  $i^{\text{th}}$  column of the matrix  $C_j$  lies in  $i^{\text{th}}$  backward Oseledets subspace  $S_i^-$  at  $j = (N_{0I} + N_{IF} + N_{FE})l$ , since under the backward evolution, a generic vector in this subspace  $S_i^-$  converges asymptotically to the  $i^{\text{th}}$  CLV.

The backward dynamics is then performed on the upper-triangular coefficient matrices  $U_j$  over the backward transient interval via (4.12) for  $N_{FE} + N_{IF}$  number of steps for  $j = (N_{0I} + N_{IF} + N_{FE} - 1)l, \dots, (N_{0I} + N_{IF})l, \dots, (N_{IF})l$ , utilizing the inverses of  $R_{j-l,j}$  which were computed and stored in the course of forward evolution over the respective time interval. This step gives us the set of upper triangular matrices  $U_j$  over the time interval  $[F, E]$  which is the backward transient interval and also over the interval  $[I, F]$  over which the relation (4.10) between BLVs and CLVs is valid with the matrices  $U_j$ . Thus over the interval  $[I, F]$ , we obtain the BLVs that are stored in the matrices  $B_j$ , and

using the upper triangular coefficient matrices  $U_j$ , we also obtain the CLVs as columns of  $C_j = B_j U_j$ .

For an more in-depth discussion of the algorithm and its convergence, we refer to [39, 52] and [67].

### 4.2.3 Data-based algorithm to calculate the Lyapunov vectors

We now describe the problem of computing the LVs when we cannot observe the full system in time, i.e., we only have access to partial and noisy observations of some of the components of the state  $x_j$  instead of the full trajectory. As discussed above in section 4.2.2, Ginelli’s algorithm requires a reference trajectory or the initial condition at  $t_j = 0$  which can be integrated forward to obtain that reference trajectory. The above procedure cannot be carried out directly when the initial condition is unknown and only partial observations of the state over time are available. This is because of the exponential divergence of nearby trajectories for chaotic systems.

Under the condition that the true trajectory may only be observed partially and indirectly through noisy measurements of state-dependent physical quantities, nonlinear filtering aims to obtain optimal estimates of the state that are in proximity to the true trajectory, or more precisely, the posterior probability distribution called the filtering distribution or simply the filter. [See, e.g., 17, 7, 36, for reviews and further references.] Most common numerical filtering algorithms compute Monte Carlo approximations of the filtering distributions, and the mean of the filter – called the analysis mean – is an optimal estimate of the true state.

When a numerical filter performs reasonably well, we expect the analysis mean to be sufficiently near the true state. The filter may also be used to give an uncertainty associated with the analysis mean, most commonly in terms of the covariance of the filter. One of the important factors affecting this uncertainty is the observational uncertainty and we focus on this aspect in this chapter. Some of the other factors affecting the filter performance include observational frequency, the sparsity of the observations, and the dynamical characteristics of the system itself.

We propose to use the analysis mean  $x_j^a$  obtained over time as an approximation of the state  $x_j$  of the true underlying trajectory, in order to compute the approximations of the true Lyapunov vectors and Oseledets’ subspaces. This leads to our proposed modification of Ginelli’s algorithm discussed in section 4.2.2: the (pseudo-)trajectory we use in this algorithm now consists of the analysis means  $\{x_j^a\}$  at times  $t_j$  obtained from a filtering algorithm. In order to get the tangent linear propagator  $\mathcal{M}_{j,j+1}$  over the time interval  $(t_j, t_{j+1})$ , equations (4.1)-(4.2) are solved with  $x_j^a$  as the initial condition for (4.1) and with  $B_j$  as the initial condition for (4.2). The other steps are exactly the same as described in the previous section 4.2.2.

This allows us to compute the LVs and the sub-spaces defined by them from the estimated trajectory obtained from any general data assimilation method. However, this comes with a caveat that nonlinear filtering algorithms result in an estimated trajectory that is not a dynamical trajectory of the model itself, i.e., there is no initial condition such that if integrated forward in time contains the filter analysis means obtained over time. But it is close to the true state over time which can be quantified by the  $l_2$  error over time or other popular metrics such as RMSE. We used ensemble Kalman filter (EnKF) as our choice of filtering algorithm to obtain the analysis mean trajectory. We apply this method to two models L63 and L96. Further details of EnKF and the data assimilation experimental setup are given in subsection 4.2.6.

#### 4.2.4 Lyapunov Vectors from perturbed trajectory

Any filter-based trajectory violates the criteria of being the dynamical trajectory of the system, as illustrated in figure 4.2. Even when the deviations  $e_j^a$  of the analysis from the true state, given by  $e_j^a = x_j^a - x_j$  are small, it is not *a priori* clear how they may affect errors in the computed LVs, through their effect on the Jacobian matrix in equation (4.2). This motivates us to investigate the stability of the LVs from a more general perspective.

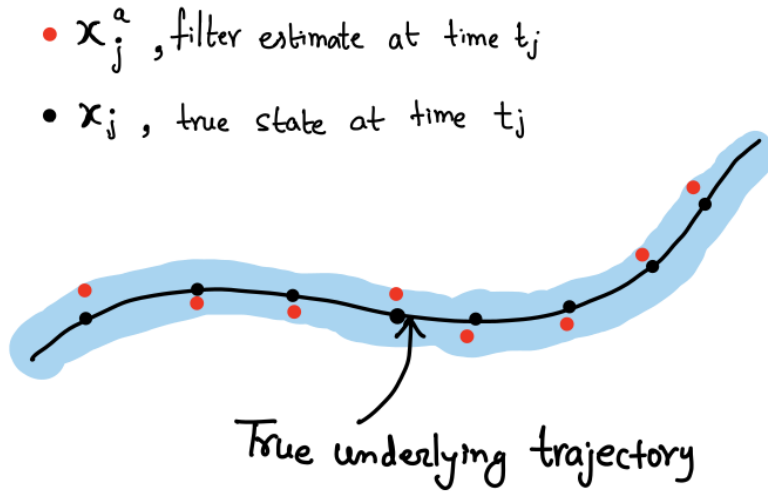


Figure 4.2: Schematic picture illustrating the idea of how filter estimated trajectory is only an approximation of the true underlying trajectory.

To systematically investigate the stability of numerically calculated LVs about the reference trajectory, we generate a perturbed or noisy version of the state  $x_j$  at each point on the trajectory by adding random perturbation to the true underlying trajectory in the following way,

$$\tilde{x}_j = x_j + \epsilon_j, \quad \text{with} \quad \epsilon_j \sim \mathcal{N}(0, \sigma^2 \mathbf{I}_n), \quad (4.13)$$

where  $\epsilon_j$  is randomly sampled from a standard normal distribution of covariance  $\sigma^2 \mathbf{I}_n$  at time  $t_j$  and  $\tilde{x}_j$  is then used as state estimates in the computation of the BLVs and CLVs at the respective time. We refer to the obtained trajectory as a perturbed orbit for the respective noise level  $\sigma$  in future discussions. We show the results for  $\sigma \in \{0.1, \dots, 0.5, 1.0, \dots, 5.0\}$  for the Lorenz-63 model and for  $\sigma \in \{0.1, \dots, 0.5\}$  for the Lorenz-96 model.

We compute all the BLVs and CLVs about the true and the perturbed trajectories over a common interval using the same procedure mentioned in section 4.2.2. The above notion of sensitivity to the perturbations in the underlying dynamical trajectory allows us to think of the analysis mean trajectory as a perturbed trajectory obtained from the true trajectory where the perturbations follow the unknown error statistics of the difference between the true state and analysis mean over time. We perform sensitivity analysis for Lorenz-63 and Lorenz-96 for  $n = 10, 20$  and  $40$ , which we introduced earlier chapters in section 2.2. As far as we know, this question of the sensitivity of BLVs and CLVs to perturbations in the space of trajectories has not been investigated either numerically or mathematically. We note that the results about their Hölder continuity with respect to initial conditions are quite distinct from the question of continuity with respect to perturbations of the whole trajectory [32, 63, 5].

## 4.2.5 Models

We apply our analysis to the compute the LV of use Lorenz-63 and Lorenz-96 [58, 61], which we introduced earlier in chapter 2, section 2.6. We briefly describe these two model here for the sake of ease. The standard Lorenz-63 model is given by the following set of equations

$$\frac{dx}{dt} = \sigma(y - x), \quad \frac{dy}{dt} = x(\rho - z) - y, \quad \frac{dz}{dt} = xy - \beta z, \quad (4.14)$$

with the parameters  $(\sigma, \rho, \beta) = (10, 28, 8/3)$  for which the system exhibits chaos and has the well-known butterfly-shaped attractor.

Lorenz-96 is a  $n$  dimensional nonlinear dissipative model with a constant external forcing term. It mimics the dynamics of a meteorological scalar variable along the latitude. The model is given by the evolution of a set of  $n$  ordinary differential equations given below:

$$\frac{dX_k}{dt} = X_{k-1}(X_{k-2} - X_{k+1}) - X_k + F \quad (4.15)$$

where  $X_k$  is the  $k^{th}$  component of the  $n$ -dimensional state and with periodic boundary conditions  $X_{k+n} = X_k$ . We choose  $n=10, 20$  and  $40$ . For the specific value of forcing  $F = 8$  for  $n = 40$  dimensions, it is a chaotic system with 13 positive Lyapunov exponents and

has a Kalpan-Yorke dimension close to 28.4.

## 4.2.6 Details of computing Lyapunov Vectors from filtered and perturbed trajectories

For the Lorenz-63 model in equation (4.14), we start with a random initial condition and integrate for a long transient up to  $t = 500$  to reach the attractor. We then choose this point on the attractor as the initial condition for the true orbit which is generated by numerically integrating the ODE for a total time  $T = 350$  with  $\delta t = 0.002$  using the Runge-Kutta 4<sup>th</sup> order scheme and storing the state every  $\Delta t = 0.01$ . To generate observation for the data assimilation experiment, we choose to observe only the  $Y$  coordinate at every  $\Delta t = 0.01$  with noisy observations given by  $y_j = Hx_j + \eta_j$  with  $H = [1, 0, 1]$  and  $\eta_j \sim \mathcal{N}(0, \mu^2)$ . We show the results for  $\mu = 0.1, 0.3, 0.5, 0.7, 0.9$ . In order to start the assimilation, we use an arbitrary initial distribution  $\mathcal{N}(x_0 + 6 \times \mathbf{1}_3, 2.0 \times \mathbf{I}_3)$  at  $t_{j=0}$ , where  $\mathbf{1}_3$  is a vector with all entries as 1, from which we generate  $N = 25$  ensemble members to initialize the EnKF algorithm. We then assimilate the previously generated observations  $y_j$  performing 35000 assimilation steps. We didn't use any inflation and localization for this simple case of Lorenz-63 ODE. Neglecting the first 5000 assimilation steps, we perform the computation of both BLV and CLV over a smaller interval excluding the transient intervals as mentioned below.

For the computation of LVs using the algorithm explained in section 4.2.2, the forward transient  $[0, I]$ , the sampling interval of interest  $[I, F]$ , and the backward transient  $[F, E]$  are all chosen to be equal to 100 which was found to be sufficient for the convergence to BLVs. The QR decomposition is performed every  $l = 1$  step. For the choice of initial perturbations  $B_0$ , we use the standard orthonormal basis vectors and integrate them forward in time using the tangent linear equations described in section 4.2.4.

For the case of the Lorenz-96 model, we perform data assimilation for 40-dimensional model, where we assimilate 20 observations taken at the evenly indexed components of the full state every time step of  $\Delta t = 0.05$ . We then generate the observations using  $y_j = Hx_j + \eta_j$  with  $H$  being the appropriate matrix to observe the alternate coordinates and  $\eta_j \sim \mathcal{N}(0, \mu^2 \mathbf{I}_{20})$  for the observation noise statistics. The initial condition for the filter was chosen to be the distribution  $\mathcal{N}(x_0 + 5.0 \times \mathbf{1}_3, 1.0 \times \mathbf{I}_{40})$  which is biased, as may happen often in practice. We implement covariance localization using the localization radius  $r = 4$  when performing the analysis as it is necessary when we have partial observations with ensemble sizes smaller than the system dimensions [17]. We vary  $\mu$  as a parameter to study the dependence on observation noise in the reconstruction of the trajectory by relating it to the *RMSE* of the obtained analysis trajectory.

To study the dimension dependence of sensitivity in L96, we perform similar numerical

experiments in dimensions equal to 10, 20 and 40. We generate a long trajectory of length  $T$ , where for dimension 10 and 20, we choose  $T = 600$  whereas for 40 we choose  $T = 1000$ , all obtained using the solver time step of  $\delta t = 0.01$  while saving the states every  $\Delta t = 0.05$ . For dimensions 10 and 20, we chose the forward transient interval  $[0, I]$  and backward one  $[F, E]$  to be of length 200 whereas for dimension 40, we choose them to be 400, and the interval of interest  $[I, F]$  is of length 200. The QR decomposition is performed every  $l = 5$  steps. The lengths of forward and backward transients were chosen based on the convergence of BLVs for the respective systems.

## 4.2.7 Comparison metrics

We compare the individual LVs and the Oseledets subspaces obtained using the perturbed or the assimilated trajectories with those obtained from the true underlying trajectory. To understand the sensitivity of the individual LVs with respect to the perturbation strength  $\sigma$ , we use the angle between the  $i^{\text{th}}$  Lyapunov vector computed from the original trajectory and the perturbed trajectory. Note that we compute the Lyapunov vectors over a length of trajectory (of course discretely sampled) and at each of the phase space points of that trajectory, we compute the cosine of the angle between the LVs of the true trajectory and LVs of the approximate trajectory. We then plot the median of the angle computed over the sampling interval along with the error bars which represent the 25<sup>th</sup> and 75<sup>th</sup> percentile of the distribution of the angle obtained.

A subspace of dimension  $k$  in  $R^n$  admits an infinite number of possible basis vectors. Thus, even though the individual LVs from the perturbed trajectory and those from the true trajectory may differ, the Oseledets subspaces spanned by them may still be “similar.” In order to quantify this “similarity” and understand the sensitivity of the Oseledets subspaces, we use the principal angles as described in detail below.

Principal angles [43] are defined as a sequence of minimum angles between two unit vectors corresponding to each of the two subspaces, such that the unit vectors are chosen to minimize the angle between them while being orthogonal to all the previous unit vectors obtained in their respective subspaces. Mathematically, for two subspaces  $\mathcal{P}$  and  $\mathcal{Q}$ , the principal angles are defined by an  $m$ -tuple of angles  $\theta(\mathcal{P}, \mathcal{Q}) = [\theta_1, \theta_2, \dots, \theta_m]$ , where  $m = \min(\dim(\mathcal{P}), \dim(\mathcal{Q}))$  and  $\theta_k$  is defined by,

$$\cos(\theta_k) = \max_{p_k \in \mathcal{P}} \max_{q_k \in \mathcal{Q}} |p_k^T q_k| \quad \text{subject to } \|p_k\| = \|q_k\| = 1, \quad p_k^T p_i = q_k^T q_i = 0, \quad i < k. \quad (4.16)$$

When the set of angles between two subspaces are all small, they are closely aligned and have a strong degree of similarity. On the other hand, if the angles are large, this implies that the subspaces are more dissimilar and have less overlap. The two subspaces are

identical (when they have the same dimension) or one is fully contained in the other (when they have different dimensions) if and only if all the principal angles are zero.

In order to investigate the quality of subspaces we obtain from the approximate LVs, we study  $i \leq k$  principal angles between the subspace spanned by the first  $k$  BLVs computed from the true and the perturbed trajectory. The aim is to understand whether a set of  $k$ -BLVs have a common lower dimensional subspace. These results are discussed in section 4.3.4. Similar to the relative angle of BLVs and CLVs, we compute the principal angles for different points along the trajectory and plot the median along with 25<sup>th</sup> and 75<sup>th</sup> percentile for the confidence intervals, see, e.g., figure 4.5).

There are several methods for calculating the principal angles between two subspaces, including the singular value decomposition (SVD) and the QR decomposition. If the two orthogonal bases  $p_i$  and  $q_i$  are arranged along the columns of two matrices  $P$  and  $Q$  respectively, we let  $P^T Q = U \Sigma V^T$  denote the singular value decomposition of  $P^T Q$ . Then the cosines of the principal angles are the diagonal elements of the matrix  $\Sigma$ .

### 4.3 Results and discussion

We now discuss the results of using the algorithms described above, by comparing the true LVs and the LVs obtained from the assimilated and the perturbed trajectories for both Lorenz-63 and Lorenz-96 systems. In section 4.3.1, we first describe the results obtained from using assimilated trajectory for different values of observation noise strength  $\mu$ . We plot the median of the angle between the true vectors and the ones obtained from the assimilated trajectory along with a confidence interval denoting the 25<sup>th</sup> and 75<sup>th</sup> percentile, which is obtained from all the points at which we compute the LV in the time interval  $[I, F]$  along the trajectory. In section 4.3.2 and 4.3.3, we present the results of the exploration of the sensitivity of the BLVs and the CLVs to the noise strength of the perturbations in the underlying true trajectory.

For Lorenz-96 in 40 dimensions, we discuss, the approximations of the Oseledets subspaces obtained using the assimilated and the perturbed trajectories by plotting the principal angles between  $k$ -dimensional subspaces spanned by the first  $k$  BLVs obtained from the true and the approximate trajectory for  $k = 2, 5, 10, 15$ , and 20. To further understand the quality of recovered subspaces from the approximate trajectories, we compare them against principal angles between randomly generated  $k$ -dimensional subspaces in section 4.3.4. We find that the angles for randomly generated subspaces are significantly larger than those between the true and perturbed Oseledets' subspaces.



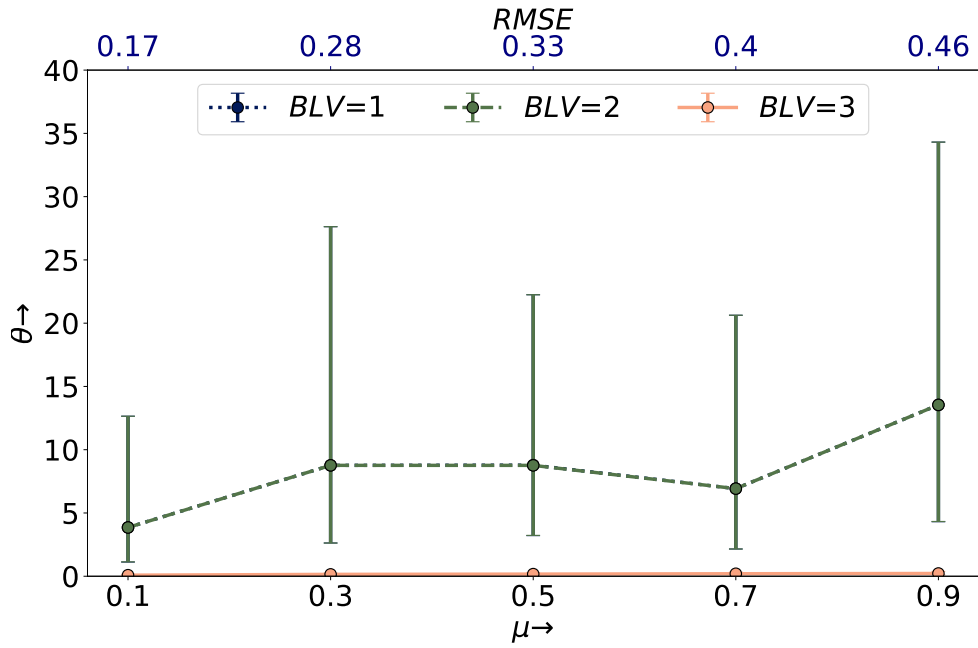
### 4.3.1 BLVs and CLVs computed from assimilated trajectories

Figure 4.3 shows the angles between the Lyapunov vectors obtained from the true trajectory with those from the assimilated trajectory for Lorenz-63 model. In this case, we only observe the  $y$ -coordinate. Left panel shows the angles for the BLV while the right one is for the CLV. We first note that since BLV are orthonormal, two of these angles are necessarily equal which happen to be those between the second and third BLV and they are quite small even for the largest observational noise strength we have used. In addition we note that the median of angle between the first - most unstable - BLV is also within 15 degrees and does not increase rapidly with the observation noise strength  $\mu$ . The CLVs show a similar behaviour but being non-orthonormal basis, all the three angles are distinct, with the angle between the third - most stable - CLV being the smallest. The top axis in these plots shows the RMSE averaged over the whole time interval of interest  $[I, F]$  of length 100 in this case. We later discuss the relation of these results to the case of perturbed trajectory discussed in detail in section 4.3.2.

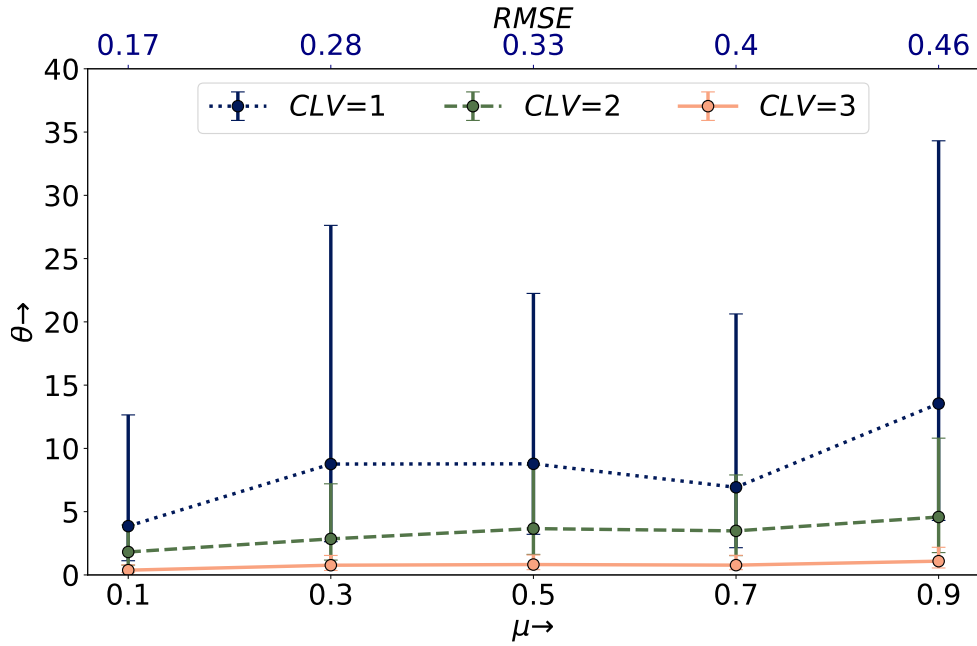
For Lorenz-96 in 40 dimensions, we performed assimilation by observing 20 alternate coordinates. We compute all the 40 BLVs and CLVs from the assimilated trajectories for different observation noise strengths  $\mu$  and the LV index versus the error in the angle. In figure 4.4 see that apart from the first few most unstable and the last few most stable LVs, the angles between the true and approximate LVs are quite large, being greater than  $45^\circ$ . Even the angles for the first - most unstable - LVs are larger than  $20^\circ$ . This indicates that the LVs obtained from the assimilated trajectories provide a very poor approximation of the true LVs. This behavior is quite distinct from the low-dimensional Lorenz-63 model for which the assimilated trajectory could be used to obtain a significantly better approximation of the true LVs, as discussed above.

Even though the individual LVs are not approximated well, it may be possible that the subspaces spanned by these vectors may have significant overlap, which is exactly the question of approximation of Oseledets' subspaces. We now investigate this question, using principle angles (PA) between these subspaces.

In figure 4.5, we plot the  $n$  principal angles between  $k$ -dimensional Oseledets' subspaces obtained from the assimilated trajectory and the true Oseledet's subspaces, for  $k \in \{2, 5, 10, 15, 20\}$ , for three different values of observational noise. We see that with an increase in the subspace dimension  $k$ , the number of principal angles which are smaller than a certain threshold, say  $20^\circ$ , increase almost linearly with  $n$ . For example, for  $k = 15$  and  $\mu = 1.0$ , there are around 11 angles less than  $20^\circ$ . This indicates within the 15-dimensional Oseledets' subspace defined by the first 15 approximate LVs, there is a 11-dimensional subspace (not necessarily Oseledets' space) which is within  $20^\circ$  of the true 11-dimensional subspace. With the increase in  $\mu$ , the relative angle increase systematically for all the P.A. The increase is more prominent for a higher P.A. index.



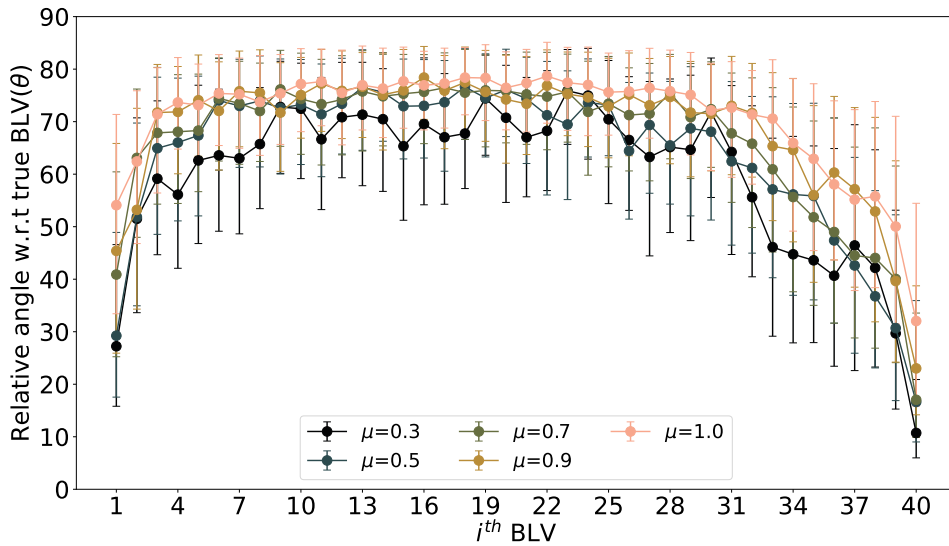
(a) Angle  $\theta$  between the true BLVs and those recovered from the analysis trajectory



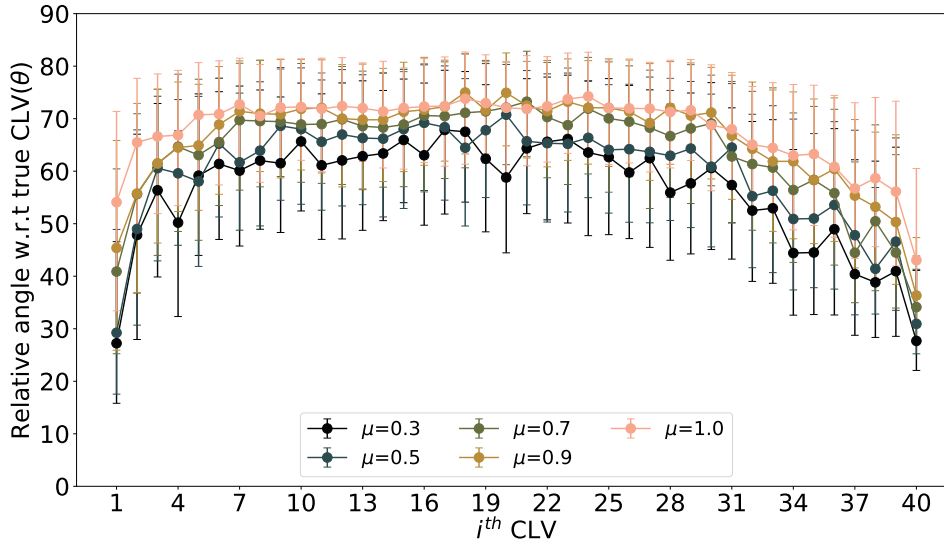
(b) Angle  $\theta$  between the true CLVs and those recovered from the analysis trajectory

Figure 4.3: The figure shows the angle  $\theta$  (in degree) between the true LVs and those recovered from the analysis trajectory, for the BLVs (top) and the CLVs (below) for the Lorenz-63 model, for different levels of observational noise  $\mu$  (bottom axis) along with the corresponding RMSE (top axis) of the analysis trajectory. The dots represent the median and the error bars represent the 25<sup>th</sup> and 75<sup>th</sup> percentiles.

To summarize, recovering individual vectors from assimilated trajectories is not possible except for the first few and the last few LVs. But embedded within any high dimensional



(a) Angle  $\theta$  between the true BLVs and those recovered from the analysis trajectory



(b) Angle  $\theta$  between the true CLVs and those recovered from the analysis trajectory

Figure 4.4: The figure shows the angle  $\theta$  (in degree) between the true LVs and those recovered from the analysis trajectory, for the BLVs (left) and the CLVs (right) for the 40-dimensional Lorenz-96 model. Different lines are for different observational noise levels  $\mu$ . The dots represent the median and the error bars represent the 25<sup>th</sup> and 75<sup>th</sup> percentiles.

Oseledets' subspace, there are lower dimensional subspaces which are close to the true subspaces. In order to understand this behaviour more clearly, we now study the dependence of this approximation on the strength of perturbation of the trajectories.

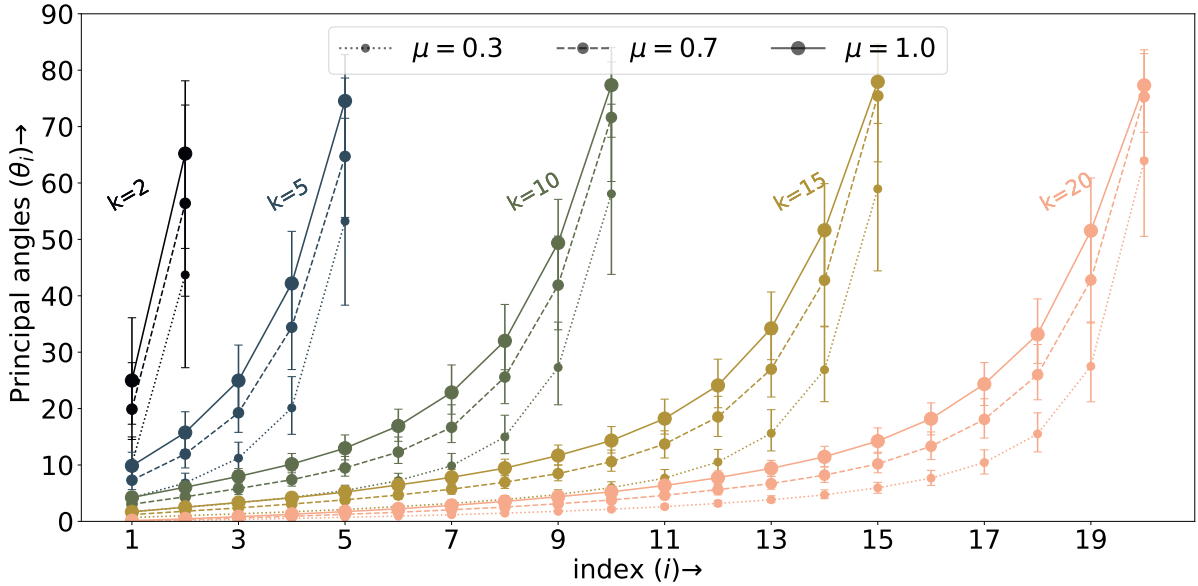


Figure 4.5: The set of principal angles between  $k$ -dimensional Oseledets' subspace and the corresponding subspace recovered from the analysis trajectory for different values of observational noise level  $\mu = 0.3, 0.7, 1.0$  for  $k = 2, 5, 10, 15, 20$  for Lorenz-96 in 40 dimensions.

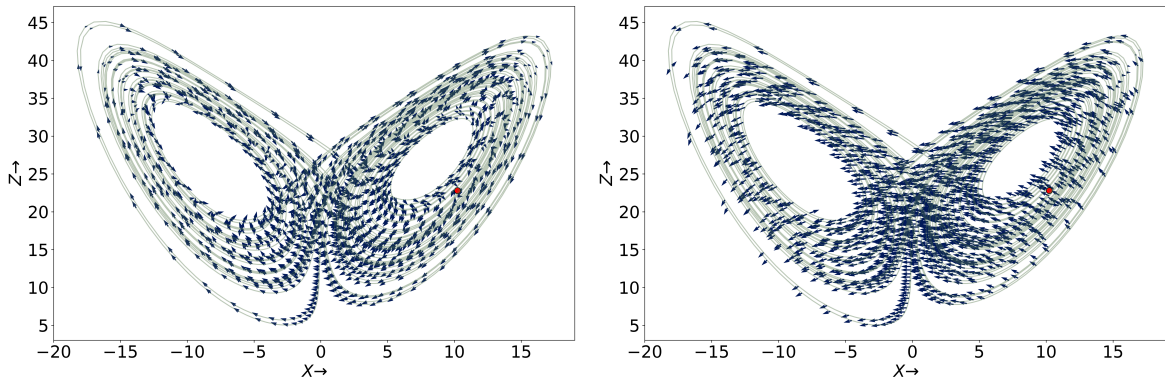


Figure 4.6: The first and the last component of the 1<sup>st</sup> and 3<sup>rd</sup> CLV plotted in XZ coordinates along the trajectory for Lorenz-63.

### 4.3.2 Dependence on perturbation strength for Lorenz-63

We first show in figure 4.6 the geometrical structure of the 1<sup>st</sup> and 3<sup>rd</sup> CLV by plotting the first and the last component of the respective vectors in the XZ-plane. We observe that the orientation of the vectors seems to change continuously as one moves along the trajectory. In general, the mutual angle between any two CLVs change along the trajectory in the phase space. It was found recently in [15] that these geometrical features such as the mutual angles between CLVs contain additional information, unlike FLVs and BLVs which are orthogonal basis vectors of the tangent space. Whenever the trajectory jumps from the right wing of the attractor to the left and vice-versa, the first two CLVs become parallel

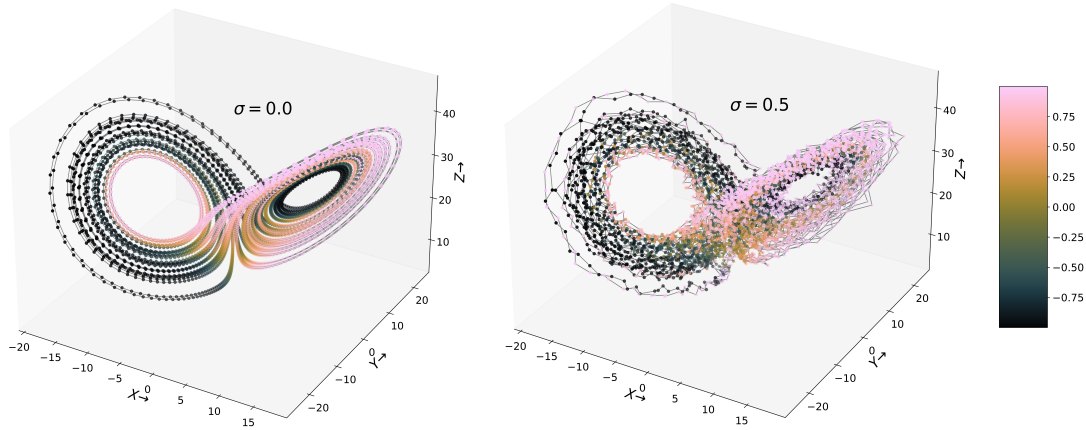


Figure 4.7: Attractor of Lorenz-63 with color indicating the cosine of the angle between the 1<sup>st</sup> and the 2<sup>nd</sup> CLV for true trajectory ( $\sigma = 0.0$ ) (left) and perturbed trajectory for  $\sigma = 0.5$  (right).

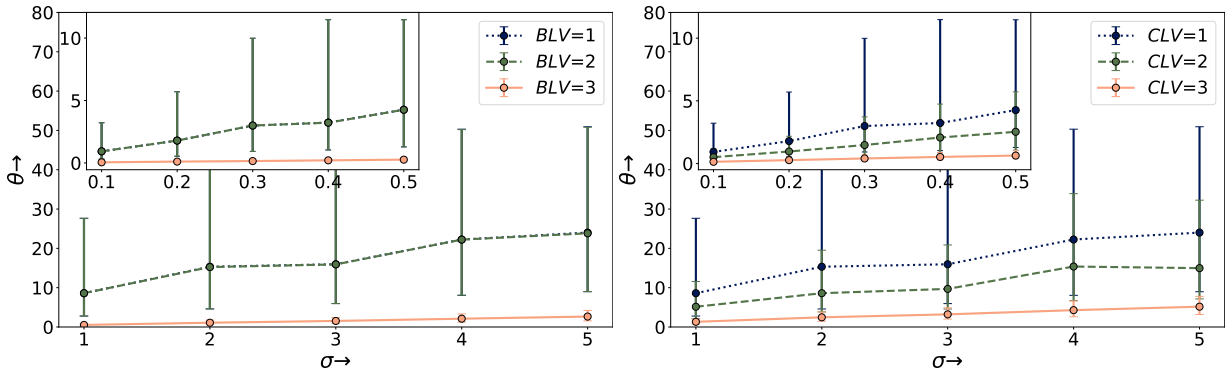


Figure 4.8: The angle  $\theta$  between the true and LV from the perturbed trajectory for different perturbation strength  $\sigma$  for Lorenz-63. The left and right panels show the results for the BLVs and the CLVs respectively. The dots represent the median and the error bars represent the 25<sup>th</sup> and 75<sup>th</sup> percentiles and the different line types for different LVs.

and anti-parallel respectively. In figure 4.7, we use colors to plot the cosine of the angle between the first two CLVs over the attractor. We also plot the same for the perturbed trajectory for perturbation strength  $\sigma = 0.5$ , which has no point on the attractor with probability 1, but the angle between the first two CLVs computed from the perturbed trajectory still captures this geometrical information, although we get a fuzzy picture of attractor using the perturbed trajectory as we increase the perturbation strength  $\sigma$ .

To study the dependence on perturbation strength  $\sigma$ , we now plot in figure 4.8 the angle between the LVs about the true trajectory and the perturbed trajectory as a function of the noise strength  $\sigma$ , for the three Lyapunov vectors. The relative angle between the true and the perturbed BLV and CLV increases gradually as we increase  $\sigma$  with a constant slope. The 1<sup>st</sup> and 2<sup>nd</sup> BLV have the same rate, whereas the rates are different for the 1st and 2nd CLVs. We also plot the absolute error in the exponents computed from the

perturbed trajectory and the original trajectory. The errors in the exponents obtained are of the order 0.2, this tells us that they are almost unaffected by the perturbation when computed from perturbed trajectories.

### 4.3.3 Dependence on dimension and perturbation strength for Lorenz-96

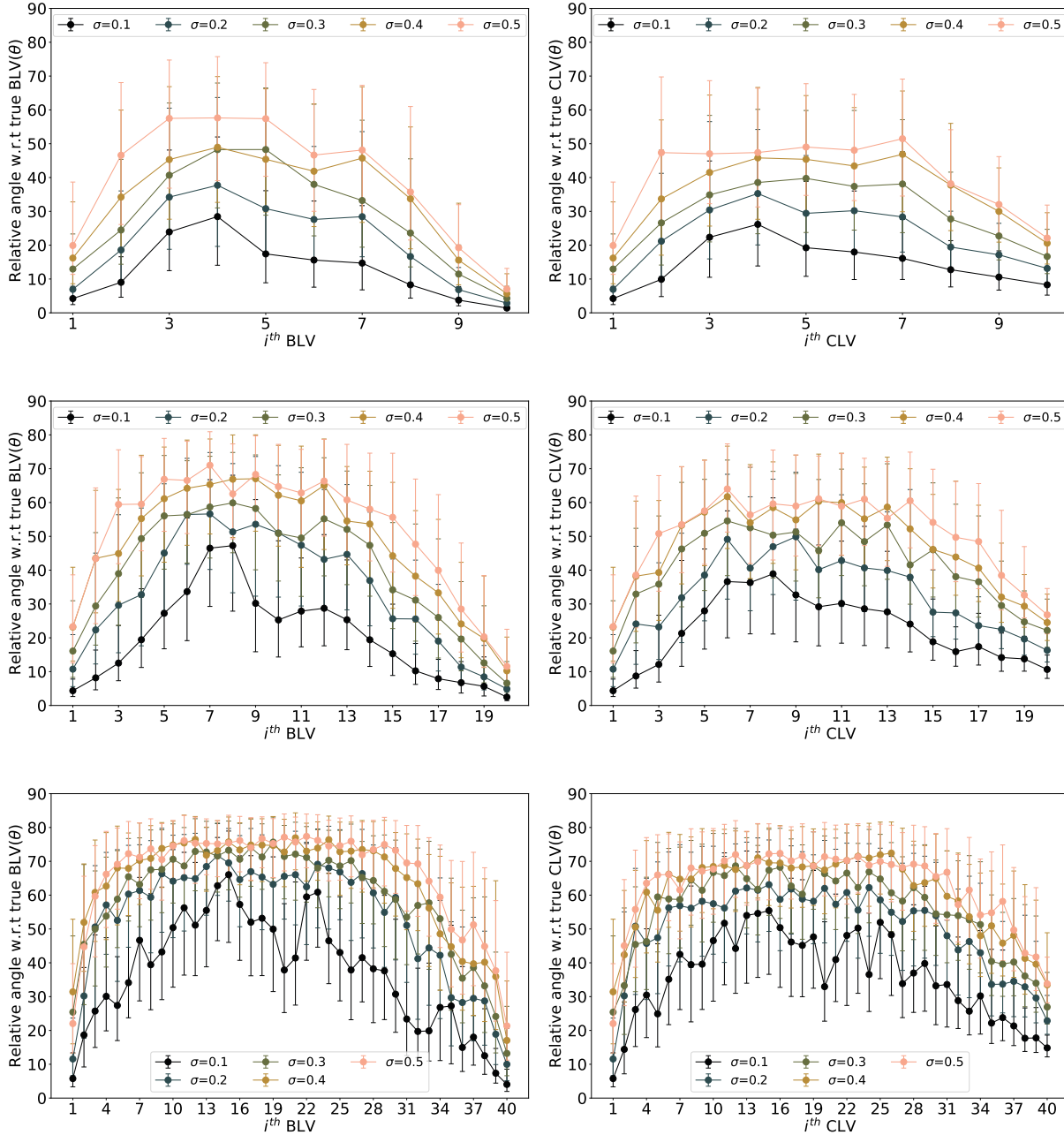


Figure 4.9: Angle  $\theta$  between the individual true and recovered BLVs (left) and CLVs (right) for different perturbation strength  $\sigma$  for Lorenz-96 in 10, 20 and 40 dimensions, in the top, middle, and bottom rows, respectively.

We now describe the results for the sensitivity of the LVs for Lorenz-96. In figure 4.9, we plot the relative acute angle between the LVs from the true and perturbed trajectory for both BLVs and CLVs for Lorenz-96 for dimensions  $n = 10, 20$ , and 40. Similar to results in the previous section 4.3.1 about LVs from assimilated trajectories, we observe that even with a small amount of noise strength, the individual vectors quickly misalign from the true vectors, which is quite different compared to Lorenz-63. Only the first few most unstable vectors and the few most stable vectors corresponding to the two opposite ends of the Lyapunov spectrum have significant projection along the true vectors. For the intermediate BLVs the angles approximately lie in the interval  $[45^\circ, 90^\circ]$ . The CLVs follow a similar picture as the BLVs for the few most unstable directions and for the stable directions. The CLVs for  $n = 10$  and 20 seem to have smaller errors than their BLV counterparts. The error in angle also increases systematically with increasing perturbation strength  $\sigma$ . The effect of dimension for this extensively chaotic system is clearly evident by the degrading sensitivity of both the BLVs and CLVs we double the dimension from  $n = 10$  to 20 and then to 40.

When applying Lorenz-96 in various dimensions, we observe increased sensitivity as the dimensionality expands. Both backward and covariant Lyapunov vectors exhibit high sensitivity, resulting in minimal errors for only a few unstable and stable vectors at the extremes.

In figure 4.10, we also plot the exponents for the true and the perturbed trajectories for different values of observational noise strength  $\sigma$ . The inset shows the absolute errors in the exponents obtained from the perturbed trajectories from the true exponents. We notice that these absolute errors are small of the order of 0.1 for  $n = 20$  and 40 and 0.2 for  $n = 10$  which suggests that the exponents themselves are not as sensitive as the BLVs and the CLVs. The relative absolute errors in the exponents do not seem to follow any trend for different values of  $\sigma$ . This shows that the Lyapunov spectrum is quite robust to the perturbations in the underlying trajectory.

#### 4.3.4 Oseledets' subspaces spanned by the LVs for different perturbation strengths

We now move towards understanding Oseledets' subspaces recovered from the perturbed trajectories instead of the individual vectors by computing the principal angles between the respective subspaces obtained from the true and the perturbed trajectories. Figure 4.11 shows the PA for a few  $k$ -dimensional subspaces for  $k = 2, 5, 10, 15, 20$  over the sampling interval  $[I, F]$  for  $n = 40$ . The behaviour is very similar to the case of using assimilated trajectory that was discussed in section 4.3.1 in figure 4.5. In particular, there are a majority of small principal angles within  $30^\circ$  for  $k = 5$  onwards, suggesting that the

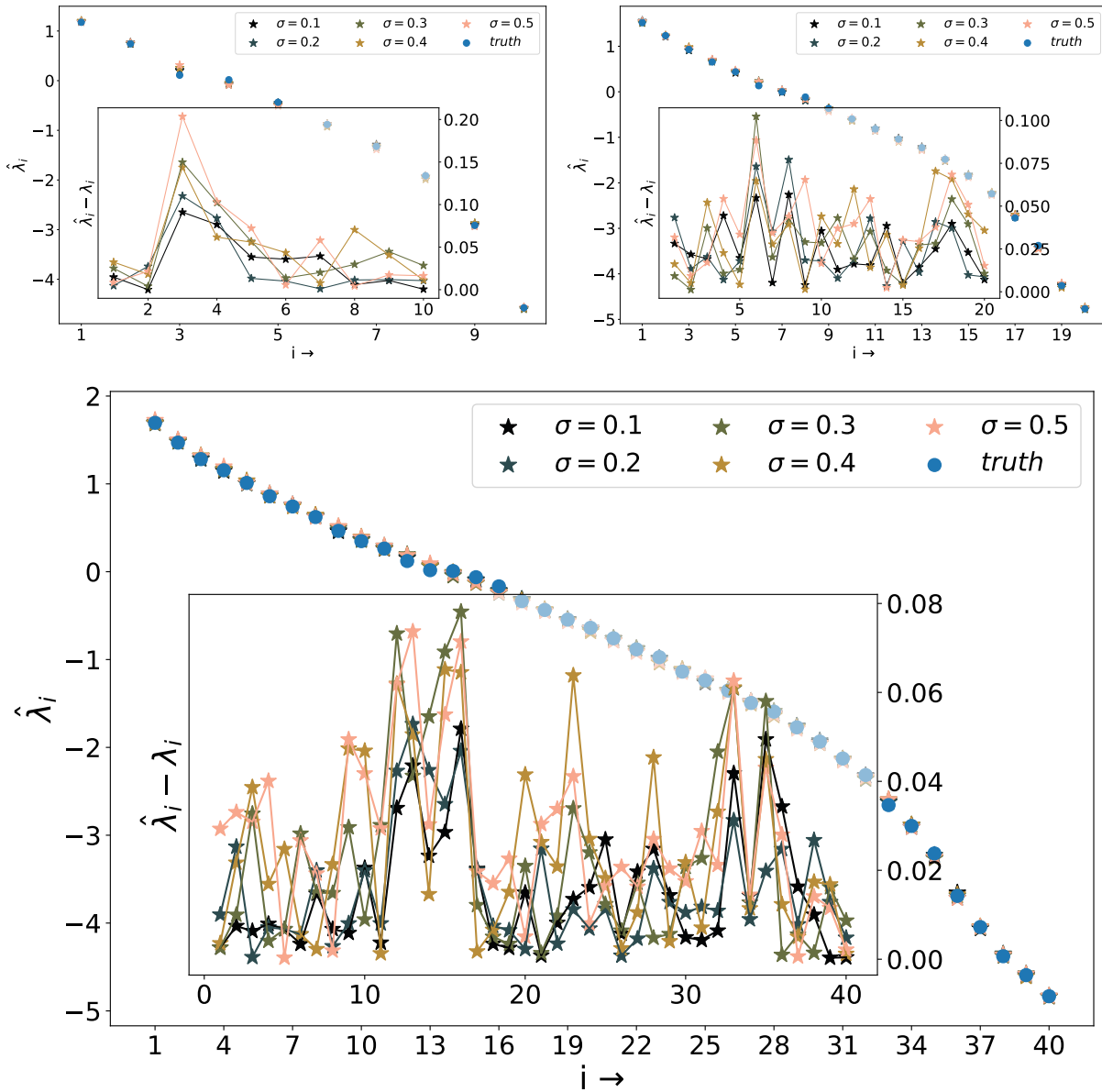


Figure 4.10: The Lyapunov exponents computed from perturbed trajectories for Lorenz-96 in 10(top left), 20 (top right) and 40 (bottom) dimensions. The inset shows relative absolute errors from the exponents of the unperturbed trajectory.

subspaces computed from the perturbed trajectory have significant overlap with the subspaces spanned by the true BLVs. As earlier, the angles increase with increasing perturbation strength  $\sigma$ .

The merits of studying the principal angles are revealed in the fact that when the angle between the individual BLVs from the perturbed trajectory is quite different, the unstable subspaces computed from them do have some similarity with the subspace spanned by the unstable vectors. When using subspaces instead of actual vectors, the BLVs from perturbed or approximate trajectories might still present some merit in capturing the unstable or stable Oseledets' subspaces. To emphasize this point further, we also plot the



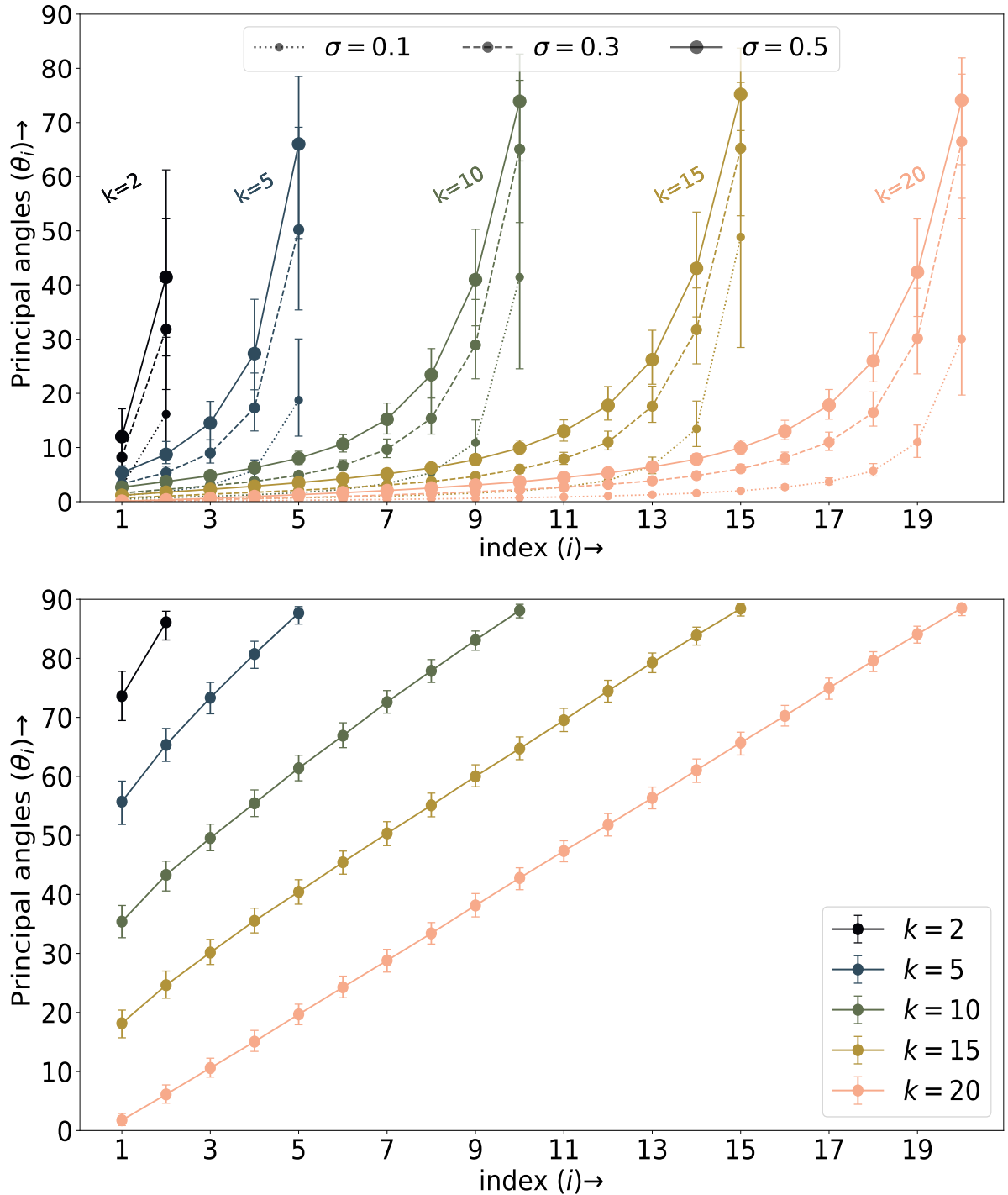


Figure 4.11: The top and the bottom panel shows principal angles (PA) between the  $k$ -dimensional Oseledet's subspace and the corresponding subspace recovered from perturbed trajectory for  $k = 2, 5, 10, 15$  and  $20$ . The line styles denote different  $\sigma$  values used, the dots represent the median of the  $i^{\text{th}}$  PA computed over the number of points in the sampling interval with the bars are 25<sup>th</sup> and 75<sup>th</sup> percentile. For comparison, the bottom panel shows the median of PA between a  $k$ -dimensional Oseledet's subspace and a random subspace of the same dimension using 100 realizations of the random subspaces. The bars represent the median with the 25<sup>th</sup> and 75<sup>th</sup> percentile.

principal angle between a  $k$ -dimensional random subspace and the Oseledet's subspace from the true trajectory. We see the number of obtained P.A. which are smaller than  $30^\circ$  are very small compared to the ones obtained from the assimilated and perturbed trajectories and these angles increase linearly with index which is quite distinct from the case of the approximate Oseledets' subspaces.

## 4.4 Summary

This paper focuses on two key questions pertaining to the computation of Lyapunov vectors when complete knowledge of the true underlying trajectory for a dynamical system is not available. Lyapunov vectors are state-dependent and their numerical computations using the commonly used algorithms by [88, 39] require the underlying trajectory or the initial condition along with the model for the system from which the trajectory can be generated. These algorithms rely crucially on a long trajectory for the convergence to the Lyapunov vectors. This poses a major challenge for chaotic systems since even a small error in initial conditions leads to a totally different trajectory.

The first question that we address is how to use partial and noisy observations in order to compute an approximation of the Lyapunov vectors. We propose a methodology for this purpose, combining the algorithm of [39] with the EnKF. Specifically, using EnKF for filtering, we use Ginelli's algorithm [39] with the filter mean as an approximate pseudo-trajectory and the model linearization between the assimilation times. We demonstrate the efficacy of the proposed idea in the context of low dimensional systems by applying it to Lorenz-63 model using only  $y$ -coordinate observations. On the other hand, for high-dimensional chaotic ODE like Lorenz-96 in various dimensions above 10, the results show that apart from a few most stable and most unstable directions, the Lyapunov vectors are very sensitive and cannot be approximated well using the filter mean as an approximate trajectory.

In order to understand the errors and biases in the recovered Lyapunov vectors, we naturally investigate the second question: how does the perturbation strength affect the LVs obtained from approximate trajectories. In particular, we explore the sensitivity of the numerical computation of both the BLVs and the CLVs to general perturbations in the underlying trajectory. In small dimensions, the results for Lorenz-63 show that the recovered vectors are quite close to the true ones, even for significant perturbation strength. This naturally explains the efficacy of recovering the LVs from filter estimates. On the other hand, the results for Lorenz-96 suggest that most of the vectors, except the most stable and most unstable, are highly sensitive to the perturbations for higher-dimensional dynamical systems. This is consistent with a very similar conclusion for LVs obtained from the filter estimates. In addition, using Lorenz-96 in different dimensions, we find that

this sensitivity grows with the number of dimensions.

Using Lorenz-96 for different dimensions, we also find that the sensitivity is dependent on the dimension of the system. We find that for both the backward and covariant Lyapunov vectors were found to be very sensitive and consequently, only the first few most unstable vectors and the last few most stable vectors had small errors.

The Lyapunov vectors span nested subspaces called Oseledets' spaces. Thus, even in cases where the individual Lyapunov vectors recovered from a perturbed trajectory are not a good approximation of the true LVs, we investigate whether the Oseledets' subspaces are approximated well. In order to quantify this approximation, we studied the principal angles between the recovered and the exact Oseledets' subspaces. Our results suggest that these subspaces are less sensitive compared to the individual vectors themselves with respect to the perturbations in the trajectory. We support this claim by showing that the principal angles between random subspaces are significantly larger than those between the recovered and exact Oseledets' spaces.

The sensitivity of the BLVs and CLVs was shown for perturbations generated from a simple Gaussian distribution. The effect of perturbation statistics itself could lead to different statistics. Although exact error statistics cannot be obtained in the actual application as one does not know the true underlying orbit, we can use another important object called ensemble variance which is of the same order as the error, for a reliable data assimilation system [2].

An important direction for future research would be to investigate the sensitivity of the Lyapunov vectors in different contracting and expanding regions of the phase space. This may better capture the local sensitivity accounting for the variations in the local stable and unstable subspaces over different points on the attractor of the system. Extending the analysis discussed in this work to the case where model errors are present is another interesting direction for future work. Applying a similar analysis to PDEs and high-dimensional models with multiscale dynamics and spatial structures would be highly relevant to practical problems such as those in earth sciences.



# Bibliography

- [1] H. Abarbanel. *Analysis of observed chaotic data*. Springer, 1996.
- [2] J. L. Anderson. A method for producing and evaluating probabilistic forecasts from ensemble model integrations. *Journal of Climate*, 9(7):1518 – 1530, 1996.
- [3] E. Andersson and J. Thépaut. Ecmwf’s 4d-var data assimilation system—the genesis and ten years in operations. *ECMWF Newsletter*, 115:8–12, 2008.
- [4] A. Apte, M. Hairer, A. Stuart, and J. Voss. Sampling the posterior: an approach to non-Gaussian data assimilation. *Physica D*, 230:50–64, 2007.
- [5] V. Araújo, A. I. Bufetov, and S. Filip. On hölder-continuity of oseledets subspaces. *Journal of the London Mathematical Society*, 93(1):194–218, 2016.
- [6] M. Arjovsky, S. Chintala, and L. Bottou. Wasserstein generative adversarial networks. In *International conference on machine learning*, pages 214–223. PMLR, 2017.
- [7] M. Asch, M. Bocquet, and M. Nodet. *Data Assimilation: Methods, Algorithms, and Applications*. SIAM, 2016.
- [8] R. Bellman. Dynamic programming. *Science*, 153(3731):34–37, 1966.
- [9] G. Benettin, L. Galgani, and A. Giorgilli. Lyapunov characteristic exponents for smooth dynamical systems and for hamiltonian systems; a method for computing all of them. part 1: Theory. *Meccanica 15, 9–20 (1980)*, pages 9–20, Jul 1980.
- [10] G. Benettin, L. Galgani, and A. Giorgilli. Lyapunov characteristic exponents for smooth dynamical systems and for hamiltonian systems; a method for computing all of them. part 2: Numerical application. *Meccanica 15.1*, pages 21–30, Jul 1980.
- [11] L. Bengtsson, M. Ghil, and E. Källén. *Dynamic meteorology: data assimilation methods*, volume 36. Springer, 1981.
- [12] T. Bengtsson, P. Bickel, B. Li, et al. Curse-of-dimensionality revisited: Collapse of the particle filter in very large scale systems. *Probability and statistics: Essays in honor of David A. Freedman*, 2:316–334, 2008.

- [13] M. Bocquet, K. S. Gurumoorthy, A. Apte, A. Carrassi, C. Grudzien, and C. K. R. T. Jones. Degenerate kalman filter error covariances and their convergence onto the unstable subspace. *SIAM/ASA Journal on Uncertainty Quantification*, 5(1):304–333, 2017.
- [14] R. G. Brown and P. Y. C. Hwang. *Introduction to random signals and applied Kalman filtering : with MATLAB exercises and solutions*. 1997.
- [15] E. L. Brugnago, J. A. C. Gallas, and M. W. Beims. Predicting regime changes and durations in lorenz’s atmospheric convection model. *Chaos: An Interdisciplinary Journal of Nonlinear Science*, 30(10):103109, 2020.
- [16] A. Carrassi, M. Bocquet, L. Bertino, and G. Evensen. Data assimilation in the geosciences: An overview of methods, issues, and perspectives. *Wiley Interdisciplinary Reviews: Climate Change*, 9(5):e535, 2018.
- [17] A. Carrassi, M. Bocquet, L. Bertino, and G. Evensen. Data assimilation in the geosciences: An overview of methods, issues, and perspectives. *WIREs Climate Change*, 9(5):e535, 2018.
- [18] A. Carrassi, M. Ghil, A. Trevisan, and F. Uboldi. Data assimilation as a nonlinear dynamical systems problem: Stability and convergence of the prediction-assimilation system. *Chaos: An Interdisciplinary Journal of Nonlinear Science*, 18(2), 2008.
- [19] A. Carrassi, A. Trevisan, L. Descamps, O. Talagrand, and F. Uboldi. Controlling instabilities along a 3dvar analysis cycle by assimilating in the unstable subspace: a comparison with the enkf. *Nonlinear Processes in Geophysics*, 15(4):503–521, 2008.
- [20] A. Carrassi, A. Trevisan, and F. Uboldi. Adaptive observations and assimilation in the unstable subspace by breeding on the data-assimilation system. *Tellus A: Dynamic Meteorology and Oceanography*, 59(1):101–113, 2007.
- [21] F. Cérou. Long time behavior for some dynamical noise free nonlinear filtering problems. *SIAM J. Control. Optim.*, 38:1086–1101, 2000.
- [22] N. Chandramoorthy and Q. Wang. An ergodic-averaging method to differentiate covariant lyapunov vectors: Computing the curvature of one-dimensional unstable manifolds of strange attractors. *Nonlinear Dynamics*, 104:4083–4102, 2021.
- [23] K. K. W. Cheung. A review of ensemble forecasting techniques with a focus on tropical cyclone forecasting. *Meteorological Applications, Royal Meteorological Society*, 8(3):315–332, 2006.

- [24] P. Chigansky. Stability of nonlinear filters: A survey, 2006. Lecture notes, Petropolis, Brazil.
- [25] P. Chigansky, R. Liptser, and R. Van Handel. Intrinsic methods in filter stability. *Handbook of Nonlinear Filtering*, 2009.
- [26] S. E. Cohn. An introduction to estimation theory. *J. Met. Soc. Japan*, 75:257–288, 1997.
- [27] P. Courtier and O. Talagrand. Variational assimilation of meteorological observations with the direct and adjoint shallow-water equations. *Tellus A: Dynamic Meteorology and Oceanography*, 42(5):531–549, 1990.
- [28] D. Crisan and B. Rozovskii. *The Oxford handbook of nonlinear filtering*. Oxford University Press, 2011.
- [29] M. Cuturi. Sinkhorn distances: Lightspeed computation of optimal transport. In *Advances in Neural Information Processing Systems*, 2013.
- [30] P. Del Moral and J. Tugaut. On the stability and the uniform propagation of chaos properties of ensemble kalman–bucy filters. *The Annals of Applied Probability*, 28(2):790–850, 2018.
- [31] A. Doucet and A. M. Johansen. A tutorial on particle filtering and smoothing: Fifteen years later. *Handbook of nonlinear filtering*, 12(656-704):3, 2009.
- [32] D. Dragičević and G. Froyland. Hölder continuity of oseledets splittings for semi-invertible operator cocycles. *Ergodic Theory and Dynamical Systems*, 38(3):961–981, 2018.
- [33] J. P. Eckmann and D. Ruelle. Ergodic theory of chaos and strange attractors. *Rev. Mod. Phys.*, 57:617–656, Jul 1985.
- [34] G. Evensen. The ensemble kalman filter: theoretical formulation and practical implementation. *Ocean Dynamics*, 2003.
- [35] J. Feydy, T. Séjourné, F.-X. Vialard, S.-i. Amari, A. Trounev, and G. Peyré. Interpolating between optimal transport and mmd using sinkhorn divergences. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pages 2681–2690. PMLR, 2019.
- [36] S. Fletcher. *Data Assimilation for the Geosciences: From Theory to Application*. Elsevier, 2017.

- [37] A. Genevay. *Entropy-regularized optimal transport for machine learning*. PhD thesis, Paris Sciences et Lettres, 2019.
- [38] A. Genevay, G. Peyré, and M. Cuturi. Learning generative models with sinkhorn divergences. In *International Conference on Artificial Intelligence and Statistics*, pages 1608–1617. PMLR, 2018.
- [39] F. Ginelli, H. Chaté, R. Livi, and A. Politi. Covariant lyapunov vectors. *Journal of Physics A: Mathematical and Theoretical*, 46(25):254005, jun 2013.
- [40] K. S. Gurumoorthy, C. Grudzien, A. Apte, A. Carrassi, and C. K. R. T. Jones. Rank deficiency of kalman error covariance matrices in linear time-varying system with deterministic evolution. *SIAM Journal on Control and Optimization*, 55(2):741–759, 2017.
- [41] T. Janjić, N. Bormann, M. Bocquet, J. Carton, S. E. Cohn, S. L. Dance, S. N. Losa, N. K. Nichols, R. Potthast, J. A. Waller, et al. On the representation error in data assimilation. *Quarterly Journal of the Royal Meteorological Society*, 144(713):1257–1278, 2018.
- [42] A. H. Jazwinski. *Stochastic processes and filtering theory*. Courier Corporation, 2007.
- [43] S. Jiang. Angles between euclidean subspaces. *Geom Dedicata*, 63:113–121, 1996.
- [44] R. E. Kalman. A New Approach to Linear Filtering and Prediction Problems. *Journal of Basic Engineering*, 82(1):35–45, 03 1960.
- [45] E. Kalnay. *Atmospheric modeling, data assimilation, and predictability*. Cambridge University Press, 2003.
- [46] L. V. Kantorovich. On one effective method of solving certain classes of extremal problems. In *Dokl. Akad. Nauk. USSR*, volume 28, pages 212–215, 1940.
- [47] L. V. Kantorovich. Mathematical methods of organizing and planning production. *Management science*, 6(4):366–422, 1960.
- [48] L. Kantorovitch. On the translocation of masses. *Management science*, 5(1):1–4, 1958.
- [49] A. Karimi and M. R. Paul. Extensive chaos in the lorenz-96 model. *Chaos: An Interdisciplinary Journal of Nonlinear Science*, 20(4):043105, 2010.
- [50] A. Katok and B. Hasselblatt. *Introduction to the modern theory of dynamical systems*. Number 54. Cambridge university press, 1995.



- [51] S. Kolouri, S. R. Park, M. Thorpe, D. Slepcev, and G. K. Rohde. Optimal mass transport: Signal processing and machine-learning applications. *IEEE signal processing magazine*, 34(4):43–59, 2017.
- [52] P. V. Kuptsov and U. Parlitz. Theory and computation of covariant lyapunov vectors. *Journal of nonlinear science*, 22:727–762, 2012.
- [53] B. K. W. Lahoz and R. Menard. *Data assimilation*. Springer, 2010.
- [54] K. Law, A. Stuart, and K. Zygalakis. *Data Assimilation*. Springer, 2015.
- [55] K. J. Law, D. Sanz-Alonso, A. Shukla, and A. Stuart. Filter accuracy for the lorenz 96 model: Fixed versus adaptive observation operators. *Physica D: Nonlinear Phenomena*, 325:1–13, 2016.
- [56] B. Legras and R. Vautard. A guide to liapunov vectors. In *Proceedings 1995 ECMWF seminar on predictability*, volume 1, pages 143–156, 1996.
- [57] N. Lei, K. Su, L. Cui, S.-T. Yau, and X. D. Gu. A geometric view of optimal transportation and generative model. *Computer Aided Geometric Design*, 68:1–21, 2019.
- [58] E. N. Lorenz. Deterministic nonperiodic flow. *Journal of Atmospheric Sciences*, 20(2):130 – 141, 1963.
- [59] E. N. Lorenz. Deterministic nonperiodic flow. *Journal of atmospheric sciences*, 20(2):130–141, 1963.
- [60] E. N. Lorenz. Predictability: a problem partly solved. In *Seminar on Predictability I*, pages 1–18. ECMWF, Reading UK, 1995.
- [61] E. N. Lorenz. Predictability: a problem partly solved. In *Seminar on Predictability I*, pages 1–18. ECMWF, Reading UK, 1995.
- [62] E. N. Lorenz. Designing chaotic models. *Journal of the atmospheric sciences*, 62(5):1574–1587, 2005.
- [63] C. Luo and Y. Zhao. Hölder continuity of oseledets subspaces for linear cocycles on banach spaces. *Physica Scripta*, 98(1):015203, 2022.
- [64] P. Mandal, S. K. Roy, and A. Apte. Stability of nonlinear filters-numerical explorations of particle and ensemble kalman filters. In *2021 Seventh Indian Control Conference (ICC)*, pages 307–312. IEEE, 2021.

- [65] P. Mandal, S. K. Roy, and A. Apte. Probing robustness of nonlinear filter stability numerically using sinkhorn divergence. *Physica D: Nonlinear Phenomena*, 451:133765, 2023.
- [66] C. Martin, N. Sharafi, and S. Hallerberg. Estimating covariant lyapunov vectors from data. *Chaos: An Interdisciplinary Journal of Nonlinear Science*, 32(3):033105, 2022.
- [67] F. Noethen. *Computing covariant Lyapunov vectors: A convergence analysis of Ginelli’s algorithm*. PhD thesis, Department of Mathematics of Universität Hamburg, 2019.
- [68] L. Oljača, T. Kuna, and J. Bröcker. Exponential stability and asymptotic properties of the optimal filter for signals with deterministic hyperbolic dynamics. *arXiv preprint arXiv:2103.01190*, 2021.
- [69] V. I. Oseledec. A multiplicative ergodic theorem, lyapunov characteristic numbers for dynamical systems. *Transactions of the Moscow Mathematical Society*, 19:197–231, 1968.
- [70] V. Oseledets. Oseledets theorem. *Scholarpedia*, 3(1):1846, 2008.
- [71] L. Palatella, A. Carrassi, and A. Trevisan. Lyapunov vectors and assimilation in the unstable subspace: theory and applications. *Journal of Physics A: Mathematical and Theoretical*, 46(25):254020, jun 2013.
- [72] T. N. Palmer and L. Zanna. Singular vectors, predictability and ensemble forecasting for weather and climate. *Journal of Physics A: Mathematical and Theoretical*, 46(25):254018, jun 2013.
- [73] G. Peyré, M. Cuturi, et al. Computational optimal transport: With applications to data science. *Foundations and Trends® in Machine Learning*, 11(5-6):355–607, 2019.
- [74] A. S. Reddy and A. Apte. Stability of non-linear filter for deterministic dynamics. *Foundations of Data Science*, 3(3):647–675, 2021.
- [75] A. S. Reddy, A. Apte, and S. Vadlamani. Asymptotic properties of linear filter for noise free dynamical system. *Systems & Control Letters*, 139:104676, 2020.
- [76] T. Salimans, H. Zhang, A. Radford, and D. Metaxas. Improving gans using optimal transport. *arXiv preprint arXiv:1803.05573*, 2018.
- [77] S. Särkkä. *Bayesian filtering and smoothing*, volume 3. Cambridge University Press, 2013.

- [78] Y. Sasaki. Some basic formalisms in numerical variational analysis. *Monthly Weather Review*, 98(12):875 – 883, 1970.
- [79] Q. Tong and K. Kobayashi. Entropy-regularized optimal transport on multivariate normal and q-normal distributions. *Entropy*, 23(3), 2021.
- [80] A. Trevisan, M. D’Isidoro, and O. Talagrand. Four-dimensional variational assimilation in the unstable subspace and the optimal subspace dimension. *Quarterly Journal of the Royal Meteorological Society*, 136(647):487–496, 2010.
- [81] A. Trevisan and F. Pancotti. Periodic orbits, lyapunov vectors, and singular vectors in the lorenz system. *Journal of the Atmospheric Sciences*, 55(3):390 – 398, 1998.
- [82] A. Trevisan and F. Uboldi. Assimilation of standard and targeted observations within the unstable subspace of the observation–analysis–forecast cycle system. *Journal of the Atmospheric Sciences*, 61(1):103 – 113, 2004.
- [83] R. van Handel. Hidden markov models. *Unpublished lecture notes*, 2008.
- [84] D. van Kekem. *Dynamics of the Lorenz-96 model: Bifurcations, symmetries and waves*. PhD thesis, University of Groningen, 2018.
- [85] S. Vannitsem and V. Lucarini. Statistical and dynamical properties of covariant lyapunov vectors in a coupled atmosphere-ocean model-multiscale effects, geometric degeneracy, and error dynamics. *Journal of Physics A: Mathematical and Theoretical*, 49(22):224001, may 2016.
- [86] C. Villani et al. *Optimal transport: old and new*, volume 338. Springer, 2009.
- [87] N. Wiener. *Extrapolation, interpolation, and smoothing of stationary time series: with engineering applications*. The MIT press, 1949.
- [88] C. L. Wolfe and R. M. Samelson. An efficient method for recovering lyapunov vectors from singular vectors. *Tellus A: Dynamic Meteorology and Oceanography*, 59(3):355–366, 2007.

