



Uncertainties estimation in Full Waveform Inversion using Ensemble methods

Julien Thurin

► To cite this version:

Julien Thurin. Uncertainties estimation in Full Waveform Inversion using Ensemble methods. Earth Sciences. Université Grenoble Alpes [2020-..], 2020. English. NNT : 2020GRALU003 . tel-02570602

HAL Id: tel-02570602

<https://theses.hal.science/tel-02570602>

Submitted on 12 May 2020

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

THÈSE

Pour obtenir le grade de

DOCTEUR DE L'UNIVERSITÉ GRENOBLE ALPES

Spécialité : Terre Solide (CETSOL)

Arrêté ministériel : 25 mai 2016

Présentée par

Julien THURIN

Thèse dirigée par **Romain BROSSIER**, enseignant chercheur,
Université Grenoble Alpes
et codirigée par **Ludovic METIVIER**, chercheur, CNRS

préparée au sein du **Laboratoire Institut des Sciences de la
Terre**
dans l'**École Doctorale Terre, Univers, Environnement**

Evaluation des incertitudes en inversion des formes d'ondes par méthodes d'ensemble

Uncertainties estimation in Full Waveform Inversion using Ensemble methods

Thèse soutenue publiquement le **27 janvier 2020**,
devant le jury composé de :

Monsieur ROMAIN BROSSIER

MAITRE DE CONFERENCES HDR, UNIVERSITE GRENOBLE ALPES,
Directeur de thèse

Monsieur LUDOVIC METIVIER

CHARGE DE RECHERCHE HDR, CNRS DELEGATION ALPES, Co-
directeur de thèse

Monsieur ANDREW CURTIS

PROFESSEUR, UNIVERSITE D'EDIMBOURG - ROYAUME-UNI,
Rapporteur

Monsieur ANDREAS FICHTNER

PROFESSEUR ASSOCIE, ETH ZURICH - SUISSE, Rapporteur

Monsieur YANN CAPDEVILLE

DIRECTEUR DE RECHERCHE, CNRS DELEGATION BRETAGNE PAYS
DE LOIRE, Président

Monsieur ALEXANDRE FOURNIER

PROFESSEUR DES UNIVERSITES, UNIVERSITE PARIS SORBONNE,
Examineur



Acknowledgments

Je tiens à remercier les nombreuses personnes qui m'ont aidé à naviguer sur les eaux tumultueuses du doctorat au cours de ces dernières années, faisant de mon voyage une croisière (presque) paisible. Sans vous, ce travail n'aurait pas été le même. Je n'aurais pas grandi de la sorte.

Je voudrais tout d'abord adresser mes plus sincères remerciements à mes deux directeurs de thèse, Romain Brossier et Ludovic Métivier. Vous m'avez tous-deux accueillis au sein de la famille SEISCOPE et accordé votre confiance dans la poursuite de mes recherches, surtout lorsqu'il a fallu explorer la littérature d'assimilation de données, qui nous était à tous trois étrangères. Merci pour votre patience, votre aide et votre soutien, dans cette aventure que je considère comme une franche réussite. Si c'était à refaire, j'embarquerais sans hésiter.

Je voudrais également adresser des remerciements spéciaux à Jean Virieux. Tu as toujours été de bon conseil, et j'ai adoré les discussions avec toi au détour d'un déjeuner. Merci pour ton énergie et ta bonne humeur indefectible.

Merci aux membres de mon jury, Andreas Fichtner, Andrew Curtis, Yann Capdeville, Alexandre Fournier, Nicolas Gillet et Jérémie Messud. J'ai eu le plaisir de faire la connaissance de certains d'entre vous au long de mon doctorat. Merci Andreas de m'avoir fait découvrir le merveilleux monde de l'American Geophysical Union en m'invitant en 2017 à présenter mon travail, à Nicolas de m'avoir donné de ton temps pour participer à mon comité de suivie thèse, et Alexandre pour ta sympathie lors de la SIAM GS19 et l'intérêt que tu as porté pour nos travaux.

Merci de même, à mes co-bureaux d'ISTerre, Phuong-Thu Trinh et Marwan Irnaka d'avoir égaillé mes journées, de m'avoir motivé dans les périodes difficiles, de m'avoir écouté, et de vous être confiés à moi. Votre confiance et votre amitié me sont précieuses. Merci aussi à Sylvain Fiolleau, qui même si arrivé plus tard, a su parfaitement compléter le quatuor du bureau 134.

Merci aux collègues de SEISCOPE. Serge Sambolian qui est bien sûr devenu le Forged in Fire Champion à Washington, tu es un gagnant Serge, et je te souhaite le meilleur pour la suite. Merci aussi à Philippe Le Bouteiller avec qui j'ai partagé mes aventures rocambolesques à la Nouvelle-Orléans, Hugo Sanchez le seul de l'équipe à comprendre les "Bayesian bullshit", Marco Salcedo, Wei Zhou, Yang Li, Yubing Li, Pengliang Yang, Paul Wellington, Hossein Aghamiry et le grand prof. Tavakoli. Vous avez tous été de fantastiques collègues, et j'ai adoré le temps passé avec vous en congrès et au laboratoire.

Merci à la bande du Gâteau du Vendredi, qui a prouvé qu'avec un peu de dévouement, on peut stratégiquement étirer les pauses et partir plus tôt en weekend, avec le ventre plein. Merci à Gaëlle Le Roy, Cyrielle Dollet, Noellie Bontemps, Antoine Guillemot, Dorian Soergel, Judith Marinier, Isabelle Dumont et Marguerite Mathey.

Merci aux amis-collocs, Aaron Carril et Gautier Chabert (qui je crois ne me doit plus d'argent). Ça a été un grand plaisir, d'habiter avec vous. Je regrette presque d'être parti avant que vous vous mettiez à

faire la vaisselle.

On garde les braves pour la fin. Arnaud Pladys, on s'est rencontrés en première année à la fac de Nice, et on s'est plus vraiment quitté. Ton amour des bonnes choses ne cessera jamais de m'étonner, et je suis curieux de savoir qu'elle sera ta prochaine folie. Je n'oublierai jamais les heures passées dans le BDE/chiottes aménagés à Nice avec la team "géophysique". Notre incroyable terrain au Bénin, qui a été l'un des points forts de ces 8 dernières années. Les sessions de révision à Phitem, et enfin les stages et la thèse et les congrès ensemble. J'aime à croire que cela ne se terminera pas avec la thèse, et je te souhaite, mon cher ami, seulement le meilleur. Grégoire Guillet, tu mérites aussi des remerciements tout particuliers. Tout comme toi, je n'étais pas rassuré lors de la sortie terrain où nous nous sommes rencontrés. On a traversé beaucoup de moments ensemble, dans la joie mais aussi les coups durs. Tu as été depuis le master et jusqu'à la fin de ma thèse, mon partner in crime, comme un frère. J'ai eu le privilège de partager un foyer avec toi, de travailler avec toi, et de profiter de tous les degrés de ton humour. Saches que tu es un ami, un scientifique, et une personne exceptionnelle, et que tu seras toujours dans mon coeur.

Merci à ma famille. Mes parents qui m'ont hissé sur les escaliers de la fac, et m'ont conduit jusqu'à ce jour du 27 janvier 2020 où je devins le premier docteur de la famille. J'ai eu une grande fierté à vous rendre fier, vous aussi. Ça y est je suis docteur, et pourtant on sait que ce n'était pas gagné d'avance! C'est à vous que je dois ma curiosité, ma passion, et l'amour des choses bien faites. Merci infiniment. Merci à ma très chère petite soeur, qui a traversé durant ma thèse, bien des épreuves, et qui je suis certain réussira sa propre aventure.

Merci à ma chère Louise. Toi qui as toujours été là et qui m'as montré ton infaillible et incommensurable soutien depuis notre rencontre. Tu as contribué à bien des égards à ma réussite, tant sur le plan de la recherche, que sur le plan humain. Je te suis et serai toujours, infiniment reconnaissant.

Abstract

Full Waveform Inversion (FWI) is an ill-posed non-linear inverse problem, aiming at recovering detailed pictures of subsurface physical properties, which are crucial to explore and understand Earth structures. Classically formulated as a least-squares optimization scheme, FWI yields a single subsurface model amongst an infinite possibility of solutions. With the general lack of systematic and scalable uncertainty estimation, this formulation makes interpretation of FWI's outcomes complex.

In this thesis, we propose an unconventional, scalable way of tackling the lack of uncertainty estimation in FWI, thanks to data assimilation ensemble methods. We develop a scheme combining both classical FWI and the Ensemble Transform Kalman Filter, that we call ETKF-FWI, and which is successfully applied on two 2-D test cases. This scheme takes advantage of the theoretical common-ground between least-squares optimization problems and Bayesian filtering. We use it to recast FWI in a local Bayesian inference framework, thanks to the ensemble representation. The ETKF-FWI provides high-resolution subsurface tomographic models and yields a low-rank approximation of the posterior covariance, holding the uncertainty and resolution information of the proposed solution. We show how the ETKF-FWI can be applied to qualitatively evaluate uncertainty and resolution of the solution. Instead of providing a single solution, the filter yields an ensemble of models, from which statistical information can be inferred.

Uncertainty is evaluated from the ensemble's variance, which relates to the diversity of solution amongst the ensemble members for each parameter. We show that lines of the correlation matrix are ideal to evaluate qualitatively parameters resolution, thanks to their adimensionality. While the methodology is computationally intensive, it has the benefit of being fully scalable. Its applicability is demonstrated on a synthetic benchmark. This preliminary test allows us to assess the sensitivity of the ensemble representation to the common undersampling bias encountered in ensemble data assimilation. While undersampling does not affect the image reconstruction in any way, it results in variance underestimation, which makes the whole exercise of quantitative uncertainty assessment complicated. Ensemble inflation has been used to mitigate this bias, but does not seem to be a practical solution.

A field data experiment is also discussed in this thesis. It makes it possible to test the sensitivity of the ETKF-FWI to complex noise structure and realistic physics. As it stands, the complexity of the problem reduces flexibility in the ensemble generation, and hence on the uncertainty estimate. Despite these limitations, results are consistent with the synthetic benchmark, and we are able to provide a qualitative uncertainty assessment. The field data case also allows us to evaluate the possibilities to use the ETKF-FWI on multiparameter inversion, which is still regarded as a challenging topic in FWI. The ETKF-FWI multiparameter inversion yields improved models compared with conventional ones. More importantly, it makes it possible to assess the uncertainty associated with parameters cross-talks.

Résumé

L'inversion de forme d'onde complète (FWI) est une méthode d'inversion non-linéaire qui a pour but l'obtention de modèles précis des propriétés physiques du sous-sol terrestre. Ces modèles, véritables cartes de propriétés physiques, sont indispensables pour l'exploration et l'étude des structures internes de la Terre. Généralement formulée sous la forme d'un schéma d'optimisation par la méthode des moindres carrés, la FWI compare des enregistrements sismiques observés en surface, avec des données synthétiques calculées à partir d'un modèle numérique de sous-sol. Alors qu'une infinité de modèles peut potentiellement expliquer les observations, la FWI, du fait de sa formulation, ne permet d'obtenir qu'un seul modèle du sous-sol fortement conditionné par le choix de modèle de départ. À cette ambiguïté s'ajoute la difficulté d'estimer l'incertitude de la solution, à cause du coût de calcul prohibitif de la FWI. La non-unicité de la solution et le manque de moyens d'estimation d'incertitude rend l'exploitation des modèles de FWI compliquée.

Dans cette thèse, nous proposons une méthode non conventionnelle et abordable, intégrant l'estimation d'incertitude au coeur de la solution de FWI. Notre méthode combine la FWI conventionnelle et l'assimilation de données par méthodes d'ensemble. De ce fait, elle tire avantage de la vitesse de convergence de la FWI conventionnelle, ainsi que des capacités d'estimation d'incertitude du Filtre de Kálmán d'Ensemble dit "Transform" (ETKF). Cette combinaison est permise par les fondements théoriques communs aux problèmes d'optimisation en FWI conventionnelle et au filtrage bayésien de l'ETKF. Nous utilisons ce schéma, l'ETKF-FWI, afin de transposer le problème de FWI dans le cadre de l'inférence Bayésienne locale. Au lieu d'une unique solution, l'ETKF-FWI retourne un ensemble de modèles qui permet à la fois de calculer la meilleure solution au sens des moindres carrés, mais aussi l'information d'incertitude et de résolution associée à chaque paramètre. Cette estimation d'incertitude est rendue possible par l'approximation de bas-rang de la matrice de covariance *a posteriori*, calculée à partir de l'ensemble. Les valeurs de variance permettent d'évaluer le degré de variabilité de la solution au sein de l'ensemble. La résolution est quant à elle, donnée par les termes hors diagonaux de la matrice de corrélation, qui est préférée à la matrice de covariance pour sa nature adimensionnelle.

L'application de l'ETKF-FWI à deux cas d'études (un test synthétique et une application sur données de terrain) nous permet d'évaluer la faisabilité, ainsi que les limites de notre technique. Malgré le coût de calcul important lié à la représentation d'ensemble, cette stratégie permet une implémentation complètement parallèle, la rendant avantageuse au regard des solutions existant dans la littérature. Ces tests nous permettent d'évaluer l'influence de la taille de l'ensemble sur l'estimation de la variance, en caractérisant le biais de sous-échantillonnage associé aux petits ensembles. Bien que ce biais soit classiquement corrigé grâce aux méthodes d'*inflation* d'ensemble, celles-ci ne semblent pas adaptées à l'ETKF-FWI, limitant l'estimation d'incertitude à des évaluations qualitatives. De plus, la complexité de l'application sur données de terrain impacte la création de l'ensemble initial, ce qui influence directement les capacités de l'ETKF-FWI à produire une estimation quantitative de l'incertitude.

Nous terminons par l'application de l'ETKF-FWI à une inversion de plusieurs paramètres physique

(vitesse des ondes P et densité), considéré comme un défi majeur en FWI conventionnelle. Ce test nous permet d'évaluer qualitativement les liens de corrélation et d'ambiguïté entre vitesse et densité, ainsi que leurs incertitude et résolution respectives. De plus, le modèle moyen issu de l'ETKF-FWI semble être de qualité supérieure, ce qui laisse supposer d'un possible effet de préconditionnement fourni par la covariance.

Contents

General Introduction	1
1 Uncertainty quantification for Full Waveform Inversion	7
1.1 Overview of FWI	7
1.1.1 Historical overview	7
1.1.2 Mathematical formulation of FWI	9
1.1.3 Interpretation of the gradient and Hessian operators	16
1.1.4 Link between the Hessian and the posterior covariance in Bayesian estimation	19
1.2 Uncertainty quantification in FWI: state-of-the-art	21
1.2.1 Uncertainty quantification through global optimization techniques	22
1.2.2 Local uncertainty estimation in FWI	27
1.2.3 Ideas from the Data Assimilation community	32
2 Data Assimilation	35
2.1 Elements of Data Assimilation	36
2.1.1 Defining the system state	36
2.1.2 Observations	37
2.1.3 A practical example of statistical estimator	38
2.1.4 The dynamical model - forecasting stage	41
2.1.5 The Kalman Filter	42
2.1.6 Extended Kalman Filter	46
2.2 Ensemble Kalman Filter	47
2.2.1 Ensemble representation	48
2.2.2 EnKF's forecast step	49
2.2.3 Analysis step	50
2.2.4 Ensemble Transform Kalman Filter	53
2.2.5 Maximum-Likelihood Ensemble Filter	58
2.2.6 Ensemble Kalman Inverse	60
2.3 Limits of EnKF methods	62
2.3.1 Undersampling characterization	62
2.3.2 Inbreeding	63
2.3.3 Filter divergence	64
2.3.4 Spurious Correlation	64
2.3.5 Solutions to undersampling	65
3 Combining DA and FWI	71

3.1	Proposition 1 - A dynamic formulation	72
3.2	Proposition 2 - Extending the state-space: The WRI analog	73
3.3	Proposition 3 - Extending the state-space: The EKI analog	74
3.4	Proposition 4 - Extending the state-space: adding adjoint	74
3.5	Proposition 5 - A simple adjoint scheme	75
3.5.1	FWI as a dynamic problem: defining a dynamic proxy	76
3.5.2	The ETKF-FWI scheme	77
3.5.3	ETKF-FWI sampling strategy	80
4	Synthetic application of the ETKF-FWI	85
4.1	Solving the FWI problem	85
4.2	ETKF-FWI on the Marmousi II synthetic benchmark	87
4.2.1	Synthetic benchmark setup	87
4.2.2	The ETKF-FWI setting	88
4.2.3	ETKF-FWI application with 600 ensemble members.	89
4.3	Investigating undersampling	97
4.3.1	Parameter estimate	97
4.3.2	Variance approximation	98
4.3.3	Correlation approximation	101
4.4	Mitigating undersampling	104
5	Field data application of the ETKF-FWI	110
5.1	The Valhall oil field dataset and ETKF-FWI parameterization	110
5.1.1	P-wave velocity reconstruction	113
5.1.2	P-wave velocity and density reconstruction	116
6	Conclusion and perspectives	124
	References	129

CONTENTS

General Introduction

Most of our current understanding of Earth's internal structure is coming from observations of various geophysical fields. Slight deviations of gravity values recorded at the surface can inform us about the variations of material density at depth, enabling the detection of large structures within the Earth. Earth's geomagnetic fields have also been permanently crystallized into seafloor spreading oceanic basalts. These records make possible to trace back our planet's geomagnetic history up to a few hundred million years, from which our planet's tectonic history was deduced. Finally, records of seismic wavefields under the form of seismograms have granted us with a considerable amount of knowledge on our planet's dynamic mechanisms and structure.

Even though seismograms might seem complicated to decipher given their cryptic look, they contain a fair share of information. Fortunately, we can rely on the strong ties between the physics governing wave propagation and the propagating medium itself to make sense of these recordings. In the same way the human body can be scanned with ultrasounds, it is possible to use seismic wavefield recorded at the surface, to reveal Earth's deep structure. This tomography method is a classical geophysical problem called *seismic tomography*.

Seismic tomography

Seismic tomography appeared at the beginning of the 20th century, as a powerful method to look through the opaque Earth and uncover its structure. One of the prime examples of discovery in this area came from Mohorovičić (1909), who noticed discrepancies in the recordings of the Zagreb earthquake of October 1909, when looking at seismograms recorded at various distance to the epicenter. To explain the variations in the seismograms content over the epicentral distance, Mohorovičić concluded that a shallow interface must be present in the subsurface, inducing reflections and refractions in the wavefield. The presence of this discontinuity (also known as the Mohorovičić discontinuity or the Moho) was inferred from observations and our knowledge of the laws governing wave propagation, as a solution of an *inverse problem*. This denomination comes from the fact that when solving an inverse problem, we ought to find the cause (Moho discontinuity) of an effect (reflections and refractions recorded in the seismograms).

Seismic tomography has since then grown and branched out into many different applications. The main categories can be defined, depending on whether the source generating the recordings is controlled or not. We often refer to these categories as passive and active tomography. The early work of Mohorovičić falls in the passive seismic category, as the source of the wavefield he studied came from an earthquake. This type of tomography is generally aimed at recovering large scale mantle structures (at a regional or global scale). It has to deal with a sparse and limited amount of natural seismic sources and receivers location: seismic stations are unevenly distributed at the surface, and seismic activity is mostly confined at the interface of tectonic plates (be it subductions or collision zones).

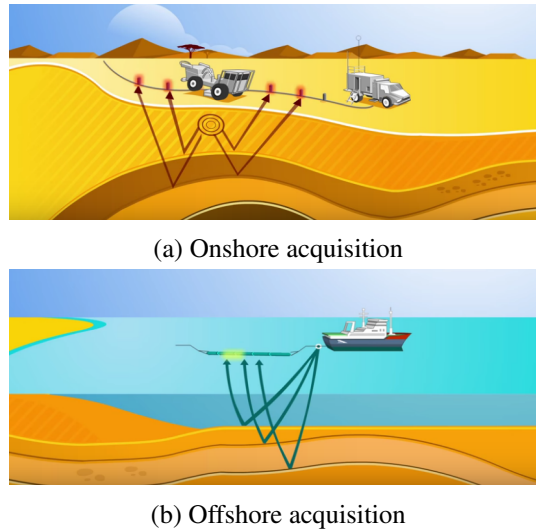


Figure 1: Schematic of an onshore and offshore setup in seismic exploration surveys. (Total).

Nevertheless, passive seismic tomography has been applied with remarkable success to regional and global scales (Aki et al., 1977; Fichtner et al., 2009; Bedle and Lee, 2009; Tape et al., 2010; Panning et al., 2010; French and Romanowicz, 2015; Bozdağ et al., 2016), granting us valuable insights on our planet's deep structures.

Active seismic tomography differs from passive seismic tomography by its voluntary nature, its use of active and controlled sources (explosion, air-gun, hammer source), its acquisition geometries, and the difference of targets. Due to the lack of energy and acquisition scale, active seismic tomography is generally focused on shallow crustal targets, ranging from a few meters to a few kilometers depth. It has been prominently used in geophysical exploration, as a means to get precise images of the subsurface, for civil-engineering applications or georesources exploitation (Plessix, 2009; Sirgue et al., 2010; Plessix et al., 2012; Warner et al., 2013b; Zhu et al., 2015; Operto et al., 2015).

Seismic exploration surveys can be conducted on land (Fig. 1a), where seismic waves are generated by applying a force directly onto the ground. The sources can either be impulsive (blast of an explosion, percussion by a heavy dead-weight) or vibrating, according to specific sweeping frequency patterns (vibrating pot or trucks). Waves generated during the survey, travel through the subsurface and are recorded at the surface by arrays of evenly distributed geophones (recording devices, similar in principle to microphones, specifically designed to record seismic waves).

Seismic surveys can also take place offshore (Fig. 1b), where the acquisition devices (sources and receivers) are dragged beneath the water surface by exploration vessels. The implosive source is generated by an air gun, which air-bubbles burst acts as an acoustic seismic source. The devices used to record the pressure wavefield in the water are called hydrophones. Despite a few differences in acquisition design, the imaging principle stays virtually the same as for onshore acquisitions. Both types of acquisition are relied-upon for oil-and-gas exploitation, and tomographic models are part of the decision-making process leading to drilling into potential reservoirs. Producing reliable tomographic images is thus critical to avoid the costly consequences of dry wells at the exploitation stage, for instance.

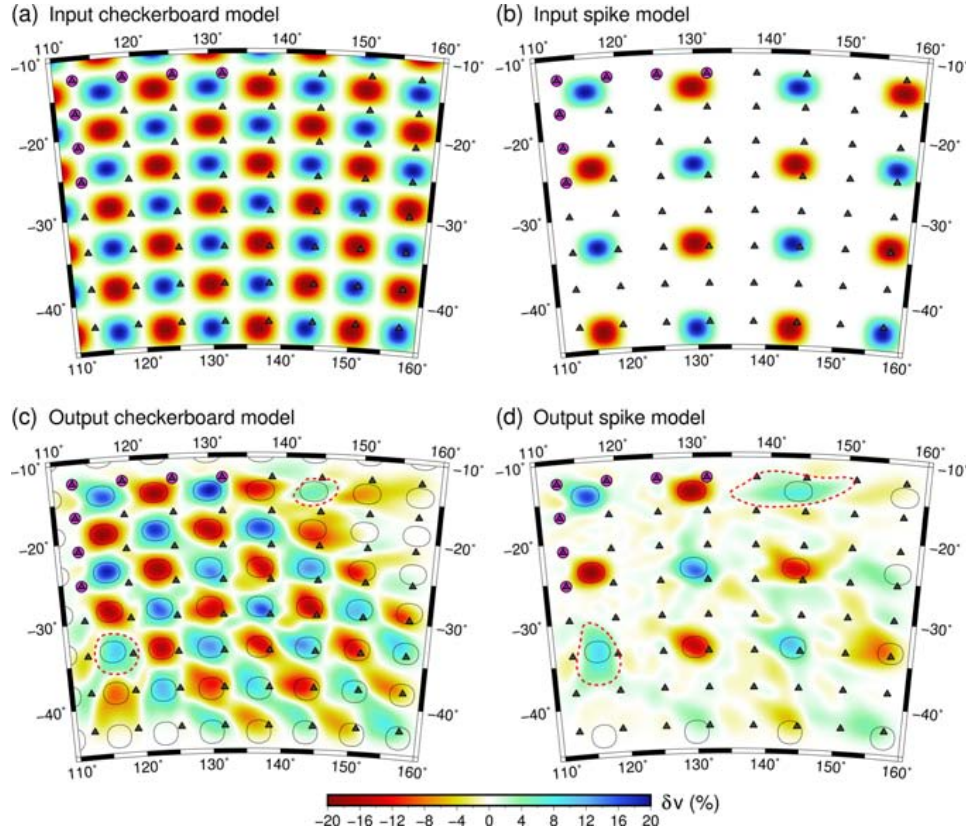


Figure 2: Example of sensitivity tests on checkerboard and spike heterogeneities. (a) Input checkerboard model; (b) input spike model; (c) output checkerboard model; (d) output spike model. Seismic sources are denoted by purple dots and receivers by black triangles. From Rawlinson and Spakman (2016)

The uncertainty problem

Over the past decades, tomography has been very successful at getting more out of seismic recordings, uncovering finer and sharper features in reconstructed tomographic images, across all scales. The techniques have also become more and more computationally expensive, and up-to-date tomographic schemes capable of the highest precision are complex. It is notably the case of *Full Waveform Inversion* (FWI), the tomography technique that is of interest in this thesis, which is often a crucial step in oil-and-gas exploration workflows.

Despite the improvements in images "quality" with modern schemes, one of the main difficulties of seismic tomography remains: assessing the reliability of tomographic models. Besides, as the complexity of the method grows, uncertainty quantification becomes more complicated to achieve, which is typically the case in FWI.

The uncertainty issue arises from the "indirect" nature of seismic tomography: the physical field (our only source of information) is measured at the surface, and it is generally not possible to check the validity of the inferred subsurface model directly. Even worst, several tomographic models could explain the observations with the same level of adequacy, meaning there are plenty of plausible solutions to the problem. Uncertainty estimation in tomographic applications thus appears essential, given its prevalence both in global seismology and exploration geophysics. In seismology, tomographic models

serve as a basis for geologic interpretations of large tectonic features and geodynamic processes, despite knowing the model might not resolve these features entirely. In exploration geophysics, interpretations of tomographic models are used to devise exploitation strategies and evaluate risks, be they financial, societal, or ecological. It is thus necessary to control and assess how trustworthy tomographic solutions are, especially when decision-making is at stake instead of committing to a single solution as it still often done nowadays. But how to check if an inferred model is representative of the truth if there is no access to the subsurface?

While it is possible to get a direct measurement of the subsurface in the case of shallow crustal targets through well-log measurement, these observations are by nature very localized. The cost of drilling also makes them very impractical, and their uses for in-situ verifications are very restricted. Therefore, if exterior means and in-situ measurements can not solve the uncertainty problem, we must seek a solution within the tomographic formalism itself.

One of the most practical methods for quality control of tomography results, fall in the category of *sensitivity analysis*. The principle is rather simple: given the acquisition geometry of a real tomographic problem, synthetic seismograms are computed in a known heterogeneous model. From a blank model, the goal is to recover these perturbations with the synthetic dataset previously computed. The quality of the recovered model is supposed to inform us of the quality we can expect from the real tomographic problem. For convenience, the known heterogeneous model generally has a specific structure, either containing spikes of positive and negative perturbations, or a "checkerboard" pattern of positive and negative anomalies (Fig. 2). However, these tests are prone to suffer from biases introduced by the chosen heterogeneities' scale and structure (Rawlinson and Spakman, 2016) and are not ideal for uncertainty estimation and quality control (despite their popularity).

Sensitivity analysis does not seem to be the definitive answer to the uncertainty problem. Instead, uncertainty estimation should ideally be based on the Bayesian inference framework for inverse problems (Tarantola, 2005). With the Bayesian inference framework, it is possible to explore the plausible solutions and assess their uncertainty, but this comes at a cost: when the scale of the problem grows large, as in FWI, uncertainty estimation based on this mathematical formalism becomes computationally prohibitive, such that no systematic uncertainty estimation method has been developed.

Nonetheless, some state-of-the-art methods for uncertainty estimation have been developed to approximate some elements of the Bayesian formalism. These methods come at the cost of expensive, additional computations to assess the model output quality and are often intrinsically not scalable, creating a computational bottleneck that prevents their systematic applications. These state-of-the-art methods will be reviewed in detail later in Chapter 1.

Instead of following that trend, this work ought to define a new solution to the lack of uncertainty estimation, by recasting our problem in a Bayesian formalism. It has to remain compatible with the requirements of large scales inverse problems by being scalable, but also integrate uncertainty estimation as part of the solution of our tomographic problem.

Seeking solutions elsewhere

As we seek affordable and systematic uncertainty estimation for FWI (and to a broader extent, any seismic tomographic technique), we can look toward another Earth-sciences community that has achieved great success in that avenue. For that, we turn to the Numerical Weather Prediction (NWP) community, which has developed, since the sixties, methods to perform uncertainty estimation for large scale problems,

similar in size with what is typically encountered in FWI. These methods are generally denoted as Data Assimilation (DA) and have since then been applied in various other research and engineering fields.

The reason we are interested in these methods in the first place is that they can provide systematic uncertainty estimation, even on large scale problems. This is achieved thanks to the core design and philosophy of DA, which uses uncertainty when inferring and forecasting physical systems state.

Thus, by integrating mathematical insights from DA into our FWI problem, we might be able to define a framework that can be successful at both inferring subsurface parameters and assessing the uncertainty of recovered solutions. To that extent, we are interested in statistical ensemble DA methods, which we believe, are the more suited for this task and have proven to be very efficient for large scale problems requiring expensive computation. These methods will be carefully reviewed later in this manuscript (chap. 2). Ensemble DA methods naturally allow representing uncertainty through the generation of multiple realizations of solution (the *ensemble* of solutions), which is a desirable feature for FWI.

Outline of the manuscript

In this manuscript, I propose a review of FWI, current uncertainty estimation methods, and ensemble DA methods, from which an original application was derived. This new methodology termed ETKF-FWI (for Ensemble Transform Kalman Filter) constitutes the main scientific contribution of this thesis work. To present this contribution, the manuscript will be organized as follows.

Chapter 1 ought to draw a complete picture of FWI, first from its historical perspectives, then by introducing its mathematical formulation as a numerical optimization problem. Later, the relationship between numerical optimization and uncertainty will be explained, as it is a natural step to understand the challenge at stake: uncertainty estimation in FWI.

The second part of this chapter will be devoted to reviewing the work that has been done by others, in estimating uncertainty in FWI, be it by the mean of global optimization, or by exploiting the mathematical objects underlying local optimization strategies. Finally, insights from the DA community on uncertainty estimation are introduced to pitch the primary goal of this thesis work: establishing a DA framework for uncertainty estimation in FWI.

In **Chapter 2**, DA theory is thoroughly reviewed. From the basic concepts of Bayesian filtering that underlie DA methods, the chapter will progress toward standard DA assimilation tools. Starting by detailing the most straightforward Bayesian estimator, the chapter expands on the Kalman filter and its popular developments based on ensemble approximations (notably the Ensemble Transform Kalman Filter). This finally leads us to discuss the typical biases and limitations that the ensemble approximation entails, along with the methods commonly employed to mitigate them.

The first two theoretical chapters bring us to **Chapter 3**, where we explore the ways of bridging together FWI and DA. I present five propositions, from which one has been selected to conduct tests in the subsequent chapters. Finally, the scheme that we came-up with and named ETKF-FWI is detailed at the end of the chapter.

Chapters 4 focuses on the application of the ETKF-FWI to the Marmousi synthetic benchmark. It discusses both the numerical implementation of the scheme, along with strategies regarding ensemble generation and analysis of the filter's output. Thanks to the straightforward nature of this synthetic benchmark, ensemble approximation biases are investigated, and an attempt at mitigating them is presented.

Chapter 5 takes on the previous application, but this time on a field dataset. This application allowed us to test the applicability of the ETKF-FWI on a more complicated test case. While it allowed verifying the observations made on the synthetic benchmark, it also allowed performing a multiparameter inversion to evaluate parameters cross-talk, which remains one of the difficulties in FWI.

Conclusions and perspectives are given in the last chapter of this thesis.

Chapter 1

Uncertainty quantification for Full Waveform Inversion

Contents

1.1 Overview of FWI	7
1.1.1 Historical overview	7
1.1.2 Mathematical formulation of FWI	9
1.1.3 Interpretation of the gradient and Hessian operators	16
1.1.4 Link between the Hessian and the posterior covariance in Bayesian estimation	19
1.2 Uncertainty quantification in FWI: state-of-the-art	21
1.2.1 Uncertainty quantification through global optimization techniques	22
1.2.2 Local uncertainty estimation in FWI	27
1.2.3 Ideas from the Data Assimilation community	32

1.1 Overview of FWI

1.1.1 Historical overview

Full Waveform Inversion (FWI) is a seismic tomography technique, aiming at interpreting seismic wave recordings, to characterize subsurface properties. As wavefields behavior and evolution are imposed by the physical properties of their propagating medium, it is possible to infer those physical parameters through inverse problem-solving. Seismic tomography applications cover a broad spectrum of scales and targets, and are commonly used for regional to global scale in the academic community (Aki et al., 1977; Fichtner et al., 2009; Bedle and Lee, 2009; Tape et al., 2010; Panning et al., 2010; French and Romanowicz, 2015; Bozdağ et al., 2016) and for crustal-scale exploration industrial applications (crustal-scale imaging, reservoir monitoring, and civil engineering targets, Plessix, 2009; Sirgue et al., 2010; Plessix et al., 2012; Warner et al., 2013b; Zhu et al., 2015; Operto et al., 2015).

FWI was first formulated by Lailly (1983); Tarantola (1984), as an attempt to bridge the gap between Claerbout (1971)'s migration imaging principle and the travelttime tomographic imaging principle. While migration imaging is mostly responsible for investigating the high-wavenumber content of subsurface

images (sharp details), and travelttime tomography is focused on the low-wavenumber content (overall kinematic), FWI aims at building broad-wavenumber models. It does so, by considering that all arrivals in a seismogram are governed by the same type of wavefield-medium interactions (Devaney, 1984), described by the wave equation.

FWI is thus formulated as a data fitting procedure, where entire observed seismograms are compared with corresponding synthetics, computed by solving the wave equation. Early applications of FWI proved to be too computationally demanding compared to the resources available at the time, even when limited to 2-D cases (Gauthier et al., 1986; Cary and Chapman, 1988; Crase et al., 1990; Jin et al., 1992; Lambaré et al., 1992). Their success were also heavily conditioned by data availability, mostly limited at that time to short-offset seismic reflections surveys. With this type of data, the sensitivity to large wavelengths structure is fairly poor, making FWI applications challenging unless a very accurate starting model is picked for the inversion. Even though the methodology was promising, the numerical cost of FWI and the difficulties associated to short-offset data prevented a significant adoption, and research stalled for a while. FWI was then push forward again in the nineties, with developments in frequency-domain FWI (Pratt, 1990; Pratt and Worthington, 1990; Pratt and Goult, 1991; Pratt et al., 1996, 1998; Pratt, 1999). These developments allowed to acknowledge the importance of long-offsets and transmission data to reconstruct large-scale structures FWI, demonstrated in 2-D cross-hole acquisition by Pratt and Worthington (1990). The frequency-domain formulation, which can be advantageous computationally-wise in 2-D, associated with advances in hardware capacity and new acquisition design, put FWI research back on track.

FWI is now a mature imaging technique, and the research interests have shifted. While early research focused on making the concept applicable and understanding FWI's imaging power, most of the current research is now focused on making FWI more robust and able to tackle more complex problems, especially so in the context of industrial exploration. For instance, FWI on land datasets is known to be particularly challenging, as it requires numerical solvers that can adequately represent complex topography and strong elastic effects that affect the wavefield in shallow subsurfaces environments (Trinh et al., 2019). Efforts are thus made to honor as much as possible the physics of wave propagation in elastic medium, while keeping numerical solver affordable and efficient.

Current research is also attempting to alleviate the inherent ill-posedness of FWI (see 1.1.2), to make its general formulation more robust. As a matter of example we can think of alternative misfit-functions that ought to relax the FWI problem, such as the envelope misfit (Bozdağ et al., 2011), instantaneous phase misfit (Fichtner et al., 2008; Maggi et al., 2009; Bozdağ et al., 2011; Lee and Chen, 2013; Tejero et al., 2015) or cross-correlation based misfit (Luo and Schuster, 1991; Tromp et al., 2005; van Leeuwen and Mulder, 2010). A new way of measuring the misfit between seismic traces have been proposed following the optimal transport paradigm (Engquist and Froese, 2014; Métivier et al., 2016, 2019). The robustness of the technique can also be ensured by using appropriate regularization terms when solving the FWI problem (van Leeuwen and Herrmann, 2013; Warner et al., 2013a; Warner and Guasch, 2014; Aghamiry et al., 2019).

Finally, in the past few years, interest has also grown over the topic of uncertainty quantification in FWI. While the importance of uncertainty estimation has always been promoted in Tarantola's visionary work (Tarantola, 2005), this aspect of the methodology has mostly been left aside through the evolution of FWI. The lack of systematic uncertainty quantification is a big issue that we ought to address in this thesis work: without uncertainty estimation, exploitation, and interpretation of tomographic models is an unsound exercise.

In the following chapter, I will introduce the physical and mathematical concepts underlying FWI.

This theoretical introduction will be followed by a review of uncertainty quantification in FWI to provide the reader with a complete overview of the current research landscape.

1.1.2 Mathematical formulation of FWI

Conventional FWI is formulated as a least-squares minimization problem between recorded wavefield data d_{obs} and synthetics d_{cal} computed from a discrete physical model m . It is an inverse problem, for which one tries to infer information about the mechanical properties of the subsurface, from information contained in observed data. In the case of FWI, the considered observations are either represented by seismograms in the time domain (Fichtner, 2011; Virieux et al., 2017), or by a complex wavefield values in the frequency domain (Pratt et al., 1998; Pratt, 1999; Virieux et al., 2009). By reducing the least-squares distance through model updates, accounting for all phases and amplitudes information, one expects to retrieve a high-resolution, plausible subsurface model, which should explain the data. Mathematically, it is formulated as finding the minimum of the *regularized* misfit functional defined as

$$\mathcal{C}(m) = \frac{1}{2} \|d_{cal}(m) - d_{obs}\|^2 + R(m, m_p), \quad (1.1)$$

where $\|\cdot\|^2$ is the least-squares norm, and m_p is a *prior model* used in the regularization term R . In general, the regularization is introduced either by specifying a prior model (which means the FWI minimization must find the good balance between the data misfit, and the distance from this prior model m_p), or by applying a smoothing term on the FWI model updates.

Before reviewing the various inversion strategies that can lead to solution of equation 1.1, it is required to introduce the *forward problem*, necessary to compute $d_{cal}(m)$ and evaluate the data misfit.

The forward problem: modeling accurate wavefields Any inverse problem requires first and foremost the ability to solve a forward problem. Here, the forward problem refers to solving

$$d_{cal} = \mathcal{H}(m) \quad (1.2)$$

where d_{cal} is the output given the input parameters m and the forward operator \mathcal{H} , which express the action of a physical system on m . This way, \mathcal{H} can be viewed as a (linear or non-linear) mapping from the model space to the data space such that

$$\mathcal{H}(m) : \mathbb{R}^n \rightarrow \mathbb{R}^d, \quad (1.3)$$

or in the case of complex valued observations:

$$\mathcal{H}(m) : \mathbb{R}^n \rightarrow \mathbb{C}^d, \quad (1.4)$$

where d is the number of discrete observations and n the number of model parameters. In the case of FWI, this forward operator requires to solve the full wave equation, that describes the behavior of a propagating wavefield u in any given medium m . In the case of seismic waves, it is common practice to describe the evolution of wavefield under the linear elasticity regime, written as:

$$\begin{cases} \rho(m) \partial_{tt} u_i(m, t) &= \partial_j \sigma_{ij}(m, t) + s_i(m, t), \\ \sigma_{ij}(m, t) &= C_{ijkl}(m) \epsilon_{kl}(m, t). \end{cases} \quad (1.5)$$

In equation 1.5, we assume Einstein's notation convention (summation over indices). u_i is the i^{th} component of the particle motion, σ_{ij} is the stress tensor, ϵ_{kl} is the strain tensor and C_{ijkl} is the stiffness tensor (with $i, j, k, l \in [x, y, z]$). The medium density is defined by ρ and s_i is the i^{th} component of an external volumetric force applied to the medium.

While the system of equations 1.5 describes compressional and shear wave propagation in elastic media, it can be interesting to consider its acoustic approximation, as the acoustic wave equation is generally cheaper to compute, while being relatively accurate in specific contexts. This is the case, for instance, in most of marine seismic exploration where the acoustic approximation of the wave equation is usually favored, as only compressional waves can be recorded in the water column.

In the acoustic approximation, the stiffness tensor reduces to the bulk modulus $\kappa(m) = \rho(m)V_p^2(m)$ where V_p is the compressional wave velocity. which leads to the acoustic approximation

$$\rho(m)\partial_{tt}u(m, t) = \kappa\nabla^2u(m, t) + s(m, t), \quad (1.6)$$

where ∇^2 is the Laplacian operator. Note that in the acoustic case, the physical field $u(m, t)$ corresponds to the (scalar) pressure field across the model m at any given time t , and the source term also reduces to a scalar field.

The acoustic wave equation can also be expressed under its time-harmonic equivalent, by considering the Fourier transform of the pressure field, which yields the wave equation in the Fourier domain

$$-\omega^2\rho u(m, \omega) = \kappa\nabla^2u(m, \omega) + s(m, \omega), \quad (1.7)$$

where ω is the angular frequency. This frequency-domain formulation can have several advantages in the frame of FWI: seismic attenuation is easily accounted for by using a complex-valued velocity V_p (Toksöz and Johnston, 1981), multi-sources can be computed efficiently if a direct solver is used to solve the corresponding linear system (Li et al., 2019). Finally, the time-harmonic formulation can allow significant dimensionality reduction of the data-space.

As we will solely focus on 2-D marine geophysics case studies, we only consider the frequency domain acoustic wave equation (which accounts for attenuation) to compute the physical field used in our FWI scheme. This choice of formalism has been driven by the low computational cost it entails for the numerical experiments conducted in Chapters 4 (where the numerical modeling scheme will be detailed) and 5 .

The acoustic wave equation is expressed under the following compact form

$$\mathbf{B}(m, \omega)u(m, \omega) = s(m, \omega), \quad (1.8)$$

where $\mathbf{B}(\omega, m)$ is a complex-valued "impedance" matrix (Marfurt, 1984). Note that in our forward problem, the wavefield depends on the modeling frequency: solving the forward problem yields a steady-state pressure field for each monochromatic frequency ω .

Under this compact form, computing the synthetic wavefield measurements $d_{cal}(\omega, m)$ comes down to

$$\begin{aligned} u(\omega, m) &= \mathbf{B}^{-1}(\omega, m)f(\omega), \\ d_{cal}(\omega, m) &= \mathbf{E}u(\omega, m), \end{aligned} \quad (1.9)$$

where \mathbf{E} is linear observation operator that extracts the values of u at receivers' locations. Finally, by noticing that the monochromatic synthetic data are only dependent on model parameters m and frequency ω , we can define a frequency-dependent forward operator $\mathcal{H}_\omega(m)$ such that

$$\mathcal{H}_\omega(m) = \mathbf{E}\mathbf{B}^{-1}(\omega, m)s(\omega) = d_{cal}(\omega, m). \quad (1.10)$$

Global solution to the FWI problem A global optimization method is a technique that ought to locate the global minimum of a continuous, possibly non-convex, misfit functional

$$\mathcal{C}(m) : \Omega \subset \mathbb{R}^n \rightarrow \mathbb{R}, \quad (1.11)$$

where Ω is defined as the *solution space*, discretized over n parameters. The interests of using such methods lie in their capacity to handle non-convex misfit functions (such as illustrated in Figure 1.1), and therefore, to deal with the occurrence of local-minima in Ω (which typically arise in ill-posed inverse problems such as FWI). With this type of approaches, one will attempt to sample adequately the solution space, rather than finding a point-localized solution in Ω : the goal is to explore and "map" the misfit function to evaluate all of its minima and find the optimum. By exploring the solution space, global search approaches are also a natural candidate for uncertainty estimation.

They require the evaluation of the misfit function in numerous points of Ω to draw an accurate description of said "global" map (the solution space). In the case of FWI, this misfit function evaluation requires to compute $d_{cal}(m)$ through the evaluation of the forward operator.

The performances of global search approaches are thus closely tied with the size of Ω , and therefore, the number of discrete points n that are used to represent the system: as the number of parameters increases, the number of samples required to accurately represent Ω increases also, which limits the applications of global searches to $n \sim 10^1$ to 10^3 .

Despite having shown little success in FWI applications (Martin et al., 2012; Bardsley et al., 2014; Biswas and Sen, 2017; Sajeve et al., 2017b), these methodology are still an active topic of research. They rely on various stochastic sampling strategies such as: Markov chain Monte Carlo (MCMC) methods (Metropolis-Hastings sampling (Metropolis et al., 1953), Gibbs sampler (Geman and Geman, 1987), Hamiltonian MCMC (Duane et al., 1987) and reversible jump MCMC (Green, 1995)) or Genetic algorithms (Mitchell, 1998) to name a few. An extensive review of global search approaches applied to FWI will be given later in this chapter.

Local solution to the FWI problem When the problem size becomes intractable for stochastic sampling, the classic alternative is to use local optimization techniques. Local optimization aims at using local information of the misfit function (such as its gradient and its curvature) to "navigate" toward one of the function's minimum. With these approaches, the convergence is always driven toward the closest minimum, hence, the choice of initial model strongly defines these schemes performances. In a way, defining the starting model m_0 defines the subspace $\mathcal{A} \subset \Omega$ in which the misfit function is minimized, de-facto reducing the number of possible solutions as the new solution space can be a small fraction of Ω . It also defines the solution one can expect to retrieve once the closest minimum is reached. The global optimum can only be reached with an appropriate starting point if the misfit function is not globally convex.

Given any initial model m_0 such that $\mathcal{C}(m_0) \in \mathcal{A}$, each subsequent solution in the sequence $m_k = (m_0, m_1, \dots, m_k)$ is considered to be located in the vicinity of the previous step (and thus in \mathcal{A}). The model update can be written as a sum of a model perturbation Δm_k and m_k such that:

$$\begin{aligned} m_{k+1} &= m_k + \Delta m_k \\ &= m_k + \alpha_k d_k, \end{aligned} \quad (1.12)$$

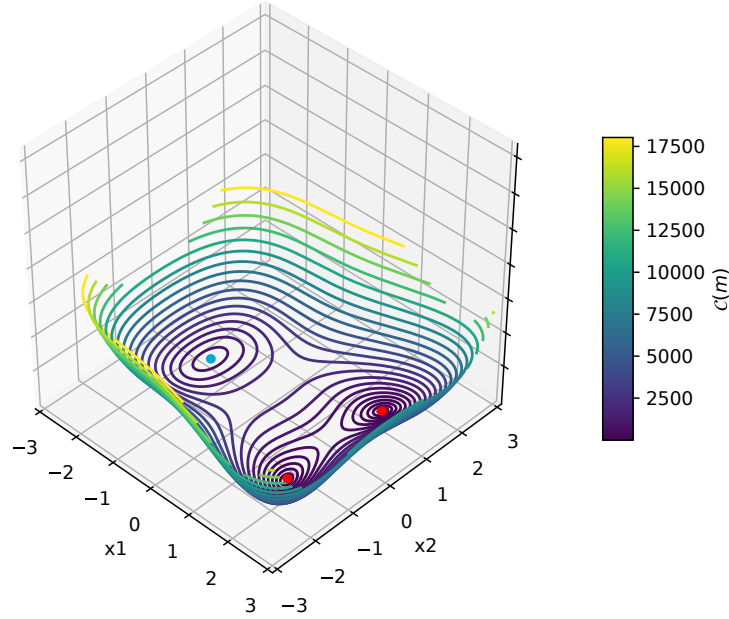


Figure 1.1: Multimodal misfit function $\mathcal{C}(m) : \Omega \subset \mathbb{R}^2 \rightarrow \mathbb{R}$. Exactly three minima are visible on this function, two of which are global (this specific optimization problem has a non-unique solution). Red dots denote the two global minima while a blue dot denotes the third minimum. Global search approaches are specially fitted for this type of settings, as they can help to identify the possible solutions to the minimization problem.

where d_k is the local *descent direction* of the misfit function at iteration k , and α_k is the *step length* in the direction of descent (how much m_k moves along d_k). These two parameters are computed at each iteration to ensure that each m_{k+1} provides a better fit to the data:

$$\mathcal{C}(m_{k+1}) < \mathcal{C}(m_k). \quad (1.13)$$

and at the same time that each iteration converges sufficiently fast toward the minimum (Wolfe, 1969).

1.1.2.1 Linesearch: the Wolfe conditions

The choice of parameter α_k is a key component to an efficient optimization strategy. Improperly chosen, it can lead to underperformances in term of convergence speed (the steps along d_k are too small) or chaotic/oscillatory behavior (oscillation around the minimum without reaching it or even moving to a different search subset) as can be seen in Figure 1.3.

Finding the optimal α_k is akin to an optimization problem on its own, in which one tries to find the optimal value α_k^* such that

$$\alpha_k^* = \underset{\alpha_k}{\operatorname{argmin}} \quad \mathcal{C}(m_k + d_k \alpha_k). \quad (1.14)$$

Minimizing this cost function iteratively with a dedicated optimization method, is referred to as solving an *exact line search* problem. In practice however, exact line search methods are only use in very special cases, such as when $\mathcal{C}(m)$ is perfectly quadratic. Instead of finding the α_k that minimizes equation (1.14), we can define conditions that are easier to satisfy, while allowing the convergence toward the

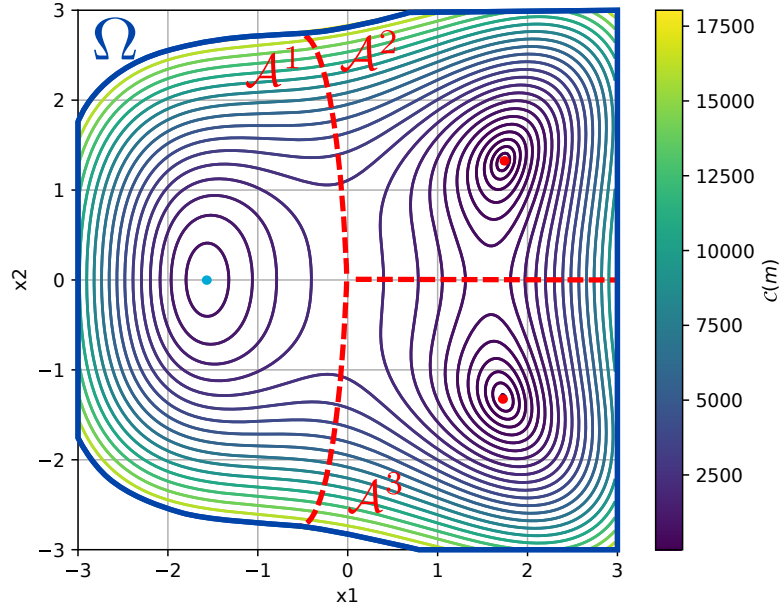


Figure 1.2: 2-D view of the multimodal misfit function $\mathcal{C}(m) : \Omega \subset \mathbb{R}^2 \rightarrow \mathbb{R}$ displayed in Figure 1.1. The domain space Ω and its three local subsets \mathcal{A}^1 , \mathcal{A}^2 and \mathcal{A}^3 associated to each of the three minima are represented. Any starting point chosen in one of the subset for the optimization method, will lead to the closest local minimum following the "topography" of $\mathcal{C}(m)$.

closest minimum. This type of approach to determine alpha are referred to as *inexact line search* methods (Nocedal and Wright, 2006).

One common set of conditions are known as the *Wolfe conditions*, which is a set of two constraints defined as

$$\begin{aligned} \mathcal{C}(m_{k+1}) &\leq \mathcal{C}(m_k) + c_1 \alpha_k d_k^T \nabla \mathcal{C}(m_k) \\ d_k^T \nabla \mathcal{C}(m_{k+1}) &\geq c_2 d_k^T \nabla \mathcal{C}(m_k) \end{aligned} \quad (1.15)$$

where the constants c_1 and c_2 are chosen such that $0 < c_1 < c_2 < 1$. The first inequality is known as the first Wolfe condition (or Armijo condition) and is a linear decrease criterion. This first condition acts as an upper-bound by imposing a linear minimization rate. The second Wolfe condition, also known as the *curvature* condition, ensures that the step along the descent direction is not too small. It imposes that the slope at $k + 1$ is greater than the slope at k times a constant c_2 and excludes small step lengths that would slow down the convergence. The curvature condition acts as a lower-bound to the line search algorithm. Taken together, the Wolfe conditions ensure that the optimization converges to the closest minimum at an optimal speed, provided c_1 and c_2 are set adequately. The behavior of the line-search algorithm is highly dependent on the values of c_1 and c_2 , and readers might refer to Nocedal and Wright (2006) for a complete review and guidelines.

We can cite the *backtracking method* (Ortega and Rheinboldt, 1970) as a popular example of line-search algorithm (Nocedal and Wright, 2006). In this thesis work, the line-search procedure used is a

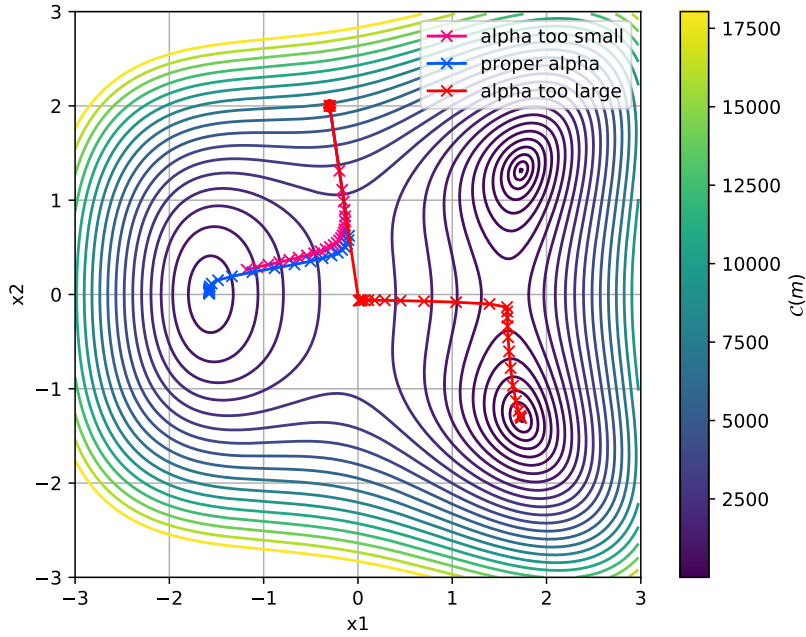


Figure 1.3: Alpha selection for steepest descent optimization strategy. The starting point of the optimization corresponds to subdomain $\mathcal{A}^1 \subset \Omega$, therefore the minimum to be recovered is the blue one (see Figure 1.2). Values of α_k have been maintained fixed for each of the $k = 25$ minimization iterations, for demonstration purposes. If alpha is fixed with a too small value, 25 iterations are not enough to drive the convergence to the minimum (pink). When α_k is set appropriately (blue), the convergence is ensured in 25 iterations. If α_k is too large, the first iteration sets us in the subset \mathcal{A}^3 and the recovered minimum is not the desired one.

bracketing method, implemented in the SEISCOPE Optimization Toolbox (Métivier and Brossier, 2016). Alternatively to the line-search method, one can define the optimal step-length according to a *trust region* method (Nocedal and Wright, 2006). However, trust region methods are deemed more complicated to implement, therefore line-search algorithm are generally preferred for their simplicity.

1.1.2.2 Defining the descent direction

The descent direction has been defined as a vector d_k that allows a minimization of the data misfit from a step k to $k + 1$, along its direction. A first order Taylor expansion of $\mathcal{C}(m_k)$ yields

$$\mathcal{C}(m_k + \Delta m_k) \approx \mathcal{C}(m_k) + \Delta m_k^T \nabla \mathcal{C}(m_k) + \mathcal{O}(\Delta m_k^2), \quad (1.16)$$

where T designate the transpose operator. As the step length is conventionally chosen to be strictly positive ($\alpha_k > 0$), the gradient descent condition in equation 1.13 gives

$$d_k^T \nabla \mathcal{C}(m_k) < 0, \quad (1.17)$$

hence the vector d_k is a descent direction only if its scalar product with the gradient of $\mathcal{C}(m_k)$ is negative.

First order methods - tangent linear approximation The first choice of descent direction that can be used to perform such minimization procedure, is the steepest descent direction, defined by the opposite

of the gradient of the misfit-function. One can think of gradient-descent optimization as a ball rolling along the steepest slope of a surface: the ball would represent the optimized model, and the surface would be the misfit function. This direction is given by

$$d_k = -\alpha \nabla \mathcal{C}(m_k), \quad (1.18)$$

which corresponds to a minimisation along a tangent-linear approximation of the misfit function. While its formulation is rather simplistic, this method is fairly limited by its linear assumptions and a low convergence rate (at best linear for a linear optimization problem), which prevents its application on practical FWI cases.

Second order methods - quadratic approximation To go beyond the tangent-linear approximation of the gradient method, we can develop equation 1.12 to the second order Taylor expansion, yielding,

$$\mathcal{C}(m_k + \Delta m_k) \approx \mathcal{C}(m_k) + \Delta m_k^T \nabla \mathcal{C}(m_k) + \frac{1}{2} \Delta m_k^T \nabla^2 \mathcal{C}(m_k) \Delta m_k + \mathcal{O}(\Delta m_k^3), \quad (1.19)$$

where $\nabla^2 \mathcal{C}(m_k)$ is the second-order derivative of the misfit function, also called the Hessian operator. For the sake of notation simplification, we will denote the gradient and the Hessian of the misfit function at m_k by \mathbb{G} and \mathbb{H} respectively. The minimizer of the quadratic approximation gives the second order descent direction. It is obtained as a stationary point of the quadratic approximation gradient. Therefore we can compute the optimum value of Δm_k that is a root of equation 1.19 derivative:

$$\nabla \mathcal{C}(m_k + \Delta m_k) = 0 \quad \text{and thus} \quad \mathbb{G} + \mathbb{H} \Delta m_k = 0, \quad (1.20)$$

hence the perturbation that minimizes $\mathcal{C}(m_k + \Delta m_k)$ is given by

$$\Delta m_k = -\mathbb{H}^{-1} \mathbb{G}. \quad (1.21)$$

While the interpretation of the gradient descent was described as "following the steepest descent direction", the second order method, or Newton's method, has a different geometrical interpretation. Because we consider the second order Taylor expansion of the misfit function, we can express it under the form of a parabola

$$y = 1 + x + \frac{1}{2} x^2. \quad (1.22)$$

Where the gradient method yield a tangent-linear approximation, the Newton's method yield a parabollic approximation of the misfit function (Figure 1.4). The minimization point correspond to the minimum of the fitting parabola: in a linear least-squares setting, a single iteration is sufficient to find the global optimum (as opposed as a linear convergence rate for the gradient method).

Thanks to their faster convergence rate, second-order methods ("Newton-type") are favored in the context of FWI. Because each cost function evaluation requires solving a costly forward problem, we are bound to choose the methods that have the faster convergence rate, and thus rule-out first-order methods. However, computation of the Hessian operator is generally non-trivial: it is a computational and memory burden, as the Hessian operator is a matrix of $n \times n$ parameters. It is impossible to build, yet to inverse in practical FWI applications and Newton's method is out of reach in most realistic settings (Pratt et al., 2008; Fichtner and Trampert, 2011a; Métivier et al., 2013).

To circumvent this limitation and still benefit from the effect of the Hessian (which encapsulate curvature information), we usually rely on *quasi-Newton* optimization methods (Nocedal and Wright,

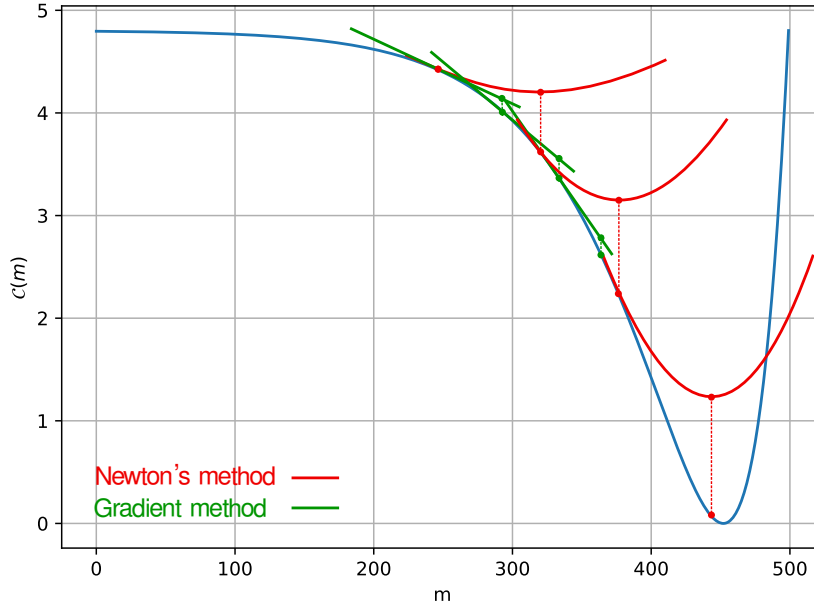


Figure 1.4: Geometrical representation of first and second order minimization methods. Second order methods demonstrate a faster convergence rate thanks to the local quadratic minimizer.

2006) which allows approximating the effect of the inverse Hessian iteratively and very cost-effectively, without explicit computation of \mathbb{H} or \mathbb{H}^{-1} . Instead of building the Hessian explicitly (either exactly or by using a finite-difference approximation), the quasi-Newton method builds an approximation of the Hessian based on the gradient's evolution across the optimization steps.

In essence, quasi-Newton express the model perturbation (or the *Newton step*) as

$$\Delta m_k = -\mathbb{B}_k^{-1} \mathbb{G}, \quad (1.23)$$

where \mathbb{B}_k is positive-definite matrix chosen to approximate \mathbb{H} , and is generally initialized to be diagonal, and is updated at each subsequent iterations. The various quasi-Newton algorithm differ in the way they update of \mathbb{B}_{k+1} .

In this work, we relied on the l-BFGS algorithm (for *limited memory*-Broyden–Fletcher–Goldfarb–Shanno). Instead of storing an $n \times n$ Hessian approximation, this algorithm reduce this computational burden by evaluating a low-rank approximation of \mathbb{B}_k that allows to recursively evaluate equation 1.23 from a limited number of vectors of size n , reducing the memory requirements of the BFGS method.

In this work, all of the FWI applications were carried-out using SEISCOPE Optimization Toolbox's (Métivier and Brossier, 2016) implementation of the L-BFGS method. Readers may refer to Nocedal and Wright (2006) for a complete overview of quasi-Newton and limited-memory quasi-Newton methods, and details on their practical implementations.

1.1.3 Interpretation of the gradient and Hessian operators

We have seen that local-optimization strategies are the go-to option for practical FWI applications, as the global-optimization options are too expensive for high-dimensional problems. However, while global

optimization methods are intrinsically linked with a probabilistic evaluation of the misfit function, the uncertainty information in local search methods is less evident at first glance.

To better understand the complex link between the misfit function, resolution, tradeoffs and uncertainty of the recovered model, it is necessary to precisely expose the two main operators that drive the minimization, namely the gradient and the Hessian operator. Both are directly tied to the misfit functional and are then the direct elements controlling the reconstruction of subsurface models through FWI, and their uncertainty. I review these operators in the following subsections, which are deeply rooted in Pratt et al. (1998) development of the Gauss-Newton and Newton methods for FWI in the frequency domain.

Gradient interpretation In frequency domain FWI, the computation of the gradient comes down to computing the partial derivative of the misfit function with respect to the subsurface model parameter m

$$\mathbb{G} = \frac{\partial \mathcal{C}(m_k)}{\partial m_k} = \Re(\mathbf{J}^T \delta d^*) = \Re \left[\left(\mathbf{E} \frac{\partial u}{\partial m} \right)^T \delta d^* \right], \quad (1.24)$$

where \mathbf{J} denotes the $n \times d$ Frechet derivative matrix (the Jacobian operator), \Re denotes the real-part and δd^* is the conjugate of the data difference

$$\delta d = d_{cal}(m_k) - d_{obs}, \quad (1.25)$$

yielding the expression of the misfit function for complex valued wavefield data

$$\mathcal{C}(m) = \frac{1}{2} \delta d^T \delta d^*. \quad (1.26)$$

Aranging equation 1.8 with respect to the parameter $m_{k,i}$ yields an expression of the Jacobian as a function of the forward operator

$$\mathbf{B} \frac{\partial u}{\partial m_{k,i}} = - \frac{\partial \mathbf{B}}{\partial m_{k,i}} u, \quad (1.27)$$

or

$$\frac{\partial u}{\partial m_{k,i}} = \mathbf{B}^{-1} s_i^\dagger, \quad (1.28)$$

where s_i^\dagger correspond to the i^{th} virtual source $s_i^\dagger = - \frac{\partial \mathbf{B}}{\partial m_{k,i}} u$. I will now reffer to $\frac{\partial u}{\partial m_{k,i}}$ as the *partial derivative wavefield* at $m_{k,i}$. This show that one can express the computation of the Frechet derivative as solving a forward problem, with the complex-valued impedance matrix and a new *source term* (as in equation 1.9). This virtual source term is the product of the predicted wavefield u and the partial derivative of the impedance matrix. This operator holds the radiation pattern of the model (i.e., how the wavefield interacts with the medium given the link between physical parameters and the angle of incidence of the wavefield). In other words, this parameter contains information about the scattering properties of the medium m_k . Thus the source term in equation 1.28 can be viewed as the predicted wavefield u , scattered by the parameter $m_{k,i}$.

Considering

$$\frac{\partial u}{\partial m} = \left[\frac{\partial u}{\partial m_1}, \frac{\partial u}{\partial m_2}, \dots, \frac{\partial u}{\partial m_{k,i}} \right] = \left[\mathbf{B}^{-1} s_1^\dagger, \mathbf{B}^{-1} s_2^\dagger, \dots, \mathbf{B}^{-1} s_i^\dagger \right] \quad (1.29)$$

or

$$\mathbf{J} = \mathbf{B}^{-1} \mathbf{S} \quad (1.30)$$

where \mathbf{S} contains all the virtual sources the gradient can be expressed as

$$\mathbb{G} = \Re(\mathbf{J}^T \delta d^*). \quad (1.31)$$

The gradient is consequently the natural direction of the model update, as it tends to locally adjust the parameters in m that are responsible for the data mismatch. This can further be highlighted if we express the gradient as

$$\mathbb{G} = \Re((s_i^\dagger)^T v) = \Re(u^T [\frac{\partial \mathbf{B}^T}{\partial m_{k,i}}] v), \quad (1.32)$$

where it clearly appears that the gradient results from a zero-lag gross correlation of the predicted and the backpropagated wavefield, projected over the radiation pattern operator (scattering properties of the medium), highlighting the missing point-scatterers in m_k . This formulation corresponds to the adjoint-state formulation in Plessix (2006).

Hessian interpretation Given that the Hessian operator is the second partial derivative of the misfit function, it can be expressed as

$$\mathbb{H} = \nabla \mathbb{G} = \mathbf{J}^T \mathbf{J} + \left(\frac{\partial \mathbf{J}}{\partial m} \right)^T (\delta d \dots \delta d), \quad (1.33)$$

such that it is expressed as a sum of two terms.

At the vicinity of the global minimum, where the gradient gives the direction of the global solution, the Hessian provides information on the local curvature and convexity of $\mathcal{C}(m_k)$. It is therefore a critical piece of information in the framework of uncertainty estimation as it is a measure of the interdependency between parameters and holds resolution measure (Pratt et al., 1998; Fichtner and Trampert, 2011b; Fichtner and van Leeuwen, 2015)

The first term is generally designated as *approximated Hessian* (or Gauss-Newton approximation of the Hessian) when considered on its own. By the nature of the $\mathbf{J}^T \mathbf{J}$ product, it corresponds to the zero-lag cross-correlation of the partial derivative wavefields, between model parameters m_i and m_j . The Gauss-Newton Hessian is a symmetrical operator, which diagonal terms contains the zero-lag autocorrelation of the partial derivative wavefields. Because this term contains the squared amplitude of the partial derivative wavefields, it naturally contains a measure of the geometrical spreading that affects the scattered wavefields. Thus, the Hessian can be used as a scaling operator to account for the geometrical spreading in \mathbb{G} and improve the balance of model updates (by "boosting" the update in areas of the model that are poorly sampled by the wavefield). The off-diagonal terms of the Gauss-Newton Hessian are given by the intercorrelations of the partial derivative wavefields of parameters m_i and m_j . They provide information about the links between different model parameters and register band limited effects and illumination. In the context of FWI, these terms act as a convolution operator applied to a point localised model parameter m_i . The maximum expected sharpness of the recovered image is intrinsically encoded in the off-diagonal terms of the Hessian operator. In the same fashion we can account for the geometrical spreading with the gradient preconditioning, the Hessian can also be seen as a de-convolution operator, to refocus the point localized model parameters which helps "de-blur" the recovered model. In the case of a multi-parameter inversion, it also helps with the scaling between the various physical parameter updates.

The second-order terms $\left(\frac{\partial \mathbf{J}}{\partial m} \right)^T (\delta d \dots \delta d)$, contain the product of second partial derivative wavefields with data residual. It accounts for double scattering (Pratt et al., 1998), which the gradient and

Gauss-Newton Hessian neglect. Because it contains second-order partial derivative wavefields, which requires numerous forward modeling solve to build, this term is generally not accounted for in FWI. This approximation is reasonable if the problem is not too strongly non-linear, or if the data-misfit is small (Tarantola, 1987; Pratt et al., 1998). Nonetheless, Métivier et al. (2017) have shown the interest of considering these second order terms in case of very strong multiscattering effects in the data (that can arise in very shallow subsurface environments or in civil engineering applications).

1.1.4 Link between the Hessian and the posterior covariance in Bayesian estimation

From the description of the gradient and Hessian operators, it is clear that both play a key role in the resolution and tradeoffs we can expect to recover from an FWI application. While the gradient can be computed explicitly at a reasonable cost, it does not have any straightforward uses for uncertainty estimation. The Hessian operator, on the other hand, is the key operator to perform uncertainty assessments and quality controls. To emphasize the importance of \mathbb{H} in uncertainty estimation, I will recall the Bayesian formulation of the linearized least-squares problem and show how statistics play a role when the model parameters of inverse problems are modeled as random processes. Consider the forward problem

$$d_{cal}(m) = \mathcal{H}(m) + \eta \quad (1.34)$$

where we add η as a noise vector. In the *Bayesian* formalism, both vectors m and η are considered as random processes, with associated probability distributions. The solution of the inverse problem in this setting is a probability distribution that updates the beliefs on m by incorporating information from d_{cal} .

Let us first consider as an example the following regularized deterministic inverse problem

$$m^* = \underset{m}{\operatorname{argmin}} \frac{1}{2} \|d_{cal}(m) - d_{obs}\|^2 + \frac{1}{2} \|m - m_p\|^2, \quad (1.35)$$

where the regularization term in equation 1.35 corresponds to the Tikhonov regularization (Tikhonov and Arsenin, 1977), and is used to mitigate the ill-posedness of the inverse problem. It does so by imposing constraints on the optimization (for instance, m_p could be a smoothness constraint to prevent overfitting the observations).

Assuming d_{obs} is affected by additive Gaussian noise and assuming the prior distribution to be Gaussian (prior knowledge about the structure of m), the posterior distribution of m knowing d_{obs} is given by the probability density function (PDF) $p_{d_{obs}}(m)$ such that

$$p_{d_{obs}}(m) \propto \exp\left(-\frac{1}{2} \|d_{cal}(m) - d_{obs}\|^2 - \frac{1}{2} \|m - m_p\|^2\right). \quad (1.36)$$

The PDF gives the probability for the random variable m to take any value in \mathcal{A} . It is closely tied to the local optimization problem, as the maximum probability values are given by minimizers of the misfit function $\mathcal{C}(m)$.

The link with the Hessian operator can be demonstrated by considering the maximum value of the PDF, the Maximum A Posteriori (MAP). By separating the "data" and "prior" term of the misfit function in equation 1.35, we can express $\mathcal{C}(m)$ as

$$\mathcal{C}(m) = \mathcal{C}^p(m) + \mathcal{C}^d(m) \quad (1.37)$$

where the prior model misfit $\mathcal{C}_p(m)$ is given by

$$\mathcal{C}_p(m) = \frac{1}{2}(m - m_p)^T \mathbf{P}_p^{-1}(m - m_p) \quad (1.38)$$

with \mathbf{P}_p being the prior covariance (uncertainty on the prior estimate).

The data misfit is given by

$$\mathcal{C}^d(m) = \frac{1}{2}(d_{cal}(m) - d_{obs})^T \mathbf{R}^{-1}(d_{cal}(m) - d_{obs}), \quad (1.39)$$

where \mathbf{R} is the measurement noise matrix (measurement uncertainty). In the linear quadratic assumption, $\mathcal{C}^d(m)$ can be expressed as

$$\mathcal{C}^d(m) = \frac{1}{2}(d_{obs} - d_{cal}(m_p) - \mathbf{H}(m - m_p))^T \mathbf{R}^{-1}(d_{obs} - d_{cal}(m_p) - \mathbf{H}(m - m_p)), \quad (1.40)$$

where \mathbf{H} is the linearized version of the forward map defined in equation 1.10. Accordingly, the gradient of the misfit function is given by

$$\nabla \mathcal{C}(m) = \mathbf{P}_p^{-1}(m - m_p) - \mathbf{H}^T \mathbf{R}^{-1}(d_{obs} - d_{cal}(m_p) - \mathbf{H}(m - m_p)). \quad (1.41)$$

The solution of $\nabla \mathcal{C}(m) = 0$ yields the MAP point of the PDF, which gives the optimal state estimate \hat{m}

$$\begin{aligned} m^* &= m_p + (\mathbf{P}_p^{-1} + \mathbf{H}^T \mathbf{R}^{-1} \mathbf{H})^{-1} \mathbf{H}^T \mathbf{R}^{-1}(d_{obs} - d_{cal}(m_p)) \\ &= m_p + \mathbf{P}_p \mathbf{H}^T (\mathbf{H} \mathbf{P}_p \mathbf{H}^T + \mathbf{R})^{-1} (d_{obs} - d_{cal}(m_p)) \end{aligned} \quad (1.42)$$

Note that the solution of equation 1.42 is the solution to the Best Linear Unbiased Estimator (see 2.1.3). From here, it is possible to make the connection between the posterior covariance matrix \mathbf{P} given by the BLUE equations (the covariance defining the Gaussian PDF in equation 1.36) and the Hessian operator in the vicinity of \hat{m} (Press et al., 1986; Draper and Smith, 1998; Tarantola, 2005). In the vicinity of the minimum, the expression of the Hessian of $\mathcal{C}(m)$ is given by

$$\mathbb{H}(m^*) = \mathbf{P}_p^{-1} + \mathbf{H}^T \mathbf{R}^{-1} \mathbf{H}, \quad (1.43)$$

while the posterior covariance of the Bayesian least-squares solution in equation 1.42 is given by the BLUE equation

$$\mathbf{P} = (\mathbf{I} - \mathbf{P}_p \mathbf{H}^T (\mathbf{H} \mathbf{P}_p \mathbf{H}^T + \mathbf{R})^{-1} \mathbf{H}) \mathbf{P}_p. \quad (1.44)$$

By applying the Sherman-Morrison-Woodbury formula (also known as the matrix inversion lemma)

$$\mathbf{P} = \mathbf{P}_p - \mathbf{P}_p \mathbf{H}^T (\mathbf{H} \mathbf{P}_p \mathbf{H}^T + \mathbf{R})^{-1} \mathbf{H} \mathbf{P}_p = (\mathbf{P}_p^{-1} + \mathbf{H}^T \mathbf{R}^{-1} \mathbf{H})^{-1} \quad (1.45)$$

we can demonstrate the equivalence between the posterior of the Bayesian least-squares solution, and the Hessian in the vicinity of the solution

$$\mathbf{P} = \mathbb{H}^{-1}(m^*). \quad (1.46)$$

The posterior covariance matrix of the Bayesian inverse problem yields information on the state estimate. It provides a measure of the uncertainty we have about the estimate \hat{m} given the prior and data uncertainty. The covariance is an $n \times n$ positive semi-definite symmetrical matrix, that is generally diagonally dominant. The diagonal of the covariance matrix contains the absolute variance values of the solution

for each of the \hat{m}_i parameters. The off-diagonal terms contain the cross-covariances between the model parameters \hat{m}_i and \hat{m}_j . The covariance matrix is thus holding all of the uncertainty information that would be desirable to estimate in FWI application. However, because neither the Hessian operator nor its inverse can be computed explicitly in the case of high-dimensional spaces, uncertainty estimation in FWI is made very challenging.

In this section, the general framework for inverse problem has been developed. The deterministic least-squares local inversion has been presented, along with the gradient and Hessian operator and their respective roles in the inversion framework. I have also presented how the Hessian operator is holding the uncertainty estimation, and how we can make a connection between the Bayesian Inference Framework via the MAP estimator. Note that this assumption only holds in the linear Gaussian case and in the vicinity of the global minimum. It is clearly identified that one should estimate the posterior covariance matrix to estimate the uncertainty of the solution. In the next section, I will review the various methodologies that have already been proposed in the literature, to underline their strength and their shortcomings.

1.2 Uncertainty quantification in FWI: state-of-the-art

In this section, I propose a review of the uncertainty estimation methods for the FWI problem, within two distinct groups: the global optimization methods, that are naturally geared to proposed uncertainty estimation, and the local optimization approaches, based on the estimate of the inverse Hessian operator in the vicinity of the solution. These have been the two main choices of methods for solving the uncertainty estimation problem in the literature. This review will be followed by a quick introduction of "alternative" methods inspired from the Data Assimilation literature, which is the main focus of this thesis. But before that, I would like to quickly address this apparent simple question...

What is even "uncertainty"?

It appears that there is no consensus on the definition of "uncertainty." There is a general definition, in the literal sense, of course (*something that is not known or certain* according to the Cambridge Dictionary). However, the lack of a precise mathematical definition is painfully apparent when one is interested in estimating or quantifying it. It is especially true in the context of inverse problem theory: how do we define uncertainty in that context? As discussed previously, uncertainty is pretty much associated with the misfit function we are dealing with, as we have shown equivalence between inverse Hessian and the covariance matrix. Nevertheless, then, which part of this misfit function is the uncertainty we are seeking to estimate? One possibility would be to evaluate the complete topology of the misfit, and locate all the possible minimum, to finally infer statistics on "which solution is the more likely." However, it appears that even when you decide which minimum is your solution, the shape of the misfit in its vicinity gives us valuable information on how certain (or not) this solution is (tradeoff and resolution information).

Both the "global" and "local" topologies of the misfit are holding uncertainty information, and even though we could argue that the "global" topology is the ideal representation of uncertainty, both global and local estimates hold a measure of uncertainty. I would like to emphasize to the readers that in this work, we have developed a local Bayesian uncertainty estimation framework, similar to local Hessian estimation methods developed in the geophysical exploration literature ((Bui-Thanh et al., 2012; Fang et al., 2018; Liu and Peter, 2019; Gineste et al., 2020)), which are also termed "uncertainty estimation methods." It is also interesting to note that the Data Assimilation community (that will be discussed

in more detail in Chapter 2) has also followed the same type of semantics for "uncertainty estimation" based on a local estimate of the Hessian (Kalmikov and Heimbach, 2014; Rao and Sandu, 2015). We will closely follow this type of framework in this thesis. Thus, "local estimation" of uncertainty is very much associated with the local estimate of the second-order derivative of the misfit in the vicinity of the solution. With that set-out, we shall review the global and local uncertainty estimation methods that have been proposed in the FWI literature thus far.

1.2.1 Uncertainty quantification through global optimization techniques

Despite being challenging to apply in the context of high-resolution (and therefore, high-dimensional solution spaces), global optimization is still an active topic of research in the FWI community. As mentioned in the previous section, the global search approaches have the advantage of providing systematic uncertainty estimation along with the solution of the problem. In theory, they are also able to sample the full solution space and map all the local minima that exist within Ω , without constraining the inverse problem with strong prior information. Additionally, the global optimization algorithms are often gradient-free, which can be particularly useful when the forward modeling problem is given as a black-box, or if the misfit function differentiation is non-trivial to compute. Thanks to this property, they benefit from lower memory requirement, algorithmic complexity and some of them are formulated as embarrassingly parallel problems.

However, there also exists significant limitations within the global optimization framework. First, the convergence is rarely guaranteed, or it requires an infinitely large amount of samples to converge. Due to this limitation, the choice of the initial sample and other hyperparameters inherent to stochastic optimization can limit the effective sampling space. The limited amount of samples one can evaluate can result in a biased overview of the solution. The stochastic properties of these methods can also be a problem, as for two identical initial conditions, the two sampled solution might be different, which make the evaluation of their performances complicated. Finally, they inherently require a large amount of forward modeling solves to evaluate the numerous samples needed to describe the solution space. While they certainly have an ideal theoretical formulation compared to local search methods, they require some tinkering and clever parameterization to be applied to FWI, in order to overcome the curse of dimensionality that arise in high-dimensional problems.

Early application of global optimization methods to FWI can be found in Sen and Stoffa (1991), with an application of simulated annealing (SA) and Stoffa and Sen (1991), with an implementation of genetic algorithm (GA). Due to computational limitations at that time, and the costly nature of these approaches, their pioneering works were conducted on 1-D velocity profiles of the earth, with at best twenty unknowns (ten layers, which velocity and density values were evaluated). Natural processes inspire both methodologies: SA is inspired by slow re-crystallization of a melted solid, and GA is based on natural selection of a population that evolve by random mutations. On top of their computational cost, they are also limited by hyperparameters that control the "crystallization" and the rate of mutations. To find the *critical temperature*, Sen and Stoffa (1991) had to perform a systematic evaluation of the temperature (systematic search approach), which would not be feasible on slightly larger inverse problems. The same goes with Stoffa and Sen (1991), as they had to define an optimal probability value to control the rate of random mutations. Though their computational cost limited both applications, GA has a clear advantage over SA, as it is easily parallelizable over the population members. Moreover, these applications have shown the importance of adequate parameterization of the problem, demonstrated with a synthetic over-determined case where the inversion results get inaccurate (Figure 1.5).

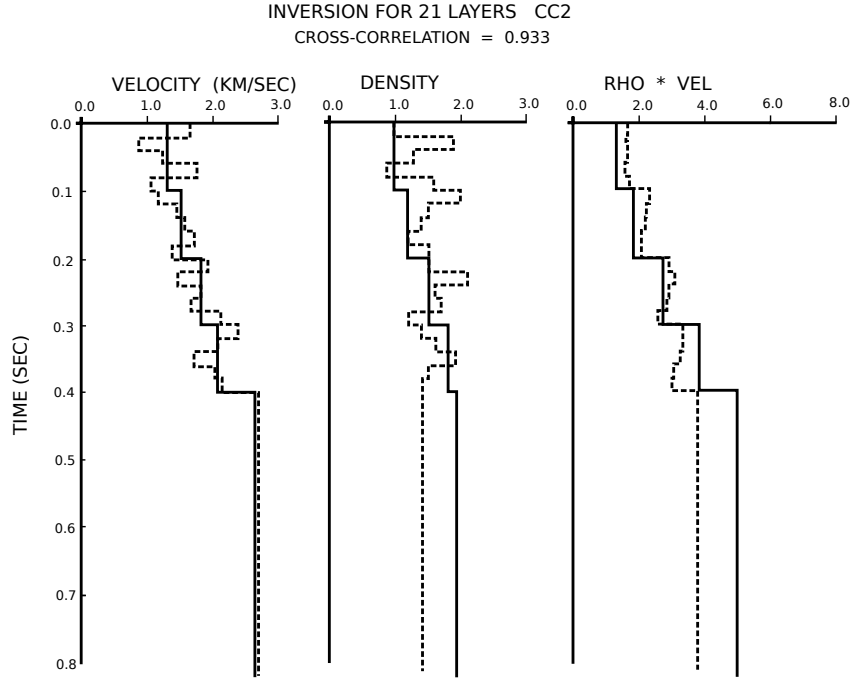


Figure 1.5: Comparison between the true (solid line) and the inverted profile (dashed line) with a cross-correlation of 0.933 (highly fitting model). This setting clearly shows how the model parameter are badly estimated when the model parameterization is not adequate (here, 21 layers instead of 10). *reproduced from Sen and Stoffa (1991)*

Two decades later, Tran and Hiltunen (2011) presented an application of SA on 2-D fully elastic FWI, with mildly varying horizontal structures. Moving from 1-D to 2-D was made possible with a clever parameterization, achieved with less than a few tenths of parameters (from ten to forty). In their study, the authors empirically showed that the number of samples required to obtain the MAP estimate was two-hundred times the number of parameters. Datta and Sen (2016), use similar technique, with a fast SA algorithm, and a model parameterization based on 2-D nonoscillatory splines functions (Figure 1.6). The authors show that slightly complex geometries can be represented with a few nodal points and defined velocities values (less than 30 nodal points are used in this study), which greatly reduce the size of the solution space.

Aleardi and Mazzotti (2016), proposed to use GA along with a Markov-Chain Monte-Carlo Gibbs sampler, to make the best use possible of the numerous forward problem solves, that are required during the GA inversion. They point out a weakness that is common to the SA and GA methods above: they are bound to produce biased posterior covariance because they are not Markov Chain Monte Carlo methods. The reason they provide biased PDF comes from their tendency to oversample the regions of the model around the minimizers (Figure 1.7). By not discarding the models that are computed during the GA inversion, Aleardi and Mazzotti (2016) shows that one can use a (MCMC) Gibbs sampler on the stored model to produce an unbiased PDF To that extent, they generate a search space from a Voronoi tessellation of the solution space, from the samples drawn during the GA search. The Gibbs sampler is run across the tessellated space to compute the unbiased PDF.

Mazzotti et al. (2016) and Sajeve et al. (2017b) later proposed an extension of the GA + Gibbs sampler, generalizing its application 2-D problems. As many of the methods presented in this section,

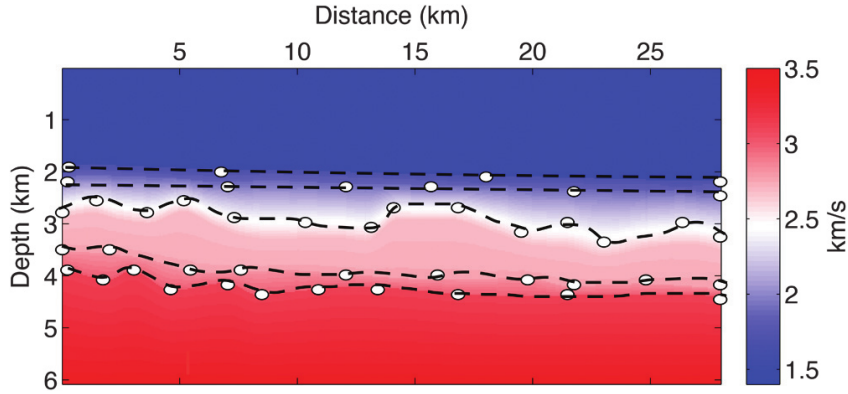


Figure 1.6: An example velocity model derived from the nonoscillatory spline function proposed in Datta and Sen (2016). The velocity model is parameterized over five interfaces. The nodes in each interface are shown with white circles whereas the dotted black line shows the interface after interpolation. The velocities between interfaces are estimated using linear interpolation. *From Datta and Sen (2016)*

their proposition relies on dimensional model reduction. Their model space is discretized over a coarse, irregularly spaced 2-D grid on which the inversion is performed. With this discretization, they were able to apply the GA + Gibbs sampler on up to 146 unknown parameters, in a fully acoustic FWI setting. The scheme is used in a very similar fashion, where the MAP estimate is given by the GA sampling, while the Gibbs sampler produces an unbiased posterior probability density function. Sajeve et al. (2017a) reviewed four classes of stochastic optimization methods: GA, SA, particle swarm optimization, and neighborhood algorithm on 1D elastic FWI. They concluded that GA was the best-suited approach, provided that the global minimum is not located in extreme values of the solution space (at the "edges" of Ω).

While the methods discussed above (from the "meta-heuristic" class) seems to be the prevalent strategy to achieve adequate posterior sampling, another class of method that is defined as purely "stochastic" has been growing in popularity for the past few years. The MCMC method (Metropolis et al., 1953; Hastings, 1970), is a class of sampling approach, akin to Monte-Carlo methods, which relies on a Markov-Chain which stationary law is the posterior PDF to be sampled. At each step, the MCMC algorithm proposes a new sample, conditional on the current one, and decide to accept or reject it according to a sampling criterion (Gallagher et al., 2009). Because this sampling strategy only takes into account the current step to produce the next sample, it does not have a memory effect and is often described as a "random walk" sampling (the sampling chain is said to satisfy the Markov property (Serfozo, 2009)). Basic MCMC algorithm is defined as a cycle of three steps:

- **Random walk in the model space.** A new sample is proposed by perturbing the current sample according to a prescribed prior probability. If Gaussian prior is used, the $k + 1$ sample perturbation will be drawn from a Gaussian PDF centered on the sample k .
- **Misfit function evaluation.** After the new sample is drawn, its data misfit is evaluated.
- **Metropolis-Hasting / Gibbs sampling.** The proposed sample is accepted or rejected, according to the probability that they can explain the data.

When a sufficient amount of samples have been accepted, the statistical property of the sampling population is computed, and typically yield the approximated posterior covariance.

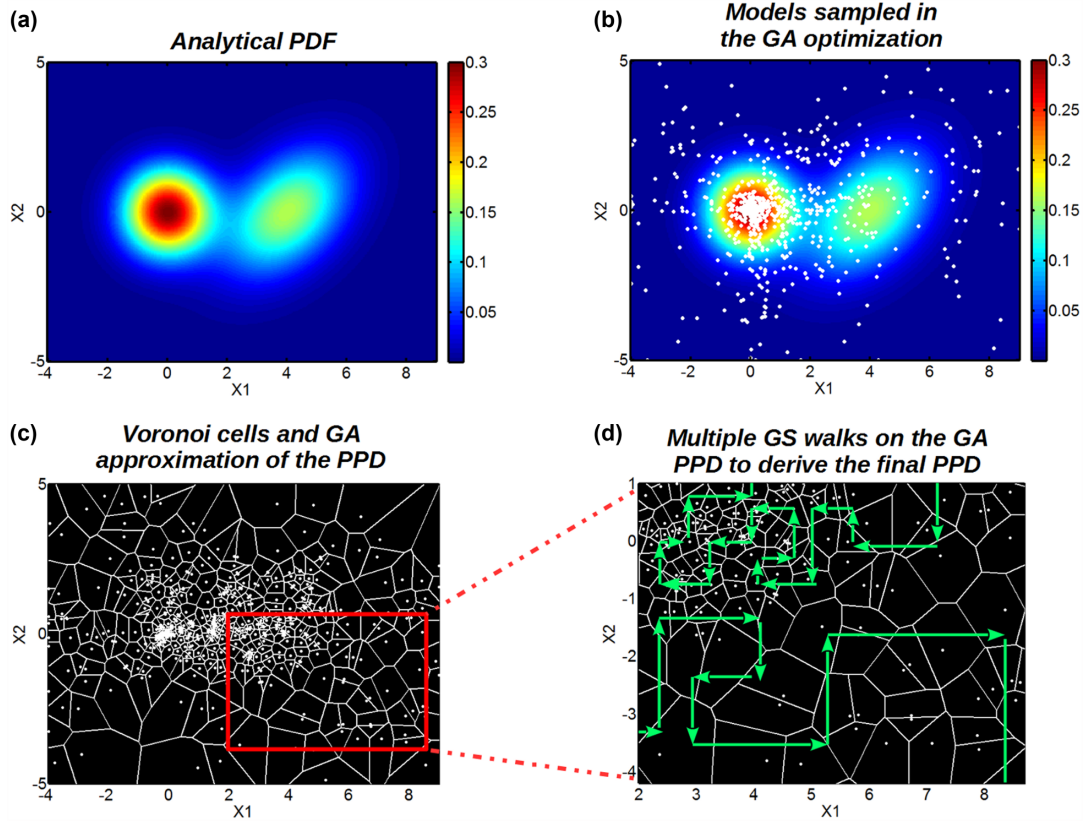


Figure 1.7: Computation of the PDF in GA + Gibbs sampler algorithm in Aleardi and Mazzotti (2016). (a) The initial analytical PDF used in the optimization, (b) the initial PDF and the 1000 samples (white dots) drawn during the GA search. The solution space has been preferentially sampled close to the two minima. (c) The solution space explored during the GA inversion is divided into Voronoi cells, with the GA samples as nuclei. The likelihood of each of the 1000 sample is associated with their respective Voronoi cell. (d) Illustration of Gibbs sampler paths that are used to draw the samples from GA that approximate the unbiased PDF *From Aleardi and Mazzotti (2016)*

While MCMC methods are popular within the geophysics community (to solve small problems), the cost of their application have been a significant limiting factor for their application in tomography, let alone FWI. Early MCMC tomographic application was proposed by Cary and Chapman (1988), on a 1-D problem, where traveltimes and waveform data were used. They focused on evaluating initial model candidates with the global search. Afterward, these candidate models were inverted with local linearized FWI. An early application on 2-D model space was performed by Vasco et al. (1993) in the context of borehole tomography, limited to traveltimes data, however. In 2-D domain, the first FWI posterior was realized by Cordua et al. (2012), also in a borehole setting (with cross-hole Ground Penetrating Radar observations), although they failed at demonstrating the applicability of the methodology on field data (where the structure of noise is not correctly characterized).

Facing difficulties to scale up the search methods, two major strategies emerged from the literature. Hybridization of meta-heuristic methods to accelerate the PDF sampling, with a subsequent MCMC, search to perform the unbiased uncertainty analysis (as presented previously) such as the MCMC algorithm proposed by Hong and Sen (2009). Some of these mixed applications have been discussed along with other meta-heuristic methods earlier in this section (Aleardi and Mazzotti, 2016; Mazzotti

et al., 2016; Sajeve et al., 2017b), but also fall in the category of MCMC samplers. The other strategy has been popularized by Bodin and Sambridge (2009), who introduced transdimensional methods to seismic tomographic imaging. The idea behind this application comes back to Sen and Stoffa (1991) observation about sub-optimal parameterization of the inverse problem. Because the number of unknown considered can significantly impact the performance of an inversion scheme, the principle of transdimensional inversion is to set the number of parameters as an unknown of the inversion problem. Depending on the target data, the transdimensional MCMC will adapt the number of parameters on which the inversion is carried. This generally allows us to consider a minimum amount of model parameters (as in the case of model reduction approaches), but also to find the optimal model discretization. The transdimensional decomposition and representation of the model space generally rely on Voronoi tessellation. At each sampling iteration of the transdimensional MCMC, the number of nuclei (and therefore of cells) can randomly be increased or decreased to fit the data better. While the original application of Bodin and Sambridge (2009) was performed with traveltime data, it has since been applied to the FWI problem in several instances: Ray et al. (2016) proposed the first application of transdimensional FWI with elastic, frequency domain data. They were able to produce a posterior sampling which MAP estimator fitted closely with previous results obtained with a deterministic method on their field-data example. Additionally, they showed the interest of the transdimensional approach by providing an evaluation of the number of interfaces to consider to fit their data adequately. Biswas and Sen (2017); Sen and Biswas (2017) have developed a two-fold approach, with transdimensional traveltime tomography, followed by an MCMC FWI performed in the best fitting model parameterization. The rationale for using traveltime tomography in the first place is the very cheap cost of misfit evaluation, which enables more robust exploration of the solution space. The outcome of the traveltime tomography MCMC constitutes a good starting point for the subsequent FWI MCMC, as it provides both the optimal starting model and the optimal model parameterization. With this application, it also appears that the Voronoi tessellation is an advantageous model parameterization. We see in Figure 1.8 that the recovered mean model is described with more discrete parameters than the sum of its parts: all the cells that are not entirely overlapping create additional discrete-points when averaged over a Cartesian grid.

The particularity of the MCMC sampling used in Biswas and Sen (2017); Sen and Biswas (2017), is the uses of Hamiltonian dynamics to drive the sampling stage effectively (Duane et al., 1987). The so-called Hamiltonian MCMC (H-MCMC), has been recently receiving much attention from the community (Fichtner et al., 2018a,b; Gebraad and Fichtner, 2018) as it seems to be a promising technique to perform fully stochastic sampling. The advantage of H-MCMC over traditional MCMC methods is that it allows using the misfit function's gradient to accelerate the sampling speed: instead of a "step-by-step" random walk on Ω , the H-MCMC sampler is affected by "momentum" and travel on Ω along paths defined by the local slope.

All the methodologies that have been presented in this segment display an evident shortcoming: the inability to produce high-resolution models within the Bayesian framework. The computational cost of global searches requires model parameterization strategies, which ultimately reduces one capacity to obtain a high-resolution model (while it is the primary appeal of FWI). As for now, the interest of global search seems limited to initial models building: exploring a coarse version of the space can help identify all the potential initial models that are fitted for linearized FWI.

Until proven wrong, it seems that the curse of dimensionality (Bellman, 2015) is too harsh of a constraint to allow global search approaches to reach the high-resolution capacity of local FWI. Trading resolution power for uncertainty estimation does not seem to be interesting in the context of FWI (but could be appropriate for traveltime tomography). Following the presentation of the global search approaches for uncertainty estimation in FWI, I will now focus on local uncertainty estimation methods

Trans-dimensional 2D FWI using RJHMC

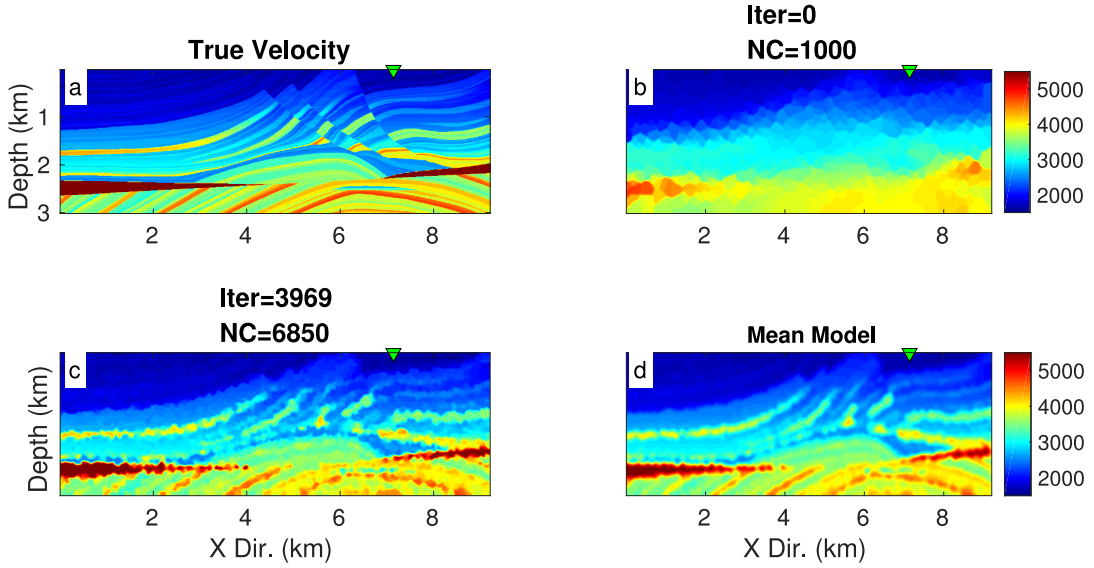


Figure 1.8: Results from transdimensional H-MCMC in Biswas and Sen (2017), at different iterations. (a) The true velocity model. (b) The initial velocity model, represented as 1000 nuclei (NC). The search in Ω is initialized at a smoothed version of the true model. (c) Sample model after 3969 evaluations, and the number of nuclei increased to 6850. (d) The final mean model, averaged over the 3969 selected samples from the transdimensional H-MCMC algorithm. Note that when average over a cartesian grid, the number of distinct values increases from the final number of nuclei ($n_{mean} > NC_{final}$ Modified from Biswas and Sen (2017)).

in the following subsection.

1.2.2 Local uncertainty estimation in FWI

The other family of uncertainty estimation is set in the local search paradigm. I have shown the clear relationship that exists between the posterior covariance matrix close to the MAP, and the local inverse Hessian operator in the vicinity of a minimum in Section 1.1.2. While this relation is only valid in the local quadratic assumption, it has allowed the development of local uncertainty estimation method. These "Hessian based" uncertainty estimation are generally performed in two steps:

1. **Performing a standard FWI.** From a good initial model m_0 , an iterative linearized least-squares optimization scheme is performed, and the model is improved with the means of a gradient or (more likely) a Newton-type optimization method.
2. **Local evaluation of the Hessian or inverse Hessian operator.** Here the Hessian or preferably its inverse is approximated one way or another, to infer the statistical properties of the inversion outcome. This step is highly dependent on the choice of m_0 , and therefore on the subset of Ω in which the misfit is minimized.

Before discussing the methods which have been proposed in the literature, we can underline an important fact: the estimated uncertainty estimation is highly dependent on the inversion hyperparameters. This

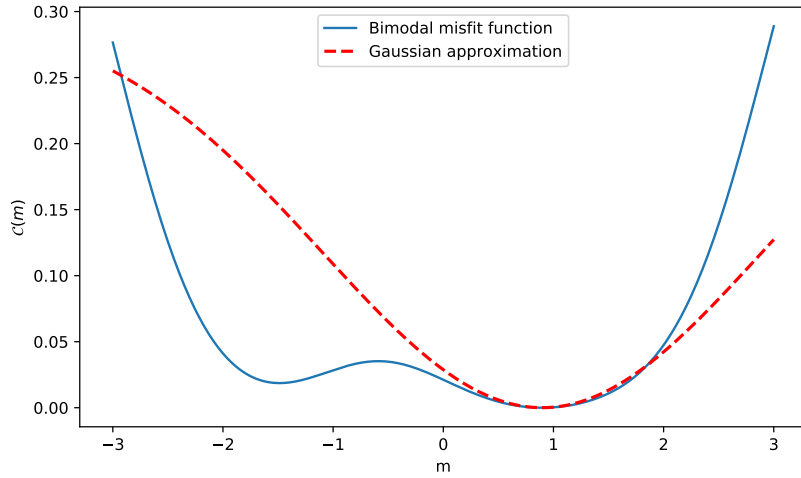


Figure 1.9: Schematic representation of the Gaussian approximation (red dashed line) of the misfit function (solid blue line). The goal of local uncertainty estimation in this context is to define the best fitting Gaussian PDF that represent the lateral variations of the misfit close to the optimum.

was already the case within the global search paradigm (where the model parameterization constrained the solution space), but while global search makes a compromise on the resolution to obtain uncertainty estimation, local approaches tend to favor high-resolution models with limited uncertainty estimation. Because the local uncertainty estimation methods are based on a strong quadratic assumption, they are limited to an evaluation of uncertainty within that approximation (Figure 1.9). Within this local paradigm, it is obvious that all of these methods will fail at evaluating the occurrence of local minima in the misfit function. At best, we can sample uncertainty information in the vicinity of the solution given by local optimization, from which we can infer information about the solution quality: correlation and cross-talk between parameters, absolute variance.

Historically, the first attempt at qualifying FWI result quality, based on second-order information of the misfit function has been proposed by Fichtner and Trampert (2011b). While their purpose was not to evaluate uncertainty *per se*, their proposition has to be acknowledged as the first direct attempt to exploit Hessian information to evaluate solutions quality in FWI. In this work, they propose an evaluation of local resolution based on Point-spread Functions (PSF), defined by the Hessian: recalling that \mathbb{H} acts as a convolutional filter, it is possible to evaluate local resolution by solving the product of \mathbb{H} with a point localized perturbation (typically a Dirac δ -function). The Hessian will spread the point-mass of the δ -function over the space, from which we can infer the resolution information contained in the Hessian. From field data experiments, they computed PSF by evaluating lines of the Hessian operator $\mathbb{H}(m, m_i)$, with m_i being the i^{th} parameter of the domain m . They observed bell-shaped PSF, centered on m_i , from which they derived a Gaussian approximation of the PSF. Because they treat the lines of the Hessian as a Gaussian PDF they can easily compute their Fourier transform, which enables to inspect model resolution in the Fourier domain. The limitations of this methodology are: 1) it is unable to provide an approximation of the inverse Hessian, 2) it can only estimate the Fourier spectrum of one line of the Hessian at a time. Based on this approach, a variant was proposed by Fichtner and van Leeuwen (2015). They also considered the effect of the Hessian to be the same as a Gaussian smoothing kernel, from which the spatial extents can be estimated with autocorrelations of the smoothed signals (Figure 1.10). Both Fichtner and Trampert (2011b) and Fichtner and van Leeuwen (2015) were derived from

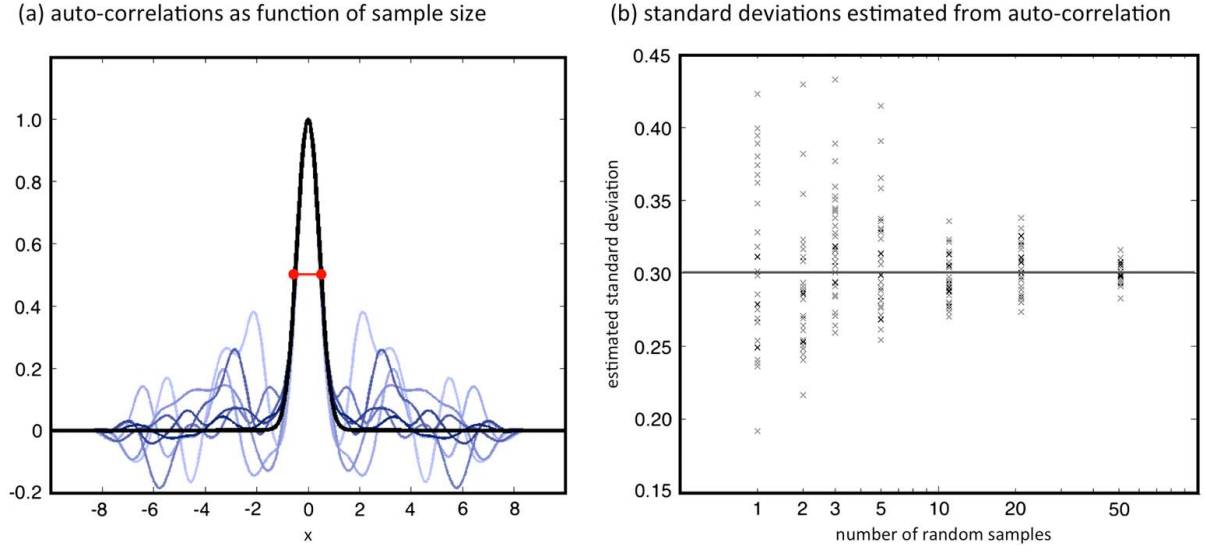


Figure 1.10: Estimation of the spatial extent of the Gaussian PDF that approximates the convolutional effect of the Hessian. *From Fichtner and van Leeuwen (2015)*

strong assumptions on the structure of the Hessian operator, and are not able to approach the posterior covariance matrix that holds uncertainty information. Instead, their problem comes down to evaluating the spatial extent of a Gaussian smoothing kernel.

As with global optimization methods, to mitigate the cost of the uncertainty estimation, a straightforward idea would be to reduce the size of the solution space, which in turns would reduce the size of \mathbb{H} and \mathbb{H}^{-1} . Following this idea, Du et al. (2012) and later Jordan (2015) have proposed a model parameterization based on B-Spline functions to achieve space reduction, which alleviates the computational burden on the estimation of the covariance \mathbf{P} . With the model size reduction, the covariance matrix can be explicitly computed, which yields an uncertainty estimate on a coarse grid. Note that even though this method has been demonstrated on a stereotomographic application, this formulation of the problem can easily be ported to FWI. However, reducing the size of the solution space in local optimization approaches does not seems to be an appealing strategy: 1) as with all local optimization strategy, the quadratic assumption of the misfit function is fairly reductive 2) the B-Spline representation would prevent reaching high-resolution results.

The second possibility is to take advantage of the data sparsity that is generally encountered in tomographic applications. As we can expect the illumination of the target subsurface to be sparse and incomplete, the Hessian operator is often rank-limited: it is, therefore, possible to project the Hessian operator on a smaller basis on which it is well defined, to compute its pseudo inverse. Bui-Thanh et al. (2013) propose to do exactly this, as the authors take direct advantage of the rapidly decaying spectrum of the Hessian matrix in their global tomographic application to approximate a low-rank Hessian. Based on previous work that established the compact nature of the Hessian for a smooth medium (Bui-Thanh and Ghattas, 2012a,b), they developed a scalable method that allow to compute a prior-preconditioned inverse Hessian, and draw samples from the posterior to evaluate pointwise statistics of the solution (Figure 1.13).

To approximate the low-rank Hessian and its inverse, they rely on matrix-free Lanczos iterations, which requires a resolution of the forward problem at each iteration. The nice perk of this methodology is that the spectrum of the Hessian can be built iteratively, and stopped once a criterion is reached to

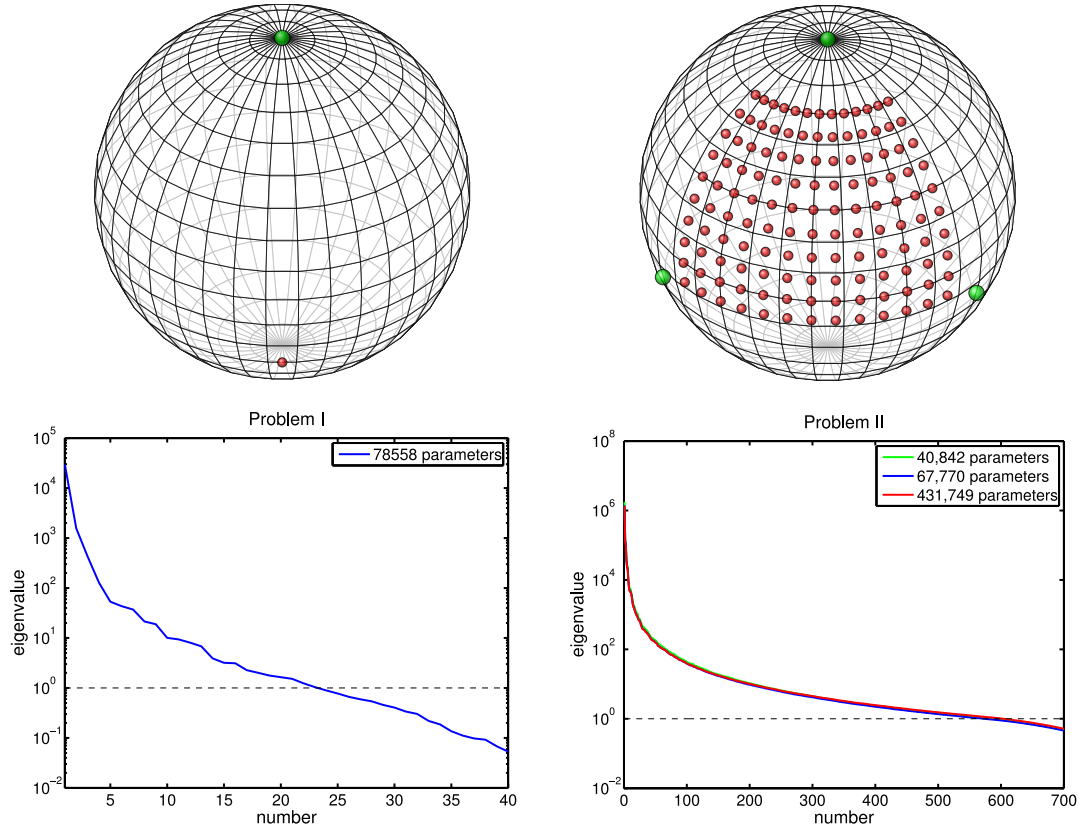


Figure 1.11: Spectrum of the Hessian operator for 2 Global FWI application. Problem I is defined with a single source (green) and receiver (red) couple. The corresponding spectrum of the Hessian displays a rapid decay and allows to be represented with less than 25 singular values. The dashed line denotes the truncation threshold. The second problem has significantly better coverage (3 sources and 130 receivers), and thus has a far better-posed Hessian matrix. As a result, it requires more singular values (at least 700) to be approximated correctly.) *Modified from Bui-Thanh et al. (2013).*

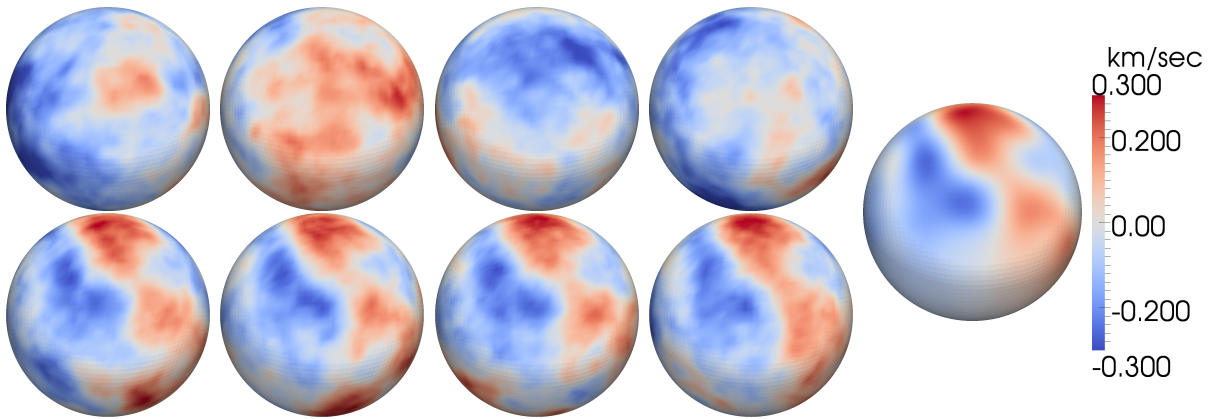


Figure 1.12: Samples drawn from the prior (top) and posterior (bottom) distributions of Problem II (see Figure 1.11). The optimal model obtained with linearized least-squares inversion is shown on the right for reference. *From Bui-Thanh et al. (2013).*

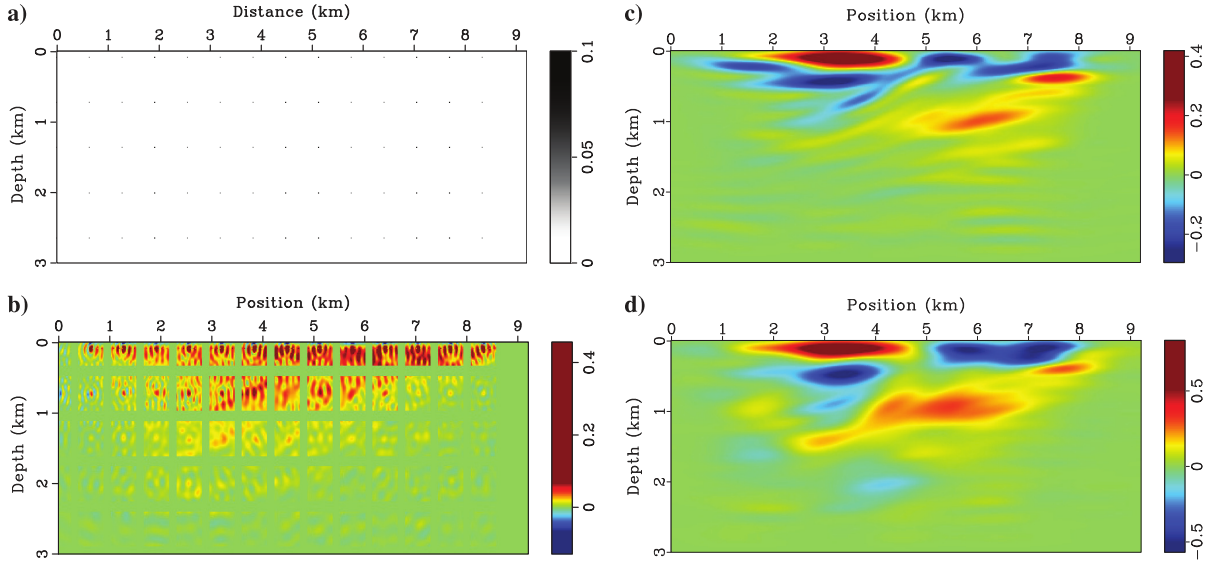


Figure 1.13: Point spread function approximation of the Hessian. (a) Set of Dirac δ -function from which the PSFs are evaluated. (b) Local subsets of the PSFs from the convolution between the Hessian and the δ -functions. (c) Product of a random vector with the true Hessian. (d) Product of the same random vector with its PSF approximation. *From Zhu et al. (2016).*

avoid computation of lower-order singular values. They reduce the uncertainty estimation problem to solving forward and adjoint simulation of the wavefield to build an approximation of the Hessian, once the MAP is reached. Because the spectral content of their prior-preconditioned Hessian matrix is strictly dependent on the acquisition, their method is not sensitive to the number of discrete grid points on which the model is represented. This leads them to obtain an approximate posterior over $n = 430000$ model parameters.

Following the same leading idea of approximating the Hessian with its truncated version, Zhu et al. (2016), Eliasson and Romdhane (2017), and Liu and Peter (2019) proposed similar approaches. Instead of using the Lanczos method to approximate the Hessian, they rely on the randomized singular value decomposition (rSVD; Liberty et al., 2007; Halko et al., 2011) which is also designed to approximate matrices of low-rank. In their development, Zhu et al. (2016) directly re-adapt the formulation of Bui-Thanh et al. (2013). As with the matrix-free Lanczos method, the rSVD requires to perform Hessian-vector product at each iteration. Based on the assumption that the convolutional properties of the Hessian do not change too rapidly across the model space, they approximate the Hessian with PSFs. By interpolation of these few PSFs, they were able to reduce the required number of Hessian-vector products to perform the rSVD and speed up their uncertainty estimation algorithm.

Additionally, it is interesting to note that the low-rank Hessian formalism can be extended to different flavors of FWI. For example, Fang et al. (2018) have applied this low-rank approximation strategy on the Wavefield Reconstruction Inversion (WRI). WRI is a waveform inversion method that relies on a specific regularization term controlled by the wave equation, which effectively changes the nature of the misfit functional. They have demonstrated that thanks to the relaxation afforded by the WRI, which further convexify the misfit function, their application is particularly adequate to the Gaussian approximation of the posterior distribution. In their proposition, they derived a factored formulation of the Gauss-Newton Hessian which allows drawing samples from the posterior without additional forward problem solves (once the square root of the Gauss-Newton Hessian has been explicitly build and stored).

This goes to demonstrate that with an appropriate formulation, provided the alternative misfit function can be locally approximated under a quadratic form, the Hessian-based approaches can be used for uncertainty estimation. Alternatives to the Euclidean norm, such as the envelope misfit function (Bozdağ et al., 2011) or the recently proposed Wasserstein distance or Graph-space distance (in the context of Optimal Transport optimization (Engquist and Froese, 2014; Métivier et al., 2016; Yang et al., 2018; Métivier et al., 2019)) could greatly benefit from this kind of development and offer uncertainty estimation on top of their reduced sensitivity to local minima.

I have introduced local uncertainty estimation method in this short review, which are all based on low-rank approximation of the Hessian matrix to produce uncertainty estimation. As with the global search approaches, all of these methodologies are ultimately limited by one's ability to harness their computational cost and solve the many forward simulations required to solve both the optimization problem and the subsequent uncertainty estimation scheme. One of the main difficulty for the application of these approaches is the intrinsic iterative nature of the Lanczos and rSVD methods, which prevent local uncertainty estimation from being implemented in a fully scalable way.

As a final note of this review, I would like to re-emphasize the very local nature of the Hessian-based uncertainty estimation schemes: In the methods above, the small singular values of the Hessian are omitted in the construction of the pseudoinverse Hessian. These values correspond to model parameters that are either in the nullspace of our inverse problem or are poorly constrained by the data (Kalmikov and Heimbach, 2014). Because the spectrum of the inverse Hessian represents parameters uncertainty, the omitted singular values of the Hessian correspond to substantial parameter uncertainty. Thus, the truncated pseudoinverse Hessian produced by these methods only measures uncertainty in the data-constrained subdomain, and discard any information related to the nullspace.

From this review, it is evident that the uncertainty estimation literature is still in its infancy in the domain of FWI. While global and local uncertainty estimation methods have their advantages and limitations, they are the two favored ways of conducting uncertainty estimation nonetheless. Still, their computational intensive nature and limited scalability have prevented a large adoption of these cutting-edge methodologies: uncertainty estimation is not systematic in FWI.

1.2.3 Ideas from the Data Assimilation community

While the FWI community is developing cutting-edge techniques for uncertainty estimation, the Data Assimilation (DA) community had concurrently a long-lasting history of uncertainty estimation paired with complex geophysical systems study. Their methods have demonstrated a particular ability to solve inverse problems with a large number of degrees of freedom, high degree of complexity, and data sparsity while integrating uncertainty quantification within their inverse problem-solving schemes. Their systematic uncertainty estimation would be highly desirable in the context of FWI.

Generally, the overall goal of DA in geophysical applications is to characterize the state of a dynamic system through time, which can be subjected to non-linear dynamics by combining sparse observations and numerical models. DA has notably been successfully implemented at operational scales in areas such as numerical weather forecasting, oceanography, reservoir characterization, and climatology (Rodell et al., 2004; Navon, 2009; Cosme et al., 2010; Lee et al., 2016).

Because DA tools can handle large-scale non-linear problems, as it is the case with the FWI problem, we might be able to take advantage of the DA formalism to bring a new look at uncertainty estimation in FWI. Applying *ensemble-based* Data Assimilation to geophysical tomography has already started being investigated. Indeed Jin et al. (2008) propose using the Ensemble Kalman Filter (see 2.2) to solve 1D

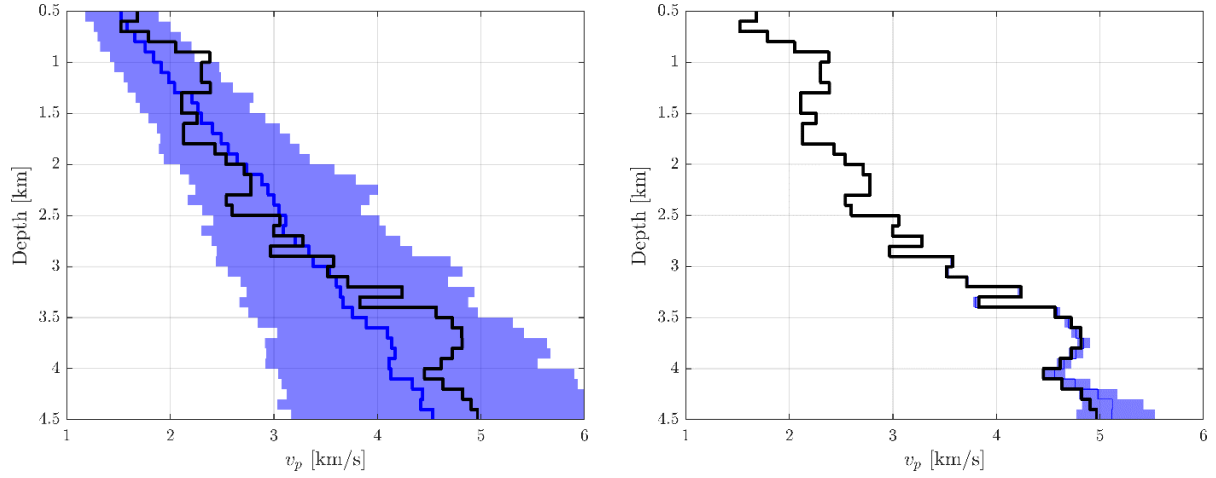


Figure 1.14: 1-D FWI application solved with the I-EnKS. The prior (left) and posterior (right) ensembles are displayed. The ensemble spread (shaded blue), and its mean (solid blue line) are the solution yielded by this ensemble method. The true model is depicted as a solid black line. This inversion yield a significant uncertainty reduction in the upper part of the domain, while the deeper part (poor illumination) displays greater uncertainty. Note also that in this application, the I-EnKS has been able to recover the true values of the velocity structure in most of the domain. *From* Gineste et al. (2019).

prestack FWI. Gineste and Eidsvik (2017) and later Gineste et al. (2019) proposed to use the Ensemble Kalman Smoother (Evensen and van Leeuwen, 2000) and the Iterative Ensemble Kalman Smoother (Bocquet and Sakov, 2014) to inverse 1D velocity profiles. Their results pictured in Figure 1.14 showed that with a very broad non-informative prior, the Iterative Ensemble Kalman Smoother was able to solve their synthetic inverse problem, while also providing a confidence interval for their solution. Liu and Grana (2018) propose to use the Ensemble Kalman Smoother to inverse jointly elastic and petrophysical rock properties in the context of reservoir monitoring. This also goes to show that the DA framework is well suited for multiparameter inversion, which is a requirement for serious FWI applications.

This review set the stage for the core work behind this thesis: an original DA-FWI scheme adapted to take advantage of both worlds. By combining a classical FWI Newton-type solver and ensemble filtering, we ought to perform FWI with an embedded uncertainty estimation solution. The linearized least-squares part is expected to provide a quick convergence rate to the MAP, along with high-resolution models. The DA part, on the other hand, is expected to provide uncertainty estimation in a Bayesian framework, based on the local quadratic assumption of the local search. This will make for a new breed of local optimization, taking advantage of the uncertainty estimation framework of DA tools.

Conclusion

In this chapter, we have introduced the general aspects of the FWI problem. Both the forward and inverse problems have been defined. It allowed us to present how the forward problem is solved in the frame of finite-difference frequency-domain FWI. Then, we gave a complete overlook of the global and local scheme one can use to solve the FWI problem. The advantages and limitation of both paradigms have been highlighted, and special attention has been given to the resolution of FWI as a linearized least-square problem. From there, the definition of the gradient and the Hessian have been given, and their link with

uncertainty estimation established. This theoretical section was followed by a comprehensive review of existing uncertainty quantification methods in FWI. We highlighted the characteristics of most of them and highlighted a few fundamental limitations that prevent their uses as a systematic solution to the uncertainty problem. Finally, a short introduction was given on DA methods, and their take on uncertainty estimation, as a general trend, but also by presenting a few novel application to FWI, leading us to the original contribution of this thesis work.

In the next chapter, I will provide an extensive theoretical development of the methodologies underlying DA tools, for the readers that might be unfamiliar with the DA literature. From the simple Kalman Filter, and its relation the Bayesian least-squares problem, I will move to more sophisticated schemes such as the Ensemble Transform Kalman Filter, that will be implemented along with FWI later in Chapter 3.

Chapter 2

Data Assimilation

Contents

2.1	Elements of Data Assimilation	36
2.1.1	Defining the system state	36
2.1.2	Observations	37
2.1.3	A practical example of statistical estimator	38
2.1.4	The dynamical model - forecasting stage	41
2.1.5	The Kalman Filter	42
2.1.6	Extended Kalman Filter	46
2.2	Ensemble Kalman Filter	47
2.2.1	Ensemble representation	48
2.2.2	EnKF's forecast step	49
2.2.3	Analysis step	50
2.2.4	Ensemble Transform Kalman Filter	53
2.2.5	Maximum-Likelihood Ensemble Filter	58
2.2.6	Ensemble Kalman Inverse	60
2.3	Limits of EnKF methods	62
2.3.1	Undersampling characterization	62
2.3.2	Inbreeding	63
2.3.3	Filter divergence	64
2.3.4	Spurious Correlation	64
2.3.5	Solutions to undersampling	65

At first glance, the common-ground between seismology, seismic tomography, and Data Assimilation (DA) might seem to be limited. DA is a general framework that is commonly used to predict and update one's belief on a physical system, and it is deeply rooted in probability theory, statistics, and even control theory in some instances. However, the shrewd eye will see through the complicated jargon and terminology of DA, an apparent connexion with inverse problem theory.

In this chapter, I wish to establish this connexion clearly: I will present the theoretical foundations that lead to the development of the most common tools in DA: from a basic estimator to the well established Kalman Filter and some of its variants. As stated in the previous chapter, this work investigates the

possibilities to couple DA and FWI together as a unified framework for uncertainty quantification; therefore, a comprehensive introduction to the DA literature is needed to establish a clear connection between these two fields.

This chapter owes much to the textbooks of Evensen (2009); Fletcher (2017), both of which are wonderful introductions to the broad topic of DA. The work of Harlim and Hunt (2005) and Hunt et al. (2007) also deserves special recognition, as it shines with clarity and exposes complex principle with a simplicity that can only be appraised. Finally, I want to acknowledge the outstanding work of Labbe (2016) whose book allowed me to build solid intuitions on the concepts of Bayesian filtering and get into the realm of DA through the right door.

2.1 Elements of Data Assimilation

I will begin by introducing the nomenclature on some key objects that are necessary to grasp the DA literature and concepts, notably the system state, the forecast, the observation vector, and the analysis state.

2.1.1 Defining the system state

The goal of DA is to characterize the *state* of a system m at any given discrete time k . In the vast majority of cases, the objects of study in DA applications are continuous systems (the Earth's atmosphere and ocean, the position of an object in a continuous space, or the properties of the subsurface, for instance). As with numerical optimization, it is practical to represent the system in the "model" space: a discrete representation of the continuous system over n points that allows its numerical manipulation. In DA, the discrete version of the continuous system is generally defined as the *state vector*

$$m_k \in \mathbb{R}^n, \quad (2.1)$$

where \mathbb{R}^n is sometimes referred to as the *state space*, and the subscript k expresses the state vector "at time k ".

The readers accustomed to DA literature might be surprised by this choice of nomenclature instead of the classical x_k that has been defined in the seminal review by Ide et al. (1997). The choice of m_k over x_k is purposely made to underline the link between numerical optimization applied to FWI and DA. This choice will also become obvious later with the theoretical developments in Chapter 3, where DA and FWI will be closely intertwined.

The state vector can hold two types of *state variables*. The state variables that are directly measurable with a sensor are called *observed variables*, while those that are inferred from the observed variables are called *hidden variables*. Let us consider a simple example to illustrate the differences between observed and hidden variables (that example will be handy to explain other concepts later in this section): we consider an object moving at an unknown constant velocity v along a 1-D axis, which position x is measured periodically. The position is directly measured: it is an observed variable. The object's velocity can be inferred from consecutive measurements and is, therefore, a hidden variable.

The system is defined as the moving object; its position and velocity define its state. From here, we can design several DA problems where one can estimate the evolution of:

- the objects' position,

- the objects' velocity,
- the objects' position and velocity altogether.

To these problems, we can associate three different state parameterizations:

$$m_k = [x_k], \quad m_k = [v_k] \quad \text{and} \quad m_k = \begin{bmatrix} x_k \\ v_k \end{bmatrix}. \quad (2.2)$$

While the first problem would be a pure estimation or filtering problem, the second and third are conceptually identical to solving inverse problems in which the velocity would be inferred from successive measurements of position. Now that the state vector has been defined, we can move on to the first obvious source of information on the system state: observations. For the remainder of this chapter, we consider the joint state vector (third case) as a simple example to illustrate various concepts of DA.

2.1.2 Observations

We define an observation vector $y_k^o \in \mathbb{R}^d$, which is the counterpart of our state vector,

$$y_k^o \in \mathbb{R}^d, \quad (2.3)$$

where the superscript o stands for "observation" and \mathbb{R}^d is generally referred to as the *observation space*. As with the state variables, the observation can cover a broad range of measurements of the system. The observation vector can contain measures of different physical attributes and come from very different sensors (various sampling rates, precision, units). Furthermore, observations can also be localized in the system space (*in-situ* measurements) or be non-local (measures at the boundaries of the system), and the length of y_k^o might change over time.

The relation that ties the state and the observation space is contained in an observation operator $\mathbf{H}_k : \mathbb{R}^n \rightarrow \mathbb{R}^d$ that yields

$$y_k^o = \mathbf{H}_k m_k. \quad (2.4)$$

In most realistic applications however, the data vector might be affected by observation errors ϵ such that equation (2.4) becomes

$$y_k^o = \mathbf{H}_k m_k + \epsilon_k. \quad (2.5)$$

Such noise can stem from the sensor precision, the level of background noise, or general uncertainty surrounding the physical state of the measurement device (uncertainty on its location, for instance). The measurement noise is often modeled as a random variable which mean $\bar{\epsilon}_k$ and covariance \mathbf{R}_k are given by

$$\bar{\epsilon}_k = \mathbb{E}[\epsilon_k], \quad (2.6)$$

and

$$\mathbf{R}_k = \mathbb{E}[(\epsilon_k - \bar{\epsilon}_k)(\epsilon_k - \bar{\epsilon}_k)^T], \quad (2.7)$$

where $\mathbb{E}[\cdot]$ designates the expectation (the average value of the random variable given an infinite amount of realizations). It is common practice to consider the *unbiased* observation errors (Asch et al., 2016), defined as a zero-mean Gaussian random variable following $\epsilon_k \sim \mathcal{N}(0, \mathbf{R}_k)$. Coming back again at our toy example, the discrete observation operator \mathbf{H}_k is given by

$$\mathbf{H}_k = \begin{bmatrix} 1 & 0 \end{bmatrix}, \quad (2.8)$$

and the observation at each timestep k are given by

$$y_k^o = \mathbf{H}_k m_k + \epsilon_k. \quad (2.9)$$

With the introduction of both the state vector and the observation vector, we are left with two different pieces of information to estimate the true system state. From an "initial guess" on the system state, we can produce an *analysis*: the estimate of the system state based on a clever combination of observational data and initial (*background*) information. Common sense would tell us that the true system state lies somewhere "in-between," and should be a tradeoff between observation and background information: if not, there is something utterly wrong with one of them! The next subsection details how this "clever" combination is made and introduce the theory behind the DA analysis.

2.1.3 A practical example of statistical estimator

Based on our toy example, let us design a practical estimation problem. At a fixed time k , I believe the object's position to be x^b (where b stands for background), but the observation x^o differs from my belief. Recalling that the solution should lie somewhere "in-between," the most naive solution to the problem would be to consider a linear combination of my initial guess and the observation

$$x^a = x^b + K(x^o - x^b), \quad (2.10)$$

where x^a denote the *analysis state* that I wish to estimate and K is a scaling factor satisfying $0 \leq K \leq 1$. To compute the "optimal" value of K , it is useful to consider the (unknown) true state x^t (where the superscript t stands for "true") and compute the different errors that are associated to this estimation problem. From

$$x^a - x^t = x^b - x^t + K(x^o - x^t - x^b + x^t), \quad (2.11)$$

we define the following errors

$$\alpha = x^a - x^t, \quad \beta = x^b - x^t \quad \text{and} \quad \epsilon = x^o - x^t, \quad (2.12)$$

where α is the analysis error, and β and ϵ are respectively, the background and the measurement error. We obtain

$$\alpha = \beta + K(\epsilon - \beta). \quad (2.13)$$

This yields

$$\mathbb{E}[\alpha] = \mathbb{E}[\beta] + K(\mathbb{E}[\epsilon] - \mathbb{E}[\beta]) \quad (2.14)$$

that can be reduced to

$$\mathbb{E}[\alpha] = \bar{\alpha} = 0 \quad (2.15)$$

provided that errors are zero-mean. Recalling the variance of a random variable can be expressed as $\sigma^2 = \mathbb{E}[(\epsilon - \bar{\epsilon})^2]$, we can express the analysis error variance as

$$\begin{aligned} \sigma_a^2 &= \mathbb{E}[\{\beta + K(\epsilon - \beta)\}^2] \\ &= \sigma_b^2 + 2K\mathbb{E}[\beta(\epsilon - \beta)] + K^2\mathbb{E}[(\epsilon - \beta)^2] \\ &= \sigma_b^2 + 2K\mathbb{E}[\beta\epsilon] - 2K\sigma_b^2 + K^2(\sigma_o^2 + \sigma_b^2) - 2K^2\mathbb{E}[\epsilon\beta] \end{aligned} \quad (2.16)$$

where σ_a , σ_b and σ_o are respectively the analysis, background and observation variances. If we consider that the forecasting error has no correlation links with the observation error, we can rewrite equation (2.16) as

$$\sigma_a^2 = \sigma_b^2 - 2K\sigma_b^2 + K^2(\sigma_o^2 + \sigma_b^2). \quad (2.17)$$

To find the best estimate of the system state, we need to find the value of K that minimizes the analysis variance, in turn minimizing the analysis error. The minimum variance is given by

$$\frac{d\sigma_a^2}{dK} = 0 = 2K(\sigma_o^2 + \sigma_b^2) - 2\sigma_b^2 \quad (2.18)$$

and therefore the optimal scaling factor K^* that minimizes the analysis error is solely given by the ratio of observation error over forecast error

$$K^* = \frac{\sigma_b^2}{\sigma_o^2 + \sigma_b^2} = \frac{1}{1 + \sigma_o^2/\sigma_b^2}, \quad (2.19)$$

The optimal state estimate given measurement and forecast error is given by

$$x^a = x^b + \frac{1}{1 + \sigma_o^2/\sigma_b^2}(x^o - x^b). \quad (2.20)$$

This estimator is classically called BLUE, for Best Linear Unbiased Estimator, which yields the optimal weighting for the linear combination of two independent pieces of information (and considering their respective errors to be zero-mean noise) (Asch et al., 2016). The BLUE can be intuitively interpreted with respect to the estimator errors:

- if we know the observations to be superior to the numerical model ($\sigma_o^2 \ll \sigma_b^2$), $x^a \approx x^o$
- if we know the numerical model to be superior to the observations ($\sigma_b^2 \ll \sigma_o^2$), $x^a \approx x^b$
- if both are strictly equivalent ($\sigma_b^2 = \sigma_o^2$), then x^a is nothing more than the arithmetic mean between the measured and forecast information.

Note that from here, we can also express the analysis error covariance in terms of the gain matrix K

$$\begin{aligned} \sigma_a^2 &= \frac{\sigma_b^2 \sigma_o^2}{\sigma_o^2 + \sigma_b^2} \\ &= (1 - K^*)\sigma_b^2 \end{aligned} \quad (2.21)$$

Where we can see that the analysis yields an update of our prior belief in the Bayesian paradigm (we update both our belief on the state of the system, but also the uncertainty on the state estimate).

In the case where both the background and measurement errors can be modeled as unbiased Gaussian noise ($\beta \in \mathcal{N}(0, \sigma_b^2)$ and $\epsilon \in \mathcal{N}(0, \sigma_o^2)$), the BLUE solution is also the minimum mean square error (MMSE) solution. The analysis is thus given in terms of the mean and variance of the posterior Gaussian distribution.

$$\mathcal{N}(x^a, \sigma_a^2) = \mathcal{N}(x^b, \sigma_b^2) \times \mathcal{N}(x^o, \sigma_o^2) \quad (2.22)$$

which makes the analysis a "Gaussian update" illustrated in Figure 2.1. Indeed, the product of two

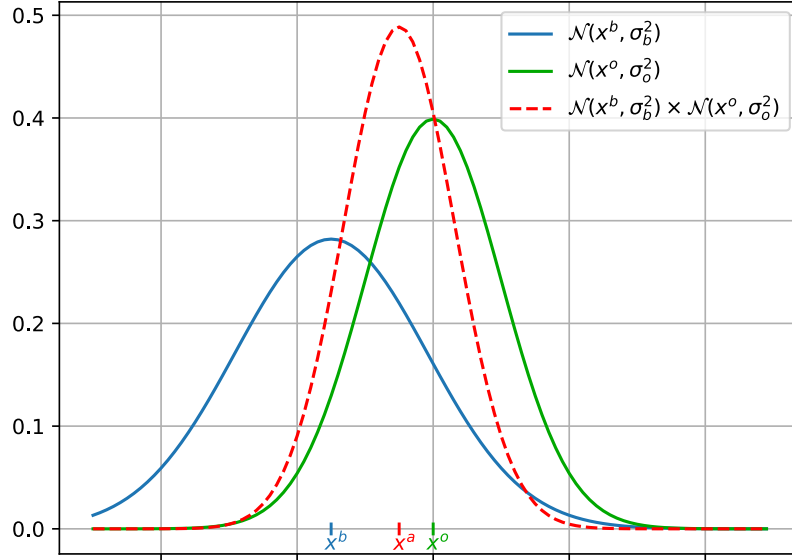


Figure 2.1: Numerical example of the BLUE in the Gaussian case. Both the background (blue) and the observation (green) are modeled as Gaussian random variables which PDFs are denoted by solid lines. The posterior distribution of the analysis (also a Gaussian) is obtained from the product of the background and observation PDFs and is denoted by the red dashed line. Note that because the observation is more reliable than the background state (smaller variance), the analysis is closer to the measurement than the initial guess.

independent gaussian with respective means x^b and x^o and variances σ_b^2 and σ_o^2 yields

The mean

$$\begin{aligned}
 x^a &= \frac{\sigma_b^2 x^o + \sigma_o^2 x^b}{\sigma_b^2 + \sigma_o^2}, \\
 &= \left(\frac{\sigma_b^2}{\sigma_b^2 + \sigma_o^2} \right) x^o + \left(\frac{\sigma_o^2}{\sigma_b^2 + \sigma_o^2} \right) x^b, \\
 &= K x^o + (1 - K) x^b, \\
 &= x^b + K(x^o - x^b),
 \end{aligned} \tag{2.23}$$

and variance

$$\sigma_a^2 = \frac{\sigma_b^2 \sigma_o^2}{\sigma_b^2 + \sigma_o^2}$$

which are both consistant with the equations of the BLUE.

This gives us the general solution for a linear Gaussian estimation problem. The principle of DA is to generalize the BLUE to dynamical systems, such that DA is defined by Asch et al. (2016) as "An analysis that combines time-distributed observations and a dynamic model". The final step toward defining DA is to introduce the dynamical model.

2.1.4 The dynamical model - forecasting stage

Knowing and predicting how the system behaves provides a great deal of information, which is important to consider when making state estimation. Thereby, the DA formulation integrates numerical forecasting models, or forecasting operators, at the core of the system state estimation, in order to complement the observations. These forecasting operators are generally built on the equations (or their approximation) that govern the evolution of the system with time: in the case of our 1-D moving object example, they would be Newton's equations of motion. Applied to the state vector at time k , they produce a forecast state at time $k + 1$ following

$$m_{k+1}^f = \mathbf{F}_k(m_k) \quad (2.24)$$

where the superscript f denotes the *forecast* and $\mathbf{F}_k : \mathbb{R}^n \rightarrow \mathbb{R}^n$, is the forecast operator.

In the (frequent) case where the model cannot adequately predict the evolution of the system state, equation (2.24) becomes

$$m_{k+1}^f = \mathbf{F}_k(m_k) + \eta_k \quad (2.25)$$

where η_k designates the *process noise*. Process noise can arise from the mismatch between the physical world and its mathematical representation. They can also appear when the correct set of equations governing the system's evolution may not be known or may only be approximated. The process noise is defined as a random variable which mean $\bar{\eta}_k$ is given by

$$\bar{\eta}_k = \mathbb{E}[\eta_k], \quad (2.26)$$

and covariance \mathbf{Q}_k by

$$\mathbf{Q}_k = \mathbb{E}[(\eta_k - \bar{\eta}_k)(\eta_k - \bar{\eta}_k)^T]. \quad (2.27)$$

\mathbf{Q}_k is generally referred to as the *process noise* and is often modeled as a zero-mean Gaussian random variable, which can be denoted as $\eta_k \sim \mathcal{N}(0, \mathbf{Q}_k)$.

Recalling the 1-D moving object example with state vector $m_k = \begin{bmatrix} x_k \\ v_k \end{bmatrix}$, and assuming a constant velocity, the position after a time Δt is given by:

$$x_{k+1} = v_k \Delta t + x_k. \quad (2.28)$$

Considering the forecast to be invariant with time, the discrete forecast operator \mathbf{F} is given by

$$\mathbf{F} = \begin{bmatrix} 1 & \Delta t \\ 0 & 1 \end{bmatrix}, \quad (2.29)$$

which allows us to express the forecast as

$$m_{k+1}^f = \mathbf{F}m_k + \eta_k \quad (2.30)$$

I stated earlier that knowing the laws that govern the system's evolution with time is already a great deal of information. However, we can also see the limits of using only a numerical model to estimate the state of a real physical system. Even in the context of our straightforward model: if the initial position is not known perfectly, the filtered position could always be over or under-evaluated. The error would be much more severe if the constant velocity is not the good one, as the state estimate would get worse with time. Even worse, if the equations governing the system are more complicated than our basic constant velocity model, the forecasting operator likely loses connection with the real system pretty quickly. This

is especially true in the context of NWP, where the chaotic nature of the atmosphere makes it difficult to predict accurately with the sole use of numerical modeling.

The solution to this problem takes the form of DA tools that integrate both the time-dependent observations with a numerical forecasting model, in order to produce the best state estimate, at any given time k . With this introductory matter aside, we can now focus on the precursor of statistical DA tools: the Kalman Filter.

2.1.5 The Kalman Filter

Historically regarded as the first statistical DA tool, the KF was proposed by Rudolph E. Kalman in his seminal paper titled "*A new approach to linear filtering and prediction problems*" (Kalman, 1960). This work was closely followed by a joint publication with Richard S. Bucy (Kalman and Bucy, 1961), but was initially rejected: one of the referees commented, "*it cannot possibly be true*" (Grewal and Andrews, 2001). A couple of years later, the KF was successfully applied to estimate and control the circumlunar trajectory of the Apollo space capsule, sparking a vivid interest for the tool. Since the sixties, the KF has been a prevalent tool in adaptive filtering for signal processing and control theory (Chen, 2003). The main idea of the KF is the generalization of the BLUE to dynamic systems, which is made possible with a straightforward algorithm.

Initialization

- Initialize the background state
- Initialize the background uncertainty

Forecast

- Project the system state ahead in time with the forecast operator
- Adjust the state uncertainty to account for the forecasting errors

Analyse

- Get observations and their associated uncertainty/precision information
- Compute the residual between the forecast state and observations
- Compute the analysis state, the best tradeoff between all pieces of information

Coming back at our object tracking problem, the KF can provide an estimate position for each k , the discrete time-steps at which we observe the linear system. To simplify the notation, we consider both the observation and forecast operator to be invariant with respect to time.

Initialization We consider the initialization of the filter in the first place. As with the BLUE, it starts with a prior belief at step $k=0$. Therefore we have

$$m_0 \quad \text{and} \quad \mathbf{P}_0, \quad (2.31)$$

where m_0 is the background state and \mathbf{P}_0 is the background uncertainty ($\mathbf{P}_0 = \mathbb{E}[\beta_0(\beta_0)^T]$ with β_0 the background error). At that stage, because we typically do not have better information than a simple initial guess, the background uncertainty \mathbf{P}_0 should be large to reflect the lack of information.

Forecast Having defined the forecast operator in equation (2.25), we can predict the state of the system at the next timestep $k = 1, \dots, K$ with

$$m_k^f = \mathbf{F}m_0 + \eta_0. \quad (2.32)$$

By applying the forecast operator, we introduce the following forecasting error

$$\begin{aligned} \epsilon_k^f &= m_k^f - m_k^t, \\ &= \mathbf{F}m_0 - \mathbf{F}m_0^t - \eta_0, \\ &= \mathbf{F}(m_0 - m_0^t) - \eta_0, \\ &= \mathbf{F}\beta_0 - \eta_0, \end{aligned} \quad (2.33)$$

and thus we need to update our belief accordingly by formulating the forecast error covariance

$$\begin{aligned} \mathbf{P}_k^f &= \mathbb{E}[\epsilon_k^f (\epsilon_k^f)^T], \\ &= \mathbb{E}[(\mathbf{F}\beta_0 - \eta_0)(\mathbf{F}\beta_0 - \eta_0)^T], \\ &= \mathbf{F}\mathbb{E}[\beta_0(\beta_0)^T]\mathbf{F}^T + \mathbb{E}[\eta_0(\eta_0)^T], \\ &= \mathbf{F}\mathbf{P}_0\mathbf{F}^T + \mathbf{Q}_k. \end{aligned} \quad (2.34)$$

The forecast error covariance is computed by applying the forecast operator to the background error covariance, to which we add the covariance term that rules the process noise.

Analysis From the forecast, the BLUE gives us the analysis equation, that can be applied as soon as observations are available. The analysis state is given by

$$m_k^a = m_k^f + \mathbf{K}_k(y_k^0 - \mathbf{H}m_k^f), \quad (2.35)$$

where \mathbf{K}_k is the Kalman gain matrix. In DA, it is common to refer to $(y_k^0 - \mathbf{H}m_k^f)$ as the *innovation* term, which measures the discrepancies between the observation and the predicted observation. The Kalman gain matrix is the optimal weight that minimizes the analysis (posterior) error covariance. Note that the Kalman gain matrix is by construction positive definite. From our development of the BLUE, its formulation is given by

$$\mathbf{K}_k = \mathbf{P}_k^f \mathbf{H}^T (\mathbf{H} \mathbf{P}_k^f \mathbf{H}^T + \mathbf{R})^{-1}, \quad (2.36)$$

which corresponds to equation (2.19) where variances have been replaced by covariance matrices, which in turn, explains the introduction of the observation operator \mathbf{H} .

Let us express $\mathbf{S}_k = \mathbf{H} \mathbf{P}_k^f \mathbf{H}^T + \mathbf{R}$ where $\mathbf{H} \mathbf{P}_k^f \mathbf{H}^T$ act as a projection of the forecast error from the state space, to the measurement space as the "measurement uncertainty". It allows expressing the Kalman gain as

$$\begin{aligned} \mathbf{K}_k &= \mathbf{P}_k^f \mathbf{H}^T \mathbf{S}_k^{-1}, \\ &\approx \frac{\text{prediction uncertainty}}{\text{measurement uncertainty}} \mathbf{H}^T, \end{aligned} \quad (2.37)$$

which goes to emphasize the optimal weight nature of the Kalman gain in the update equation as

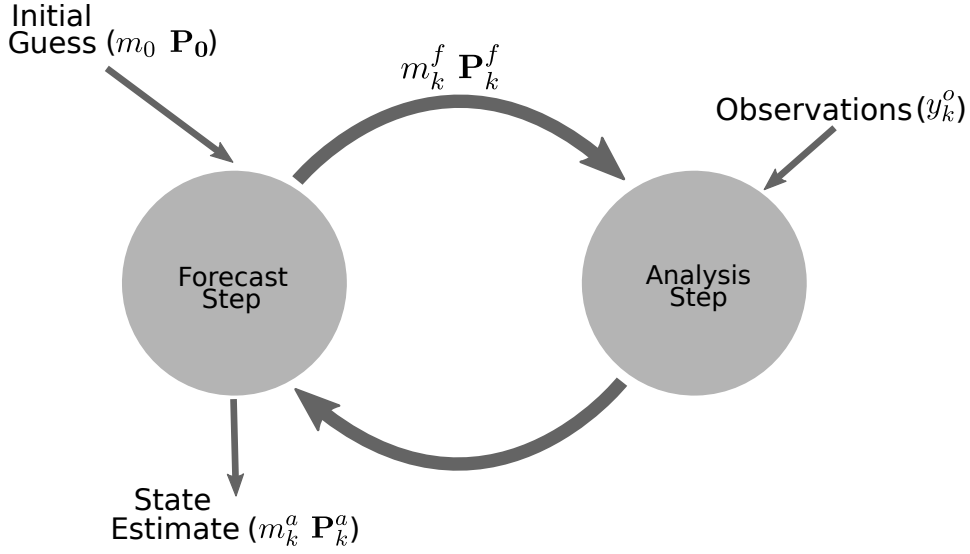


Figure 2.2: Schematic representation of the KF algorithm. From an initial guess/initial conditions, the sequential filter is initialized to produce forecast and analysis cycles. After each cycle, the output of the filter is the state estimate and its associated uncertainty. *Modified from* (Labbe, 2016).

$$\begin{aligned}
 m_k^a &= m_k^f + \mathbf{P}_k^f \mathbf{H}^T \mathbf{S}^{-1} (y_k^o - \mathbf{H} m_k^f), \\
 &\approx m_k^f + \frac{\text{prediction uncertainty}}{\text{measurement uncertainty}} \mathbf{H}^T (y_k^o - \mathbf{H} m_k^f).
 \end{aligned} \tag{2.38}$$

Finally, the last step of the KF analysis is to compute the posterior error covariance \mathbf{P}_k^a given by

$$\mathbf{P}_k^a = (\mathbf{I} - \mathbf{K}_k \mathbf{H}) \mathbf{P}_k^f, \tag{2.39}$$

which is nothing more than the multivariate version of equation (2.21). Note that due to the positive definite nature of the Kalman gain matrix, the analysis covariance should always be a smaller fraction of the forecast covariance. If the analysis uncertainty is greater than the forecast uncertainty, it is generally an indication of filter degeneracy (often caused by numerical errors). With this short example, we have demonstrated that from a rough initial guess, it is possible to predict the future state of the system and to correct it when observations are available. From there, the analysis of the previous step becomes the starting point for the next cycle (Figure 2.2), and each new cycle benefits from the variance reduction of the previous analysis. This is what allows starting from a large prior uncertainty: the filter is expected to mitigate the initial lack of knowledge of the system over time (provided the system is behaving purely linearly and that the forecast operator is designed correctly)

Once all the developments are out of the way, the KF cycles reduce to a few equations

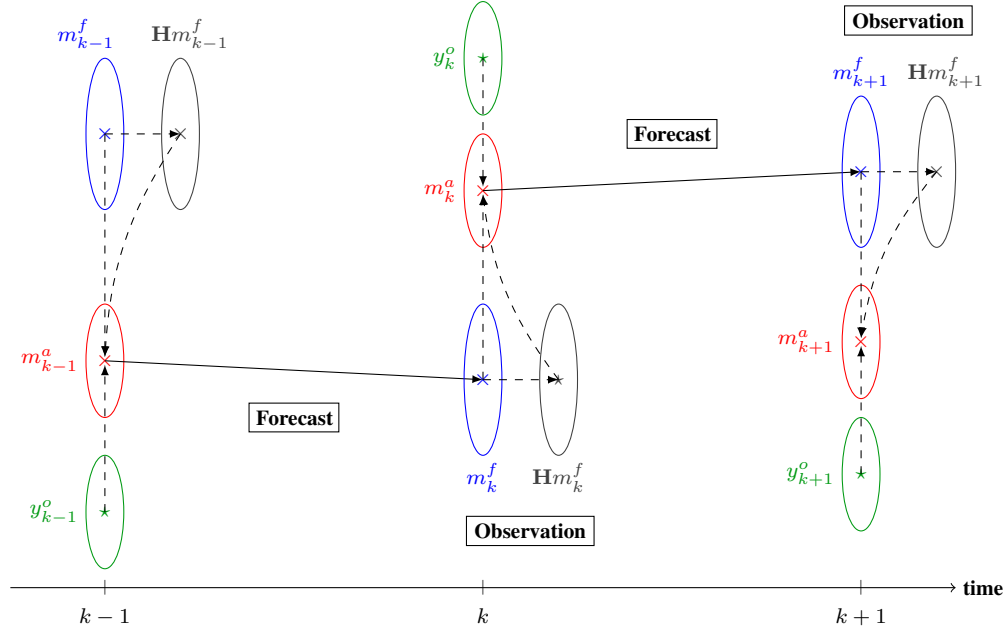


Figure 2.3: Schematic representation of the KF algorithm. Crosses denote state vectors, stars denote measurement vectors, and ellipses represent covariances. Blue denotes the forecast state, red the analysis, green the observed data, and grey the forecast data. The dashed lines are indicative of the information used to produce the analysis state. Note that here, because the observation error covariance (green) is lower than the forecast error covariance (blue), the analysis is closer to the data vector than the forecast vector.

Forecast

$$m_k^f = \mathbf{F}_k m_{k-1}^a + \eta_k,$$

$$\mathbf{P}_k^f = \mathbf{F}_k \mathbf{P}_{k-1}^a \mathbf{F}_k^T + \mathbf{Q}_k,$$

Analysis

$$\mathbf{K}_k = \mathbf{P}_k^f \mathbf{H}_k^T (\mathbf{H}_k \mathbf{P}_k^f \mathbf{H}_k^T + \mathbf{R})^{-1},$$

$$m_k^a = m_k^f + \mathbf{K}_k (y_k^o - \mathbf{H}_k m_k^f),$$

$$\mathbf{P}_k^a = (\mathbf{I} - \mathbf{K}_k \mathbf{H}_k) \mathbf{P}_k^f,$$
(2.40)

yielding a very powerful filter, shining by its simplicity and ease of implementation.

A "sequential view" of the KF is illustrated in Figure 2.3 for a 1-D problem. This way, we can see that after each analysis, the filter defines the best tradeoff between the forecast (blue) and the measurement (green) information.

Additionally, the Gaussian assumption of the KF allows expressing the analysis as a Bayesian least-squares minimization problem (its "variational" formulation). This comes down to minimizing the following cost function

$$\mathcal{C}(m_k) = \frac{1}{2} \|\mathbf{H}_k m_k - y_k^o\|^2 + \frac{1}{2} \|m_k - m_k^f\|^2, \quad (2.41)$$

that can also be expressed as

$$\mathcal{C}(m_k) = \frac{1}{2}(y_k^o - \mathbf{H}m_k)^T \mathbf{R}^{-1}(y_k^o - \mathbf{H}m_k) + \frac{1}{2}(m_k - m_k^f)^T \mathbf{P}_k^{f-1}(m_k - m_k^f), \quad (2.42)$$

and again, we can see that the analysis is a balancing act between the observation term (first term), and the forecasting term (second term).

From the derivation of the BLUE, we have seen that the KF is an optimal tool when dealing with linear dynamics and Gaussian noises. However, its abilities to tackle non-linear problems are relatively limited. Typically, the filter will fail whenever the forecast model has a strongly non-linear response, or if the observation operator is non-linear. To mitigate this shortcoming, one of the earliest proposition took the form of the Extended Kalman Filter (EKF) (Jazwinski, 2007).

2.1.6 Extended Kalman Filter

As its name suggests, the EKF allows to "extend" the KF formalism, to tackle weakly non-linear problems. The idea of the EKF date back to the early papers of Kopp and Orford (1963); Cox (1964) and was formally established by Jazwinski (2007)'s first edition in 1970. The idea stems from consideration of the non-linear case for the KF, where the forecast state and the observations equations become

$$\begin{aligned} m_k^f &= \mathcal{F}(m_{k-1}^a) + \eta_k, \\ y_k^o &= \mathcal{H}(m_k^f) + \epsilon_k, \end{aligned} \quad (2.43)$$

where \mathcal{F} is a non-linear forecasting operator, and \mathcal{H} is a non-linear observation operator. Using this set of equations in the KF would violate its underlying Gaussian assumption (a Gaussian function through a non-linear process is not Gaussian anymore), and therefore an optimal solution cannot be derived for the non-linear setting. We can instead rely on a local tangent-linear approximation of \mathcal{F} and \mathcal{H} , at the current state

$$\begin{aligned} \mathbf{F} &= \left. \frac{\partial \mathcal{F}(m_{k-1}^a)}{\partial m} \right|_{m_{k-1}^a} \\ \mathbf{H} &= \left. \frac{\partial \mathcal{H}(m_k^f)}{\partial m} \right|_{m_k^f} \end{aligned} \quad (2.44)$$

where \mathbf{F} and \mathbf{H} are the Jacobians of the forecasting and observation operator, yielding respectively the discrete forecast operator and the discrete observation operator (Jazwinski, 2007; Labbe, 2016). This gives the following algorithm for the EKF

Forecast

$$\begin{aligned}\mathbf{F} &= \left. \frac{\partial \mathcal{F}(m_{k-1}^a)}{\partial m} \right|_{m_{k-1}^a} \\ m_k^f &= \mathbf{F}_k m_{k-1}^a + \eta_k \\ \mathbf{P}_k^f &= \mathbf{F} \mathbf{P}_k^f \mathbf{F}^T + \mathbf{Q}_k\end{aligned}$$

Analysis

(2.45)

$$\begin{aligned}\mathbf{H} &= \left. \frac{\partial \mathcal{H}(m_k^f)}{\partial m} \right|_{m_k^f} \\ \mathbf{K}_k &= \mathbf{P}_k^f \mathbf{H}^T (\mathbf{H} \mathbf{P}_k^f \mathbf{H}^T + \mathbf{R})^{-1} \\ m_k^a &= m_k^f + \mathbf{K}_k (y_k^o - \mathcal{H}(m_k^f)) \\ \mathbf{P}_k^a &= (\mathbf{I} - \mathbf{K} \mathbf{H}) \mathbf{P}_k^f\end{aligned}$$

In practice, the linearized \mathbf{F} and \mathbf{H} are often too inaccurate to produce reliable forecasts and observations. Their non-linear versions computed from an appropriate numerical integration scheme are generally preferred (Labbe, 2016). For this reason, the Jacobians are only used to compute the error covariances and the Kalman Gain matrix. If the problem allows it, the Jacobians can be computed analytically. Otherwise, they are approximated with a numerical solver. The EKF allows formulating a linear approximation to the non-linear problem, which requires re-evaluating the Jacobians of \mathbf{F} and \mathbf{H} at each iteration.

Unfortunately, both the KF and EKF are limited to studies of small-scale systems, with at best a few hundred parameters. When the size of the state space increases, the covariance matrices \mathbf{P}^f and \mathbf{P}^a (of size $n \times n$) become untractable, and other strategies must be used. As with FWI, in operational NWP and other DA applications, it is not uncommon to deal with state spaces in the range of $10^{6 \sim 9}$ degrees of freedom, which prevents storing and manipulating the covariance matrices. To overcome this limitation, it is possible to rely on rank-limitation strategies, such as ensemble methods, to limit the computational burden of DA applications.

2.2 Ensemble Kalman Filter

As discussed in the previous section, using the KF to study large-scale systems can be problematic, especially when it comes to manipulating covariance matrices. To overcome this limitation, Evensen (1994) proposed a Monte-Carlo approximation of the KF, based on an *ensemble* representation of the system state: the *Ensemble Kalman Filter* (EnKF). He showed that for an infinite number of ensemble members and considering a linear-Gaussian setting, his methodology could yield the exact KF solution. Therefore, with a limited number of ensemble members, it can produce an approximation of the KF solution. In the non-linear case, however, the state approximated by the ensemble can differ from the KF filter solution (Mandel and Beezley, 2009; Le Gland et al., 2009; Mandel et al., 2011). Despite representativity errors, its successful applications in geophysics have made it one of the prevalent DA methodology, and its limitations have been overlooked. Thanks to its success, variants of the EnKF are currently being developed at an operational level in several meteorological centers on up to 10^9 degrees of freedom (as it is the case for the MOGREPS global assimilation system ran at the Met Office (United-Kingdom) or the ICON global domain model ran at the Deutscher Wetterdienst (Germany) as

part of their numerical weather prediction routines). Note that even though the EnKF original formulation required the observation operator to be linear, the EnKF has, since then, been successfully applied with non-linear observation operators (Evensen, 2003; Hunt et al., 2007).

As the KF and EKF, the EnKF is a Bayesian sequential filter, which is also based on alternating forecast and analysis steps. The ensemble representation allows computing a reduced-rank analysis, which in turn, allows a significant cost reduction in high-dimensional problems. In this section, I will detail the EnKF formalism, list some of the popular variants of this filter, and discuss its practical implementation.

2.2.1 Ensemble representation

The overarching idea of Evensen's EnKF is to approximate the system state (its state estimate and its error covariance matrices), with an ensemble \mathbf{m}_k of N_e state vectors $m_k^{(i)} \in \mathbb{R}^n$,

$$\mathbf{m}_k = \{m_k^{(1)}, m_k^{(2)}, \dots, m_k^{(N_e)}\} \quad (2.46)$$

where \mathbf{m}_k is a $n \times N_e$ matrix, whose column contains all the individual state vectors and where typically $N_e \ll n$. From the ensemble representation, provided its repartition is Gaussian, the state estimate and uncertainty are represented by the first and second-order Gaussian moments (mean and variance). Consequently, it is possible to compute the relevant metrics involved in the original KF formulation from the ensemble (as illustrated with a 2-D case in Figure 2.4), which makes the extension of the KF scheme to large scale problems affordable.

The system state estimate (first Gaussian moment) is given by

$$\bar{m}_k = \frac{1}{N_e} \sum_{i=1}^{N_e} m_k^{(i)}. \quad (2.47)$$

We denote the *perturbation* matrix \mathbf{M}_k , which columns contain the deviation to the mean

$$\mathbf{M}_k = [m_k^{(1)} - \bar{m}_k, m_k^{(2)} - \bar{m}_k, \dots, m_k^{(N_e)} - \bar{m}_k]. \quad (2.48)$$

This matrix is at best of rank $N_e - 1$ by construction, as the sum of all of its column equates to zero. The state uncertainty is given by the second Gaussian moment that we can approximate by

$$\mathbf{P}_{e,k} = \frac{1}{N_e - 1} (\mathbf{m}_k - \bar{m}_k)(\mathbf{m}_k - \bar{m}_k)^T = \frac{1}{N_e - 1} \mathbf{M} \mathbf{M}^T. \quad (2.49)$$

where the subscript e denote the "ensemble" covariance matrix (or "empirical" covariance estimated from the ensemble). Computing the Gaussian moments of the distribution only requires storing N_e state vectors, which is very interesting when $N_e \ll n$.

This ensemble representation strategy, which yields a low-rank approximation of the system state, is the basis of the EnKF. From here, we can study how this approximation carries over the forecast and analysis steps.

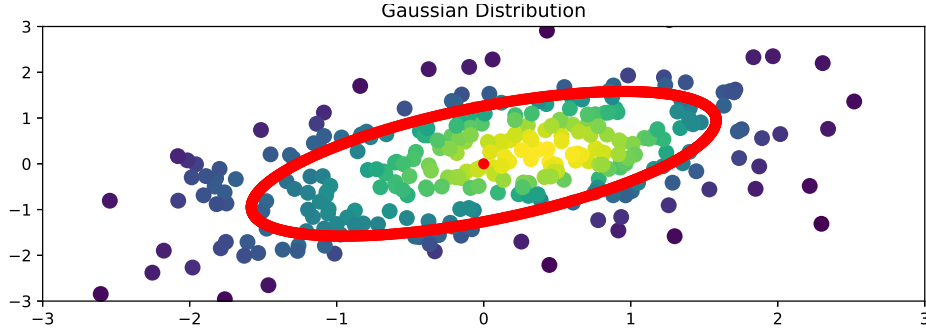


Figure 2.4: Schematic representation of the ensemble strategy. Here, the true mean and covariance (red dot and red ellipse) can be estimated from the ensemble of points. Note that this sampling strategy is not interesting in this 2-D example, as $N_e \gg 2$. However, it becomes advantageous as soon as $N_e \ll n$.

2.2.2 EnKF's forecast step

Considering first the forecast step in the ensemble formalism. Following the KF's forecast equation, the ensemble forecast from step $k - 1$ to k is given by

$$m_k^{f(i)} = \mathcal{F}_k(m_{k-1}^{(i)}) + \eta_k^{(i)} \quad i = 1, 2, \dots, N_e, \quad (2.50)$$

where the forecast operator \mathcal{F}_k can be non-linear, and $\eta_k^{(i)}$ is the process-noise which is generally modeled as a multivariate Gaussian random vector ($\eta_k^{(i)} \in \mathcal{N}(0, \mathbf{Q}_k)$), with \mathbf{Q}_k being the forecast error covariance matrix. As with the ensemble covariance matrix, it is possible to generate an ensemble of random noise vector $\{\eta_k^{(1)}, \eta_k^{(2)}, \dots, \eta_k^{(N_e)}\}$, which empirical covariance $\mathbf{Q}_{k,e}$ tends toward \mathbf{Q}_k when the size of the ensemble goes to infinity. In practice, however, we generally have limited knowledge on \mathbf{Q}_k and on the natural process errors that affect the system. In that case, one generally reduces equation (2.50) as

$$m_k^{f(i)} = \mathcal{F}_k(m_{k-1}^{(i)}) \quad i = 1, 2, \dots, N_e. \quad (2.51)$$

While omitting the process noise greatly simplify the forecasting step, this simplification tends to under evaluate the forecast covariance. This underestimation can be mitigated with an *inflation* procedure that artificially increases the ensemble spread after the forecast step, to avoid an over-confidence in the forecast. This procedure is detailed in subsection 2.3.5.

The advantage of the EnKF is that the forecast covariance can be approximated from the forecasted ensemble, such that

$$\mathbf{P}_{k,e}^f = \frac{1}{N_e - 1} \mathbf{M}^f (\mathbf{M}^f)^T, \quad (2.52)$$

where each column of $\mathbf{M}_k^{f(i)} = [m_k^{f(i)} - \bar{m}_k^f]$. In practice, the ensemble covariance is never explicitly computed due to computational and memory limitations. The implications of this approximation over the analysis step are detailed in the next subsections.

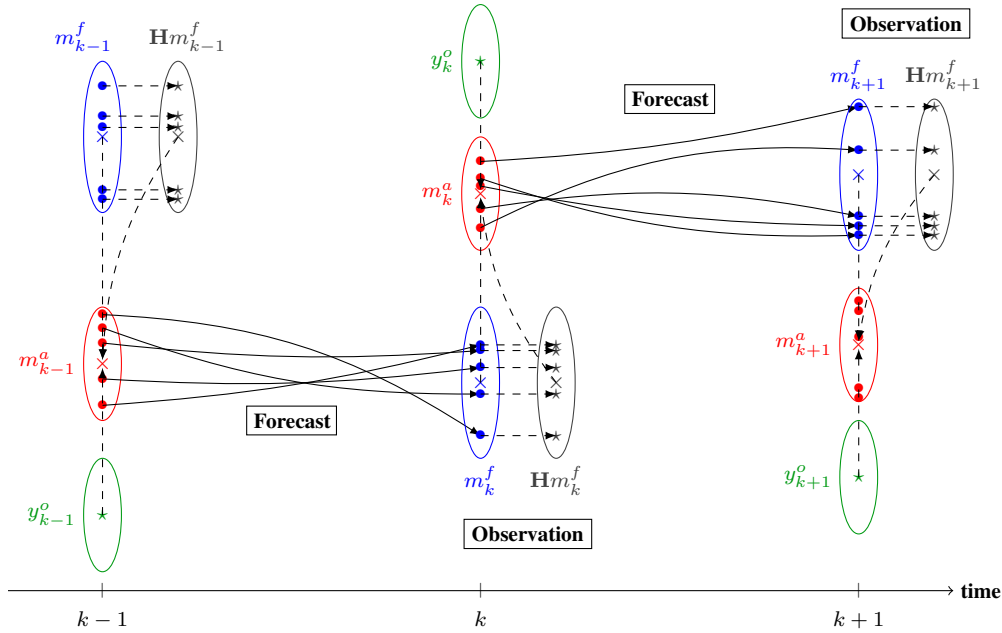


Figure 2.5: Schematic representation of the EnKF algorithm. Bullets represent the ensemble members, crosses denote ensemble means, stars denote measurement vectors, and ellipses represent ensemble covariances. Blue denotes the forecast state, red the analysis, green the observed data, and grey the forecast data. The dashed lines are indicative of the information used to produce the analysis state.

2.2.3 Analysis step

As stated before, the goal of the EnKF is to provide a low-rank approximation of the KF equations, without building explicitly any covariance matrices. In its original formulation, Evensen (1994) proposed to consider only the case where the observation operator is linear and directly plugged the empirical covariance $\mathbf{P}_{k,e}^f$ in the KF analysis equation. In the subsequent developments, the subscript k has been omitted to ease the readability.

Evensen's original analysis scheme is based on the formulation of \mathbf{P}^a in terms of \mathbf{P}^f as in equation (2.16). In the multivariate case, equation (2.16) becomes

$$\begin{aligned} \mathbf{P}^a &= \mathbb{E}[(\alpha)(\alpha)^T], \\ &= (\mathbf{I} - \mathbf{KH})\mathbf{P}^f(\mathbf{I} - \mathbf{H}^T\mathbf{K}^T) + \mathbf{K}\mathbf{R}\mathbf{K}^T, \\ &= \mathbf{P}^f - \mathbf{KHP} - \mathbf{P}^f\mathbf{H}^T\mathbf{K}^T + \mathbf{K}(\mathbf{HP}^f\mathbf{H}^T + \mathbf{R})\mathbf{K}^T, \\ &= (\mathbf{I} - \mathbf{KH})\mathbf{P}^f. \end{aligned} \tag{2.53}$$

However, Burgers et al. (1998) and Houtekamer and Mitchel (1998) pointed-out that Evensen's formulation of \mathbf{P}_e^a does not coincide with \mathbf{P}^a unless observations are treated as random variables with covariance \mathbf{R} . Evensen's original development thus lacked an equivalent to the term $\mathbf{K}\mathbf{R}\mathbf{K}^T$, making his formulation of \mathbf{P}_e^a converges toward

$$\mathbf{P}_e^a = (\mathbf{I} - \mathbf{KH})\mathbf{P}^f(\mathbf{I} - \mathbf{KH})^T, \tag{2.54}$$

which does not satisfy the BLUE equation of the covariance analysis.

From this demonstration, careful research went into formulating the ensemble analysis that satisfies the BLUE equation while retaining the significant cost reduction offered by the ensemble approximation.

Stochastic analysis

To mitigate the shortcoming of the original EnKF formulation, Burgers et al. (1998); Houtekamer and Mitchel (1998) proposed a new EnKF scheme in which the observations are perturbed with Gaussian random noise, by defining an ensemble of N_e perturbed observations $\mathbf{y}_\epsilon^o = [y_\epsilon^{o(1)}, y_\epsilon^{o(2)}, \dots, y_\epsilon^{o(N_e)}]$ with

$$y_\epsilon^{o(i)} = y^o + \epsilon^{(i)}, \quad i = 1, 2, \dots, N_e. \quad (2.55)$$

The stochastic EnKF analysis ensemble $\mathbf{m}^a = [m^{a(1)}, m^{a(2)}, \dots, m^{a(N_e)}]$ is expressed as

$$m^{a(i)} = m^{f(i)} + \mathbf{K}_e [y_\epsilon^{o(i)} - \mathbf{H}(m^{f(i)})], \quad (2.56)$$

with the Kalman gain defined as

$$\mathbf{K} = \mathbf{P}_e^f \mathbf{H}^T (\mathbf{H} \mathbf{P}_e^f \mathbf{H}^T + \mathbf{R})^{-1}. \quad (2.57)$$

Note that this formulation implies that the state estimate is given by

$$\bar{m}^a = \bar{m}^f + \mathbf{K}_e [\bar{y}_\epsilon^o - \mathcal{H}(\bar{m}^f)], \quad (2.58)$$

where $\bar{y}_\epsilon^o = y^o$ corresponds to the unperturbed observations. With these developments, the authors ensured that the relation between \mathbf{P}_e^a and \mathbf{P}_e^f in the EnKF, is the same as the relation between \mathbf{P}^a and \mathbf{P}^f in the original KF scheme. The analysis ensemble covariance is given by

$$\begin{aligned} \mathbf{P}_e^a &= \mathbb{E}[(\mathbf{m}^a - \bar{m}^a)(\mathbf{m}^a - \bar{m}^a)^T] \\ &= \frac{1}{N_e - 1} \mathbf{M}^a (\mathbf{M}^a)^T \\ &= (\mathbf{I} - \mathbf{K}_e \mathbf{H}) \mathbf{P}_e^f + \mathcal{O}(N_e^{-1/2}) \end{aligned} \quad (2.59)$$

where $\mathbf{M}^a = (\mathbf{m}^a - \bar{m}^a)$. With the addition of perturbed observations, their version of the analysis ensemble covariance satisfies the BLUE. They also showed that this formulation could be expanded to account for non-linear observations: introducing the non-linear observation operator \mathcal{H} , the general formulation of the stochastic EnKF analysis is given by the following set of equations:

$$\begin{aligned} y^{f(i)} &= \mathcal{H}(m^{f(i)}) \\ \bar{y}^f &= \frac{1}{N} \sum_{i=1}^{N_e} y^{f(i)} \\ \mathbf{Y}^f &= (\mathbf{y}^f - \bar{y}^f) \\ \mathbf{K}_e &= \frac{1}{N_e - 1} \mathbf{M}^f (\mathbf{Y}^f)^T \left(\frac{1}{N_e - 1} \mathbf{Y}^f (\mathbf{Y}^f)^T + \mathbf{R} \right)^{-1} \\ m^{a(i)} &= m^{f(i)} + \mathbf{K}_e (y_\epsilon^{o(i)} - y^{f(i)}). \end{aligned} \quad (2.60)$$

This formulation solved the two major problems of the KF: the stochastic EnKF satisfies the BLUE even for non-linear forecast and observation operators and is compatible with large systems study thanks to its covariance matrix-free formulation.

While the stochastic EnKF relies on approximating \mathbf{R} with an ensemble of perturbed observations, alternatives "deterministic" EnKFs, have been formulated to account for \mathbf{R} explicitly. The interest of deterministic methods is that they eliminate the risks of sampling error for \mathbf{R} . They are also generally more accurate and stable than the stochastic EnKF (Whitaker and Hamill, 2002; Tippett et al., 2003), and also exhibit better performances (Whitaker and Hamill, 2002; Sakov and Oke, 2008).

Deterministic EnKFs are of particular interest in this study, as they have been preferred to stochastic EnKF in our numerical experiments, due to the advantages mentioned above. I propose a short review of the deterministic EnKFs in the following subsections.

Deterministic analysis

As mentioned previously, the first motivation for the derivation of deterministic EnKFs was to avoid sampling error introduced by the perturbed observations, while also satisfying the BLUE equation (2.53). Following Whitaker and Hamill (2002), these two requirements can be met by expressing the EnKF analysis equations as a two-step update:

$$\bar{\mathbf{m}}^a = \bar{\mathbf{m}}^f + \mathbf{K}(\bar{\mathbf{y}}^o - \mathbf{H}\bar{\mathbf{m}}^f), \quad (2.61)$$

$$\mathbf{M}^a = \mathbf{M}^f + \hat{\mathbf{K}}(\mathbf{Y}^o - \mathbf{H}\mathbf{M}^f), \quad (2.62)$$

where \mathbf{K} is the regular Kalman gain, and $\hat{\mathbf{K}}$ is a *reduced gain* matrix used to update the perturbation matrix. The stochastic EnKF presented in the previous subsection can be expressed as a special case of this formulation where $\mathbf{K} = \hat{\mathbf{K}}$ and \mathbf{Y}^o are a set of random perturbations drawn according to the measurement noise distribution \mathbf{R} (as in equation (2.55)). The need for perturbed observation then appear clear when considering the special case where $\mathbf{K} = \hat{\mathbf{K}}$ and $\mathbf{Y}^o = 0$, which can be shown to yield

$$\mathbf{P}_e^a = (\mathbf{I} - \mathbf{K}\mathbf{H})\mathbf{P}_e^f(\mathbf{I} - \mathbf{K}\mathbf{H})^T, \quad (2.63)$$

and therefore does not satisfy the BLUE analysis, as seen previously (Burgers et al., 1998).

Whitaker and Hamill (2002), propose to set the reduced gain $\hat{\mathbf{K}}$ so that it satisfies

$$(\mathbf{I} - \hat{\mathbf{K}}\mathbf{H})\mathbf{P}_e^f(\mathbf{I} - \hat{\mathbf{K}}\mathbf{H})^T = (\mathbf{I} - \mathbf{K}\mathbf{H})\mathbf{P}_e^f, \quad (2.64)$$

which has a solution under the form

$$\hat{\mathbf{K}} = \mathbf{P}_e^f \mathbf{H}^T [(\sqrt{\mathbf{H}\mathbf{P}_e^f \mathbf{H}^T + \mathbf{R}})^{-1}]^T [\sqrt{(\mathbf{H}\mathbf{P}_e^f \mathbf{H}^T + \mathbf{R})} + \sqrt{\mathbf{R}}]^{-1}. \quad (2.65)$$

Therefore, the perturbation update can be expressed as

$$\mathbf{M}^a = (\mathbf{I} - \hat{\mathbf{K}}\mathbf{H})\mathbf{M}^f. \quad (2.66)$$

This deterministic analysis is referred to as an ensemble *square root* filter (EnSRF), where the analysis covariance is obtained by updating its square root \mathbf{M}^a . Contrarily to the stochastic EnKF, there are several implementations of deterministic EnKFs, even though they are algebraically equivalent. Their differences stem from the non-unique ways of computing \mathbf{M}^a as a *transformation* of \mathbf{M}^f under the form

$\mathbf{M}^f \mathbf{T}$ or $\hat{\mathbf{T}} \mathbf{M}^f$, where \mathbf{T} and $\hat{\mathbf{T}}$ are *transformation matrices* that satisfy the BLUE analysis equations (Sakov and Oke, 2008). For a complete review of EnSRFs methods and other deterministic alternatives, readers might refer to Tippett et al. (2003); Sakov and Oke (2008); Fletcher (2017).

For the remainder of this section, I will review additional forms of deterministic EnKFs:

- the Ensemble Transform Kalman Filter (ETKF) proposed by Bishop et al. (2001), as it is the scheme we chose to perform numerical experiments in this study.
- the Maximum Likelihood Ensemble Filter (MLEF) proposed by Zupanski (2005), which non-linear formulation and iterative nature make it an interesting bridge with numerical optimization methods.
- the Ensemble Kalman Inverse (EKI) proposed by Iglesias et al. (2013a), which is a stochastic optimization method based on statistical DA principles.

2.2.4 Ensemble Transform Kalman Filter

The ETKF is a deterministic EnKF scheme that was initially proposed by Bishop et al. (2001). Its derivation starts by recognizing that at any time, the state estimate error covariance can be reduced to a product of the covariance' square roots

$$\mathbf{P}_e = \frac{1}{N_e - 1} \mathbf{M} \mathbf{M}^T. \quad (2.67)$$

Therefore, the most cost-effective way of computing the analysis ensemble, as suggested in the ensemble square root filters, is to update the mean and the perturbation matrix separately. Considering equation (2.67) we can re-write the Kalman gain matrix as

$$\begin{aligned} \mathbf{K} &= \frac{1}{N_e - 1} \mathbf{M}^f (\mathbf{M}^f)^T \mathbf{H}^T \left[\frac{1}{N_e - 1} \mathbf{H} \mathbf{M}^f (\mathbf{M}^f)^T \mathbf{H}^T + \mathbf{R} \right]^{-1} \\ &= \frac{1}{\sqrt{N_e - 1}} \mathbf{M}^f (\mathbf{M}^f)^T \mathbf{H}^T \frac{1}{\sqrt{N_e - 1}} \left[\frac{1}{\sqrt{N_e - 1}} \mathbf{H} \mathbf{M}^f (\mathbf{M}^f)^T \mathbf{H}^T \frac{1}{\sqrt{N_e - 1}} + \mathbf{R} \right]^{-1}. \end{aligned} \quad (2.68)$$

Defining $\mathbf{A} = \frac{1}{\sqrt{N_e - 1}} \mathbf{H} \mathbf{M}^f$, we can recognize the matrix identity

$$\mathbf{A}^T (\mathbf{A} \mathbf{A}^T + \mathbf{R})^{-1} = (\mathbf{I}_{N_e} + \mathbf{A}^T \mathbf{R}^{-1} \mathbf{A})^{-1} \mathbf{A}^T \mathbf{R}^{-1}, \quad (2.69)$$

and thus the Kalman gain's expression becomes

$$\mathbf{K} = \frac{1}{N_e - 1} \mathbf{M}^f \left[\mathbf{I}_{N_e} + \frac{1}{N_e - 1} (\mathbf{M}^f)^T \mathbf{H}^T \mathbf{R}^{-1} \mathbf{H} \mathbf{M}^f \right]^{-1} (\mathbf{M}^f)^T \mathbf{H}^T \mathbf{R}^{-1}. \quad (2.70)$$

Replacing the Kalman gain in equation (2.53) by the form we have derived in equation (2.70) allows us to express the analysis covariance as

$$\begin{aligned}
 \mathbf{P}_e^a &= \frac{1}{N_e - 1} (\mathbf{I}_{N_e} - \mathbf{K}\mathbf{H}) \mathbf{M}^f (\mathbf{M}^f)^T, \\
 &= \frac{1}{N_e - 1} (\mathbf{M}^f - \mathbf{K}\mathbf{H}\mathbf{M}^f) (\mathbf{M}^f)^T, \\
 &= \frac{1}{N_e - 1} \left\{ \mathbf{M}^f - \frac{1}{N_e - 1} \mathbf{M}^f \left[\mathbf{I}_{N_e} + \frac{1}{N_e - 1} (\mathbf{M}^f)^T \mathbf{H}^T \mathbf{R}^{-1} \mathbf{H} \mathbf{M}^f \right]^{-1} (\mathbf{M}^f)^T \mathbf{H}^T \mathbf{R}^{-1} \mathbf{H} \mathbf{M}^f \right\} (\mathbf{M}^f)^T, \\
 &= \frac{1}{N_e - 1} \mathbf{M}^f \left\{ \mathbf{I}_{N_e} - \frac{1}{N_e - 1} \left[\mathbf{I}_{N_e} + \frac{1}{N_e - 1} (\mathbf{M}^f)^T \mathbf{H}^T \mathbf{R}^{-1} \mathbf{H} \mathbf{M}^f \right]^{-1} (\mathbf{M}^f)^T \mathbf{H}^T \mathbf{R}^{-1} \mathbf{H} \mathbf{M}^f \right\} (\mathbf{M}^f)^T.
 \end{aligned} \tag{2.71}$$

This formulation is particularly interesting in its way of handling the analysis update. We can first notice that this expression is entirely free of covariance matrices manipulation: instead of \mathbf{P}_e^f , its factorized form is preferred. We can also note that the inner term can be expressed as a symmetrical operator $\tilde{\mathbf{P}}^a = \mathbf{T}\mathbf{T}^T$, which dimensions are $N_e \times N_e$,

$$\mathbf{P}_e^a = \frac{1}{N_e - 1} \mathbf{M}^f \tilde{\mathbf{P}}^a (\mathbf{M}^f)^T. \tag{2.72}$$

Recalling that the ensemble covariance matrices are at best of rank $N_e - 1$, $\tilde{\mathbf{P}}^a$ allows carrying the analysis over the ensemble subspace $\mathcal{S} \in \mathbb{R}^{N_e}$, on which \mathbf{P}_e^f and \mathbf{P}_e^a are well defined. This projection over the arbitrary subspace spanned by the ensemble allows a notable cost reduction, which makes the ETKF very affordable. $\tilde{\mathbf{P}}^a$ is nothing more than the "effective analysis covariance" spanning the ensemble subspace (Hunt et al., 2007). In his review, Tippett et al. (2003) points out that the computational complexity of the ETKF is $\mathcal{O}(N_e^3 + nN_e^2 + dN_e^2)$ and is thus linearly dependent on the size of the state and observation spaces. This low-rank reduction based on the ensemble representation is similar to the basis reduction strategies discussed in Chapter 1, which makes this formalism very interesting for uncertainty estimation in FWI.

The second advantage of the ETKF, is that it allows us to consider non-linear observation operators in a very straightforward manner. This property comes from the fact that in equation (2.71), each time the linearized observation operator appears, it is next to the perturbation matrix \mathbf{M}^f . Because $\mathbf{H}\mathbf{M}^f$ is the first Taylor expansion of $\mathcal{H}(\mathbf{m}^f) - \mathcal{H}(\bar{\mathbf{m}}^f)$, we can linearly approximate $\mathbf{H}\mathbf{M}^f$ as a matrix \mathbf{Y}^f (Harlim and Hunt, 2005) such that

$$\mathbf{Y}^f = \mathcal{H}(\mathbf{m}^f) - \bar{\mathbf{y}}^f \tag{2.73}$$

where $\bar{\mathbf{y}}^f$ is given by

$$\bar{\mathbf{y}}^f = \frac{1}{N_e} \sum_{i=1}^{N_e} \mathcal{H}(\mathbf{m}^{f(i)}). \tag{2.74}$$

By replacing occurrences of $\mathbf{H}\mathbf{M}^f$ by \mathbf{Y}^f in equation (2.71), the analysis formula becomes

$$\begin{aligned}
\mathbf{P}_e^a &= \frac{1}{N_e - 1} \mathbf{M}^f \left\{ \mathbf{I}_{N_e} - \frac{1}{N_e - 1} \left[\mathbf{I}_{N_e} + \frac{1}{N_e - 1} (\mathbf{Y}^f)^T \mathbf{R}^{-1} \mathbf{Y}^f \right]^{-1} (\mathbf{Y}^f)^T \mathbf{R}^{-1} \mathbf{Y}^f \right\} (\mathbf{M}^f)^T, \\
&= \frac{1}{N_e - 1} \mathbf{M}^f \left[\mathbf{I}_{N_e} + \frac{1}{N_e - 1} (\mathbf{Y}^f)^T \mathbf{R}^{-1} \mathbf{Y}^f \right]^{-1} (\mathbf{M}^f)^T, \\
&= \mathbf{M}^f \left[(N_e - 1) \mathbf{I}_{N_e} + (\mathbf{Y}^f)^T \mathbf{R}^{-1} \mathbf{Y}^f \right]^{-1} (\mathbf{M}^f)^T.
\end{aligned} \tag{2.75}$$

Finally the effective ensemble analysis covariance is expressed as

$$\tilde{\mathbf{P}}^a = \mathbf{T} \mathbf{T}^T = \left[(N_e - 1) \mathbf{I}_{N_e} + (\mathbf{Y}^f)^T \mathbf{R}^{-1} \mathbf{Y}^f \right]^{-1}. \tag{2.76}$$

The covariance analysis reduces to finding the appropriate square root of $\mathbf{T} \mathbf{T}^T$ such that $\mathbf{M}^a = \mathbf{M}^f \mathbf{T}$ (which has non-unique solutions). In his original formulation, Bishop et al. (2001) proposed to compute the SVD of the $N_e \times N_e$ operator $\mathbf{T} \mathbf{T}^T = \mathbf{U} \mathbf{\Gamma} \mathbf{U}^T$ and to perform the update as $\mathbf{M}^a = \mathbf{M}^f \mathbf{U} \mathbf{\Gamma}^{1/2}$. However, his version of the ETKF introduces some biases to the analysis (Wang and Bishop, 2003): the sum over the columns of \mathbf{M}^a yields a non-zero vector, which deviates the ensemble from the analysis mean. It also tends to introduce a bias on the variance repartition. Leeuwenburgh et al. (2005) demonstrated on a scalar model by assimilating a single observation that this formulation returned an ensemble perturbation in which all entries but-one were zeros. Consequently, if the number of observations is insufficient, Bishop et al. (2001)'s original scheme could result in zero-entry vectors in the ensemble perturbation matrix. This induces a tendency to generate outliers that concentrate most of the variance information (Sakov and Oke, 2008). The variance repartition bias is problematic because it tends to reduce the number of contributing perturbations to the representation of variance, which in turn can result in an underestimation of the total variance. In other words, Bishop et al. (2001)'s analysis scheme introduces rank-deficiency in the ensemble, such that it cannot sample the covariances adequately and might introduce instabilities along with the DA cycles.

To correct for this bias, Wang et al. (2004) and Ott et al. (2004) proposed to use a strictly symmetrical transformation operator, such that a better balance in the variance repartition is ensured while preserving the mean. To that extent, they introduced the transform operator as $\mathbf{T} = \mathbf{U} \mathbf{\Gamma} \mathbf{U}^T \mathbf{C}$, with \mathbf{C} being an orthogonal matrix that preserves the mean, satisfying

$$\mathbf{C} \mathbf{C}^T = \mathbf{I}_{N_e}, \quad \mathbf{C} \mathbf{1} = \mathbf{1} \tag{2.77}$$

where $\mathbf{1}$ is a vector of size N_e which entries are ones (Sakov and Oke, 2008). In practice, this matrix is often set as the identity matrix.

The analysis scheme of this mean-preserving ETKF (*spherical simplex* ETKF (Ott et al., 2004; Sakov and Oke, 2008)) can be resumed as follow:

$$\begin{aligned}
 y^{f(i)} &= \mathcal{H}(m^{f(i)}), \\
 \bar{y}^f &= \frac{1}{N} \sum_{i=1}^{N_e} y^{f(i)}, \\
 \mathbf{Y}^f &= (\mathbf{y}^f - \bar{y}^f), \\
 \tilde{\mathbf{P}}^a &= ((N_e - 1)\mathbf{I} + \mathbf{Y}^{fT} \mathbf{R}^{-1} \mathbf{Y}^f)^{-1} = \mathbf{U} \mathbf{\Gamma} \mathbf{U}^T, \\
 \mathbf{T} &= \mathbf{U} \mathbf{\Gamma}^{-1/2} \mathbf{U}^T \mathbf{I}_{N_e}, \\
 \mathbf{M}^a &= \sqrt{(N_e - 1)} \mathbf{M}^f \mathbf{U} \mathbf{\Gamma}^{-1/2} \mathbf{U}^T \mathbf{I}_{N_e}, \\
 \bar{m}^a &= \bar{m}^f + \mathbf{M}^f \tilde{\mathbf{P}}^a (\mathbf{Y}^f)^T \mathbf{R}^{-1} (y^o - \bar{y}^f), \\
 \mathbf{m}^a &= \bar{m}^a + \mathbf{M}^a.
 \end{aligned} \tag{2.78}$$

Additionally, as with the regular KF, because all underlying assumptions of Gaussianity hold for the ETKF, it can be expressed under the following variational formulation (Hunt et al., 2007)

$$\mathcal{C}(m) = \frac{1}{2} (y^o - \mathcal{H}(m))^T \mathbf{R}^{-1} (y^o - \mathcal{H}(m)) + \frac{1}{2} (m - \bar{m}^f)^T \mathbf{P}_e^{f-1} (m - \bar{m}^f), \tag{2.79}$$

which further underlines the link with least-squares optimization schemes discussed in Chapter 1: in the Bayesian linear setting, the ETKF analysis is equal to the KF analysis, which has been shown to be equivalent to Newton's method of optimization (Humpherys et al., 2012) in the least-squares optimization setting.

To solve this minimization problem and compute the analysis state, the singularity of \mathbf{P}_e^f in equation (2.79) must be dealt with, as it is not invertible under that formulation. We can see however, that the column space \mathcal{S} of \mathbf{P}_e^f is the same column space spanning \mathbf{M}^f , the *ensemble subspace*, on which $(\mathbf{P}_e^f)^{-1}$ is well-posed. Because $(m - \bar{m}^f)$ also lies in the forecast ensemble perturbation subspace, the minimization of $\mathcal{C}(m)$ is well-posed over \mathcal{S} . Therefore, the appropriate coordinates system has to be used in order to project this minimization over the subspace \mathcal{S} .

To do so, there are two possible approaches: using the singular vectors of \mathbf{M}^f (Ott et al., 2004) or directly use the columns of \mathbf{M}^f to project the problem in the ensemble subspace (Hunt et al., 2007). We favor the second methodology that does not requires to solve the singular value decomposition of \mathbf{M}^f , which is generally more affordable.

We change the coordinate system by letting

$$m = \bar{m}^f + \mathbf{M}^f w \tag{2.80}$$

where $w \in \mathbb{R}^{N_e}$ is a weight vector to be determined. This is equivalent to consider that the deviation from the mean of the system state is given by a linear combination of the weight vector w and the forecast ensemble perturbations (Harlim and Hunt, 2007). We can re-write the variational cost function as

$$\begin{aligned}
 \mathcal{C}(\bar{m}^f + \mathbf{M}^f w) &= \frac{1}{2} (y^o - \mathcal{H}(\bar{m}^f + \mathbf{M}^f w))^T \mathbf{R}^{-1} (y^o - \mathcal{H}(\bar{m}^f + \mathbf{M}^f w)) + \\
 &\quad \frac{N_e - 1}{2} (\mathbf{M}^f w)^T [\mathbf{M}^f (\mathbf{M}^f)^T]^{-1} (\mathbf{M}^f w)
 \end{aligned} \tag{2.81}$$

which can be simplified as

$$\mathcal{C}(w) = \frac{1}{2}(y^o - \mathcal{H}(\bar{m}^f + \mathbf{M}^f w))^T \mathbf{R}^{-1}(y^o - \mathcal{H}(\bar{m}^f + \mathbf{M}^f w)) + \frac{N_e - 1}{2} w^T w \quad (2.82)$$

This change of coordinate system (on \mathcal{S}) allows solving the minimization problem over a N_e dimensional space, rather than a n dimensional space. From there, finding the weight vector w^a that minimizes equation (2.84), is equivalent to finding the $\bar{m}^a = \bar{m}^f + \mathbf{M}^f w^a$ that minimizes equation (2.79).

Following the linearization in equation (2.73), we can linearly approximate the measurement operator by considering

$$\mathcal{H}(\bar{m}^f + \mathbf{M}^f w) = \bar{y}^f + \mathbf{Y}^f w. \quad (2.83)$$

Then, the cost function is expressed as

$$\mathcal{C}(w) = \frac{1}{2}(y^o - \bar{y}^f + \mathbf{Y}^f w)^T \mathbf{R}^{-1}(y^o - \bar{y}^f + \mathbf{Y}^f w) + \frac{N_e - 1}{2} w^T w, \quad (2.84)$$

which minimum is given by equating its gradient to zero

$$\nabla \mathcal{C}(w) = (N_e - 1)w - (\mathbf{Y}^f)^T \mathbf{R}^{-1}(y^o - \bar{y}^f + \mathbf{Y}^f w) = 0. \quad (2.85)$$

Equation 2.85 has a solution under the form

$$w^a = \tilde{\mathbf{P}}^a (\mathbf{Y}^f)^T \mathbf{R}^{-1}(y^o - \bar{y}^f) \quad (2.86)$$

where $\tilde{\mathbf{P}}^a = ((N_e - 1)\mathbf{I} + \mathbf{Y}^{fT} \mathbf{R}^{-1} \mathbf{Y}^f)^{-1}$. In the state space, the ensemble is given by

$$\begin{aligned} \bar{m}^a &= \bar{m}^f + \mathbf{M}^f w^a, \\ \mathbf{P}_e^a &= \frac{1}{N_e - 1} \mathbf{M}^f \tilde{\mathbf{P}}^a (\mathbf{M}^f)^T. \end{aligned} \quad (2.87)$$

We finally obtain the same algebraic solution to the ETKF problem, as derived earlier. In this study, we chose to implement the ETKF rather than other alternatives for several reasons:

- Its deterministic nature that guaranteed better performance over the stochastic EnKF.
- Its variational formulation lets us establish a clear link with least-squares minimization.
- Its ease of implementation.
- Its low computational and memory burden.

In the next subsection, I briefly introduce the Maximum-Likelihood Ensemble Filter (Zupanski, 2005), which directly used the variational formulation of the ETKF to formulate an iterative EnKF variant. On top of its advantageous iterative nature and the possibility to formally consider non-linear observations, the MLEF is an important stepping stone of the EnKF literature, as it is the basis for several other iterative methods such as the Ensemble Randomized Maximum Likelihood filter (Chen and Oliver, 2012) and the Iterative Ensemble Kalman filter (Sakov et al., 2012).

2.2.5 Maximum-Likelihood Ensemble Filter

One of the limitations of the ensemble filters lies in the nature of their observation operator. The analysis of the EnKF and ETKF guarantee an optimal solution only if the observation operator is strictly linear (although we have seen that non-linear observation can be considered). When the observation operator is strongly non-linear, the behavior of these filters might be compromised. To deal with non-linearities, we have seen previously that we typically consider a first-order Taylor expansion of the non-linear operator, in place of using a linearized operator. Instead of this approximation, the Maximum-Likelihood Ensemble Filter (MLEF) of Zupanski (2005) propose a fully non-linear development, and then determine a solution in the reduced space, with a projection similar to the transformation performed in the ETKF.

As its name suggests, the MLEF optimal state estimate corresponds to the maximum of the posterior PDF (Zupanski et al., 2008). The MLEF thus differs from the methods we have seen so far that are focused on variance reduction (note that both state estimates should be strictly equivalent in the linear Gaussian case). While other ensemble methods provide the state estimate in a single iteration, the MLEF integrates an inner iterative minimization loop inside of its algorithm, which requires using a local-optimization method. While this methodology was not tested in the frame of this thesis work, the iterative nature of the MLEF makes it appealing to solve non-linear problems such as FWI.

The MLEF starts with defining a state vector m_0 , which represents the initial condition of the DA application. As with previous cases, this state vector lies in the state space \mathbb{R}^n and is associated to an ensemble of *initial perturbation* p_0^i , $i = 1, \dots, N_e$ spanning the ensemble subspace \mathcal{S} . Zupanski (2005) define these perturbations as the square-root error covariance, which is an $n \times N_e$ matrix, and makes it in essence, very similar to the ensemble perturbation matrix \mathbf{M} previously defined. From here, the initial ensemble is given by

$$m_0^{(i)} = m_0 + p_0^{(i)}, \quad i = 1, \dots, N_e \quad (2.88)$$

In a similar fashion than other deterministic EnKFs, the goal of the analysis is to evaluate the square root $(\mathbf{P}^a)^{1/2}$, by computing $p^{a(i)}$, $i = 1, \dots, N_e$. The forecast of the MLEF is expressed as

$$p_{k+1}^f{}^{(i)} = \mathcal{F}(m_k^{a(i)}) - \mathcal{F}(m_k^a) = \mathcal{F}(m_k^a + p_k^{a(i)}) - \mathcal{F}(m_k^a), \quad (2.89)$$

which ultimately corresponds to the forecast of a perturbed ensemble, to which we subtract the forecast state vector $\mathcal{F}(m_k^a)$: the forecast is only about projecting the ensemble of perturbations. Once the forecast is obtained, the goal of the analysis focuses on maximizing the posterior's likelihood. As with the ETKF, this is done by minimizing the BLUE cost function. From now on, we drop the time index for readability.

$$\mathcal{C}(m) = \frac{1}{2}(y^o - \mathcal{H}(m))^T \mathbf{R}^{-1}(y^o - \mathcal{H}(m)) + \frac{1}{2}(m - m^f)^T (\mathbf{P}^f)^{-1}(m - m^f) \quad (2.90)$$

where m^f is the forecasted state vector, and $\mathbf{P}^f = (\mathbf{P}^f)^{1/2}[(\mathbf{P}^f)^{1/2}]^T$. The covariance square root $(\mathbf{P}^f)^{1/2}$ appears to be at best of rank $N_e - 1$ and therefore a basis reduction strategy can be used to apply the MLEF. The state estimate update Δm^a is defined as

$$\Delta m^a = m^a - m^f \quad (2.91)$$

and lie in the ensemble subspace $\mathcal{S} \in \mathbb{R}^{N_e}$, which is defined by the column space of $\mathbf{P}^{f^{1/2}} = \{p^{f(1)}, p^{f(2)}, \dots, p^{f(N_e)}\}$. The analysis can be expressed as a linear combination between the forecast

perturbations and a coefficient vector $w \in \mathcal{S}$ and thus the state estimate update is expressed as

$$\Delta m^a = w_1 p^{f(1)} + w_2 p^{f(2)}, \dots, w_{N_e} p^{f(N_e)}. = \mathbf{P}^{f1/2} w \quad (2.92)$$

The covariance square-root act as a projection operator (similar to the ETKF) from the ensemble subspace back to the state space. The cost function increment is given by

$$\Delta \mathcal{C}(m) = \mathcal{C}(m + \Delta m) - \mathcal{C}(m) \quad (2.93)$$

for $\Delta m \in \mathcal{S}$, and since minimizing $\Delta \mathcal{C}(m)$ is equivalent to minimizing $\mathcal{C}(m)$ (Zupanski et al., 2008) the cost function can be expressed as

$$\begin{aligned} \mathcal{C}(m + \Delta m) &= \mathcal{C}(m) + (\Delta m)^T (\mathbf{P}^f)^{-1} (m - m^f) \\ &\quad - [\mathcal{H}(m + \Delta m) - \mathcal{H}(m)]^T \mathbf{R}^{-1} [y^o - \mathcal{H}(m)] \\ &\quad + \frac{1}{2} (\Delta m)^T (\mathbf{P}^f)^{-1} (\Delta m) \\ &\quad + \frac{1}{2} [\mathcal{H}(m + \Delta m) - \mathcal{H}(m)]^T \mathbf{R}^{-1} [\mathcal{H}(m + \Delta m) - \mathcal{H}(m)], \end{aligned} \quad (2.94)$$

and for a differentiable observation operator, we can write

$$\begin{aligned} \mathcal{C}(m + \Delta m) &= \mathcal{C}(m) + (\Delta m)^T (\mathbf{P}^f)^{-1} (m - m^f) \\ &\quad - (\Delta m)^T \left[\frac{\partial \mathcal{H}}{\partial m} \right]^T \mathbf{R}^{-1} [y^o - \mathcal{H}(m)] \\ &\quad + \frac{1}{2} (\Delta m)^T (\mathbf{P}^f)^{-1} (\Delta m) \\ &\quad + (\Delta m)^T \left[\frac{\partial \mathcal{H}}{\partial m} \right]^T \mathbf{R}^{-1} \left[\frac{\partial \mathcal{H}}{\partial m} \right] (\Delta m) \\ &\quad + \mathcal{O}(\|\Delta m\|^3), \end{aligned} \quad (2.95)$$

which is equivalent to the second order Taylor's expansion of $\mathcal{C}(m)$ in the vicinity of m (Zupanski et al., 2008). From there, it is possible to define \mathcal{Z} , a non-linear analog to the ETKF observation perturbation matrix \mathbf{Y} , by defining

$$\begin{aligned} \mathcal{Z}(m) &= [z^1(m), z^2(m), \dots, z^{N_e}(m)] \\ z^i &= \mathbf{R}^{-1/2} [\mathcal{H}(m + p^{(i)}) - \mathcal{H}(m)], \end{aligned} \quad (2.96)$$

from which we can define the gradient and Hessian of $\mathcal{C}(m)$

$$\nabla \mathcal{C}(m) = (\mathbf{P}^f)^{-1/2} (m - m^f) - [\mathcal{Z}(m)]^T \mathbf{R}^{-1/2} (y - \mathcal{H}(m)) \quad (2.97)$$

$$\nabla^2 \mathcal{C}(m) = \mathbf{I}_{N_e} + [\mathcal{Z}(m)]^T \mathcal{Z}(m), \quad (2.98)$$

where the gradient is a N_e dimensional vector and the Hessian is a $N_e \times N_e$ matrix, as the minimization is carried in the ensemble subspace. In Zupanski et al. (2008), the authors demonstrate that their formulation of the analysis can be expressed as a generalization of the Conjugate-Gradient and the BFGS optimization methods, with the added benefit of not requiring explicitly differentiable operators. The final step of the MLEF is to update the ensemble perturbations such that

$$(\mathbf{P}^a)^{1/2} = (\mathbf{P}^f)^{1/2} [\mathbf{I}_{N_e} + [\mathcal{Z}(m^a)]^T \mathcal{Z}(m^a)]^{-1/2}, \quad (2.99)$$

which satisfies

$$\mathbf{P}^a = (\mathbf{P}^f)^{1/2} [\mathbf{I}_{N_e} + [\mathcal{Z}(m^a)]^T \mathcal{Z}(m^a)]^{-1} (\mathbf{P}^f)^{T/2}. \quad (2.100)$$

The MLEF analysis is very similar to what we had derived for the ETKF equations, except that it allows considering a non-linear observation operator. With minor modification to the misfit function, Fletcher and Zupanski (2006) also showed that the MLEF could handle log-normal PDFs, which makes it work even in non-Gaussian context (though it is unable to tackle multimodal PDFs).

While the MLEF has proven to be a great candidate for NWP applications (Carrassi et al., 2008), its implementation complexity made it less attractive than the ETKF, for the prospective work we ought to perform in this thesis work. Nonetheless, I believe that this method, which makes a clear bridge between statistical DA methods and iterative local optimization approaches, is a good indicator that both DA and FWI methodologies can be bridged together to produce more robust FWI tools. The MLEF ability to perform well under non-linear regimes while delivering a systematic uncertainty quantification would be highly desirable for FWI applications.

Before moving on to the next section, I would like to briefly introduce the work of Iglesias et al. (2013a), which seems to have stayed under the radar of the DA community. While I have not found the time to play with their ideas, their work has grasped my attention, as it constitutes a true *missing link* between statistical DA and optimization: the Ensemble Kalman Inverse (EKI).

2.2.6 Ensemble Kalman Inverse

The main idea underlying the EKI is to use the EnKF formalism, and adapt it as an iterative minimization method, that share the same kind of restrictions than local optimization methods. Recalling the fundamental inverse problem introduced in Chapter 1, the EKI ought to find the model parameters m from a set of observations y^o , related by

$$y^o = \mathcal{H}(m) + \epsilon, \quad (2.101)$$

where, as in Chapter 1, $\mathcal{H} : \mathbb{R}^n \rightarrow \mathbb{R}^d$ is the forward map from the state space to the measurement space (equivalent to the model and data space from our early definition of inverse problems) and $\epsilon \sim \mathcal{N}(0, \mathbf{R})$ is a random noise vector with $\epsilon \in \mathbb{R}^d$. As the solution of the inverse problem is sought within the EnKF framework, they consider a regularized problem with the misfit function under the form of

$$\mathbf{C}(m) = \frac{1}{2}(y^o - \mathcal{H}(m))^T \mathbf{R}^{-1}(y^o - \mathcal{H}(m)) + \frac{1}{2}(m - \bar{m})^T \mathbf{P}_{prior}^{-1}(m - \bar{m}), \quad (2.102)$$

where the second term is a regularization term introducing prior information on m . The effect of this regularisation is to introduce a finite-dimensional subset \mathcal{A} in which the solution of the inverse problem is sought (Iglesias et al., 2013a). The EKI is able to solve both dynamic and static inverse problems but for the sake of simplicity and to keep notations concise, I will focus only in the static case, where there is only one instance of d_{obs} .

To perform the minimization, the EKI requires to augment the space, in order to create an artificial dynamic system on which the EKI evolves. This state augmentation resembles the joint state-observation space that is considered in the Ensemble Adjustment Filter of Anderson (2001). The extended space is defined as \mathbb{R}^{n+d} along with the artificial dynamics mapping $\Upsilon : \mathbb{R}^{n+d} \rightarrow \mathbb{R}^{n+d}$ as

$$\Upsilon(z) = \begin{pmatrix} m \\ \mathcal{H}(m) \end{pmatrix} \quad \text{for} \quad z = \begin{pmatrix} m \\ y \end{pmatrix}. \quad (2.103)$$

where z is the augmented vector of length $n + d$ containing both the state parameters m and the observations y . The artificial dynamics of the EKI is defined as

$$z_k = \Upsilon(z_{k-1}), \quad (2.104)$$

and the observation related to the artificial dynamics are as follows

$$y_k = \mathbf{O}(z_k) + \epsilon_k \quad (2.105)$$

where $\mathbf{O} : \mathbb{R}^{n+d} \rightarrow \mathbb{R}^d$ is defined as a projection operator such that $\mathbf{O} = (O, \mathbf{I}_d)$ with \mathbf{I}_d being an Identity matrix of size d , and $\epsilon_k \sim \mathcal{N}(0, \mathbf{R})$.

At each iteration of the EKI, an ensemble of models spanning the subspace \mathcal{A} are updated by combining the artificial dynamics, and a perturbed version of the single observation via the EnKF analysis equation. Perturbed observations are used to drive the search for the solution inside of \mathcal{A} as it allows to "move around" the subspace. From the extended state vector, the state estimate is given by

$$\bar{m}_k = \frac{1}{N_e} \sum_{i=1}^{N_e} m_k^{(i)} = \frac{1}{N_e} \sum_{i=1}^{N_e} \mathbf{O}^\perp z_k^{(i)}. \quad (2.106)$$

where $\mathbf{O}^\perp : \mathbb{R}^{n+d} \rightarrow \mathbb{R}^n$ is the reciprocal projection operator mapping the extended space back to the state space and is given by $\mathbf{O}^\perp = (\mathbf{I}_n, 0)$. Finally, given an ensemble of size N_e , each iteration of minimization is given as

Forecast

$$\begin{aligned} \mathbf{z}_k^f &= \Upsilon(\mathbf{z}_{k-1}^a) \\ \bar{z}_k^f &= \frac{1}{N_e} \sum_{i=1}^{N_e} z_k^{f(i)} \\ \mathbf{P}_k^f &= \frac{1}{N_e - 1} (\mathbf{z}_k^f - \bar{z}_k^f)(\mathbf{z}_k^f - \bar{z}_k^f)^T, \end{aligned} \quad (2.107)$$

Analysis

$$\begin{aligned} \mathbf{K}_k &= \mathbf{P}_k^f \mathbf{H}^T (\mathbf{H} \mathbf{P}_k^f \mathbf{H}^T + \mathbf{R})^{-1} \\ z_k^{a(i)} &= z_k^{f(i)} + \mathbf{K}_k \left[y_k^{(i)} - \mathbf{O}(z_k^{f(i)}) \right], \end{aligned}$$

where $y_k^{(i)}$ are instances of randomly perturbed observation generated $y_k^{(i)} = y^o + \epsilon_k^{(i)}$. Once the iteration is performed, the state estimate is computed according to equation (2.106), to check if the convergence criterion has been met. Once this optimization step has been performed, the Euclidean distance $\|y^o - \mathcal{H}(\bar{m}_k^a)\|^2$ should be reduced. The proposition in Iglesias et al. (2013a) is then followed by several demonstrations that show

- The solution to the EKI is strictly equivalent to the least-squares problem' solution under linear and Gaussian assumptions.
- The EKI is a derivative-free, iterative minimization technique.
- The EKI is able to tackle non-linear inverse problems.

Note that the original paper did not thoroughly investigate the EKI capacities in terms of uncertainty estimation. It seems, however, that despite being able to yield an accurate state estimate, the uncertainty estimation is biased, as the ensemble members steadily collapse toward the solution in the misfit space, losing statistical information (Iglesias et al., 2013b; Chada et al., 2018). Nonetheless, as with the MLEF, I believe this type of method requires careful investigation, as it may be a solution to FWI lack of systematic uncertainty estimation methods.

In this section, I have mainly talked about the advantages of ensemble Kalman filtering methods and focused on their theoretical foundations to underline the characteristics they share with numerical optimization. It is now due time to discuss their limiting factors and drawbacks.

2.3 Limits of EnKF methods

The statistical DA tools that we have presented up to now are all variants or generalization of the BLUE to specific problems. While the KF generalizes the BLUE for linear dynamical problems, the ensemble methods introduce a generalization for large scale problems, yielding an exact solution for the linear Gaussian case. They also allow for weakly non-linear forecast and observation operators, but in this case, their state estimate is no-longer optimal and does not yield the BLUE solution (Le Gland et al., 2009).

As they are all tied to the BLUE, these tools are strictly limited, in theory, to Gaussian forecast and measurement errors. In that sense, they are bound to explore a subspace \mathcal{A} of the solution space Ω , the same way local optimization methods are limited by quadratic assumptions (see Chapter 1). Hence, the first significant limitation of EnKFs is that they are not fit to quantify uncertainty in the non-Gaussian case (though their success in NWP would tend to advocate otherwise). Despite their intrinsic statistical nature, ensemble DA methods are closer to local optimization approaches than from global optimization methods.

The second limitation of ensemble methods lies in their low-rank approximations. While the ensemble representation is very advantageous to carry computations over the ensemble subspace, small ensembles are bound to generate errors and biases in EnKF applications. When an ensemble is too small to represent the system in a statistically meaningful manner, we say it is affected by *undersampling*.

2.3.1 Undersampling characterization

One of the primary challenges of any ensemble DA method is to ensure that the ensemble stays statistically meaningful over cycles while $N_e \ll n$. In that context, maintaining a representative ensemble becomes challenging, and the filter might be affected by undersampling (Guzzi, 2015).

I have stated earlier that the forecast error was generally poorly known, and therefore is often neglected in practical applications. This causes a general tendency to underestimate the forecast error in DA systems, and this problem is exacerbated when the ensemble is too small (Furrer and Bengtsson, 2007). Indeed, while the forecast errors are underestimated cycle after cycle, artificial confidence is placed in the forecast. Moreover, because the Kalman gain is given as a ratio of forecast error over measurement error, this tends to severely bias the filter. This overconfidence in the forecast means that the observations are neglected during the analysis, and the analysis update is likely not to represent the system accurately. Undersampling is also responsible for spurious correlation terms appearing in the off-diagonal terms of the covariance matrix as can be seen in Figure 2.6.

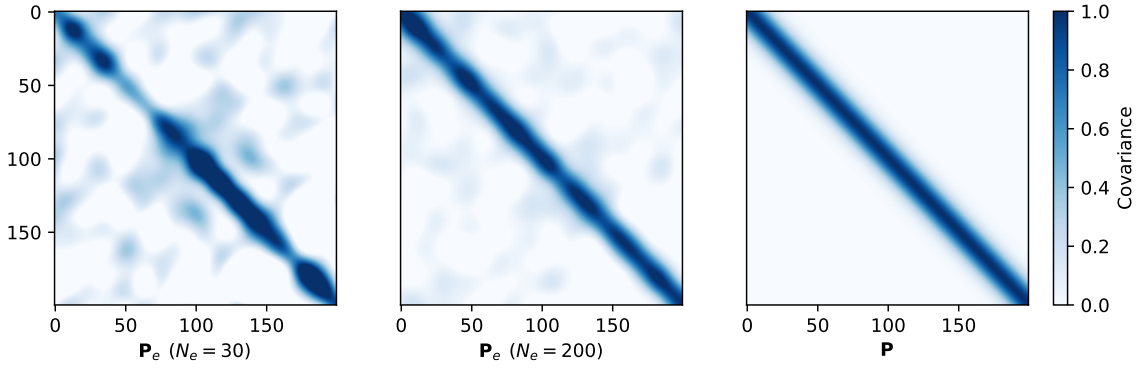


Figure 2.6: Numerical example of undersampling. The first two panels are Gaussian covariances estimated with respectively $N_e = 30$ and $N_e = 200$, for a problem with $n = 200$. The true covariance is shown in the third panel for reference. In this case, undersampling is responsible for spurious terms appearing on both diagonal and off-diagonal terms. Note that even with $N_e = n$ the sampling is not adequate to perfectly approximate the true covariance.

The obvious solution to solve the undersampling issue in EnKF approaches would be to increase significantly the number of ensemble members, which is more often than not computationally impossible. Unfortunately, as long as $N_e \ll n$, which is typically the case in ensemble filtering, ensembles are bound to be undersampled, which in turns can cause three types of problems. These three common problems of ensemble filtering are known as filter divergence, inbreeding, and spurious distant correlations.

2.3.2 Inbreeding

Inbreeding was first introduced by Houtekamer and Mitchel (1998) and describes an undersampling bias associated with variance underestimation. Inbreeding often refers to the occurrences where the analysis error covariance is underestimated during the analysis update. As mentioned previously, this is mostly due to the underestimation of forecast error. By definition, the analysis error covariance should always be smaller than the forecast error covariance (as it is given by $\mathbf{P}_e^a = (\mathbf{I} - \mathbf{KH})\mathbf{P}_e^f$) which becomes problematic as it can be underestimated (Furrer and Bengtsson, 2007). Step after step, the forecast and analysis error covariances are then reduced, which can result in the ensemble collapsing and losing all statistical meaningfulness.

The subsequent imbalance caused by inbreeding will skew the analysis in favor of the forecast, as forecast uncertainty get reduced, causing artificial overconfidence in the forecast step.

Inbreeding occurs when the ensemble subspace does not span the model subspace adequately to represent the state estimation errors (Petrie, 2008). The smaller the ensemble is, the smaller its subspace will be, which increases the chances of misrepresenting the system state. Another source of inbreeding is found in EnKF formulations that perturb the analyzed observations (such as Evensen (1994) original formulation). The sampling error caused by these perturbed observations increases the odds of inbreeding in stochastic EnKFs (Whitaker and Hamill, 2002). For this reason, the choice of deterministic EnKFs over the stochastic EnKF can be justified, as this rules-out one of the sources of inbreeding. In DA, inbreeding is known to be responsible for filter divergence (Houtekamer and Mitchel, 1998).

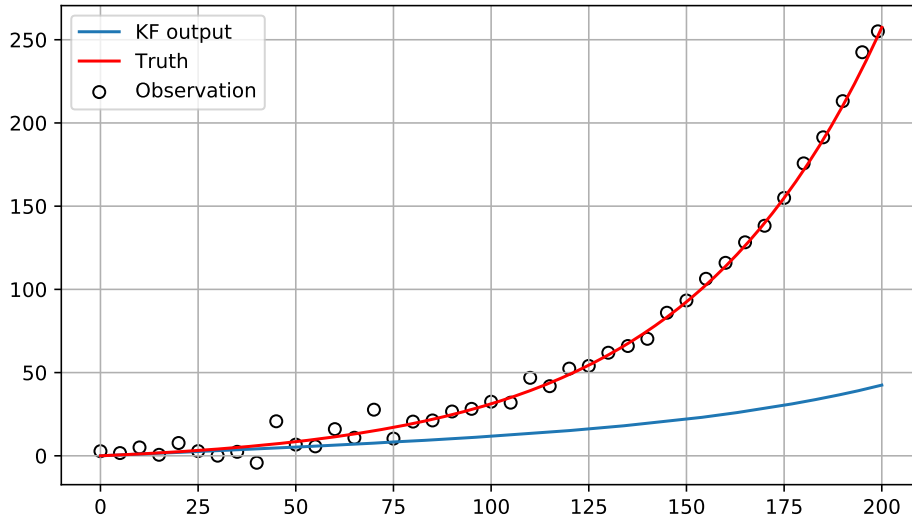


Figure 2.7: Numerical example of filter divergence on a 1-D KF. In this case, the filter divergence is the result of overconfidence in the forecast. This overconfidence results in a cutback of the analysis update, and observations are almost ignored. In this example, I purposely underestimated the process noise matrix \mathbf{Q} to demonstrate how it can affect the filter. In the case of the EnKF, this overconfidence could come from inbreeding.

2.3.3 Filter divergence

As Lorenz (1963) demonstrated in his publication titled “*Deterministic Nonperiodic Flow*,” which introduced his famous *attractor*, small changes of initial conditions in a dynamic system model can lead to instability of the solution. Therefore, when the analysis is strongly biased by inbreeding (or other sources of error such as numerical rounding errors), the analysis state might be disconnected with real system state. Subsequently, the filter becomes overconfident in an incorrect state estimate and becomes unable to correct it with observations: the EnKF is reduced to a succession of model forecasts, that are likely to diverge farther from the truth at each cycle. This is especially true for a chaotic dynamic model such as the atmospheric models used in NWP.

While inbreeding is one of the primary sources of filter divergence, Hamill et al. (2001) suggested that both variance underestimation and strong cross-covariance terms (off-diagonal of \mathbf{P}_e^f) could lead to filter divergence. It also has been shown that biases in forecast error covariance tend to reduce the spread of the ensemble after the analysis (van Leeuwen, 1999; Furrer and Bengtsson, 2007; Crystalng et al., 2011). Note that, as filter divergence denotes a filter’s inability to produce meaningful state estimates, the KF can also be subjected to divergence (Figure 2.7). This can happen when studying non-linear systems, or when the forecast covariance matrix \mathbf{P}^f has a high-condition number (Houtekamer and Zhang, 2016).

2.3.4 Spurious Correlation

The last type of bias that can be introduced by undersampling is known as *spurious distant correlations*. We have shown an effect of this type of artifact in Figure 2.6, that arise because of sampling errors. The correlations are said to be spurious when the forecast cross-covariance terms relate two parameters, while we know there is no physical relationship between them, or while they are apart a significant distance.

An example of such spurious correlation in NWP would be the forecast covariance linking together the pressure field in Brazil, with temperature evolution in Tokyo (which are likely to be unrelated due to their distance).

Spurious correlations are also defined as unphysical updates during the analysis, driven by long-distance observations, or between variables that are known to be uncorrelated or decoupled (Anderson and Anderson, 1999; Evensen, 2009; Petrie and Dance, 2010). Due to spurious correlations, any observation has the potential to detrimentally impact state variables that are remote from the measurement point in the state space. Typically, that problem arises when the spurious correlations in the forecast covariance estimate are more significant than the true correlation links between the state parameters (Hamill et al., 2001). Hamill et al. (2001) also show the apparent relation between ensemble size and spurious correlations, as larger ensemble exhibit less sensitivity to sampling noise, and therefore are less impacted by spurious correlation terms (see Figure 2.6). Sampling theory indicates that the spurious noise in \mathbf{P}_e^f (and \mathbf{P}_e^a) is proportional to $1/N_e$, which tends to indicate that adding ensemble members should quickly mitigate the spurious correlation problem, and is consistent with Hamill et al. (2000)'s observations: as with inbreeding, the effect of spurious correlations is strongly dependent on the size of the ensemble.

Because of their potential severe effects on ensemble analysis, spurious correlations and inbreeding have been a hot topic of research since the advent of ensemble methods in DA, in an attempt to prevent filter divergence in practical applications.

2.3.5 Solutions to undersampling

I introduced in this section, the two primary sources of filter divergence, which are inbreeding and spurious distant correlations. These two biases have a very different effect on the state representation and therefore call for different methods of mitigation.

Covariance inflation

As mentioned previously, inbreeding takes its origin in underestimation of covariance error in the forecast, which then results in overconfidence in the forecast state. In order to mitigate inbreeding, one can artificially increase the forecast error, to account for its almost certain underestimation. This method, called covariance inflation, was introduced by Anderson and Anderson (1999). This technique artificially "inflate" the forecast deviation to the mean (the distance between each ensemble members, encapsulated in \mathbf{M}^f), by multiplying the forecast covariance with a constant factor r , slightly larger than 1.0. The inflated perturbation matrix \mathbf{M}_i^f is given by

$$\mathbf{M}_i^f = r(\mathbf{m}^f - \bar{\mathbf{m}}^f), \quad (2.108)$$

and thus the inflated ensemble, used in the analysis becomes

$$\mathbf{m}_i^f = \bar{\mathbf{m}}^f + \mathbf{M}_i^f. \quad (2.109)$$

Note that for the ETKF, the multiplicative inflation can be conveniently implemented by modifying the formulation of $\tilde{\mathbf{P}}^a$ as

$$\tilde{\mathbf{P}}^a = \mathbf{T}\mathbf{T}^T = \left[\frac{(N_e - 1)}{r^2} \mathbf{I}_{N_e} + (\mathbf{Y}^f)^T \mathbf{R}^{-1} \mathbf{Y}^f \right]^{-1}. \quad (2.110)$$

This formulation differs from the EnKF inflation, as the ensemble is inflated in the space spanned by the observation perturbation matrix \mathbf{Y}^f . It also has the advantage of having a computational complexity in $\mathcal{O}(N_e)$ rather than in $\mathcal{O}(n \times N_e)$.

The value of r is usually chosen empirically or by trial and error: its optimal value is highly dependent on the size of the ensemble (Hamill et al., 2000). While the optimal value cannot be expressed analytically, the recommendations in the literature are calling for values between 1% and 7% inflation (Hamill et al., 2000; Anderson, 2001; Whitaker and Hamill, 2002), with r varying depending on the implementation of EnKF (Evensen's formulation, the EnSRFs or other deterministic ensemble filters). Whitaker and Hamill (2002)' study also indicates that the choice of r might be highly dependent on the nature of the dynamic forecast operator: some dynamical systems might be more sensitive to the accumulation of errors, while some might have long error decay with time. Note also that a few methodologies have proposed to tune the inflation parameter based on maximum likelihood estimation (Miyoshi, 2011a) or based on Desrozier's diagnostic (Desroziers et al., 2005; Li et al., 2009). A long-standing history of applying multiplicative inflation factors to EnKF can be found in the literature (Anderson and Anderson, 1999; Hamill et al., 2001; Whitaker and Hamill, 2002; Oke et al., 2007; Anderson, 2007, 2009; Li et al., 2009; Miyoshi, 2011a; Bocquet and Sakov, 2012).

Despite being the most prevalent method in the literature, multiplicative inflation is not the only strategy available to deal with inbreeding. Instead of inflating the ensemble by a constant factor, it is also possible to add additive-noise either before or after the analysis such that each ensemble members become

$$m^{f(i)} \leftarrow m^{f(i)} + d^{(i)}. \quad (2.111)$$

where $d^{(i)}$ are random vectors chosen so that $d^{(i)} \sim \mathcal{N}(0, \mathbf{D})$, with \mathbf{D} being an approximation of the process noise matrix \mathbf{Q} , and \leftarrow denotes the value replacement. More details are given in Mitchell and Houtekamer (2000), and examples of applications can be found in Houtekamer et al. (2005); Hamill and Whitaker (2005); Houtekamer and Mitchell (2005); Hamill and Whitaker (2011).

Both additive and multiplicative inflation have been implemented together as a "hybrid" inflation scheme by Zhang et al. (2004); Whitaker and Hamill (2012), where the ensemble is inflated following

$$\begin{aligned} \mathbf{M}^a &\leftarrow (1 - r)\mathbf{M}^a + r\mathbf{M}^f \\ \text{or} \\ \sigma^a &\leftarrow (1 - r)\sigma^a + r\sigma^f \end{aligned} \quad (2.112)$$

where σ^a is the analysis ensemble standard-deviation, and σ^f is the forecast ensemble standard-deviation. Note that in these hybrid inflation methods, the inflation is carried over the analysis ensemble rather than on the forecast ensemble. A more extensive review of covariance inflation can be found in Luo and Hoteit (2013), where the authors list a few more methodologies, having a different philosophy than additive and multiplicative inflation methods that do not seems to be adopted to the same extent. While covariance inflation methodologies tend to increase the stability of ensemble filtering methods, they do not provide a solution to the second undersampling bias. To correct for spurious distant correlations, one has to rely on *localization* methods.

Localization methods

Localization is a way to cut-off distant spurious correlation in the ensemble forecast covariance, according to a specific state-space cut-off distance. This method allows considering smaller ensembles, even though

we know undersampling will introduce strong spurious cross-covariance terms in the ensemble estimate of \mathbf{P}^f . The purpose of localization is to mask-out any unwanted distant-correlation terms in the forecast covariance with either a taper function centered around the diagonal (*covariance localization* - CL), or by excluding remote data from the parameter that is being updated during the analysis (*local analysis* - LA). Instead of muting explicitly off-diagonal terms of the covariance matrix, the local analysis approach implies that the analysis step is conducted sequentially, one state parameter at a time, accounting exclusively for a few selected neighboring data. While these two approaches tackle the localization problem very differently, Sakov and Bertino (2011) showed that both yield comparable results on "weak" assimilation (when \mathbf{R} is sufficiently high, and the analysis variance reduction is not too drastic, i.e., when the system is far from optimality). The choice between CL and LA is much more a question of scalability and ease of implementation rather than a question of performances, as both methods should be equivalent in practical applications.

It has to be noted that, while both forms of localization can reduce the spatial domain of influence of observation during the analysis, they only make sense if the observations are *local*. This is the case of in-situ measurement (that are attributed to a certain point of the state space, or rather a group of points of the state space). Non-local observation such as radiative transfer data (or waveform data in the context of FWI) cannot benefit from localization, as they are often measured at the boundary of the state space: a large chunk potentially influences them (if not all) of the state space and a "localization distance" does not make sense in this context.

Covariance localization

Introduced by Houtekamer and Mitchel (1998); Hamill et al. (2000); Whitaker and Hamill (2002), the covariance localization method (also known as covariance filtering) ought to mask-out off-diagonal correlation terms (that can be spurious or not) by taking its Hadamard product (Horn and Johnson, 2012) with a local-support correlation matrix ρ (Fig 2.8). The filtered covariance is given by

$$\mathbf{P}_l = \rho \circ \mathbf{P}_e, \quad (2.113)$$

where \circ denotes the Hadamard product (an element-wise product between two matrices), and the subscript l denotes the localized nature of the matrix. The correlation matrix is defined as a band of non-zero values centered on the leading diagonal. The correlation coefficients are set to 1 at the diagonal, and the coefficient values decay to 0 at a given distance of the diagonal (Petrie and Dance, 2010).

This correlation matrix is often built on the Gaspari-Cohn function (Gaspari and Cohn, 1999) which is a 5th order piecewise rational function such that each line of ρ is given by

$$\rho^{(i)} = \begin{cases} -\frac{1}{4} \left(\frac{z}{c}\right)^5 + \frac{1}{2} \left(\frac{z}{c}\right)^4 + \frac{5}{8} \left(\frac{z}{c}\right)^3 - \frac{5}{3} \left(\frac{z}{c}\right)^2 + 1 & 0 \leq z \leq c, \\ \frac{1}{12} \left(\frac{z}{c}\right)^5 - \frac{1}{2} \left(\frac{z}{c}\right)^4 + \frac{5}{8} \left(\frac{z}{c}\right)^3 + \frac{5}{3} \left(\frac{z}{c}\right)^2 - 5 \left(\frac{z}{c}\right) + 4 - \frac{2}{3} \left(\frac{z}{c}\right)^{-1} & c < z \leq 2c, \\ 0 & z > 2c, \end{cases} \quad (2.114)$$

where z is the euclidean distance between two state parameters or between an observation point and a state parameter location, and c is the length scale defined as the localization radius (or filtering length scale). It is defined such that when $z > 2c$ the correlation terms vanish, and is problem-dependent. Equation (2.114) yields a positive semidefinite matrix that conveniently approximates a Gaussian function, but falls at zero at a finite distance.

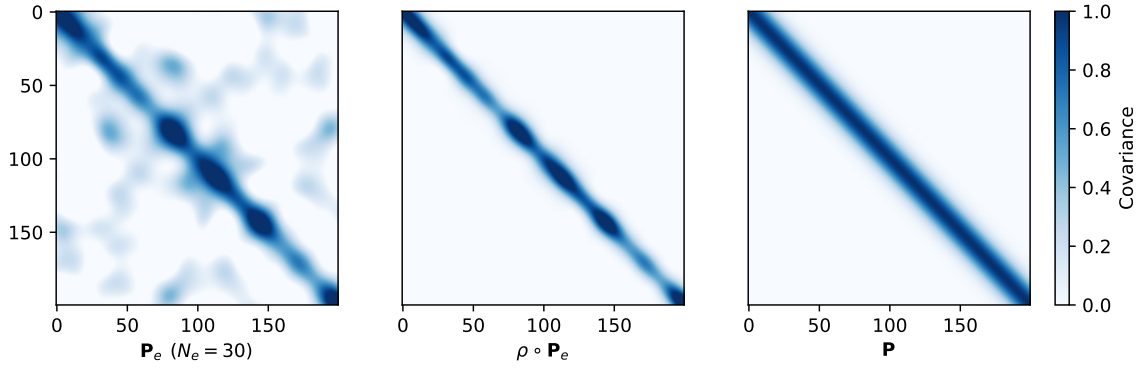


Figure 2.8: Numerical example of covariance localization over an ensemble covariance matrix \mathbf{P}_e with $N_e = 30$ (first panel). The localization is performed in the second panel, where a localization operator ρ is applied directly to the ensemble covariance matrix through Hadamard's product. The localization operator is generally defined as a matrix containing correlation functions with local support centered around the diagonal. In our ideal case, ρ was set to be a Gaussian function which decay corresponds to the true covariance matrix \mathbf{P} (third panel).

Implementation wise, the Hadamard product is generally carried over the Kalman gain computation such that the localized Kalman gain \mathbf{K}_l is given by

$$\mathbf{K}_l = \left[(\rho \circ \mathbf{P}_e^f) \mathbf{H}^T \right] \left[\mathbf{H} (\rho \circ \mathbf{P}_e^f) \mathbf{H}^T + \mathbf{R} \right]^{-1}. \quad (2.115)$$

Given that ρ has a regular structure and when the observation operator \mathbf{H} is linear, it is much more efficient to consider the localized Kalman gain as

$$\mathbf{K}_l = \left[\rho \circ (\mathbf{P}_e^f \mathbf{H}^T) \right] \left[\rho \circ (\mathbf{H} \mathbf{P}_e^f \mathbf{H}^T) + \mathbf{R} \right]^{-1}, \quad (2.116)$$

which is how Houtekamer and Mitchell (2001) chose to implement the Hadamard product.

The choice of the filtering length scale is particularly important in ensuring the stability of the filter. While the spurious correlation should be tapered-out of the covariance matrix, the true physical correlations of the system must be preserved. If the filtering length scale is too long, the localization will do a poor job at removing unwanted terms, and spurious correlation will remain in the covariance matrix. On the other hand, if the filtering length scale is set too short, meaningful correlation information might be removed from the system and not accounted for during the analysis. Defining the optimal length scale is, unfortunately, mostly based on trial and error.

Additionally, performing the localization through the Hadamard product has the benefit of increasing the rank of the covariance (Hamill et al., 2001; Oke et al., 2007), which increases the number of directions spanned by the ensemble subspace beyond $N_e - 1$. By increasing the effective ensemble size, it is expected that the analysis will yield better results (Oke et al., 2007). Finally reducing most entries of \mathbf{P}_e^f to zeros also results in making the covariance matrix highly sparse, which allows the use of sparse matrix solvers and increase the computational efficiency of the analysis step (Lorenc, 2003).

However, this localization technique is formulated with explicit covariance matrices manipulations and therefore is not applicable to schemes such as the ETKF, where manipulations of matrices square

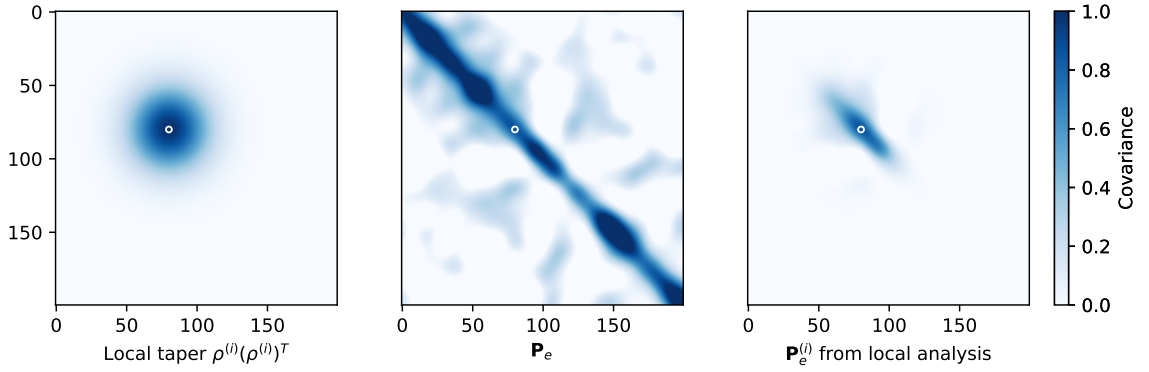


Figure 2.9: Effect of the Local Analysis on the ensemble covariance matrix. The i^{th} line of the localized covariance matrix is equivalent to the taper function applied to the ensemble covariance \mathbf{P}_e .

roots are preferred. While Petrie (2008); Petrie and Dance (2010) have attempted to formulate an approximation of the Hadamard product for the ETKF by expressing the localization as

$$\rho \rho^T \circ \mathbf{P}_e = \rho \rho^T \circ \frac{1}{N_e - 1} \mathbf{M} \mathbf{M}^T \approx \frac{1}{N_e - 1} (\rho \circ \mathbf{M})(\rho \circ \mathbf{M})^T, \quad (2.117)$$

this approximation has proven to yield poor results. To this day, no one has been able to express the localization of the ETKF in terms of the Hadamard product.

2.3.5.1 Local analysis

LA is the second prevalent method in localization and has been introduced by Evensen (2003); Anderson (2003); Ott et al. (2004). It relies on performing the analysis sequentially, one state parameter at a time, by selecting only neighboring observation points around the analyzed parameter (Fig. 2.9). The advantage of LA is that it is not scheme-dependent as CL, and has been implemented in the ETKF (Ott et al., 2004; Harlim and Hunt, 2005; Hunt et al., 2007), which has been granted the name LETKF by Hunt et al. (2007) for "Local" ETKF.

The analysis of the LETKF is based on the equations of the regular ETKF, except that y^o, \bar{y}^f and \mathbf{Y}^f are truncated to include observation in a local region around a grid point, while \bar{m}^f and \mathbf{M}^f are truncated to include only the neighboring parameters. Harlim and Hunt (2005) underline the importance of the consistency between analyzed parameters and recommend that physically close state parameters should be analyzed with a similar set of observations. If measurements are dense over the state-space, observations for each state parameters should overlap sufficiently to ensure a consistent analysis: it should not introduce discontinuities in the analysis state estimate. To further ensure smoothness of the analysis state, the effect of observations over the state-space is tapered, enforcing a decay to zero as the distance from the observation point grows. Hunt et al. (2007) remarked that this could be achieved by applying a correlation function to the inverse of the measurement noise matrix such that

$$\mathbf{R}_l^{-1} = \rho \circ \mathbf{R}^{-1}. \quad (2.118)$$

This has the effect of smoothly increasing the uncertainty over the local region defined by the function ρ until the uncertainty is set to infinity over a certain distance (from which there is no effect during the analysis).

As with the inflation parameters, localization (be it CL or LA) still requires manual tuning, and its effects on EnKFs optimality have often been neglected in practical applications (Sakov and Bertino, 2011). Nonetheless, localization is now commonly used at an operational level in NWP and is still an active topic of research (Bocquet, 2016; Bishop et al., 2017; Bocquet and Farchi, 2019), with an emphasis on localizing remote sensing observations, that are inherently non-local. Though the perspective of applying localization on non-local observations is an exciting prospect, it is difficult to say if these methods will have uses outside of the NWP and atmospheric sciences fields.

Conclusion

In this chapter, I gave an overview of the popular statistical DA methods. I introduced the theory from the ground-up, and the EnKF and some of its interesting variants were discussed. Apart from practicality issues (undersampling mitigations), most of this chapter has been focused on the BLUE and variational forms of the EnKFs: we have seen that when the conditions are right, these methods' analysis yield the solution to an inverse problem in which we try to reduce the least-squares distance between observed and forecast data.

The purpose of this chapter was both to introduce DA theory, and to underline the connexions between DA and local optimization methods. Now that these connexions have been shown, we have some reasons to believe that ensemble DA methods are an adapted option to bring systematic uncertainty estimation in geophysical imaging, and particularly to FWI applications. In the next chapter, I will evaluate the possibilities that we have established to pair DA and FWI. This includes the various options to define the model and observation spaces, but also in defining the forecasting operator for a tomographic (and therefore static) problem, or on the type of sequential approaches that can be set-up to ensure the stability of such a mixed tool. As when building a KF for a specific problem, we will see there are many filter parameterization that might be adequate to solve our problem.

Chapter 3

Combining DA and FWI

Contents

3.1	Proposition 1 - A dynamic formulation	72
3.2	Proposition 2 - Extending the state-space: The WRI analog	73
3.3	Proposition 3 - Extending the state-space: The EKI analog	74
3.4	Proposition 4 - Extending the state-space: adding adjoint	74
3.5	Proposition 5 - A simple adjoint scheme	75
3.5.1	FWI as a dynamic problem: defining a dynamic proxy	76
3.5.2	The ETKF-FWI scheme	77
3.5.3	ETKF-FWI sampling strategy	80

With the introduction of both FWI and DA's theories and the apparent links between the two frameworks (united by the BLUE formulation), we can introduce possible ways of combining DA and FWI to produce uncertainty estimation. While I have presented some of the most popular DA tools in Chapter 2, I wish to keep the DA parameterization general: once the problem has been formulated (by defining the state vector, the observations, the dynamical and observations operators), the means of solving the problem are interchangeable.

Following these DA parameterization propositions, I will review the practical implementation of the ETKF coupled with FWI, and discuss the importance of the initial ensemble building within this formulation.

Notations

In this chapter, we re-consider the FWI tomographic problem and re-introduce the following notations:

- n is the number of degrees of freedom (or model parameters), d is the length of the observation vector, and l is the number of gridpoints on which the forward modeling is discretized.
- $m \in \mathbb{R}^n$ is the subsurface model with n parameters. Note that n does not necessarily denote the number l of modeling gridpoints.
- $d_{obs} \in \mathbb{R}^d$ is the vector containing the observed data used in the FWI optimization problem.

- $d_{cal} \in \mathbb{R}^d$ is the vector containing the synthetic data used in the FWI optimization problem.
- $u \in \mathbb{R}^l$ is a vector containing the propagating wavefield, obtained by solving the wave equation for the model m . Note that in the frequency domain, $u \in \mathbb{C}^l$ instead.
- $\mathcal{M} : \mathbb{R}^n \rightarrow \mathbb{R}^l$ ($\mathcal{M} : \mathbb{R}^n \rightarrow \mathbb{C}^l$ for the frequency domain equivalent) is a generic non-linear forward problem operator (with \mathcal{M} for modeling). Given a source term s , it allows to express and compute the wavefield as $u = \mathcal{M}^{-1}(m)s$.
- $\mathcal{H} : \mathbb{R}^n \rightarrow \mathbb{R}^d$ ($\mathcal{H} : \mathbb{R}^n \rightarrow \mathbb{C}^d$ for the frequency domain equivalent) is the FWI forward problem operator defined in equation (1.2). While \mathcal{M} yields the wavefield in the entire domain, \mathcal{H} gives synthetic seismograms at receiver location such that $\mathcal{H}(m) = d_{cal}(m)$.
- $\mathbf{E} : \mathbb{R}^l \rightarrow \mathbb{R}^d$ ($\mathbf{E} : \mathbb{C}^l \rightarrow \mathbb{C}^d$ for the frequency domain case) is the linear extraction operator that gives $d_{cal} = \mathbf{E}u$, by extracting the values of the wavefield at the receiver locations.
- $\mathcal{I}_i : \mathbb{R}^n \rightarrow \mathbb{R}^n$ is the FWI non-linear operator. Applied to the model m_{k-1} , this operator yields $\mathcal{I}_i(m_{k-1}, d_{obs,k}) = m_k$ where m_k is the optimized model after i non-linear inversion iterations.

3.1 Proposition 1 - A dynamic formulation

Earlier in Chapter 2, we saw that DA was designed to study the evolution of dynamic systems, by generalizing the BLUE to time-dependent problems. This property of DA could prevent us from applying its framework to FWI, as tomographic imaging is a "static" problem: from a collection of recorded seismic data, the subsurface image is built, without any kind of time-dependency. Even though there are special cases (geoessource pumping sites, geothermal plants, or active seismic faults), in which physical properties can change within a short timespan, these changes in medium properties are generally not fast enough to consider the construction of a tomographic image as a time-dependent problem.

For this reason, we may instead parameterize our DA system based on the dynamic part of the FWI workflow: the wavefield propagation in the medium. Because the wavefield's propagation in the subsurface is formulated as a time-dependent problem (when considering the time-domain wave equation), we can define a DA problem centered on forecasting the wavefield evolution through time. Doing so, we may consider the following parameterization:

- The state vector $u \in \mathbb{R}^l$, contains the wavefield at a given timestep, discretized over l gridpoints.
- $y^o = d_{obs} \in \mathbb{R}^d$. Most of the time, the only observations available in FWI are wavefield data, which defines our data assimilation observation vector.
- $\mathbf{H} = \mathbf{E}$. The linear observation operator corresponds to the extraction operator, responsible for extracting the values of the wavefield at the receiver locations.
- The forecast is defined as $\mathcal{F}(u) = \mathcal{M}^{-1}(m)s$, corresponding to a step of the numerical integration scheme that compute the wavefield's propagation given the model m .

This way, the forecast step corresponds to propagating the wavefield from a timestep to the next: it tries to "predict" the true wavefield from the current subsurface model. The predicted wavefield would

then be corrected during the analysis, from the discrepancies between the observation vector and the forecasted wavefield u^f .

The associated misfit functional in terms of the BLUE is given by

$$\mathcal{C}(u) = \frac{1}{2}(y^o - \mathbf{H}u)^T \mathbf{R}^{-1}(y^o - \mathbf{H}u) + \frac{1}{2}(u - u^f)^T (\mathbf{P}^f)^{-1}(u - u^f). \quad (3.1)$$

which is equivalent to solving the following minimization problem

$$\mathcal{C}(u) = \min_u \frac{1}{2} \|y^o - \mathbf{H}u\|^2 + \frac{1}{2} \|u - u^f\|^2. \quad (3.2)$$

The analysis would yield an "optimal" wavefield, lying in between d_{cal} and d_{obs} . Note that even though the forecast uncertainty \mathbf{P}^f is closely tied to the model uncertainty (as the wavefield is parameterized by the medium's properties), this connection would be hard to establish. Therefore, this formulation solely allows estimating the wavefield's uncertainty and would be unable to inform us about the subsurface model parameters and uncertainties.

This "time dependent" DA-FWI formulation does not match with our goal of providing uncertainty estimation on the tomographic images.

3.2 Proposition 2 - Extending the state-space: The WRI analog

To obtain information on the subsurface model m , a solution might be to consider an augmented state-space, containing both the wavefield and subsurface model. This formulation is conceptually very close to the Wavefield Reconstruction Inversion (WRI, van Leeuwen and Herrmann, 2013), from the FWI literature. WRI is an original take on the FWI problem, which tries to find a wavefield that minimizes the least-squares misfit for both the observed data and the forward operator \mathcal{F} parameterized by the subsurface model (the wave equation is an added constraint to the minimization problem).

We can, therefore, formulate the DA-FWI problem as

- $z = \begin{bmatrix} m \\ u \end{bmatrix}$, with $z \in \mathbb{R}^{n+l}$, is the joint model-wavefield state space.
- $y^o = d_{obs}$
- $\mathbf{H} = \begin{bmatrix} 0 & 0 \\ 0 & \mathbf{E} \end{bmatrix}$, with $\mathbf{H} : \mathbb{R}^{n+l} \rightarrow \mathbb{R}^d$ The observation operator extracts values of the wavefield at receiver locations.
- $\mathcal{F}(z) = \begin{bmatrix} \mathbf{I}_m & 0 \\ 0 & \mathcal{M}^{-1}(m)_s \end{bmatrix}$, with $\mathcal{F} : \mathbb{R}^{n+l} \rightarrow \mathbb{R}^{n+l}$ where the forecast operator is also a step of the numerical integration scheme that compute the wavefield's propagation given the model m and \mathbf{I}_m is the Identity.

This formulation is equivalent to solving the following minimization problem

$$\mathcal{C}(z) = \min_z \frac{1}{2} \|y^o - \mathbf{H}z\|^2 + \frac{1}{2} \|z - z^f\|^2. \quad (3.3)$$

This time, we ought to find the augmented state-vector that minimize the data misfit, while accounting for a prior regularization term (expressed as a function of the extended space z). Thanks to the extended space, the uncertainty estimation would be given both in terms of model uncertainty and wavefield uncertainty. This would also allow accounting for the cross-covariance terms between the model and the physical field (how the variations in the model result in variation in the wavefield, and vice-versa).

While this proposition is certainly more advantageous than the previous one in terms of uncertainty characterization, it might face challenges in regards of the optimization side of things: if the initial model m_0 is too far from the true model (in the least-squares sense), the non-linearities of the FWI cost function might prevent to reach the global minimum of the cost function with this parameterization. However, iterative DA methods such that the MLEF (2.2.5) might alleviate the sensitivity to non-linearities.

3.3 Proposition 3 - Extending the state-space: The EKI analog

Instead of a regularization based on the entire wavefield, we could also derive a variant of the previous scheme based on the observables. We can express the DA problem as follows,

- $z = \begin{bmatrix} m \\ d_{cal} \end{bmatrix}$, with $z \in \mathbb{R}^{n+d}$, is the joint model-wavefield state space.
- $\mathcal{F}(z) = \begin{bmatrix} \mathbf{I}_m & 0 \\ 0 & \mathcal{H}(m) \end{bmatrix}$, with $\mathcal{F} : \mathbb{R}^{n+d} \rightarrow \mathbb{R}^{n+d}$ where the forecast operator plays a role similar to the artificial dynamics in the EKI and computes synthetic observations from the model m .
- $y^o = d_{obs}$
- $\mathbf{H} = \begin{bmatrix} 0 & 0 \\ 0 & \mathbf{I}_d \end{bmatrix}$, with $\mathbf{H} : \mathbb{R}^{n+d} \rightarrow \mathbb{R}^d$ and \mathbf{I}_d is the identity matrix with d entries.

This formulation comes-down to solving the following minimization problem

$$\mathcal{C}(z) = \min_z \frac{1}{2} \|y^o - \mathbf{H}z\|^2 + \frac{1}{2} \|z - z^f\|^2, \quad (3.4)$$

which is almost equivalent to Proposition 2, albeit the difference in observation operator.

The interest of this formulation is that it corresponds precisely to the EKI (as presented in 2.2.6). The forward operator would be used to generate the artificial dynamics, and by considering perturbed observations, we could perform the EKI iterative minimization scheme. Formulating a DA problem this way should allow investigating the solution space thanks to the ensemble of perturbed observations while benefiting from the iterative nature of the EKI. However, the convergence rate and the sensitivity toward non-linear regimes of this method have yet to be tested.

3.4 Proposition 4 - Extending the state-space: adding adjoint

While I focused on adjoint-free parameterization so far, it is possible to formulate the problem to benefit from conventional FWI adjoint formulation. This should, in theory, yield better results in terms of optimization, as the FWI adjoint-based optimization would ensure fast convergence rates,

and stability. Moreover, using quasi-Newton methods within the DA scheme would allow us to use Hessian-preconditioning to retrieve better-optimized models.

To stay in the extended space paradigm, we consider the following DA-FWI formulation

- $z = \begin{bmatrix} m \\ u \end{bmatrix}$, with $z \in \mathbb{R}^{n+x}$, is the joint model-wavefield state space.
- $\mathcal{F}(z) = \begin{bmatrix} \mathcal{I}_i(m_0) & 0 \\ 0 & \mathcal{M}^{-1}(m_i)s \end{bmatrix}$, with $\mathcal{F} : \mathbb{R}^{n+x} \rightarrow \mathbb{R}^{n+x}$. The forecast step is twofold: we first compute an optimized model m_i with respect to d_{obs} , and then the wavefield is computed in m_i (solution of the tomographic problem).
- $y^o = d_{obs}$
- $\mathbf{H} = \begin{bmatrix} 0 & 0 \\ 0 & \mathbf{E} \end{bmatrix}$, with $\mathbf{H} : \mathbb{R}^{n+l} \rightarrow \mathbb{R}^d$. As in proposition 2, the observation operator extracts values of the wavefield at receiver locations.

While similar to Proposition 2, this formulation differs from the inclusion of the FWI process within the DA scheme. The advantage over the other propositions is that the adjoint-based FWI should speed up the convergence of the model estimate. The uncertainty estimation would be given in terms of both wavefield uncertainty, and uncertainty of the FWI solution, at the expense of a large solution space (as with the other extended state-space propositions).

To alleviate the complexity of this extended-space adjoint formulation, we can go back to a simpler state-space.

3.5 Proposition 5 - A simple adjoint scheme

The last proposition I wish to present is formulated with the tomographic model at its core. Because we ultimately want to quantify the uncertainty of optimized models, we propose to define the state vector solely as the subsurface model. This way, we can benefit from the adjoint-based FWI advantages and reduce the memory requirements of the extended space.

The uncertainty produced by this method will be measuring the quality of the adjoint-based FWI solution, making this approach similar to the local uncertainty estimation methods presented in 1.2.2. The filter parameterization is given as follow:

- The state vector is defined as a subsurface model $m \in \mathbb{R}^n$.
- $\mathcal{F}(m) = \mathcal{I}_i(m_0)$, with $\mathcal{F} : \mathbb{R}^n \rightarrow \mathbb{R}^n$. During the forecast, we now solely solve the FWI problem with respect to d_{obs} .
- $y^o = d_{obs}$
- $\mathcal{H} = \mathcal{H}(m)$, with $\mathcal{H} : \mathbb{R}^n \rightarrow \mathbb{R}^d$. The observation operator comes down to computing d_{cal} from the subsurface model m . Note that this is the only non-linear observation operator, out of all the five proposition.

with a misfit function defined as

$$\mathcal{C}(m) = \frac{1}{2}(y^o - \mathcal{H}(m))^T \mathbf{R}^{-1}(y^o - \mathcal{H}(m)) + \frac{1}{2}(m - m^f)^T (\mathbf{P}^f)^{-1}(m - m^f), \quad (3.5)$$

which has the same form as the Tikhonov-regularized FWI misfit function.

We now have to take into consideration that our state-space (subsurface model) does not evolve with time; hence, we lack the dynamic axis required by the DA framework. To cope with this issue, we propose to define a proxy for the dynamic axis, based on the FWI process itself.

3.5.1 FWI as a dynamic problem: defining a dynamic proxy

While it might seem counter-intuitive, we propose to consider the FWI process itself, as a dynamic system. To back this proposition, we can draw parallels between the forward map used in "classical" DA problems and our tomographic operator:

- From the initial conditions (initial model), FWI yields an update of the physical parameters (the model "evolves" with each iteration).
- Each inversion iteration becomes the initial conditions to the next iteration (analogous to time integration schemes in forward-modeling engines).
- It is a fully deterministic process.

Therefore, provided an adequate dynamic axis can be defined, any mapping (be it forward or inverse) could potentially play the role of the forecasting operator.

The frequency continuation strategy

In the case of FWI, the most obvious choice for this dynamic axis would be to consider the hierarchy in modeling/inversion frequencies. It is common practice in FWI to consider increasingly high-frequency data, starting from the lowest frequency possible. This *frequency-continuation* strategy (also known as the *multi-scale* approach) has been highlighted by the work of Bunks et al. (1995) and is commonly employed to mitigate cycle-skipping artifacts.

Cycle-skipping is a common FWI artifact that occurs when oscillating seismic signals are shifted from more than half a period (Fig 3.1): an ambiguity arises when matching a synthetic arrival with its counterpart in the observed signal (this corresponds to a local minimum of the misfit function). We typically see this type of artifact arising when the initial model cannot represent the basic kinematics of the wave propagation through the medium (underestimated velocities, for instance).

As it can be seen in Figure 3.2, the lower frequency considered, the more the misfit function becomes convex (Bunks et al., 1995): in the low-frequency regime, the signal's period is longer, reducing the possible ambiguity between phases. Thanks to this property, we can follow the frequency hierarchy, to design a dynamic axis proxy for our DA-FWI scheme, following the multi-scale approach.

Thus, our FWI system evolves along an axis of frequency continuity, starting from the lowest to the highest frequency. Doing so, FWI becomes a dynamic process on its own, responsible for the evolution of subsurface models according to the increasing data/modeling frequency.

Out of the five propositions, we believe this last formulation is the most interesting for initial numerical tests. It gives uncertainty in terms of the model parameters, which is the goal of our research

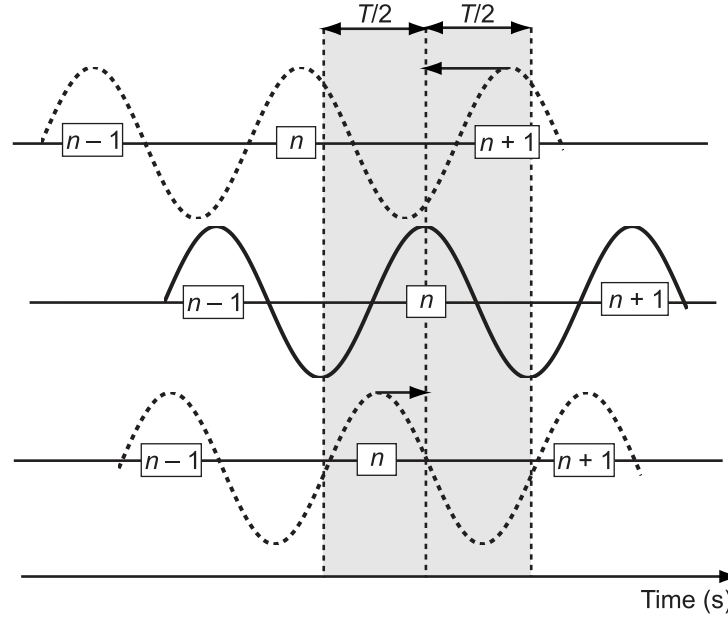


Figure 3.1: Illustration of the cycle-skipping artifact in FWI. The solid black line represents a monochromatic seismogram of period T as a function of time. The upper dashed line represents the modeled monochromatic seismograms with a time delay greater than $T/2$. In this case, FWI will update the model such that the $n + 1^{\text{th}}$ cycle of the modeled seismograms will match the n^{th} cycle of the observed seismogram, leading to an erroneous model. The bottom example will not produce cycle-skipping because the time delay is less than $T/2$. From Virieux and Operto (2009).

while allowing us to use efficient quasi-Newton solvers during the forecasting stage. Its formulation makes it easily implementable, and working with this simple state-space rather than with an extended space alleviates possible memory overburden. In the next subsection, I will explain how we can pair this parameterization with the ETKF, presented in 2.2.4.

3.5.2 The ETKF-FWI scheme

As I have stated earlier in Chapter 2, the ensemble DA scheme we have chosen for this study is the ETKF. Its straightforward implementation and efficiency, coupled with our last parameterization proposition, makes it a very simple scheme, from which we performed numerical experiments to investigate the potential of ensemble DA for uncertainty estimation in FWI.

Our ETKF-FWI scheme is defined as follows,

1. Define an optimal starting model m_0 , as the best starting model available.
2. Define the K ordered steps of the dynamic axis in frequency continuation as $k = 1, 2, \dots, K$. Note that k can contain either a single frequencies or a frequency group.
3. Generate an initial ensemble such that $\mathbf{m} \sim \mathcal{N}(m_0, P_0)$, with P_0 being the prior model covariance matrix and \mathbf{m} containing N_e ensemble members.
4. For each k until K is reached:

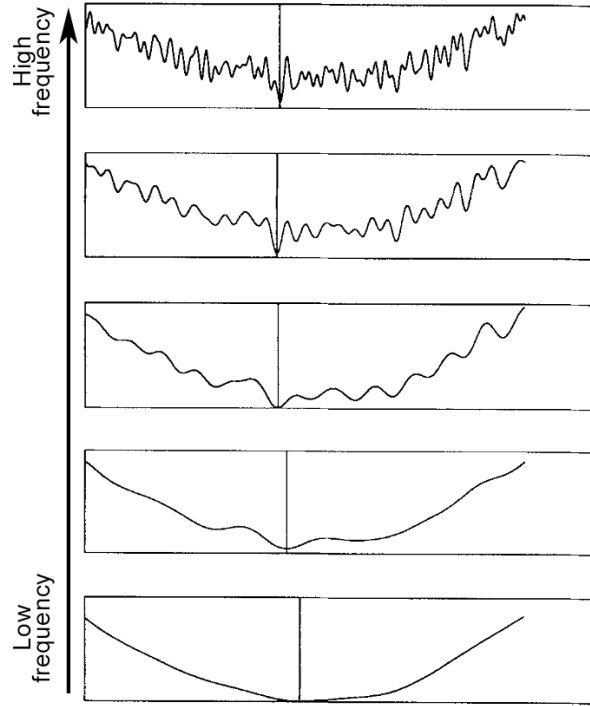


Figure 3.2: Illustration of the misfit-function dependency on the data-frequency content. When low-frequency data are considered, the misfit-function is smoother and contains few local minima. When the data frequency increase, the misfit-function gets more local minimum and is less convex. *Modified from Bunks et al. (1995).*

- (a) Compute the forecast ensemble by solving $m_k^{f(i)} = \mathcal{F}_i(m_{k-1}^{(i)}, d_{obs,k})$, by minimizing the least-squares distance to the observed data $d_{obs,k}$, corresponding to the frequency or frequency group k . Note that the forecast is set to perform an arbitrary number i of non-linear inversion iterations.
- (b) Compute the forecast observation $d_{cal}^{f(i)}$ by solving the forward problem as $d_{cal}^{f(i)} = \mathcal{H}(m_k^{f(i)})$
- (c) Compute the analysis ensemble, following the ETKF analysis scheme. The analysis ensemble is obtained by solving a minimization problem over the ensemble subspace $\mathcal{S} \in \mathbb{R}^{N_e}$, which cost is negligible compared to the rest of the scheme.

From this formulation, we can see that the forecasting step will be the computational bottleneck of this scheme, as it requires to perform as many independent FWI iterations as we have ensemble members. However, the independent nature of the forecast makes it an embarrassingly parallel problem, which allows very efficient parallelization, provided we have access to computing resources that are sufficient to fit all the FWI processes in memory. To alleviate the computational cost, we decided to work with frequency-domain FWI, as it reduces the computational requirements dramatically, to solve the forward modeling problem in 2-D. Therefore, we consider only time-harmonic monochromatic wavefield data as our observations.

In the ETKF-FWI scheme, pictured in Figure 3.3, we have combined ensemble DA and the quasi-Newton FWI methodology, into a framework that can produce both high-resolution tomographic models,

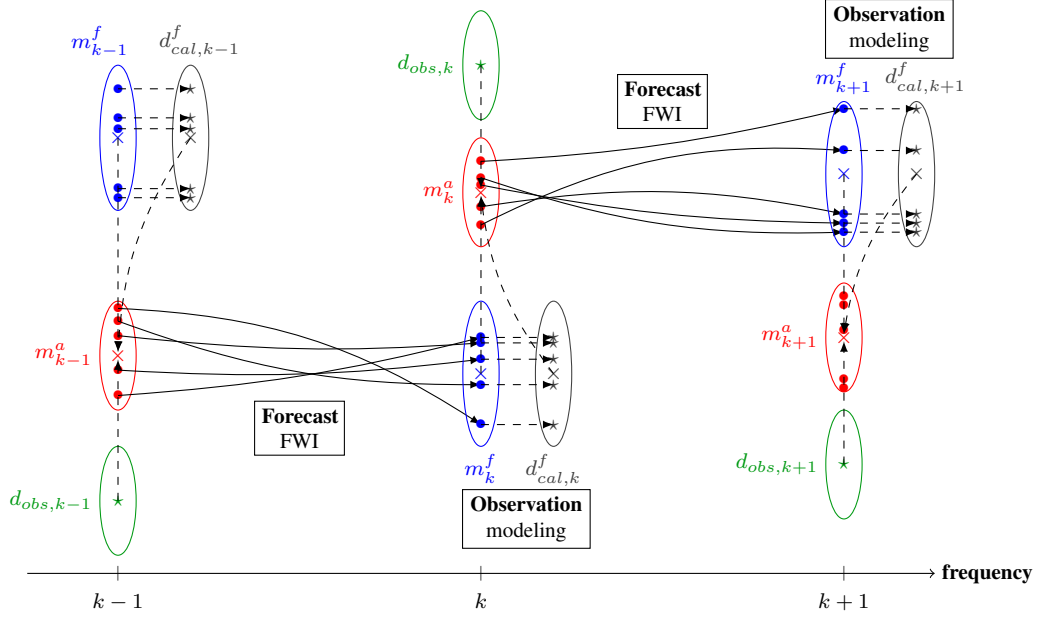


Figure 3.3: Schematic representation of the ETKF-FWI algorithm. Bullets represent subsurface models, crosses denote ensemble means, stars denote observations, and ellipses represent ensemble covariances. As previously, the blue color denotes the forecast state, the red denotes the analysis, green is for the observed data, and grey for the forecast data. The dashed lines are indicative of the information used to produce the analysis state.

and uncertainty estimation of the solution at the same time. Even though the total computational time might seem to be daunting, the ETKF-FWI method we propose is fully scalable, which is a feature lacking from the other local uncertainty analysis proposed in the literature.

The particularity of this scheme is that we essentially replaced the forward map by an inverse map, in the forecasting step (Fig. 3.3), which sets this method apart from classical DA problems. Additionally, because the analysis step is equivalent to an inverse problem, our scheme can be viewed as cycling through two different inverse problems.

Note that the observed data for both the forecast and the analysis are set to be the same $d_{obs,k}$ data at step k . This approach also deviates from common ETKF scheme where it is assumed to introduce new information during the analysis rather than relying on previous data. However, because \mathcal{I}_i is updating the ensemble of models using $d_{obs,k}$, the wavenumber content of these updates is closely tied to the frequency content of $d_{obs,k}$. The ensemble of optimized models obtained from $d_{obs,k}$ are likely to lack the higher wavenumber content to explain or "predict" higher frequency data $d_{obs,k+1}$, and may cause cycle-skipping to happen during the analysis. The ETKF-FWI can thus be broken down into two successive minimization procedures, with respect to the same objective data:

- The forecast is a *model-wise* optimization, with a quasi-Newton solver. Each of the N_e minimizations is independent of the others.
- The analysis results in an *ensemble-wise* minimization. Contrarily to the forecast, the analysis uses the information contained in the ensemble repartition, and the measurement uncertainty to balance the forecast and observations.

The model-wise optimization (forecast) is tasked with bringing together all the ensemble members close to the global minimum (provided the ensemble is built adequately and all ensemble members are sampling the same subset \mathcal{A}). After the model-wise adjoint-based inversions, the analysis perform an ensemble-wise inversion, rearranging the ensemble around the mean solution, ensuring coherency of the solution and reducing the ensemble variance.

Finally, we justify the compatibility between the underlying assumptions of the ETKF and FWI by the following: provided that all the ensemble members are located in \mathcal{A} , the local convexity of the cost function assumed to perform quasi-Newton optimization, is deemed a good first-order approximation of the Gaussian probability density function required by the ETKF. Therefore, the initial ensemble must be generated such that it will correctly sample the subset defined by the initial model $m_0 \in \mathcal{A}$. In the next subsection, I discuss the implication of the ensemble repartition and the sampling strategies that are involved in satisfying this Gaussian assumption.

3.5.3 ETKF-FWI sampling strategy

To ensure the ETKF-FWI scheme' stability, the initial ensemble must sample adequately the subspace \mathcal{A} defined by the choice of the initial model m_0 : if m_0 is bound to converge toward the minimum in \mathcal{A} , we must ensure that the whole ensemble converges toward this same minimum.

Because of the hypothesis underlying the ETKF, we must draw an initial population that has a Gaussian repartition. To ensure that we are sampling \mathcal{A} , the initial ensemble mean should be m_0 such that the ensemble can be modeled as $\mathbf{m}_0 \sim \mathcal{N}(m_0, \mathbf{P}_0)$. Thus, the initial ensemble can be built as

$$m_0^{(i)} = m_0 + \eta^{(i)}, \quad (3.6)$$

with $\eta^{(i)}$ being zero-mean multivariate gaussian perturbations, defined as $\eta^{(i)} \in \mathcal{N}(0, \mathbf{P}_0)$ with $i = 1, 2, \dots, N_e$. We are left with defining the initial covariance.

As we generally lack information regarding subsurfaces' physical properties, we could be tempted to choose the initial covariance matrix such that the initial ensemble sample as much of the solution space as possible, akin to global optimization methods. However, due to FWI non-uniqueness, and the nature of our forecasting operator, this might cause the ensemble to split over several local minima, as illustrated in Figure 3.4. Taking the misfit function presented in Chapter 1, 10 random samples were drawn (red dots) following an arbitrary Gaussian repartition centered on an arbitrary initial point $m_0 = (-1.2, 1)$. After minimization, each ensemble member converges toward their respective closest local minimum (blue dots). If outliers are present in the initial repartition, we sample several minima of the misfit, effectively splitting the ensemble into several basins of attractions. Consequently, we see the occurrence of two biases:

- Because the ensemble is sampling several local minimum, the repartition of the ensemble becomes multimodal, which violates the Gaussian hypothesis of the ETKF.
- As the ensemble is splitted, the mean of the ensemble (black cross) does not correspond to an approximation of the MMSE (the minimum at $(-1.2, 1)$).

Note that the severity of these biases is heavily dependent on the "topography" of the misfit function. If the misfit function is strongly non-convex, the odds of sampling several local minima (and thus splitting the ensemble and biasing the ensemble mean) are increasing.

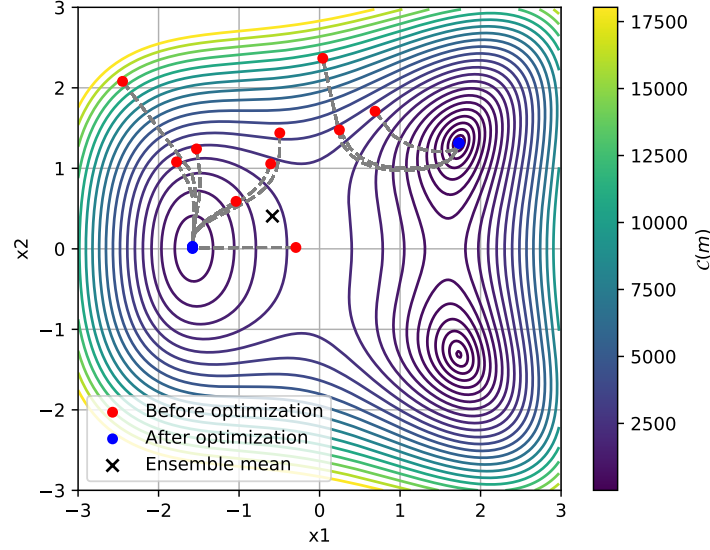


Figure 3.4: Changes in repartition under an inverse mapping forecast (as in the ETKF-FWI scheme). Ten samples are drawn (red) from a Gaussian distribution with mean $(-1.58, 0)$ and arbitrary large variances. Optimization paths are denoted with dashed grey lines and optimized model by blue dots. The ensemble mean after inversion is denoted by the black cross. In such a case where outliers are present in the initial repartition, the optimized ensemble mean is bound to be biased.

To mitigate these biases, we can either minimize the influence of outliers by increasing the number of ensemble members (such that most of the samples will be drawn within the subset \mathcal{A} , as in Figure 3.5) or by reducing the spread of the ensemble such as in Figure 3.6, reducing the odds of outliers being sampled outside of \mathcal{A} . While increasing the ensemble size slightly mitigates the sampling bias, it is not a satisfying solution: increasing the ensemble size results in a computing overburden while still producing a biased mean estimate and an ensemble splitting. The results may also largely vary if the function is strongly non-convex. Therefore, reducing the initial variance of \mathbf{P}_0 to better match the extent of \mathcal{A} seems to be a better solution, as it allows sampling the correct misfit subset.

While it would be difficult to define the appropriate variances that fit the local subset \mathcal{A} in an arbitrary inversion problem (especially so as the number of degrees of freedom gets substantially big), we can take advantage of FWI problem formulation to do so: as cycle-skipping is mostly responsible for the non-convexity of the cost function, we can generate an ensemble contained within \mathcal{A} by generating model-perturbations that do not generate cycle-skipping with respect to the lowest frequency considered.

Additionally, we have to make sure that the wavenumber content of the initial perturbations matches the local resolution imposed by the FWI resolution power for a given frequency band (Devaney, 1984; Wu and Toksöz, 1987).

With these constraints in mind, we define the generation of the initial ensemble as follows: Each ensemble member is built by taking an initial model m_0 deemed suited for convergence, to which we add a perturbation. Perturbations are generated by convolution of zero-mean, uniformly distributed random vector $\psi^{(i)}$ (with $i = 1, 2, \dots, N_e$) with a non-stationary Gaussian function \mathcal{G} which correlation length

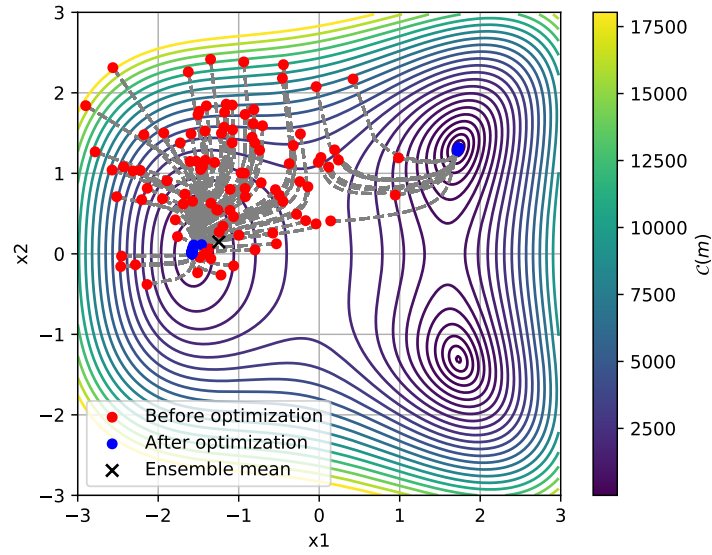


Figure 3.5: Changes in repartition under an inverse mapping forecast (as in the ETKF-FWI scheme). The random sampling was kept identical with Figure 3.4, except that 100 samples were drawn. While the position of the mean is less biased, increasing the number of samples also increases the odds of sampling the wrong local subset.

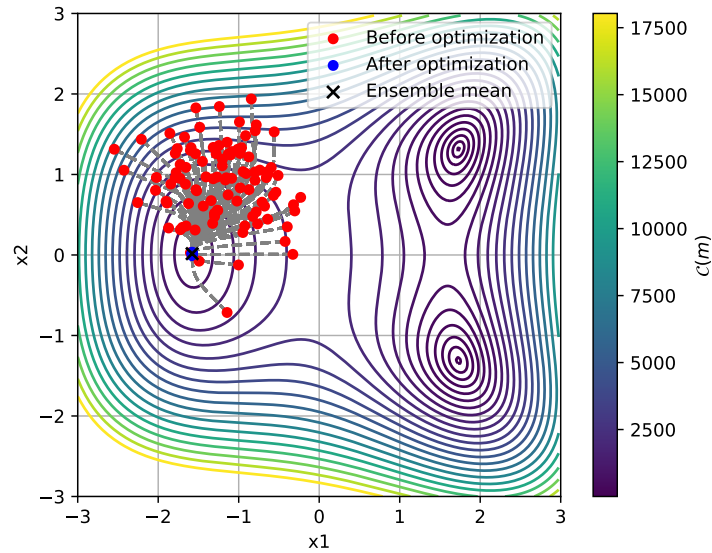


Figure 3.6: Changes in repartition under an inverse mapping forecast (as in the ETKF-FWI scheme). By changing the sampling distribution (reducing the variances of \mathbf{P}_0) we reduce the odds of sampling the wrong local subset. In this example, the 100 ensemble members are sampling the desired subset \mathcal{A} and the ensemble mean corresponds to the minimum.

	Subsurface uncertainty	Adjoint-free	Iterative	Extended state-space
Proposition 1	✗	✓	✗	✗
Proposition 2	✓	✓	✗	✓
Proposition 3	✓	✓	✓	✓
Proposition 4	✓	✗	✗	✓
Proposition 5	✓	✗	✗	✗

Table 3.1: Summary of Propositions 1 to 5.

and amplitude are varying according to the local velocity in m_0 , such that

$$m_0^{(i)} = m_0 + \mathcal{G}\psi^{(i)}, \quad i = 1, 2, \dots, N_e. \quad (3.7)$$

$\mathcal{G}\psi^{(i)}$ produces smooth perturbations which wavenumber is half of the wavefields' wavelength, corresponding to the maximum spatial frequency that can be recovered (Wu and Toksöz, 1987). The initial ensemble is then inspected with an Eikonal solver, to ensure that the initial population of models will not allow cycle skipping at our starting frequency, that could be provoked by too dramatic initial perturbations. Even though this test only allows assessing the first arrival cycle skipping, it is deemed sufficient as a first-order diagnosis of the initial ensemble quality. To further ensure favorable initial conditions, we verify that the rank of the initial ensemble is equal to N_e .

This way, we make sure that the initial sampling is adequate with respect to our starting model, and the subset \mathcal{A} we wish to sample. From there, we start the cycle of forecast and analysis steps, following the frequency continuation axis, until the highest frequency is reached.

Conclusion

In this chapter, we have seen that there can be several ways of combining DA and FWI. This is especially true when defining an arbitrary dynamic axis based on hierarchical approaches. An emphasis was put on Proposition 5, as it is the parameterization we adopted for subsequent numerical experiments on both synthetic benchmark and field-data experiments. We picked this scheme because of its ease of implementation and simple formulation, which allowed us to start experimenting with the joint DA-FWI formulation quickly. While we cannot claim that Proposition 5 is the optimal choice to produce uncertainty estimation in FWI, it provided us with numerous insights and results that make for substantial advances in this attempt at bridging DA and FWI together.

The other propositions are nonetheless interesting, especially so the extended-space formulations that are analogous to the WRI and EKI. Even though these have not been implemented and tested on practical cases, I believe they are still worth investigating, as the extended formulation might improve results' stability (by improving the misfit-convexity) or reduce the computational cost (thanks to the adjoint-free formulation for the EKI). The basic properties of each method have been summed up in Table 3.1.

From this summary, we see that Proposition 1 can be ruled out, as it is unable to inform us directly on the subsurface parameters uncertainty. However, the other characteristics such as "adjoint-free" are only indicative (being adjoint free can be viewed as an advantage in terms of implementation, but might also be a disadvantage in terms of convergence rate and stability), such that propositions 2,3 and 4 are still relevant.

Note that if we were to include additional observables to the problem (like well-log data, for instance), we could introduce additional parameterization to account for these new pieces of information.

In the next chapters, I will present the numerical experiments in which the ETKF-FWI has been used to produce uncertainty estimation. Chapter 4 will be focused on the monoparameter synthetic-benchmark, while Chapter 5 will present the results obtained on a field-data experiment, both in mono and multiparameter FWI.

Chapter 4

Synthetic application of the ETKF-FWI

Contents

4.1 Solving the FWI problem	85
4.2 ETKF-FWI on the Marmousi II synthetic benchmark	87
4.2.1 Synthetic benchmark setup	87
4.2.2 The ETKF-FWI setting	88
4.2.3 ETKF-FWI application with 600 ensemble members.	89
4.3 Investigating undersampling	97
4.3.1 Parameter estimate	97
4.3.2 Variance approximation	98
4.3.3 Correlation approximation	101
4.4 Mitigating undersampling	104

In precedent chapters, we have presented both the FWI and the DA theories and proposed a few ways of bridging the two methodologies into a unified framework. From the different parameterization propositions, we identified the possibility of combining the ETKF and the quasi-Newton FWI framework as a frequency-dependent process.

In this chapter, I will first present the technicalities involved in solving the numerical FWI problem. Then the synthetic benchmark, along with the ETKF-FWI filter parameterization, will be presented, followed by inversion results and uncertainty estimation. Finally, I will discuss the implications of undersampling in our ETKF-FWI schemes and how we attempted to mitigate its biases.

4.1 Solving the FWI problem

We have seen in Chapter 1 that an accurate numerical solver is required to compute the incident and adjoint wavefields, in order to build the FWI gradient. As stated previously, we will solely focus on solving the 2-D frequency domain, visco-acoustic wave equation for our numerical tests. In that case, computing the wavefield comes down to solving a linear system of the form $\mathbf{B}u = s$.

In our implementation of the ETKF-FWI, we solve this linear system with the open-source TOY2DAC code developed in the SEISCOPE Consortium. TOY2DAC's forward modeling engine is based on a

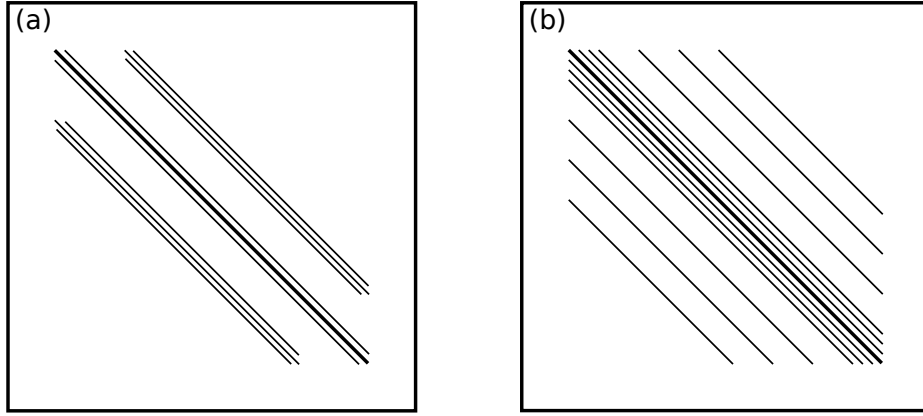


Figure 4.1: Sketch representation of the sparse impedance matrices $\mathbf{B}(\omega, m)$ in finite-difference setting. The panel a) represent the matrix structure with the mixed finite-difference stencil proposed in (Hustedt et al., 2004), containing $9 \times n$ non-zero entries. The panel b) illustrate the matrix structure under a "parsimonious" fourth-order finite-difference stencil, with $13 \times n$ non-zero entries. *Reproduced from* (Hustedt et al., 2004)

direct LU (for Lower-Upper) solver (Duff and Reid, 1983; Liu, 1992; Amestoy et al., 2000), which relies on the decomposition of the sparse impedance matrix \mathbf{B} (see Figure 4.1) into lower \mathbf{L} and upper \mathbf{U} triangular matrices,

$$\mathbf{B}u = (\mathbf{L}\mathbf{U})u = \mathbf{L}(\mathbf{U}u) = s. \quad (4.1)$$

The monochromatic stationary wavefield u is given by forward and backward substitution of a triangular set of equations

$$\begin{aligned} \mathbf{L}y &= s \\ \mathbf{U}u &= y \end{aligned} \quad (4.2)$$

This formulation is particularly interesting, as it allows quick computation of the wavefield for multiple sources once the decomposition of \mathbf{B} is stored (Marfurt, 1984). Note that even though the direct solver allows some significant computational cost reduction compared to 2-D time-domain modeling, it is not the case for 3-D applications. Due to the complexity of the forward and backward substitution being linearly linked with the size of the discretization space, it becomes a limiting factor for 3-D frequency-domain full-waveform modeling, as the number of grid points l is expected to grow rapidly (Virieux and Operto, 2009).

The TOY2DAC forward modeling solver implementation relies on an optimized finite-difference discretization strategy with a compact stencil providing accuracy, equivalent to fourth-order methods (Hustedt et al., 2004; Operto et al., 2009). The LU decomposition required within this direct solver is based on the MUMPS sparse solver (MUMPS team, 2017), which has been designed to decompose large sparse matrices efficiently and is thus particularly adapted for FWI applications as we typically have to deal with large discretized domains.

The other interest of this formulation is that once the matrix decomposition has been stored, it can be used to compute both the direct and the adjoint wavefields. We can thus use the LU-decomposed impedance matrix twice, thanks to the adjoint formulation of the gradient, which is computationally efficient.

Under the adjoint-state formalism, the gradient can be computed in two ways, depending on the expression for the source term of the backpropagated wavefield v . It can either be based on the conjugate of the data misfit and corresponds to time-reversed residuals. Or we can use the self-adjoint nature of the wave-equation and replace $[(\mathbf{B}^{-1})^T \delta d^*]$ by $[(\mathbf{B}^{-1})^{T*} \delta d]$ in equation 1.31, where $(\mathbf{B}^{-1})^{T*}$ is the adjoint forward operator. Instead of propagating the time-reversed data misfit through the forward operator, the data misfit is directly the source of the adjoint forward operator (a "time-reversed" wave equation). In any case, both formulations can be considered to formulate the gradient as they are equivalent, and both allow to compute \mathbb{G} without building the Jacobian matrix explicitly.

From the computation of the gradient, based on the adjoint-state method, the local optimization scheme used to solve the FWI problem is the l-BFGS scheme, implemented in SEISCOPE optimization tool-box (Métivier and Brossier, 2016). This second-order quasi-Newton method provides a simple and efficient inversion scheme, perfectly suited for the computationally intensive experiments we ought to perform.

Finally, provided the observed data d_{obs} and a starting model m_0 are available, we can solve the FWI problem and iteratively compute an optimized model.

4.2 ETKF-FWI on the Marmousi II synthetic benchmark

In this section, I will first present the synthetic benchmark, before presenting the ETKF-FWI application in which it is involved, and the results we obtained from our scheme.

4.2.1 Synthetic benchmark setup

The synthetic benchmark we use in this numerical experiment is the Marmousi II synthetic model (Martin et al., 2006). The model's domain width and depth are respectively $x = 16.025km$ and $z = 3.250km$ with vertical and horizontal resolutions of $dx = dz = 25m$, for a total of $n = 83300$ degrees of freedom. The true model $m_{true} \in \mathbb{R}^n$ used to simulate the observed data is pictured in Figure 4.2 along with the starting model $m_0 \in \mathbb{R}^n$, obtained by smoothing the true model with an isotropic Gaussian kernel. Note that the 500m water-depth is considered to be known, such that the water-seabed interface does not have to be smoothed. In this FWI setting, we ought to minimize the least-squares distance $\mathcal{C}(m) = \|d_{cal}(m_0) - d_{obs}(m_{true})\|^2$

Data are simulated using a fixed spread surface acquisition configuration, with 144 sources and 640 receivers evenly spaced, to mimic marine acquisition geometry, resulting in an observation vector $d_{obs}(m_{true}) \in \mathbb{C}^d$ with $d = 92160$ entries for each mono-frequency data. The total number of discrete data is equal to $92160 \times n_\omega$ where n_ω is the number of considered frequencies. In our case, we define the FWI workflow based on the successive inversion of $n_\omega = 15$ mono-frequency complex-valued data from 3 Hz to 10 Hz, with a 0.5 Hz increment between each inversion. This FWI workflow will serve as a basis for our ETKF-FWI application.

To avoid inverse crime, a complex Gaussian random noise was added to the synthetic observed data (Eikrem et al., 2019):

$$d_{noisy} = d + \frac{\|d\|}{\sqrt{r * E(\|w\|^2)}} w \quad (4.3)$$

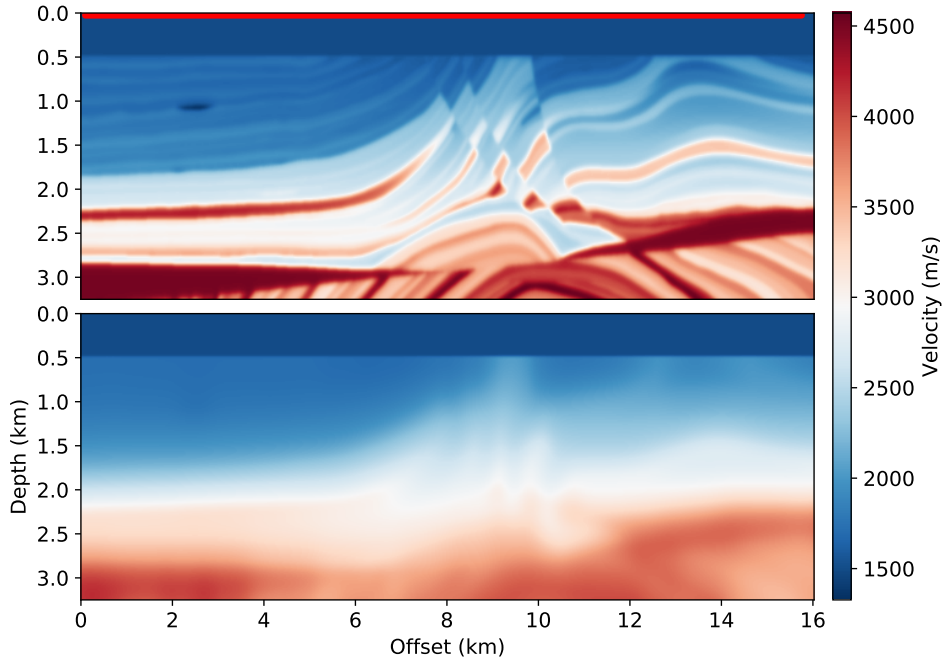


Figure 4.2: Numerical experiment setting. Top: True Marmousi II model. A red line denotes the acquisition footprint at the surface. Bottom: A smoothed version of the true model, used as a starting model m_0 .

with d_{noisy} being the noisy signal, d is the original noise-free signal, $\|\cdot\|$ denotes the Euclidean norm and E the expectation. The vector $w \in \mathbb{C}^d$ is defined as

$$w = v_1 + iv_2, \quad (4.4)$$

where $v_1, v_2 \in \mathbb{R}^d$ are vectors of normally distributed random numbers and r is defined as the signal to noise ratio, such that

$$r = \frac{\|d\|^2}{\|w\|^2}. \quad (4.5)$$

In the following experiments, we set up $r = 8$ as our reference noise value through all the successive inversion. Even though in field-data applications, we expect the noise repartition to be non-uniform (high-frequency data typically have a better SNR than low-frequency data), we kept the noise value fixed for the sake of simplicity.

4.2.2 The ETKF-FWI setting

Now that the synthetic FWI has been defined, we can set-up the ETKF-FWI according to the scheme presented in 3.5.2.

Ensemble generation: In this context, the initial model m_0 , becomes the mean of our initial ensemble, which members are built according to equation 3.7. As said in Chapter 3, the generation of the initial ensemble requires to satisfy both the resolution power of FWI at the lowest frequency (3 Hz) and must not generate first arrival cycle-skipping. An example of velocity perturbations meeting these criterions are shown in Figure 4.3.

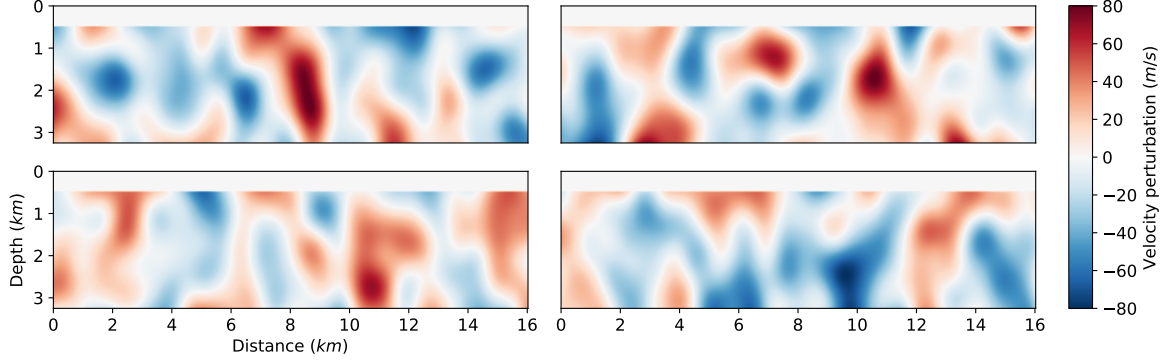


Figure 4.3: Example of initial perturbation for the synthetic Marmousi II experiment.

Defining the forecast: As the ETKF-FWI dynamic axis is defined to follow the frequency continuation strategy, we chose to define $K = 15$ ETKF-FWI cycles, as we have $n_\omega = 15$ monochromatic data. In each of the cycles, the forecast operator \mathcal{I}_n is set to perform $n = 10$ minimization iterations with the l-BFGS optimization scheme, on each of the N_e velocity models, with mono-frequency synthetic calculated $d_{cal,k} \in \mathbb{C}^d$ and noisy observed data $d_{obs,k} \in \mathbb{C}^d$ at frequency k with $k = 1, 2, \dots, K$. The cost of the methodology is thus linearly linked with the number of ensemble members and non-linear FWI iterations. In this instance, the number of forward modeling is $N_e \times 10 \times 2$, as both the incident and adjoint wavefields are computed at each iteration.

Computing the forecast data: Once the forecast state is obtained (an ensemble of optimized initial models), we compute the forecast data at the frequency k with observation operator \mathcal{H} . This generates the ensemble of synthetic data, computed in the forecast ensemble. Finally, the analysis is performed, which updates the ensemble mean, and reduces the ensemble variance.

Balancing forecast and observations: Recalling that the two factors controlling the analysis are the forecast covariance matrix $\mathbf{P}_{e,k}^f$ defined by the forecasted ensemble, and the measurement noise matrix \mathbf{R} , it appears necessary to correctly set \mathbf{R} to ensure that the analysis successfully balances the forecast ensemble and the observations. In the synthetic case, as the noise source in the observations is perfectly known, we can set-up \mathbf{R} accordingly:

$$\mathbf{R} = \mathbf{I}_d \sigma^2 = \mathbf{I}_d \left[\frac{\|d\|^2}{r * E(\|w\|^2)} \right] \quad (4.6)$$

where \mathbf{I}_d is an identity matrix of size d and σ is the standard deviation of the observation noise. By considering \mathbf{R} as a scaled identity, we consider that the noise vectors added to each receiver data are uncorrelated: this assumption comes from the lack of information about possible correlated measurement errors in seismic acquisitions. While the benefits of taking correlated noise structures into account have been highlighted (Stewart et al., 2008; Weston et al., 2014), it is not possible in our case. Note that if we could estimate the off-diagonal terms of \mathbf{R} , the computation of \mathbf{R}^{-1} during the analysis step would become computationally challenging.

4.2.3 ETKF-FWI application with 600 ensemble members.

We start with a favorable case: the ensemble size is chosen to be arbitrarily large, to mitigate the possible undersampling effects. Thus, we generate an ensemble with $N_e = 600$ ensemble members.

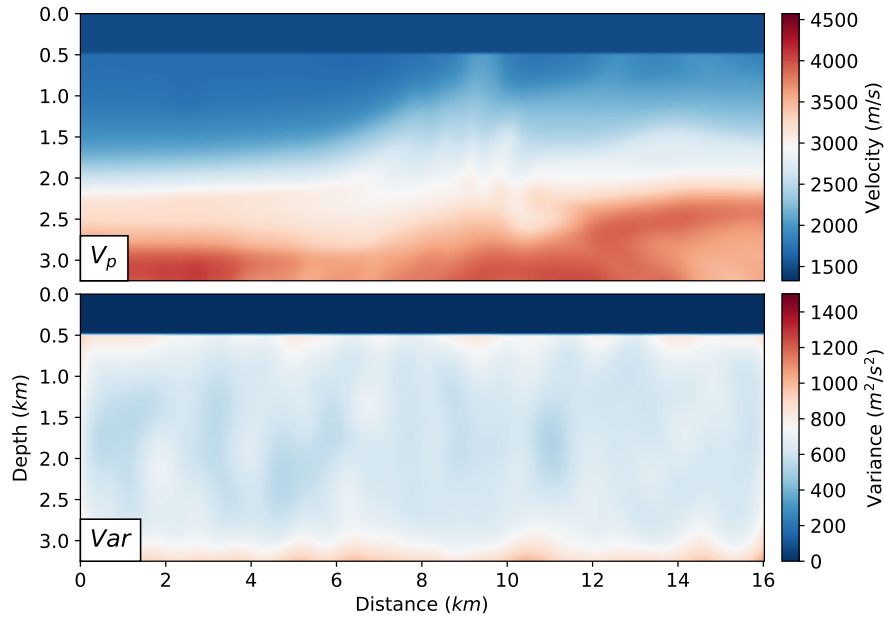


Figure 4.4: Initial ensemble mean and variance. The ensemble mean (top) corresponds to the mean velocity map. The ensemble variance (bottom) denotes the initial uncertainty, associated with the ensemble generation.

Ensemble generation and first cycle In order to explain what is occurring during each cycle, I will first present the initial ensemble along with the first forecast and its analysis. The initial ensemble mean and variance are displayed in Figure 4.4.

As mentioned previously, the initial ensemble mean corresponds to the smoothed model m_0 . Recalling that the perturbations used to generate the ensemble are modeled as $\eta^{(i)} \sim \mathcal{N}(0, P_0)$, we can indeed expect the initial ensemble mean to be m_0 . The initial variance (diagonal of \mathbf{P}_0) has been set to be nearly homogeneous across the domain to reflect our lack of initial knowledge regarding the subsurface structures.

Following the ensemble generation, the ensemble forecast is computed by applying the FWI operator \mathcal{I}_n , which minimizes the least-squares distance $\mathcal{C}(m)^{(i)} = \left\| d_{cal}(m_0^{(i)}) - d_{obs,1} \right\|^2$. The forecast mean and the forecast variance (diagonal of $\mathbf{P}_{e,1}^f$) are displayed in Figure 4.5.

The forecast ensemble mean differs from the initial ensemble mean: it contains more features of the true model, especially in the shallow part of the domain. This change is due to the convergence of the ensemble members toward their closest local minima (as every model is being optimized with respect to the observations).

In the variance map, we observe a reduction of variance in shallow areas (where most of the model update take place), but also a significant increase along sharp velocity contrasts (mainly at 1.5km depth). Because we cannot ensure that all of the N_e models will resolve the interface within the same number of iterations, the variance might increase in this specific area of the model. Thus, observing this type of behavior on sharp discontinuities corresponds to our expectations, as these features are generally hard to recover in FWI. This first forecast variance map allows us to identify the features displaying a higher degree of uncertainty.

Finally, after the forecast step, the analysis is performed, with respect to the same objective data $d_{obs,1}$. The analysis ensemble mean and variances are displayed in Figure. 4.7.

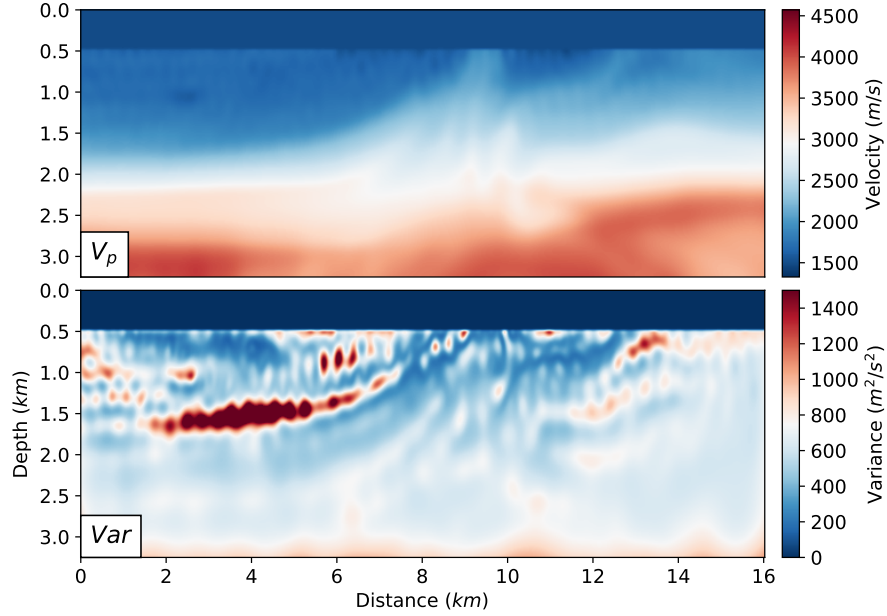


Figure 4.5: First forecast mean (top) and variance (bottom). Some features of the true model are being included in the mean model, as the ensemble progresses toward the closest minimum of the misfit function. In the variance map, we can see that some features, primarily sharp velocity contrasts, are displaying an increase of variance.

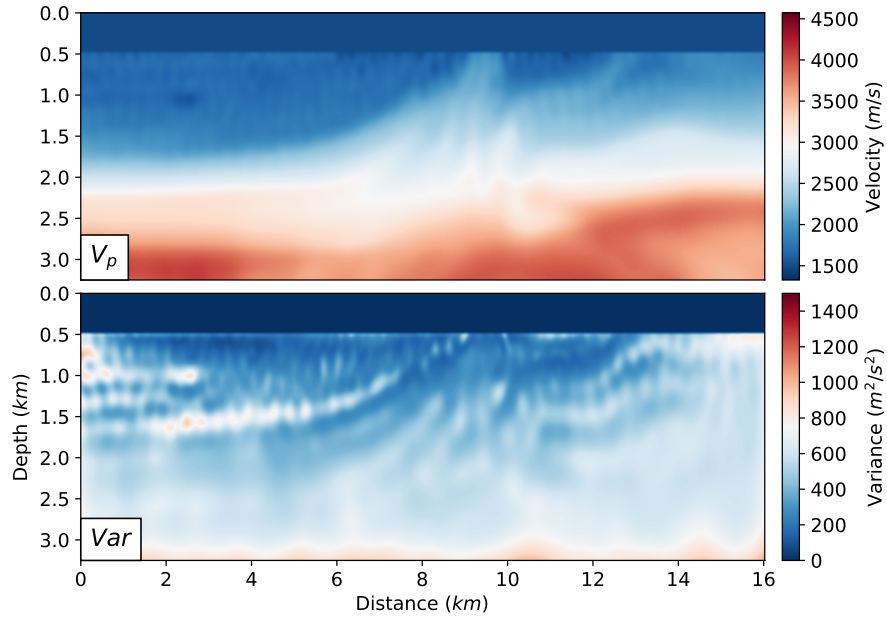


Figure 4.6: First analysis mean (top) and variance (bottom). We can see significant decrease in variance in key areas, notably on shallow reflectors.

The most striking changes are visible in the variance map; notably, the high uncertainty zones in the forecast variance have been drastically reduced. This emphasizes the role of the analysis, which is essential to rebalance the ensemble around the optimal mean (in the least-squares sense) and lower the variance of the forecast ensemble. It allows controlling the ensemble' spread all along with the cycles, mitigating the risks of the ensemble splitting over several minima. Naturally, the ensemble mean has also been updated by the analysis. While the changes are subtle (slightly sharper velocity contrasts), it is interesting to note that most of the changes have occurred along with high variance reduction areas. A map of the velocity update from the analysis step is displayed in Figure 4.7. We can see a clear connection between the changes in variance and the velocity update.

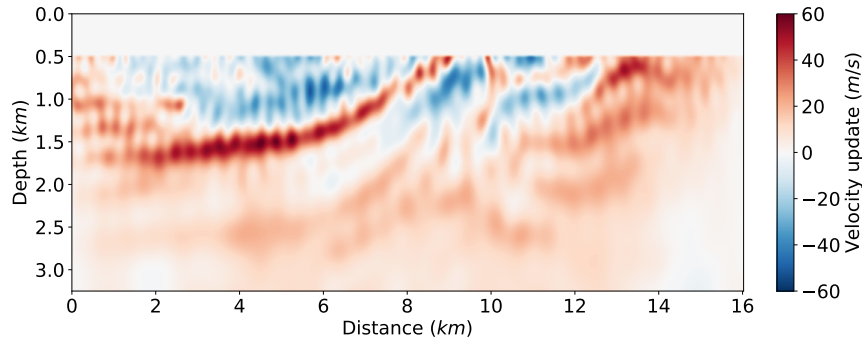


Figure 4.7: Velocity update of the first analysis step. The changes in velocity from the analysis are consistent with the variance reduction zones.

Results We have seen in detail how the forecast and analysis affect the ensemble of velocity models during the first cycle. We now directly jump to the results after the 15 cycles, from 3 Hz to 10 Hz. The final ensemble mean (and hence the final parameter estimate) and the final variance map are displayed in Figure 4.8.

The ETKF-FWI parameter estimate tends to a regular FWI result: most of the features of the true model have been correctly recovered, such that the ensemble mean can be considered to be a standard inversion result. Note that the ETKF-FWI solution is also lacking the resolving power at depth, as it can be expected due to the surface acquisition and the way waves propagate through the subsurface. When the initial ensemble is designed appropriately, its ensemble members will converge toward the global solution along with the ETKF-FWI cycles.

The added value of the ETKF-FWI approach thus lies in the evaluation of the posterior covariance matrix, which individual lines and diagonal are easily accessible, provided the ensemble has been stored. The ensemble covariance matrix allows estimating the uncertainty and resolution information of the solution, which is the main objective of this method. In figure 4.8, for instance, the final variance map allows identifying the prevalent uncertainty zones in the final model. To precisely locate "uncertain" features, local variance peaks have been plotted on both maps. Variance peaks were extracted with a maximum filter of radius 275m. The maximum filter dilates the variance map, and create local zones of homogeneous values. Peaks (or local maxima) are defined as parameters located where the variance map and the output of the maximum filter are equal.

In this example, we can see that most of the measured variance peaks are consistently located along interfaces where high-velocity layers are overlaying lower velocity layers. We might associate high

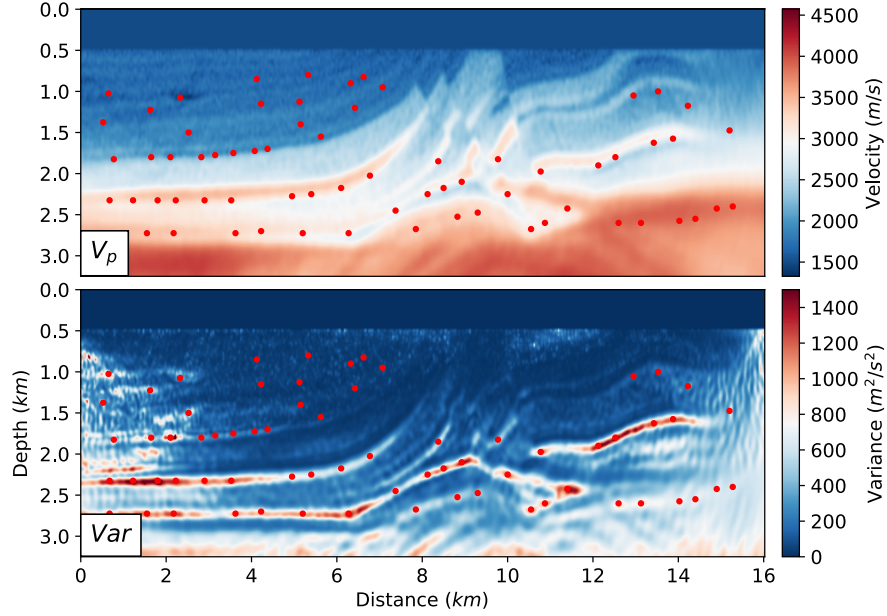


Figure 4.8: Final ensemble mean (top) and variance (bottom). Red dots denotes local maximum variance peaks in both maps.

variance at interfaces with the band-limited context of our application: band-limitation is expected to limit the ability of the optimization scheme to recover sharp discontinuities, which will tend to smooth the interfaces because of the lack of high-frequency content. Another possible source of variability in interface recovering might be the inherent velocity-depth ambiguity in reflection tomography (Yilmaz, 1993).

Additionally, we can observe high variance values toward the depth and lateral limits of the physical domain, where poor illumination is expected. In this case, where sources and receivers are located at the surface, the geometrical spreading of the wavefield prevents the sampling of the sides and bottom of the medium, hence the higher relative uncertainty.

Individual lines of the covariance matrix can also be computed to evaluate how parameters are linked with each other. However, unlike absolute variance values that can be put in perspective when compared with the initial variance, absolute covariances are more challenging to interpret. Thus, we propose to compute the lines of the correlation matrix instead of the covariance matrix.

The correlation matrix is a dimensionless operator that contain correlation coefficients from -1 to $+1$ (Feller, 2008). When the correlation coefficients tend to $+1$, it reflects a strong positive link between two parameters, implying that they share similar physical properties and are evolving similarly. Conversely, a negative correlation coefficient of -1 denotes a strong link but expresses an opposite behavior between parameters. Finally, a correlation coefficient of 0 implies the absence of a physical connection between parameters. To compute the correlation matrix, we first need to define \mathbf{D} as a diagonal matrix containing the variance terms of \mathbf{P} . The correlation matrix is then given by,

$$\mathbf{C} = (\mathbf{D})^{-1/2} \mathbf{P} (\mathbf{D})^{-1/2}. \quad (4.7)$$

\mathbf{C} is thus a dimensionless, normalized version of the covariance matrix, which diagonal terms are all equal to 1 (correlation of a parameter with itself). In the following, we evaluate several correlation maps,

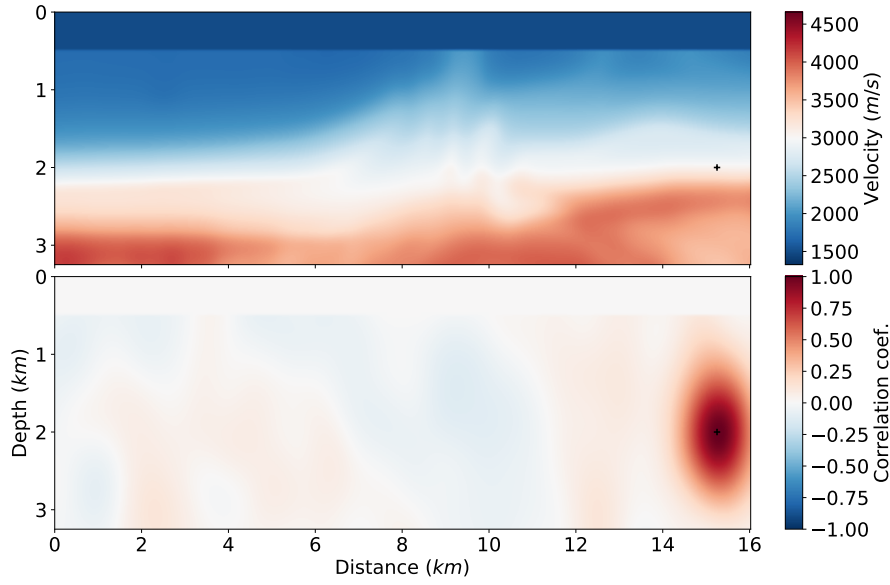


Figure 4.9: Velocity map (top) and correlation map for the velocity parameter located at $x = 15.75\text{km}$ and $z = 2\text{km}$ in the initial ensemble. The circular shape of the positive correlation zone is directly a result of the initial sampling.

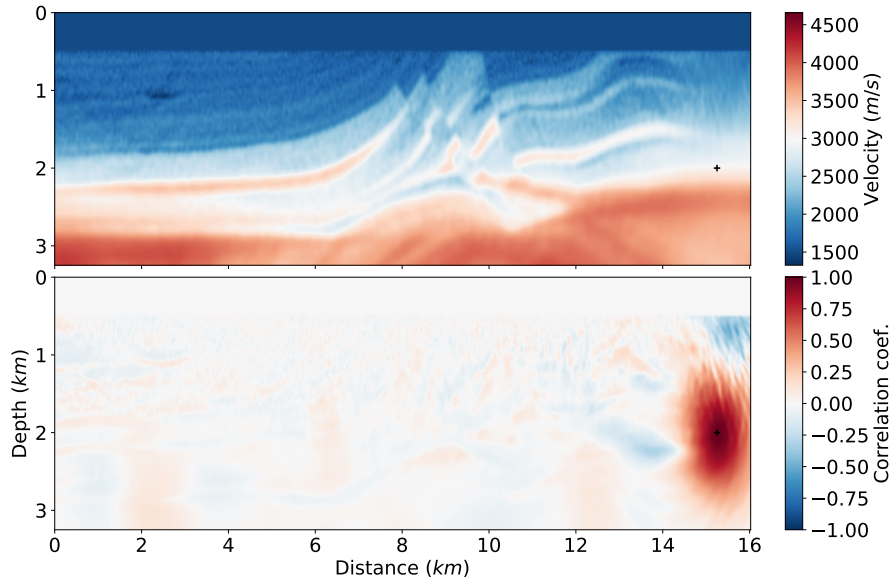


Figure 4.10: Velocity map (top) and correlation map (bottom) for the velocity parameter located at $x = 15.75\text{km}$ and $z = 2\text{km}$ in the final ensemble.

to investigate how recovered parameters are linked with their surroundings. We chose for reference a parameter outside of the optimal illumination zone, close to the lateral limit of the model. We evaluate how the local resolution changes from the initial ensemble (Fig. 4.9) to the final ensemble (Fig. 4.10).

In the initial ensemble (Fig. 4.9), the positive correlation zone around the parameter at position $x = 15.75\text{km}$ and $z = 2\text{km}$ is solely defined by the correlation length of the initial perturbations. It is

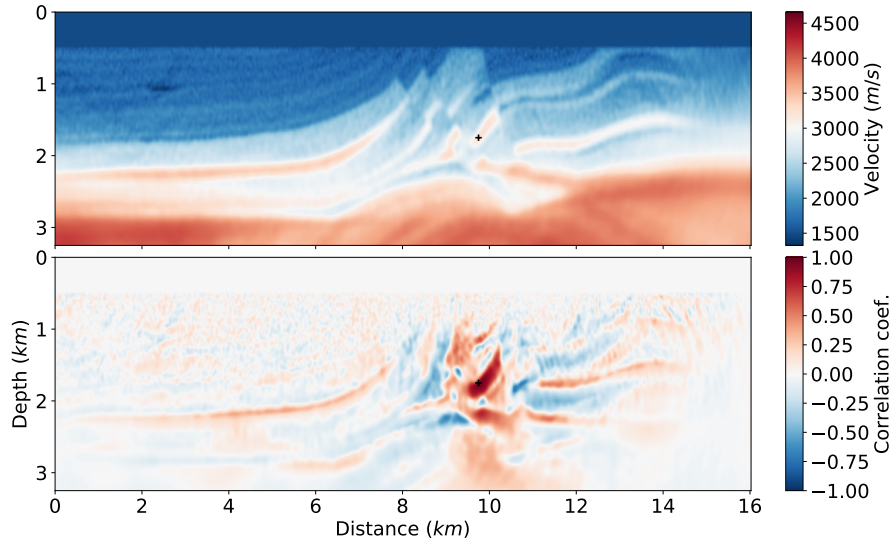


Figure 4.11: Velocity map (top) and correlation map (bottom) for the velocity parameter located at $x = 10.25\text{km}$ and $z = 1.75\text{km}$ in the final ensemble.

thus expected to observe a circular positive correlation zone around the chosen parameter. In the final ensemble, however, for the same parameter, the correlation map is almost identical (Fig. 4.10), meaning the local resolution has not been improved. This is due to the poor illumination in the lateral limits of the domain, which implies that this parameter can only be weakly constrained by the data.

In the following, we chose three other parameters to produce the following correlation maps (Fig. 4.11, 4.12 and 4.13). The parameter picked to produce Figure 4.13 is located inside a small structure within the tilted blocks of the Marmousi model. This has a strong impact on the correlation map, as the high positive correlation zone around the parameter is mostly contained within this small unit section. Note that, contrarily to the example in Figure 4.10, the span of the positive correlation zone has been reduced and is bounded the surrounding structure. This result is somehow expected as the velocity parameters inside the same geological unit should share the same physical properties; it is thus logical that a strong correlation link is found within the unit. Note also that weaker, distant positive correlations structure are present on this map with other high-velocity layers. Negative correlations are also visible in the close vicinity of this unit, as the whole structure close to the parameter influences its recovered velocity.

The third example, in Figure 4.12 has been chosen because it lies precisely on the boundary between a low-velocity layer and a high-velocity layer. In that case, the positive correlations are stretched over the interface, which is also expected in this context. Indeed, as rocks along the interface should be of the same nature, they should share the same physical properties. In this case, the vertical resolution is largely superior to the lateral resolution. Note that this parameter is also positively correlated with the overlaying interface at $z = 1.7\text{km}$ depth. This might indicate that both velocity contrasts are interdependent, or that an ambiguity exists between them in the data. We can also observe what can be associated with finite-frequency artifacts that translate to a change of correlation polarity in the vicinity of the evaluated parameter.

The final example (Fig. 4.13), has been picked to highlight the finite-frequency oscillations in the correlation maps. This parameter has been selected in the upper part of the medium, close to its lateral limit, and we can expect that it has been sampled by a limited number of up-going waves. Thus, the

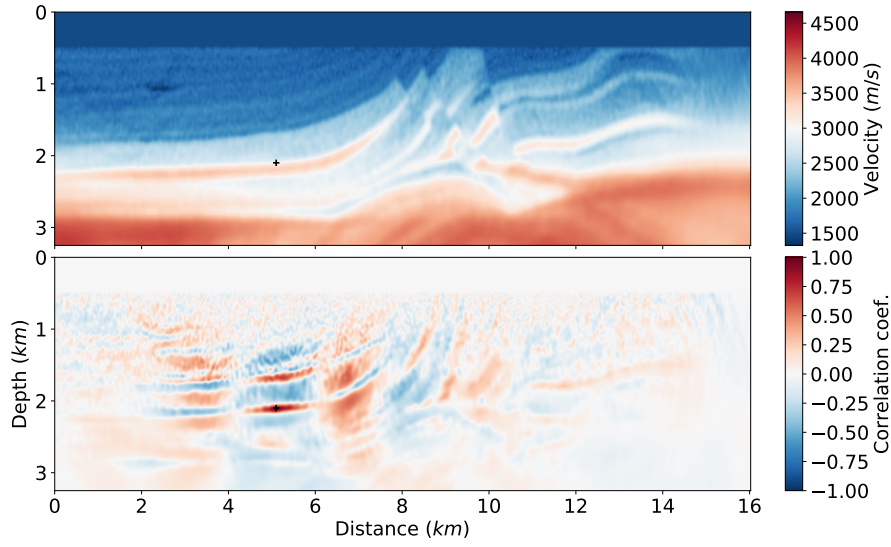


Figure 4.12: Velocity map (top) and correlation map (bottom) for the velocity parameter located at $x = 5.6\text{km}$ and $z = 2.1\text{km}$ in the final ensemble.

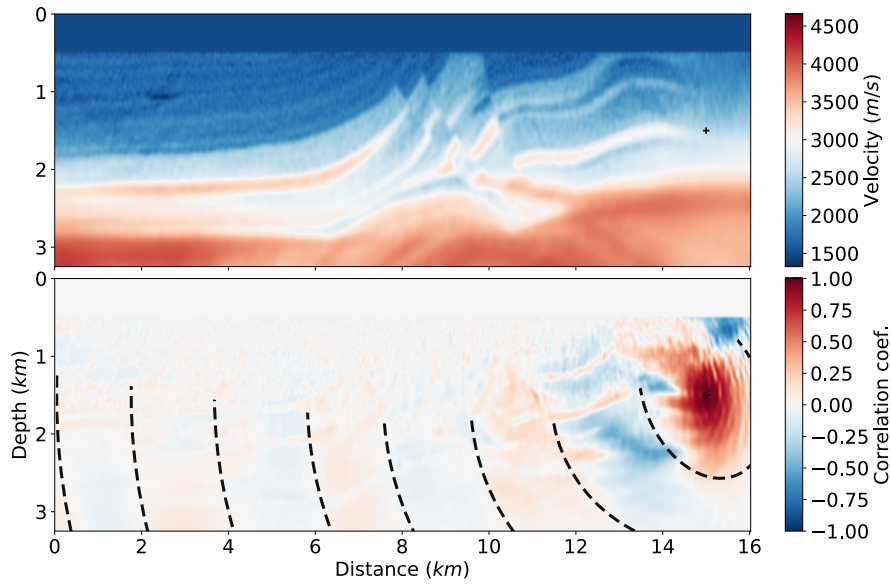


Figure 4.13: Velocity map (top) and correlation map (bottom) for the velocity parameter located at $x = 15.5\text{km}$ and $z = 1.5\text{km}$ in the final ensemble. Dashed lines denote the polarity transition in the correlation map, highlighting the finite-frequency effects.

polarity of the finite-frequency effect in the vicinity of this parameter is no more horizontal, as seen in Figure 4.12, but is tilted upward due to an incomplete illumination. We have drawn the sign transitions of these oscillations with dashed black lines on the correlation map, to highlight the possible link between this effect and the point-spread function of the parameter, caused by band-limited data.

4.3 Investigating undersampling

We have seen in the previous results, that the ETKF-FWI can provide both a parameter estimates, deemed as a satisfying FWI solution, along with uncertainty estimation, both in terms of variance and correlation between the physical parameters. However, the ensemble size of $N_e = 600$, despite seeming to be a favorable case, generates a great computational overhead. It is thus of interest to investigate how the filter behaves when the ensemble size is reduced, as we expect to face undersampling biases, as introduced in 2.3.1. In this section, we will be comparing the results from the test at $N_e = 600$ with smaller ensembles defined as $N_e = [20, 100]$.

4.3.1 Parameter estimate

We first evaluate the impact of ensemble size over the state estimate. To that extent, we generate smaller ensembles by randomly selecting pre-computed initial perturbations, from our previous $N_e = 600$ case. Doing so, we make sure that our three ETKF-FWI tests are built on very similar starting ensembles. To ensure that our comparison is meaningful, the ETKF-FWI is run with the exact same settings as in the previous section, the only difference being the ensemble size (and thus the computational cost). The results of the three tests are displayed in Figure 4.14.

From these results, it seems that the ensemble size has no visible influence over the quality of the reconstructed model. All three test cases lead to consistent parameter estimation, as all results are fairly visually comparable. This outcome might seem surprising, given that we expected undersampling to occur (and perhaps filter divergence could have happened). However, because our forecast operator is solving an optimization problem rather than a forward problem, it seems that the ensemble mean is clear of any drifting or divergence effect, commonly encountered in typical dynamic EnKF applications. Indeed, even with small ensembles, if they have been appropriately generated, all the ensemble members should converge toward the same solution of the optimization problem, preventing any unwanted divergence effects.

To complete this analysis, we have computed the root-mean-square-error (RMSE) values of the ensemble final means, with the true model as a reference:

$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (m_{n,true} - \bar{m}_n^a)^2}, \quad (4.8)$$

where $m_{n,true}$ is the n^{th} parameter of true velocity model, and \bar{m}_n^a is the n^{th} parameter of the final ensemble mean. Values of RMSE reduction between the initial model m_0 and final ensemble means are displayed in Table 4.1: RMSE reduction is not affected by the ensemble size and all parameters estimates are nearly identical, which is consistent with the results observed in Figure 4.14.

Table 4.1: Normalized root-mean-square model error (RMSE) reduction with respect to the initial model for various ensemble sizes. The RMSE values are computed between the final ensemble mean and the true model. The RMSE reduction is computed with the initial model RMSE as a reference.

N_e	20	30	40	50	60	80	100	200	300	600
RMSE reduction	15.4%	15.4%	15.6%	15.0%	15.2%	15.9%	15.7%	15.6%	15.3%	15.6%

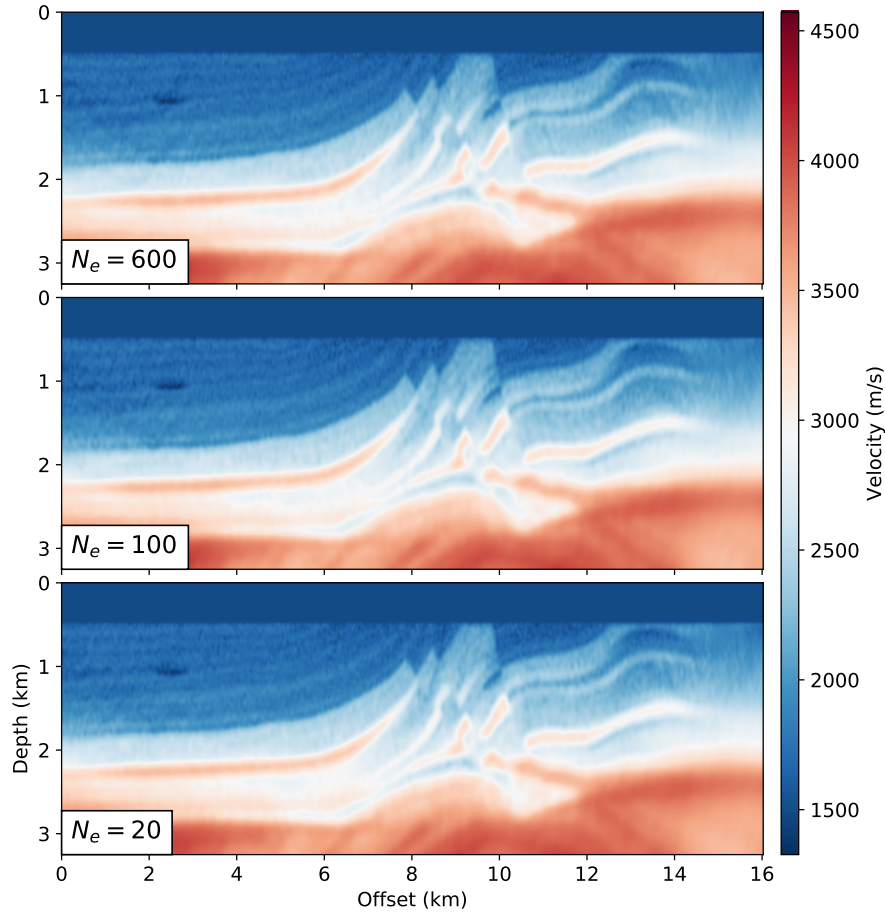


Figure 4.14: Final analyzed ensemble means for ensemble sizes $N_e = 600, 100, 20$ after 15 ETKF-FWI cycles from $3Hz$ to $10Hz$.

4.3.2 Variance approximation

We now look at the ensemble repartition to evaluate how undersampling may affect the results of our uncertainty estimation scheme. As a matter of reference, we first present the three initial variance maps in Figure 4.15. In this figure, we can see that the initial repartition is already affected by sampling biases. While the initial variance map for $N_e = 600$ is nearly homogeneous, the map for the $N_e = 20$ ensemble has significant variations of amplitude, purely due to the insufficient sampling (as it has been presented earlier in Figure 2.6).

The final variance maps after 15 ETKF-FWI cycles for the three test cases are plotted in Figure 4.17. Contrarily to the state parameter estimate, we observe a substantial lack of result consistency regarding the ensemble size. The various ensemble size tested reveals that the ensemble covariance is strongly affected by undersampling, as we expected. In this instance, undersampling seems to be responsible for variance underestimation.

Assuming that the test with ensemble $N_e = 600$ is the least sensitive to undersampling biases, we define it as our reference result to make several observations. The $N_e = 20$ case displayed in Figure 4.15 exhibit a severe underestimation of the variance values over the whole domain. In the deeper part of

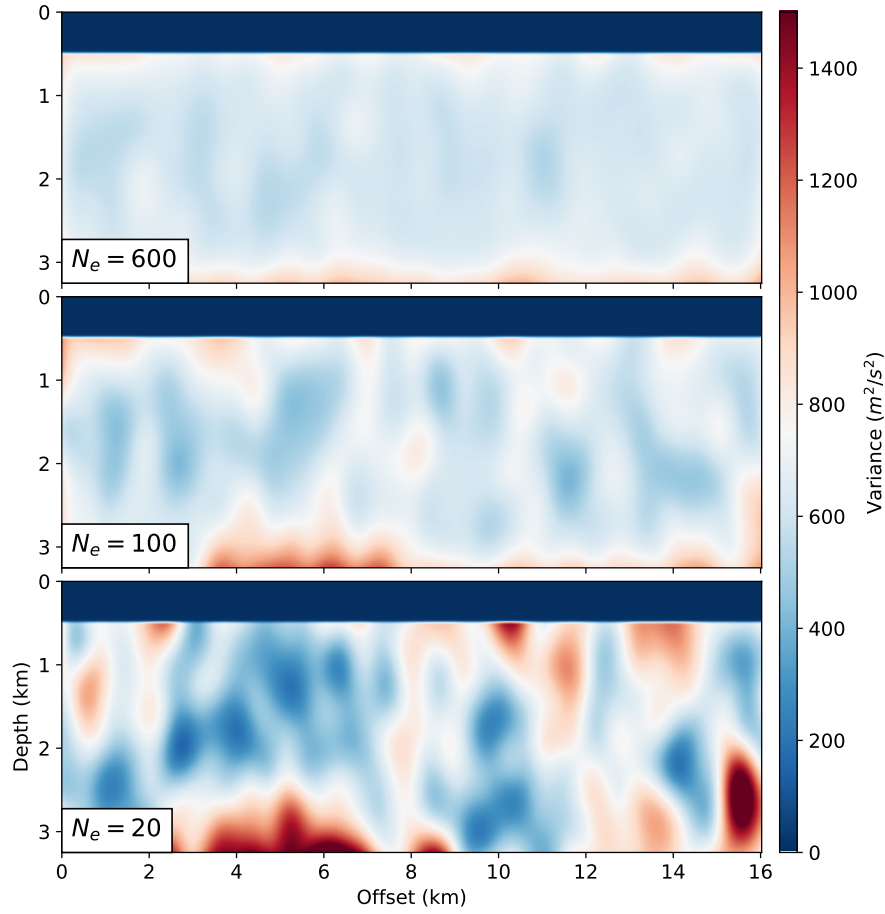


Figure 4.15: Initial ensemble variance maps for ensemble sizes $N_e = 600, 100, 20$.

the domain, we also observe non-physical oscillatory behavior, resulting from the poor initial variance approximation. In the $N_e = 100$ case, these oscillations are not visible, thanks to a better initial sampling. The variance values are nonetheless slightly underestimated. The qualitative aspect of the variance map is at least preserved, as opposed to the $N_e = 20$ case.

To better understand the results of Figure 4.15 and go beyond simple qualitative comparison, we evaluate absolute variance values from a set of ETKF-FWI realizations for $N_e = [20, 600]$. We evaluate the underestimation of variance by computing the mean variance value for each variance map. We plot the averaged variance against ensemble size in Figure 4.16. Crosses denote the ensemble sizes used to generate this curve.

As it stands, the trend in absolute variance values seems to be consistent with the variance underestimation observed in Figure 4.17. It is also worth noting that variance estimates behave almost asymptotically, which means we can hope to find a compromise between too small and too big ensembles. Although, it seems complicated to estimate this "optimal" ensemble size in advance, nor it is practical to evaluate it by trial and error. We could instead mitigate variance underestimation by using multiplicative-inflation (2.3.5), provided an adequate inflation factor can be chosen.

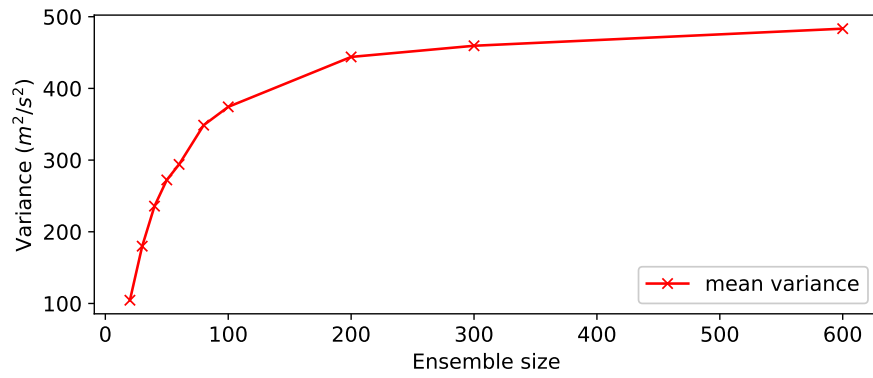
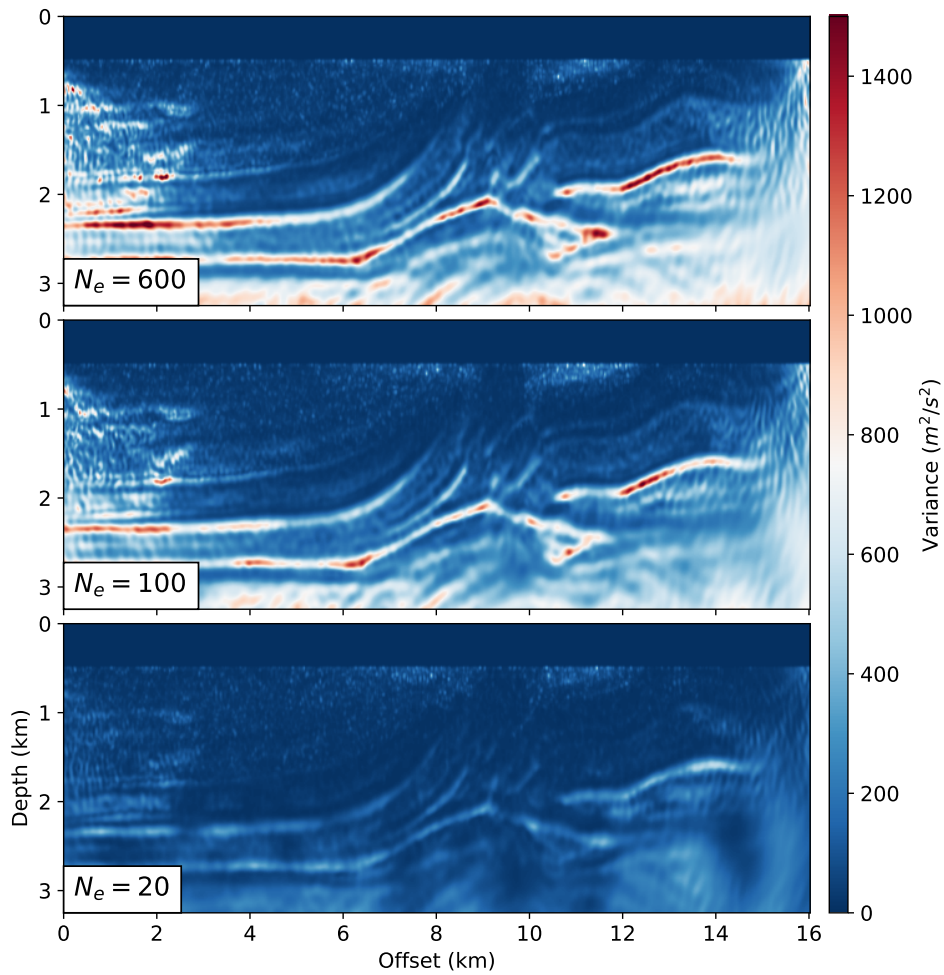


Figure 4.16: Average variance plotted against ensemble size.

Figure 4.17: Initial ensemble variance maps for ensemble sizes $N_e = 600, 100, 20$.

4.3.3 Correlation approximation

Additionally, we compute the correlation maps for the parameters presented in Figures 4.11, 4.12 and 4.13, to evaluate the effects of undersampling on the off-diagonal terms of the covariance matrix. The results for the three parameters are plotted in Figures 4.18, 4.19 and 4.20. Again, we consider the result at $N_e = 600$ to be the reference case, as it should be less affected by undersampling biases. Note that even though the off-diagonal covariance terms might be underestimated, correlation maps allow a fair comparison, as the correlation matrix is a normalized version of the covariance.

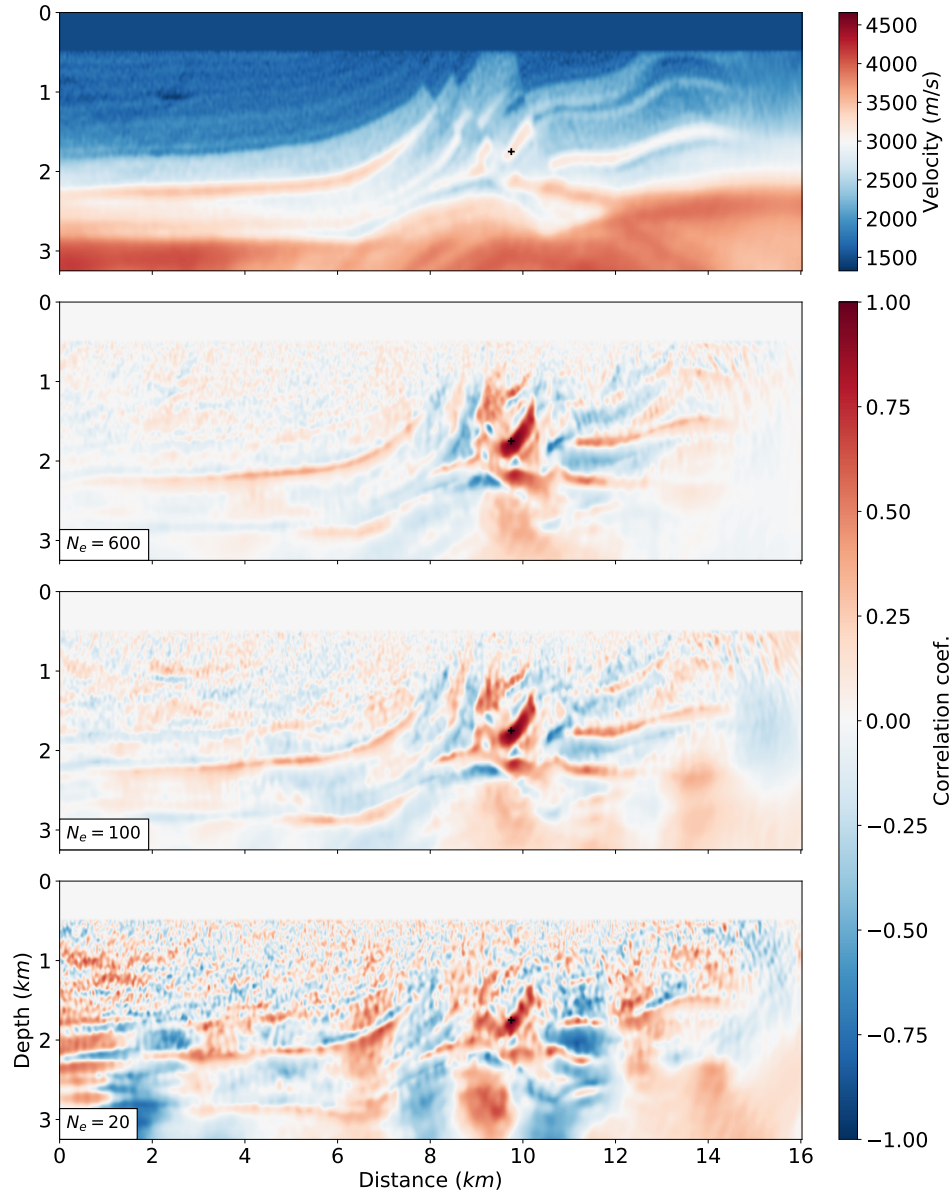


Figure 4.18: Velocity map (top) and correlation map (bottom) for the velocity parameter located at $x = 10.25\text{km}$ and $z = 1.75\text{km}$ in the final ensemble, for the three test cases.

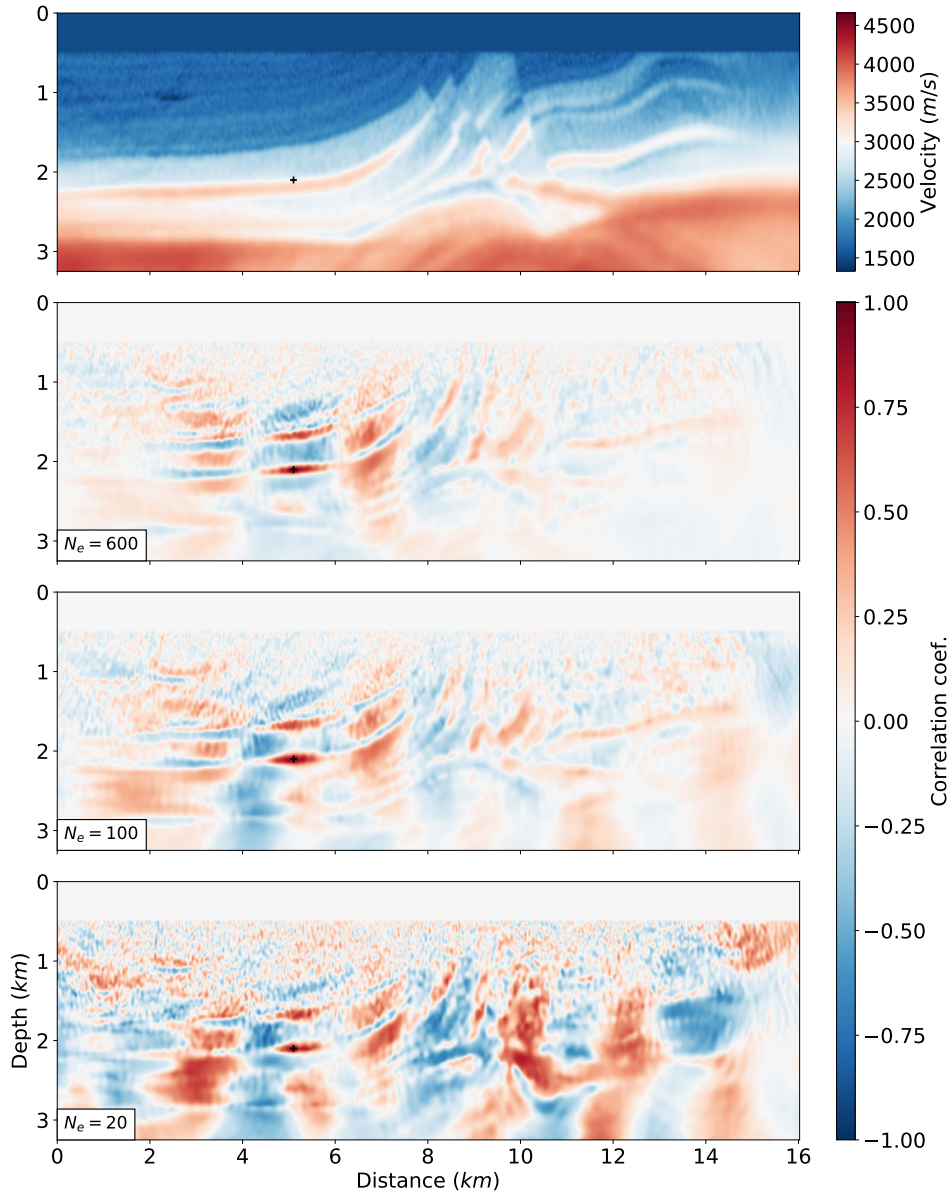


Figure 4.19: Velocity map (top) and correlation map (bottom) for the velocity parameter located at $x = 5.6\text{ km}$ and $z = 2.1\text{ km}$ in the final ensemble, for the three test cases.

On both the $N_e = 20$ and $N_e = 100$ cases, we observe undersampling biases, which are manifested by the appearance of spurious correlation terms in the correlation maps. The $N_e = 100$ case allows a good approximation of the reference case, which is consistent with the undersampling curve presented in Figure 4.16.

On the other hand, the $N_e = 20$ case suffers from serious undersampling biases: spurious correlations are well visible in the whole domain, and their amplitudes are not decaying along with the distance from the inspected parameter, as in the $N_e = 100$ and the reference cases.

However, it seems that despite the strong undersampling biases, short-range correlation terms are somehow preserved. Thus, even though the $N_e = 20$ is displaying severe undersampling, we might be

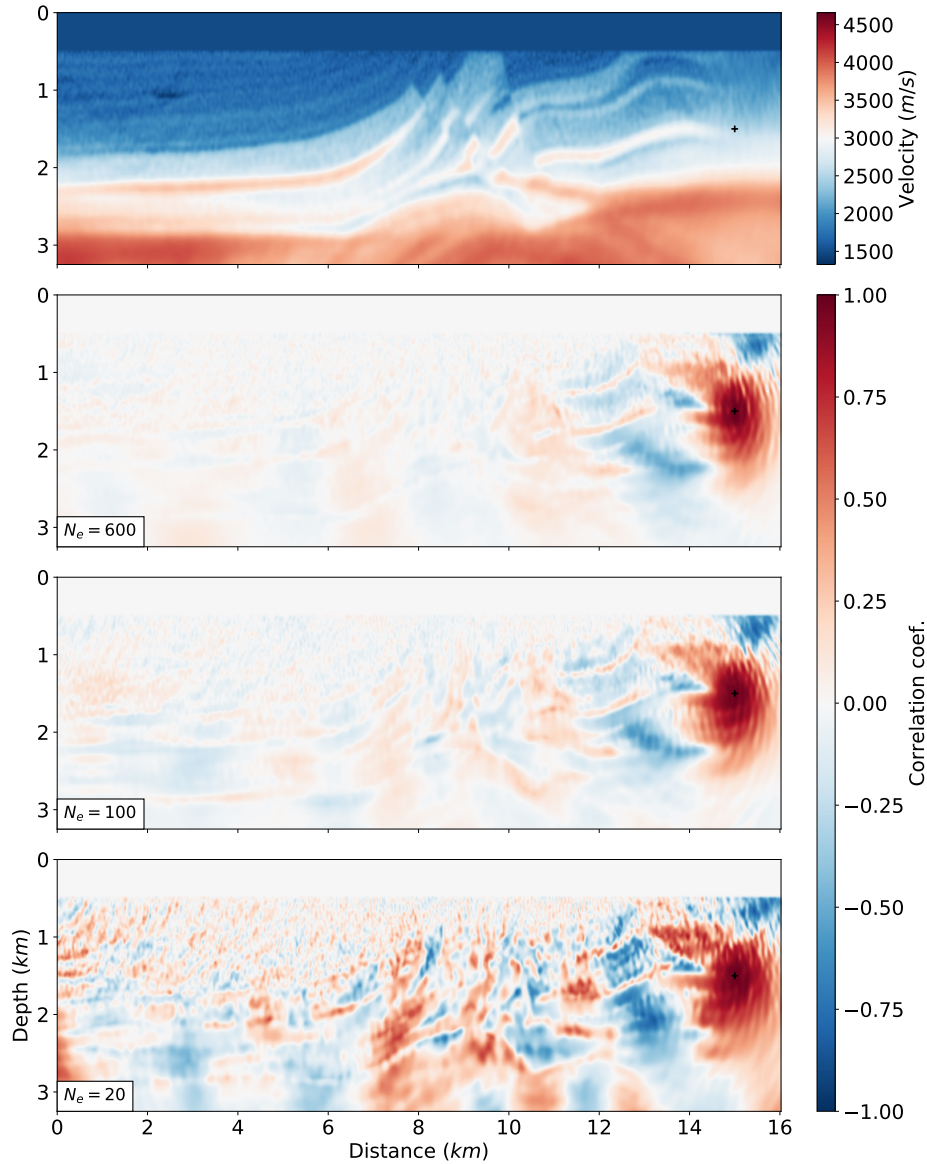


Figure 4.20: Velocity map (top) and correlation map (bottom) for the velocity parameter located at $x = 15.5\text{km}$ and $z = 1.5\text{km}$ in the final ensemble, for the three test cases.

able to extract local structural information, such as dipping angles or thickness of geological units and "local resolution" from these biased correlation maps.

I have mentioned in Chapter 2 that spurious distant correlation could be corrected with the use of Covariance Localization (2.3.5) or Local Analysis (2.3.5.1). In our case, CL via a Hadamard product seems unrealistic, as building large covariance matrices to filter-out unwanted off-diagonal terms would be computationally prohibitive. LA does not appear to be a satisfying solution either, due to the non-local nature of our observations: recorded signal travels through a significant part of the medium, and each parameter along the way directly influence the recorded signal. It would thus be difficult to design an appropriate local domain and perform local analysis without losing meaningful correlation terms arising from our observation operator $\mathcal{H}(m)$.

4.4 Mitigating undersampling

Now that the previous tests have highlighted undersampling biases, we may attempt to mitigate them with methods from the DA literature. Providing solutions to mitigate undersampling would improve the ETKF-FWI outcomes for smaller ensembles and thus yield satisfying solutions at a lower computational cost.

I have mentioned previously that neither CL nor LA is easily applicable due to the size of our problem and the nature of our observations. Therefore, we will solely focus on tackling inbreeding biases with multiplicative inflation in order to mitigate variance underestimation. To do so, we have replicated the previous comparative study, with several inflation factors, in order to evaluate the effect of multiplicative inflation on our ETKF-FWI outcomes.

To apply multiplicative inflation, we artificially inflate the forecasted ensemble according to

$$\mathbf{M}_i^f = r(\mathbf{m}^f - \bar{m}^f), \quad (4.9)$$

where r is an inflation factor. We performed several tests with $r = [1.005, 1.01, 1.015, 1.02]$ (recalling that inflation factors suggested in the DA literature are ranging from 1% to 7%).

Inflation is performed before the observation operator is applied to the ensemble members to compute forecast data. The analysis is then computed following the exact same scheme as in the previous test. Multiplicative inflation is thus easily implementable and does not add any significant computing operations.

To visualize the effect of the inflation factor over the ensemble, we first compute mean variance curves of the inflated ensembles, as in Figure 4.16. The results of this comparative inflation test are presented in Figure 4.21.

It appears that multiplicative inflation successfully increases the final variance values amongst all ensemble size, and its effect is more significant on larger ensembles. From this plot, it seems that the ensemble of size $N_e = 50$ with a multiplicative inflation factor of $r = 1.02$ should result in variance estimates close to the reference case (denoted by the black dashed line in Figure 4.21). We compare the reference variance map with the $N_e = 50$ case with, and without multiplicative inflation in Figure 4.22.

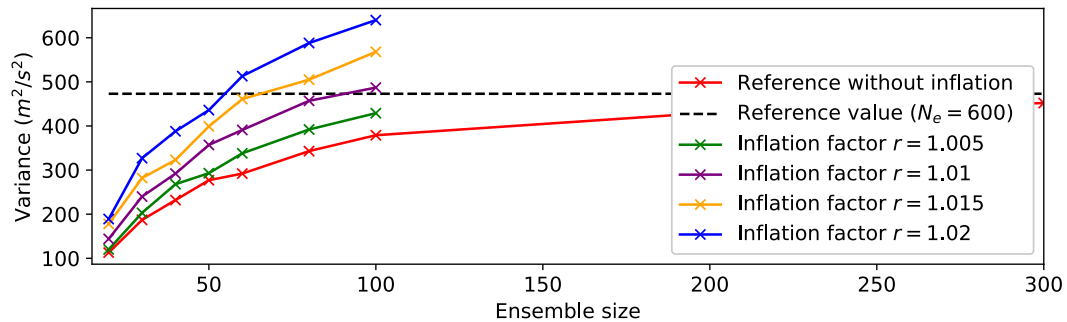


Figure 4.21: Average variance plotted against ensemble size for several multiplicative inflation values. The black dashed line represent the "optimal" reference case obtained with $N_e = 600$ and without multiplicative inflation.

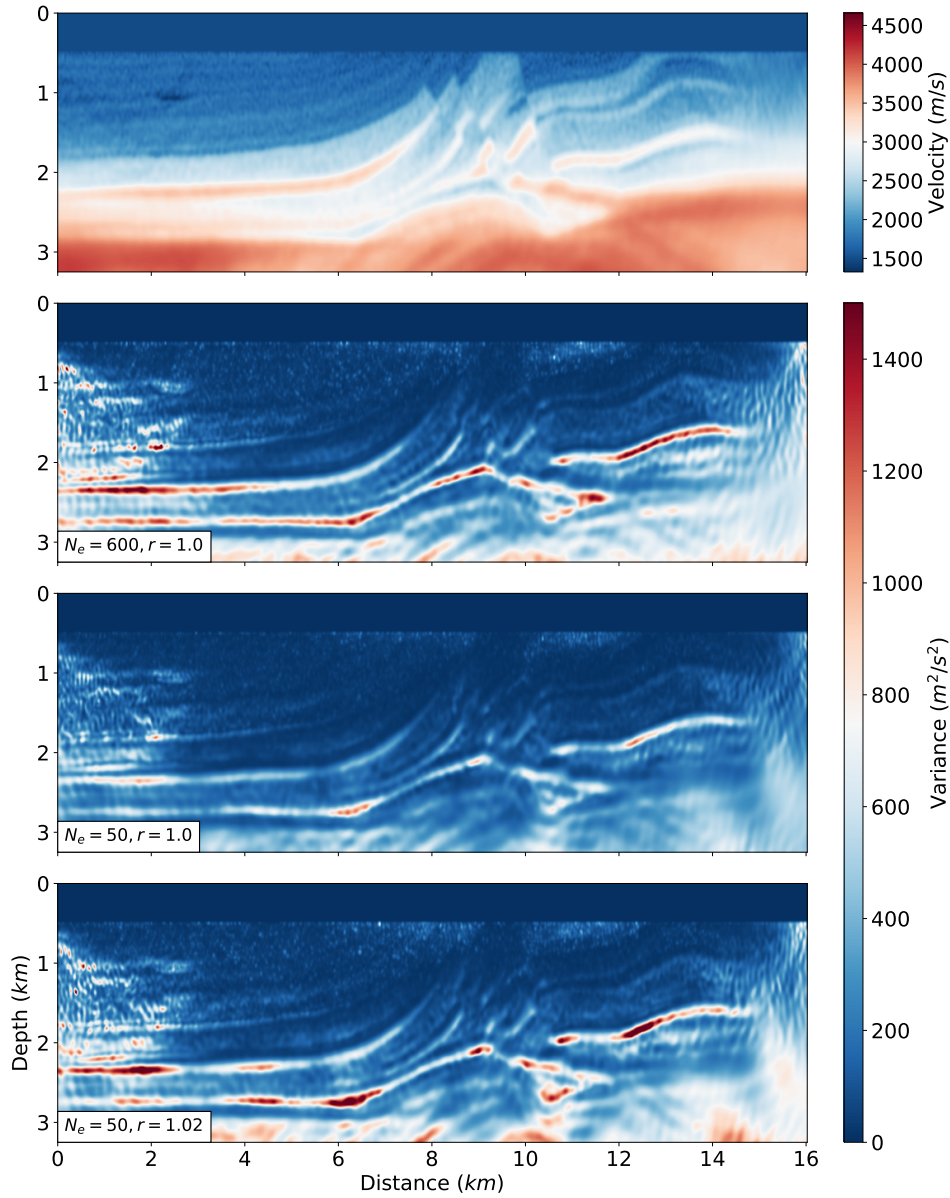


Figure 4.22: Comparison of variance maps on the $N_e = 50$ case, with and without inflation. The $N_e = 600$ case and the mean model are included for reference. With multiplicative inflation, the final variance map is close to the reference case, despite the original $N_e = 50$ case being severely undersampled.

As predicted by the variance curve in Figure 4.21, the ETKF-FWI with ensemble size $N_e = 50$ and multiplicative inflation parameter $r = 1.02$, allows us to improve the outcome of the filter, compared with the case without inflation. Significant improvements are visible at major interfaces, where the variance has been correctly estimated. Variance has also been increased in the deeper domain, highlighting deep structures, which better fits the reference case variance map.

To further inspect the effect of multiplicative inflation on estimating the covariance matrix, we also look at correlation maps with and without inflation. While spurious correlations are generally corrected

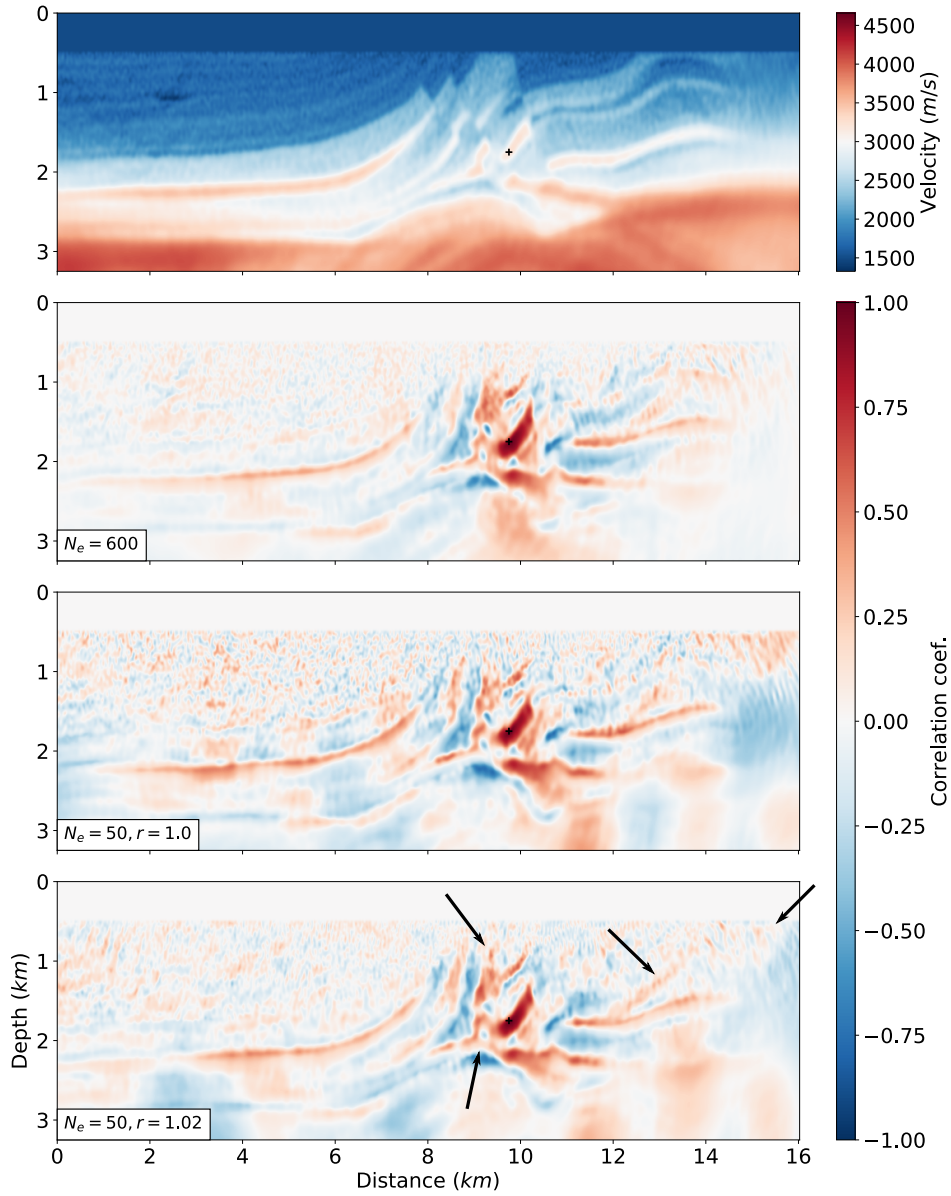


Figure 4.23: Comparison of correlation maps on the $N_e = 50$ case, with and without inflation. The $N_e = 600$ case and the mean model are included for reference. Black arrows in the bottom panel are highlighting features that have been improved in the result with multiplicative inflation.

with localization methods, we ought to verify if the improved variance estimation, allows retrieving better correlation maps (as the variance is involved in correlation computation). Results are displayed in Figure 4.23.

We can see in Figure 4.23 that improving the variance estimate with multiplicative inflation yields slight improvements in the correlation maps. Notably, we can observe a reduction of amplitude in the shallow region, away from the parameter. We have highlighted features which recovery has been improved by multiplicative inflation by black arrows. We also computed correlation maps comparisons for the previously tested parameters in Figure 4.24 and 4.25.

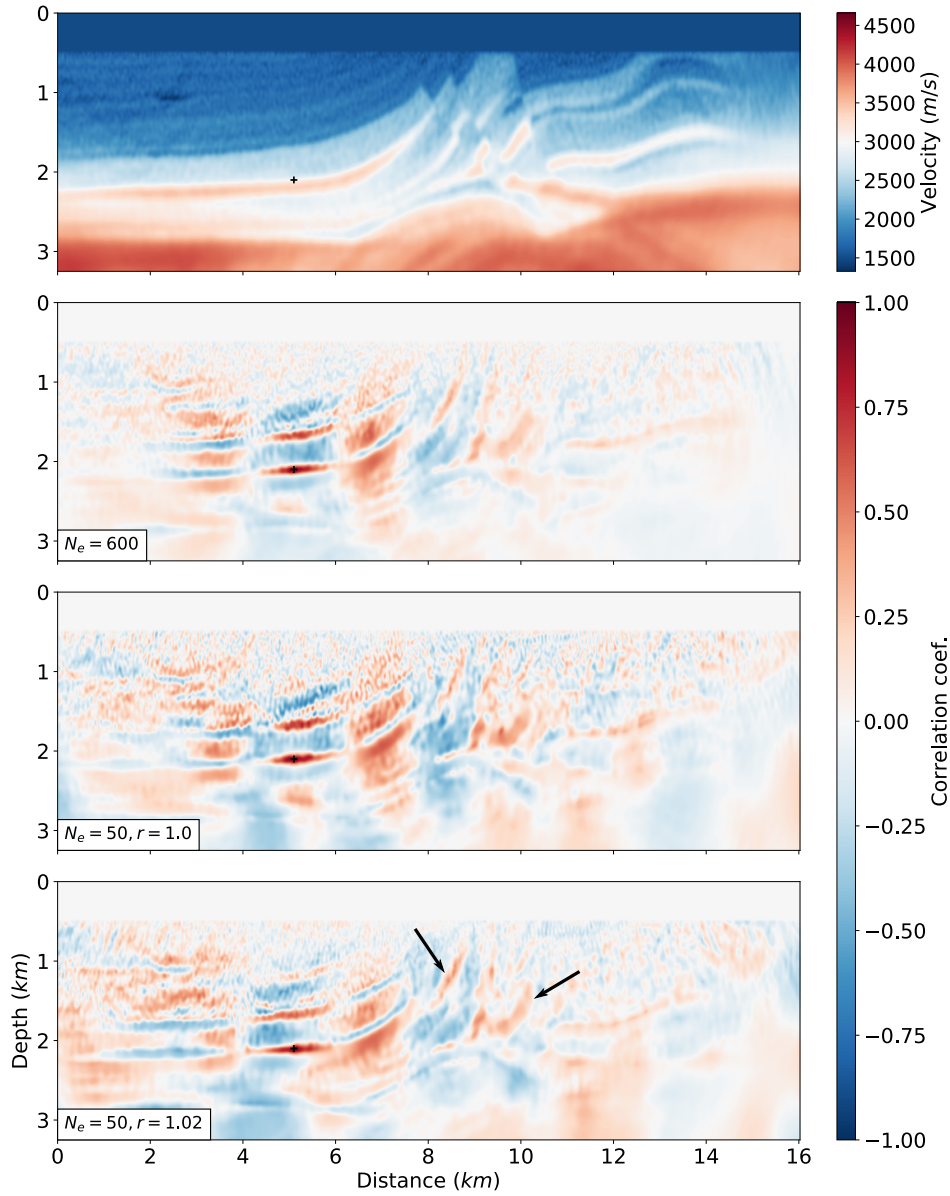


Figure 4.24: Comparison of correlation maps on the $N_e = 50$ case, with and without inflation. The $N_e = 600$ case and the mean model are included for reference. Black arrows in the bottom panel are highlighting features that have been improved in the result with multiplicative inflation.

While we observe some slight improvements over the case without inflation, all three cases are still exhibiting strong spurious correlations terms that are not corrected by inflation (as expected).

Finally, note that in this favorable case, the inflation parameter has been chosen a-posteriori after an extensive series of tests with varying ensemble size and inflation parameters. In practice, the application of multiplicative inflation would be troublesome, as we would be left with determining r purely arbitrarily or by trials and errors. Therefore, despite showing promising results, we cannot expect to apply multiplicative inflation in practical applications.

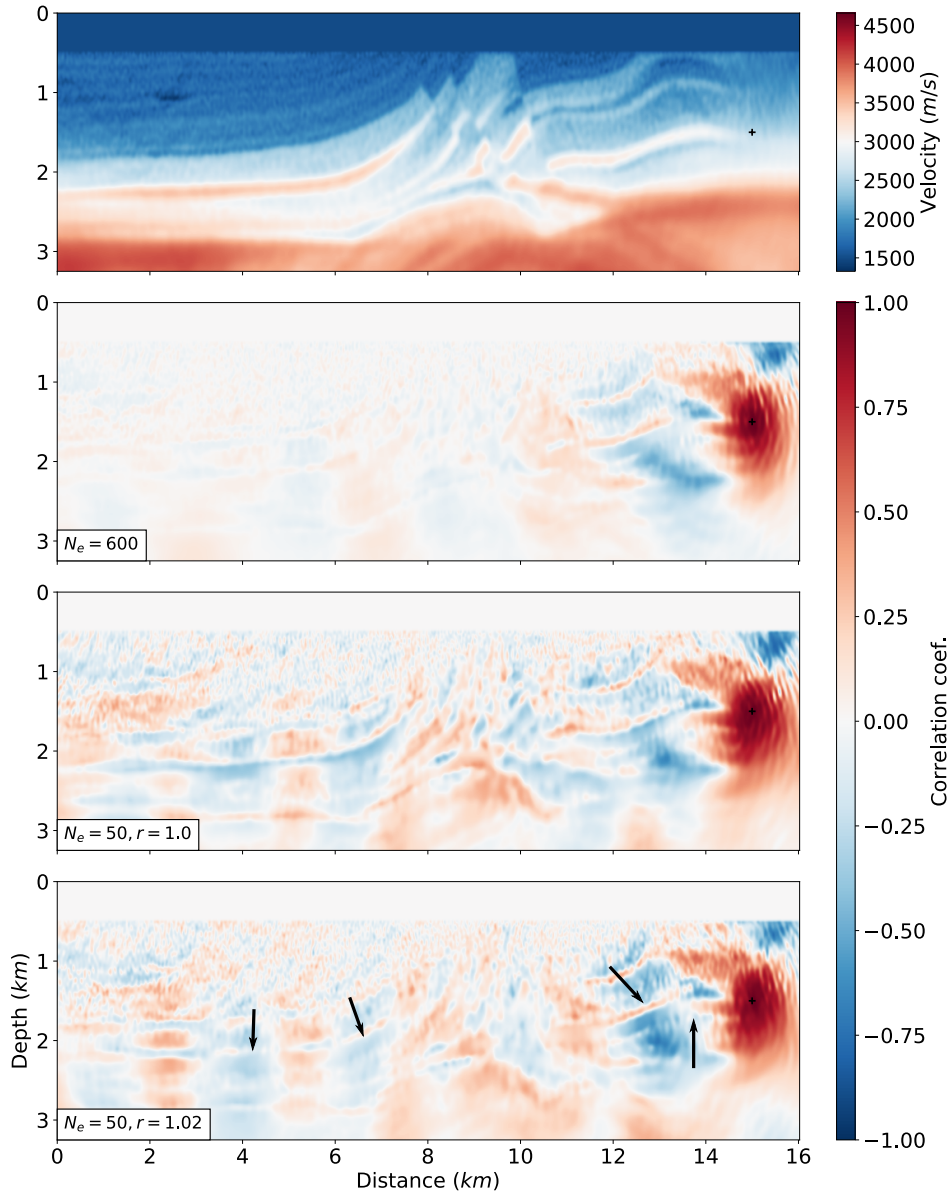


Figure 4.25: Comparison of correlation maps on the $N_e = 50$ case, with and without inflation. The $N_e = 600$ case and the mean model are included for reference. Black arrows in the bottom panel are highlighting features that have been improved in the result with multiplicative inflation.

Conclusion

In this chapter, we have presented a complete application of the ETKF-FWI on the Marmousi II synthetic benchmark. We have shown that the ETKF-FWI can yield a satisfying parameter estimate along with uncertainty estimation. We have also highlighted the importance of the analysis step in maintaining coherency and reducing uncertainty along the ETKF-FWI cycles. We have discussed how uncertainty information could be extracted from the diagonal and lines of the covariance matrix, provided the final ensemble is stored.

Variance maps provide a direct qualitative uncertainty assessment of the solution by indicating which parameters in the medium have a relatively high variance. We have seen that in this particular setting, uncertainty seems to be dominated by the geometrical spreading effect, and structural uncertainty: higher variance tends to be located at depth and on the lateral limits of the domain, and on sharp velocity contrasts. These observations are in agreement with our understanding of FWI's properties, regarding parameter recovery and illumination. Mind that the quantitative variance estimate and hence, our uncertainty estimation is limited by several factors:

- It is directly tied to our initial variance, from which it can only be compared.
- It is limited by the finite-frequency wave propagation involved in our forecast, hence uncertainty is dependent on how waves "see" the medium
- It provides information on the misfit function solely in the vicinity of the proposed solution, as we rely on a local optimization scheme.

Thus, we chose to address mainly qualitative uncertainty assessment.

We have also been able to exploit off-diagonal terms of the covariance matrix by computing their correlation counterpart, as suggested by Tarantola (2005). Correlation maps appear to be a good proxy for resolution analysis, as they indicate how parameters are linked with their surroundings. So far, the results we have observed in correlation maps are in agreement with our theoretical expectations (strong positive correlations between parameters in the same unit, for instance).

Following these preliminary results, we were able to investigate and characterize how undersampling can affect ETKF-FWI outcomes. This investigation leads us to think that state estimation should be devoid of undersampling biases. This could be explained by the inverse nature of our forecasting operator that tends to bring the ensemble members together and avoid possible ensemble divergence. Undersampling is more serious when it comes to uncertainty estimation, as a small ensemble may lead to an unreliable uncertainty measurement. Variance underestimation thus makes quantitative uncertainty estimation complicated, as we are likely to underestimate uncertainty in most cases.

Finally, we successfully mitigated variance underestimation by implementing a multiplicative inflation procedure within our ETKF-FWI scheme. Although inflation has effectively reduced this undersampling bias, it is not a satisfactory solution to our problem. Indeed, adjusting the inflation parameter to the ensemble size would probably involve trial and error, which makes it improper to be a reliable solution in our case.

To complete this thesis work, I ought to present a second ETKF-FWI application. Thus, the next chapter will be dedicated to field-data seismic exploration FWI application to demonstrate the applicability of the ETKF-FWI scheme, in a less favorable case.

Chapter 5

Field data application of the ETKF-FWI

Contents

5.1 The Valhall oil field dataset and ETKF-FWI parameterization	110
5.1.1 P-wave velocity reconstruction	113
5.1.2 P-wave velocity and density reconstruction	116

To complement the synthetic benchmarks detailed in Chapter 4, I wish to present a field-data application, to demonstrate the feasibility of our method on real cases. The results in this chapter owe much to the data processing work of Zhou (2016); Zhou et al. (2018), which greatly simplified setting up the following experiments.

The interests in applying our ETKF-FWI scheme to a field-dataset are multiples. For instance, it makes it possible to evaluate the sensitivity of the ETKF-FWI scheme to complex noise structures. It also allows assessing its robustness regarding complex inverse problems, as we will be accounting for visco-acoustic and anisotropic effects presents in the data. Finally, we will also evaluate the scheme capacity to deal with the inversion of multiple subsurface physical parameters (P-wave velocity and density) and to measure multi-parameters crosstalks.

Most of this final chapter is based on the following publication:
J. Thurin, R. Brossier, L. Métivier, Ensemble-based uncertainty estimation in full waveform inversion, *Geophysical Journal International*, Volume 219, Issue 3, December 2019, Pages 1613–1635.

5.1 The Valhall oil field dataset and ETKF-FWI parameterization

Field-data setting and domain parameterization: The region of interest of this exploration field-data experiment is the Valhall oil field, located in the southern Norwegian North sea, 300 km southwest of the city of Stavanger (Fig. 5.1), within the Central Graben area (Munns, 1985; Leonard and Munns, 1987; Barkved et al., 2010; Sirgue et al., 2010). It features 2.4 km of Tertiary succession overlaying two Upper Cretaceous oil-bearing chalk formations. The Tertiary overburden contains a low-velocity, gas-charged shale formation (referred to as a "gas-cloud" in previous studies), which severely distorts seismic acquisition (Munns, 1985) and makes seismic imaging challenging.

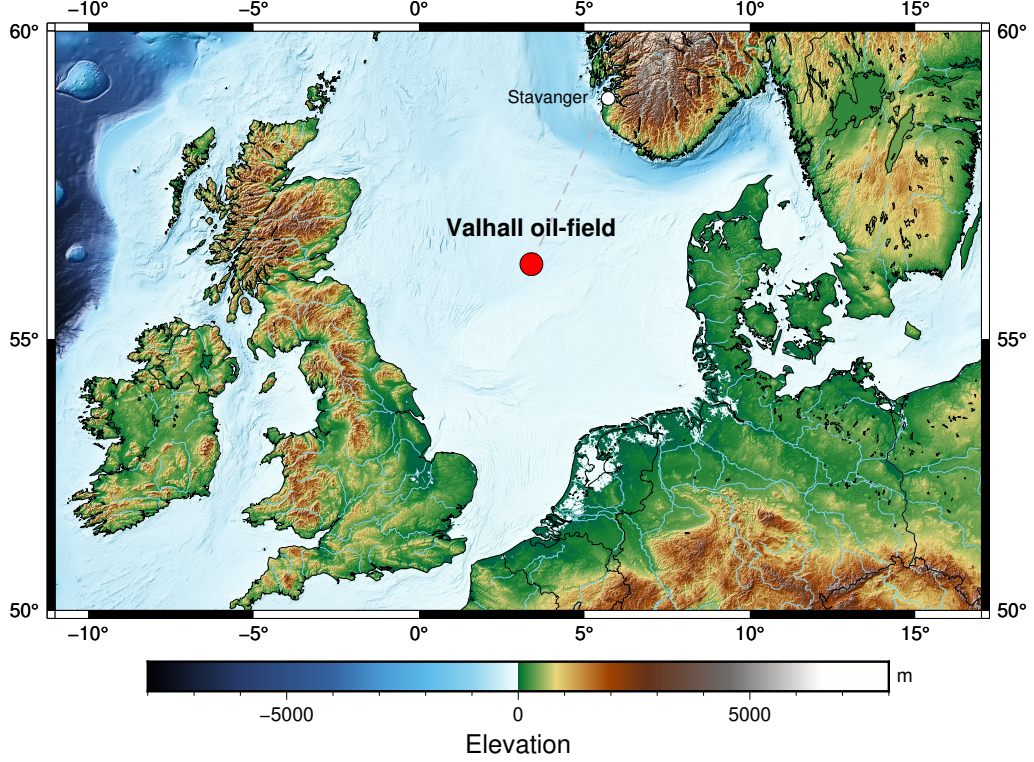


Figure 5.1: Location of the Valhall oil field in the Norwegian North sea. Distance from Stavanger has been drawn with a dashed line and is approximately 300 km.

The Valhall oil field is located in a shallow water environment (70m water column), and benefited from a deployment of thirteen, four-component Ocean Bottom Cables (OBC) in the frame of the Valhall Life of Field Seismic project (Barkved et al., 2003). On top of its impressive instrumentation, the advantage of this case study is that it has been well documented and FWI has already proven to be successful with both 2-D (Prioux et al., 2013a,b; Gholami et al., 2013; Zhou et al., 2018) and 3-D datasets (Sirgue et al., 2010; Operto et al., 2015; Operto and Miniussi, 2018). For our experiment, we consider a 2-D section with a domain width and depth of respectively $x = 16.725\text{ km}$ and $z = 5.025\text{ km}$ with a vertical and horizontal resolutions of $dx = dz = 25\text{ m}$, for a total of $l = 134469$ discrete grid points. We recall that in the monoparameter case, the number of grid points l is equal to the number of degrees of freedom n .

Observations properties: The dataset is composed of 4 components OBC recordings. From the full acquisition which contains 50,824 shots for 2320 receivers, we extract a 2D line containing 192 sources and 315 receivers (which makes each frequency data vector composed of 60480 entries), the same as the one used by (Zhou et al., 2018) and corresponds to the "Cable 13" in this study. The total number of discrete data is equal to $60480 \times n_\omega$. OBC receivers are evenly spaced (50 m) and lie fixed on the seabed (70 m depth). The selected sources are also evenly spaced (50 m) at a constant 5 m depth. In this application, we only exploit the hydrophone out of the 4 components recordings.

Defining the forecast: The ETKF-FWI scheme follows the same setup as in Chapter 4. To ensure the best-case scenario result, we work with an ensemble of $N_e = 600$ members, as the application size is of the same order of magnitude as the synthetic test case. We choose to work with $K = 6$ ordered groups of frequencies ranging from 3.56 Hz to 7.01 Hz (Table 5.1). This frequency

	Freq. 1 (Hz)	Freq. 2 (Hz)	Freq. 3 (Hz)	Freq. 4 (Hz)
$k = 1$	3.567913	3.937008	4.306102	
$k = 2$	4.060039	4.429133	4.798228	
$k = 3$	4.552165	4.921259	5.290354	
$k = 4$	5.044291	5.413385	5.782480	
$k = 5$	5.536417	5.905511	6.274606	
$k = 6$	5.536417	5.905511	6.397637	7.012795

Table 5.1: Table of the frequency groups used in Zhou et al. (2018). These frequencies are used in our frequency continuation strategy and thus define our dynamic axis proxy.

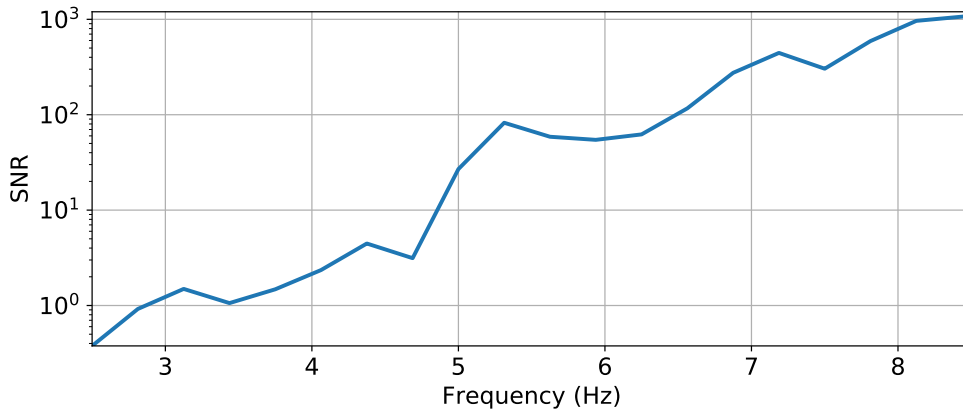


Figure 5.2: Estimated signal to noise ratio plotted against frequency.

selection strategy has been suggested in preliminary work conducted by Zhou et al. (2018) on this dataset and has proven to be adequate for this specific application. Using frequency groups rather than mono-frequency data ensures that each inversion cycle relies on redundant information, which helps to mitigate the impact of noise and inter-parameters cross-talk for multi-parameter FWI. This brings the amount of mono-frequency data pieces of each ETKF-FWI cycle to $n_\omega = 17$. In each of the cycles, the forecast operator \mathcal{I}_n is set to perform $n = 10$ minimization iterations with a preconditioned l-BFGS optimization scheme.

Defining the measurement noise matrix: Following the successful application in Chapter 4, we also define the measurement noise matrix \mathbf{R} as a scaled Identity matrix. We computed each of the monochromatic data noise variances, from measured SNR values (Fig 5.2) in the dataset, according to equation 4.3. As we are working with frequency groups rather than monochromatic data, the measurement noise matrix is block-diagonal, each block of \mathbf{R} corresponding to a monochromatic data variance.

In the following, we first present a monoparameter inversion test, analogous to the synthetic benchmark. The initial velocity model m_0 obtained with reflection tomography, along with the ensemble initial variance map, are displayed in Figure 5.3.

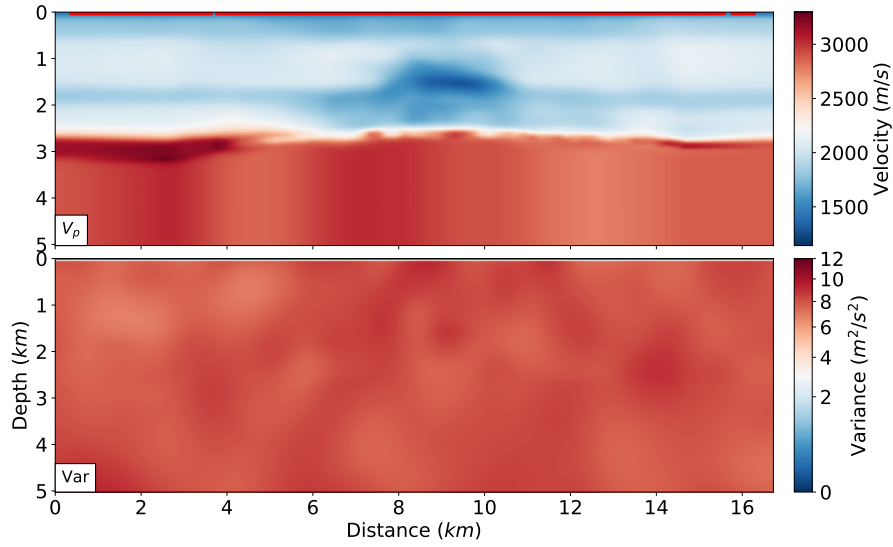


Figure 5.3: Top : Initial ensemble mean velocity model m_0 . Acquisition is denoted by a red line at the surface. Bottom : Initial variance map for $N_e = 600$.

5.1.1 P-wave velocity reconstruction

The final ensemble mean and the final variance map are displayed in Figure 5.4. As expected from the previous experiment, the final mean model provides a net increase in resolution and tends to a conventional FWI solution.

Layered structures are well defined in the top half of the domain, and from this result, we can identify what can be interpreted as hydrocarbon-charged units overlaying the anticline structure. The deep layered

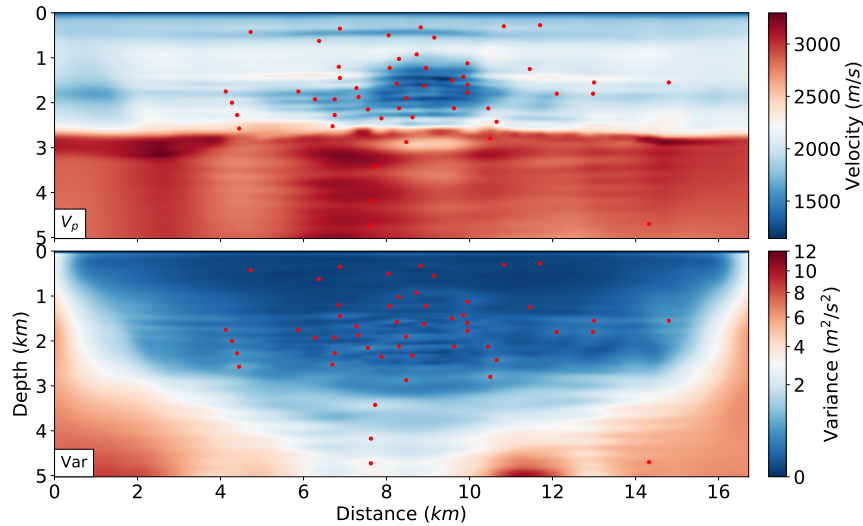


Figure 5.4: Final ensemble mean (top) and final variance map (bottom) for $N_e = 600$, after 6 ETKF-FWI cycles between 3.56 Hz to 7.01 Hz.

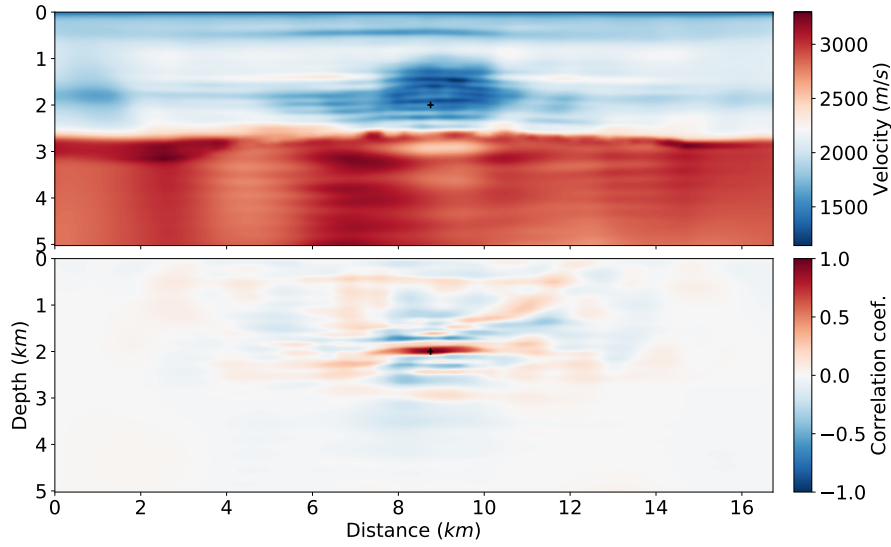


Figure 5.5: Velocity map (top) and correlation map (bottom) for the velocity parameter located at $x = 8.75 \text{ km}$ and $z = 2.0 \text{ km}$ in the final ensemble.

structures are not as sharp as the top section because of the strong impedance contrast between the upper and lower units of the medium. The strong P-wave velocity contrast between the upper and lower domain is expected to reduce the illumination power in the deeper part of the model, along with the geometrical spreading effect.

While the initial variance (Fig. 5.3) is relatively homogeneous in the entire domain (the water depth is not perturbed), the final variance displays the same two tendencies as in the synthetic case. The first order uncertainty structure is dominated by the geometrical spreading and the sharp velocity contrast between the upper and lower units at 2.5 km depth. Second to that are the variance values imposed by the reflectors estimated in the solution. Note that we use a non-linear colorscale to underline uncertainty associated with reflectors.

To repeat the procedure detailed in the synthetic application and evaluate how the variance aligns with the velocity structure, we compute maximum peak locations in the final variance map. The search radius has been reduced to 150 m because of the smoothness of the variance map. Despite the map smoothness and the thin layered structure in the final velocity model, we can confirm that local uncertainty maxima are preferentially located along structure discontinuities.

In the following, we chose four parameters to produce correlation maps computed in the final ensemble, following the same procedure detailed in Chapter 4 (Fig. 5.5, 5.6, 5.7 and 5.8).

The parameter chosen in Figure 5.5 corresponds to one of the hydrocarbon layer recovered at the center of the model. We chose this specific parameter for its key location in the medium, as the hydrocarbon layers are amongst the sharpest features we were able to recover in the final model. As in the synthetic test case, we observe positive correlations, aligned with the velocity structure. The amplitude of correlation structures decays rapidly with the distance, meaning the layer in which the parameter lies in is well constrained.

The parameter chosen in Figure 5.6 lie underneath a shallow (approximately 500m deep) velocity layer. It displays the same type of behavior as the previous parameter, as correlations are decaying with distance, and the positive correlation is aligned with the velocity structure of the reconstructed mean

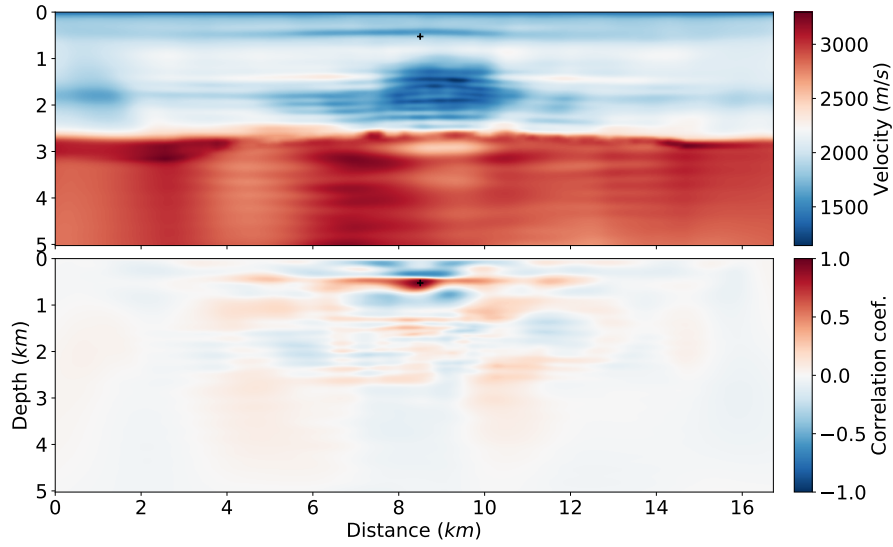


Figure 5.6: Velocity map (top) and correlation map (bottom) for the velocity parameter located at $x = 8.5 \text{ km}$ and $z = 0.525 \text{ km}$ in the final ensemble.

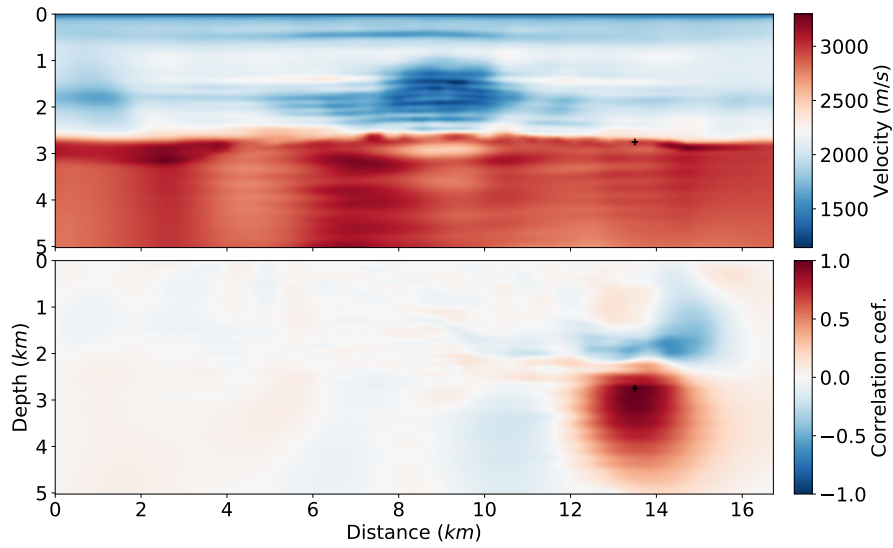


Figure 5.7: Velocity map (top) and correlation map (bottom) for the velocity parameter located at $x = 13.35 \text{ km}$ and $z = 2.75 \text{ km}$ in the final ensemble.

model.

The parameter in Figure 5.7 has been chosen to illustrate the partitioning between the upper (low-velocity) and lower (high-velocity) domains. This point lies at the limit of the illuminated domain and has thus a poor spatial resolution (large positive correlation spot around the parameter). However, despite this low resolution, we can see that there is a clear upper limit to the correlation zone, at the velocity transition zone (at approximately $z = 2.7 \text{ km}$). This is coherent to observations made in Chapter 4: the parameters sharing the same type of lithologic properties are expected to be positively correlated, hence the partitioning between the upper-lower domains of the model.

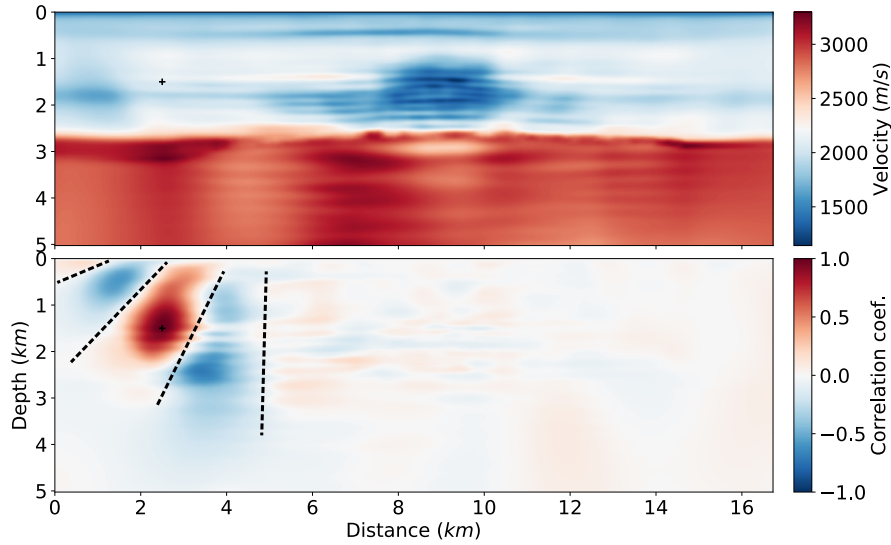


Figure 5.8: Velocity map (top) and correlation map (bottom) for the velocity parameter located at $x = 2.5 \text{ km}$ and $z = 1.5 \text{ km}$ in the final ensemble. Dashed lines denote the polarity transition in the correlation map, highlighting the finite-frequency oscillations effects.

The last parameter (Fig. 5.8) has been picked to show the occurrence of finite-frequency oscillations on this field-data application. Similarly, as in Figure 4.13, black dashed lines have been drawn to highlight this effect. The polarization(orientation) of the oscillations once again corresponds roughly to the trajectory of the body-waves that are sampling the shallow lateral bounds of the domain.

Finally, we compare the velocity model obtained from our ensemble approach to a classical FWI result in Figure 5.9.

At first glance, we can verify that the ETKF-FWI does produce a parameter estimate, close to a standard FWI solution. We can, however, notice that the resolution of the mean ETKF-FWI is slightly higher and the velocity contrasts between layers appear sharper in the ETKF-FWI result, both in the shallow and deep parts of the model. This might be due to the effect of the analysis step, which provides a correction from the estimated covariance matrix. This could have an effect similar to the one of a preconditioner which approximates the inverse Hessian operator. This is further discussed in the following multiparameter application.

We emphasize that the quality of these results is strongly linked to the initial ensemble parameterization. Modifying the initial perturbations correlation length or amplitude will result in a different outcome, or might cause instabilities if incorrectly chosen.

5.1.2 P-wave velocity and density reconstruction

In the following, we present preliminary multiparameter inversion results to show the potential of the method for uncertainty estimation in this context. Multiparameter FWI is known as a challenging problem, especially because of the presence of cross-talks between parameters (Operto et al., 2013). Recovering information about the uncertainty linked to these cross-talks is thus crucial, and might be an important benefice from strategies such as the ETKF-FWI scheme presented here. We modify the system state vector such that the columns of the ensembles contain both the velocity parameter V_p and

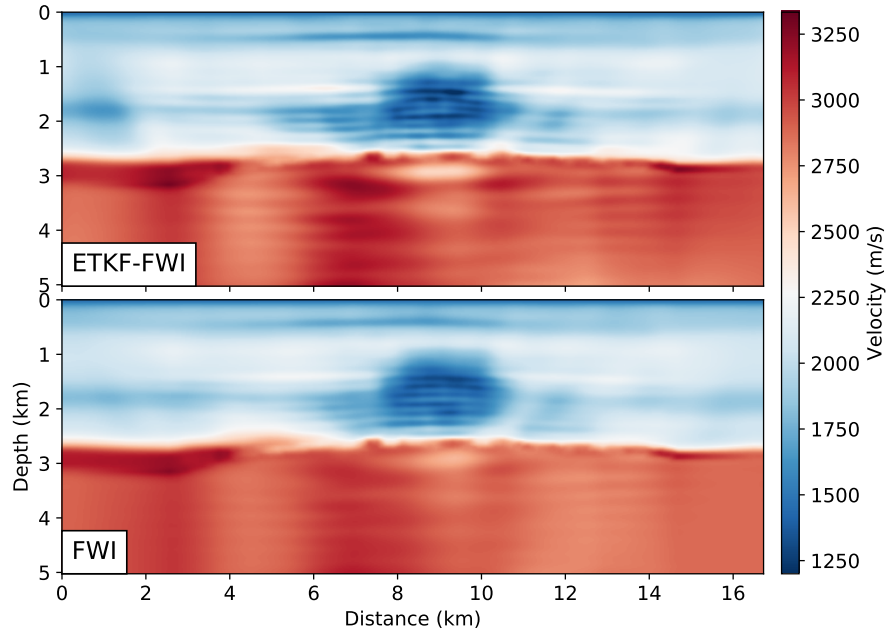


Figure 5.9: Comparison of monoparameter ETKF-FWI (top) and FWI (bottom) results with similar inversion setup (inversion parameters, regularization, acquisition geometry and data frequency groups).

the density ρ instead of the velocity alone

$$m_{V_p, \rho}^{(i)} = \begin{pmatrix} V_p \\ \rho \end{pmatrix}. \quad (5.1)$$

Considering the joint state $m_{V_p, \rho}^{(i)}$ makes it possible to take the changes of density during the forecast optimization steps into account when the analysis is performed. Note that the extension of the state vector also implies an extension of the state covariance matrix. It is expected that the cross-talk terms between V_p and ρ (off-diagonal blocks of the covariance matrix) will play a role in the Kalman Gain estimate.

The initial density perturbations are derived from the initial perturbed velocity model according to Gardner's empirical relationship (in soil only) (Gardner et al., 1974)

$$\rho = 0.31 V_p^{0.25}. \quad (5.2)$$

This way, initial ensemble members' velocity and density perturbations are physically linked.

The starting ensemble mean velocity and density models are displayed in Figure 5.10. The ETKF-FWI scheme is applied following the same setup as detailed for the monoparameter test, except for the forecast that now includes inversion of the density parameter alongside the inverted velocity. The parameter estimation after 6 ETKF-FWI cycles are shown in Figure 5.11.

The recovered velocity model is almost identical to the velocity estimate from the monoparameter case. As for the density inversion, the horizontally layered structures observed in the velocity map, are closely matching the density estimate. A lower density is seen in the central area where hydrocarbon charged layers are expected to be located.

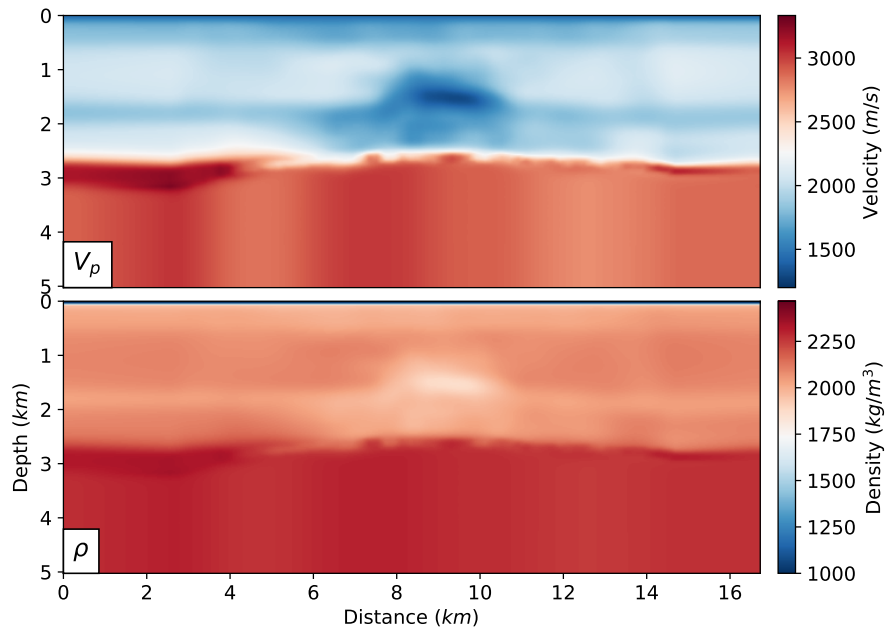


Figure 5.10: Top : Initial ensemble mean velocity model m_{0,V_p} . Bottom : Initial ensemble mean density model $m_{0,\rho}$

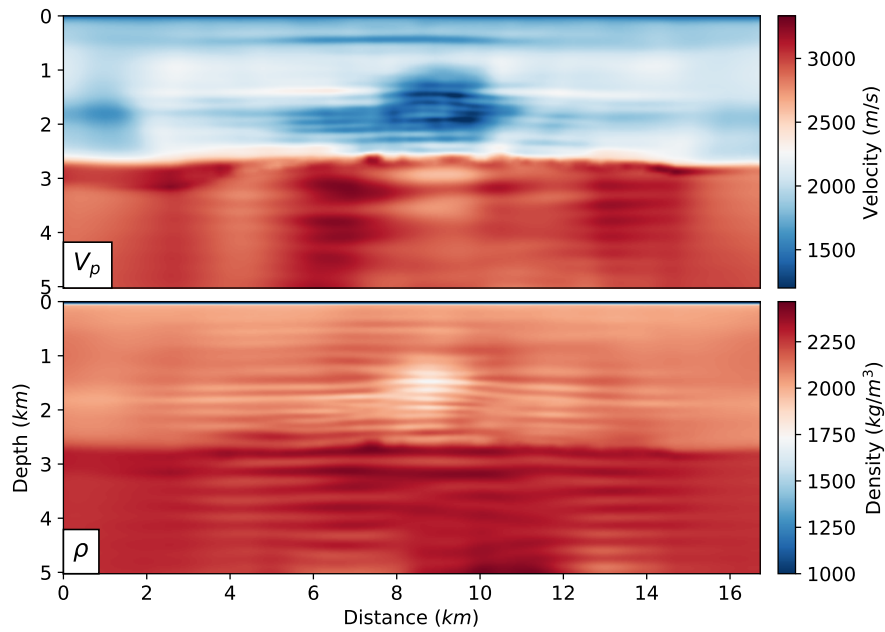


Figure 5.11: Top : Final ensemble mean velocity model m_{V_p} . Bottom : Final ensemble mean density model m_ρ

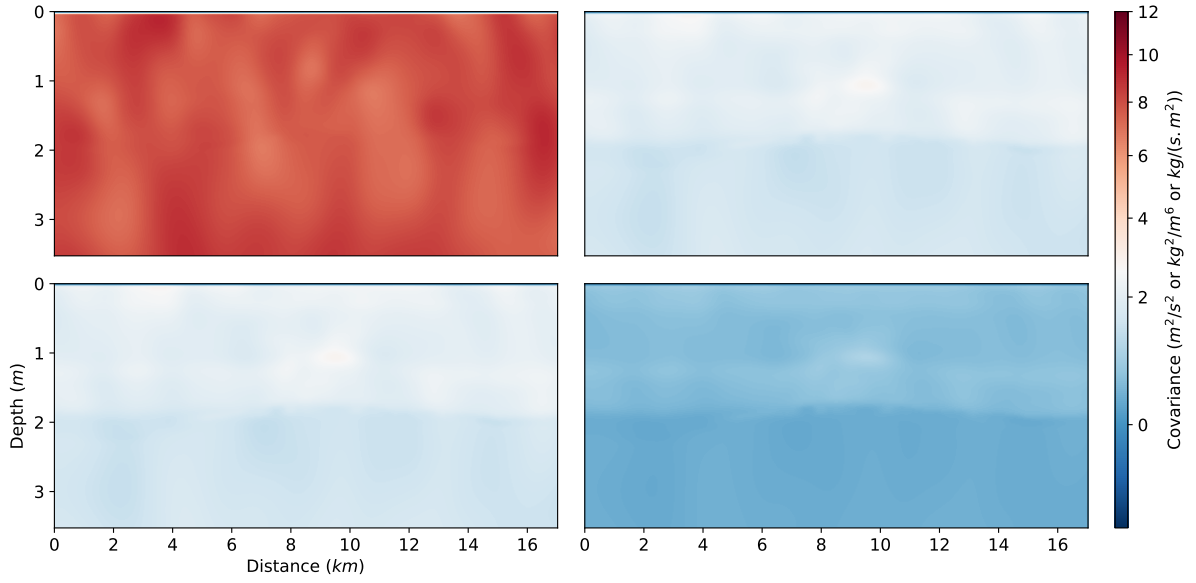


Figure 5.12: Diagonal elements of the initial joint-covariance matrix, plotted in the physical domain arranged according to their respective position in the block matrix. Top left: P-wave velocity variance in m^2/s^2 . Bottom right: density variance in kg^2/m^6 . Bottom left and top right: V_p, ρ cross-covariance maps in $kg/(s.m^2)$.

The joint covariance matrix for the multiparameter case contains four blocks. Its structure is defined by

$$\mathbf{P}_{[V_p, \rho]} = \begin{bmatrix} \mathbf{P}_{V_p V_p} & \mathbf{P}_{V_p \rho} \\ \mathbf{P}_{\rho V_p} & \mathbf{P}_{\rho \rho} \end{bmatrix}, \quad (5.3)$$

where $\mathbf{P}_{V_p V_p}$ and $\mathbf{P}_{\rho \rho}$ are the variance matrices of the marginal distribution of V_p and ρ respectively, and $\mathbf{P}_{V_p \rho}$ and $\mathbf{P}_{\rho V_p}$ are the cross-covariance blocks. Note that since $\mathbf{P}_{V_p, \rho}$ is symmetric, we have $\mathbf{P}_{\rho V_p} = \mathbf{P}_{V_p \rho}^T$ by definition. The $\mathbf{P}_{V_p V_p}$ block is expected to yield results similar to the covariance matrix in the mono-parameter case, while the $\mathbf{P}_{\rho \rho}$ block is its equivalent for the recovered density. The cross-covariance blocks are instead a measure of the link between the two parameters, and therefore makes it possible to quantify the inversion cross-talk between velocity and density. Starting with the parameter's uncertainty and cross-talk, we extract the four diagonal elements of the block joint-covariance matrix and plot them as variance and cross-covariance maps in Figures 5.12 and 5.13.

The initial variance maps are displayed in Figure 5.12. The initial velocity variance distribution tends to the monoparameter case starting distribution, while the initial density variance map is very different. This is a result of the use of Gardner's law to produce the initial density models from perturbed velocity. The cross-covariance maps are symmetric and appear to be a combination of both velocity and density variances.

The final variance maps are displayed in Figure 5.13. As in the previous results, the geometrical spreading effect is the prevalent source of uncertainty in the velocity reconstruction, while structural uncertainty is the dominant effect in the density variance map. Although the geometrical spreading is not directly visible in the density variance map, the higher variance values are located in the deeper region of the model nonetheless. The cross-covariance maps seem to indicate that the cross-talk between parameters is strongly linked to their respective uncertainties. The differences between the velocity

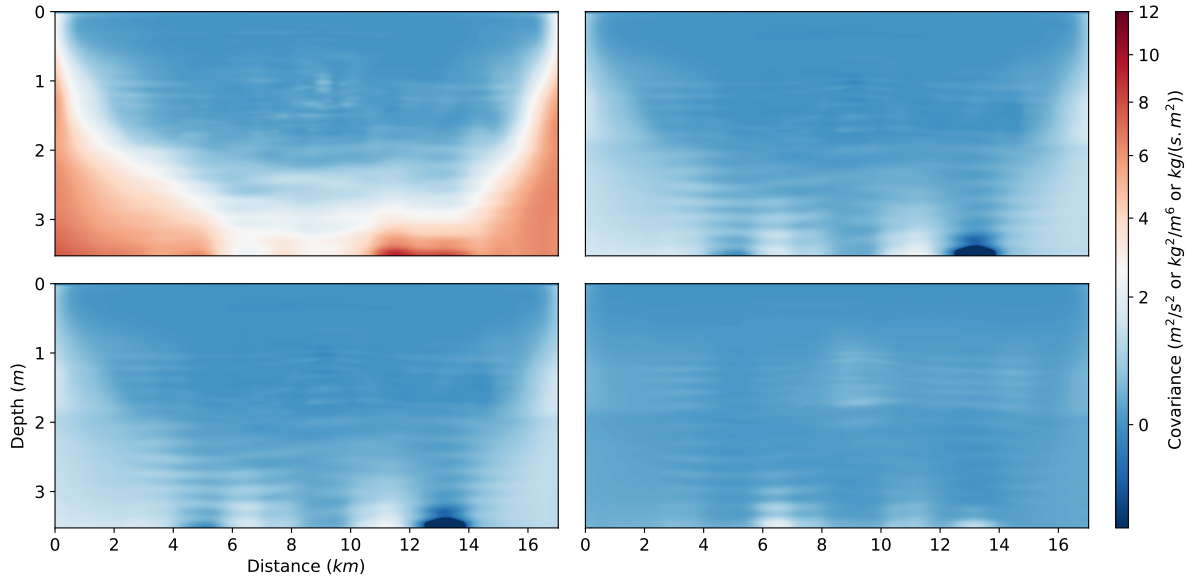


Figure 5.13: Diagonal elements of the posterior joint-covariance matrix, plotted in the physical domain arranged according to their respective position in the block matrix.

variance map and the density variance map can be linked to wave propagation theory. The prevalence of the geometrical spreading effect can be associated to the higher sensitivity of the body-waves to velocity perturbations, while the structural uncertainty in the density map could be explained by the higher sensitivity of short-offset reflected-arrivals toward density changes.

Added to the diagonal elements of the block-covariance matrix, individual parameters resolution and cross-talk terms of the block-correlation matrix are evaluated. This is achieved by extracting four corresponding lines out of the different blocks and mapping the correlation coefficients into the physical domain. This procedure is the extension of the correlation maps computation of the previous applications, to the block-diagonal structure. We choose arbitrarily a parameter located at $z = 2.0$ km; $x = 9.6$ km. Its initial correlation maps are plotted in Figure 5.14 followed by its final correlation maps in Figure 5.15.

Although the initial correlations are identical in all blocks due to the models' generation, the final correlation patterns are entirely different in the final maps. There is a sharp difference of resolution in velocity and density: velocity correlations are laterally oriented along the structure, while density correlations are oriented along a vertical axis across the domain. The resolution information is coherent with theoretical expectations as stated previously; velocity reconstruction is mostly constrained by diving waves that can explain lateral ambiguity, while density is constrained by short offset reflections arrivals, which can explain the higher vertical uncertainty.

Besides, correlation cross-talk maps allow evaluating the coupling effect between velocity and density across the whole domain. In this case, the $\{V_p, \rho\}$ and $\{\rho, V_p\}$ correlations terms are weaker than the $\{V_p, V_p\}$ and $\{\rho, \rho\}$ correlation terms. It means that for this parameter, velocity and density are well decoupled (density reconstruction does not seem to be contaminated by velocity leakage during the inversion).

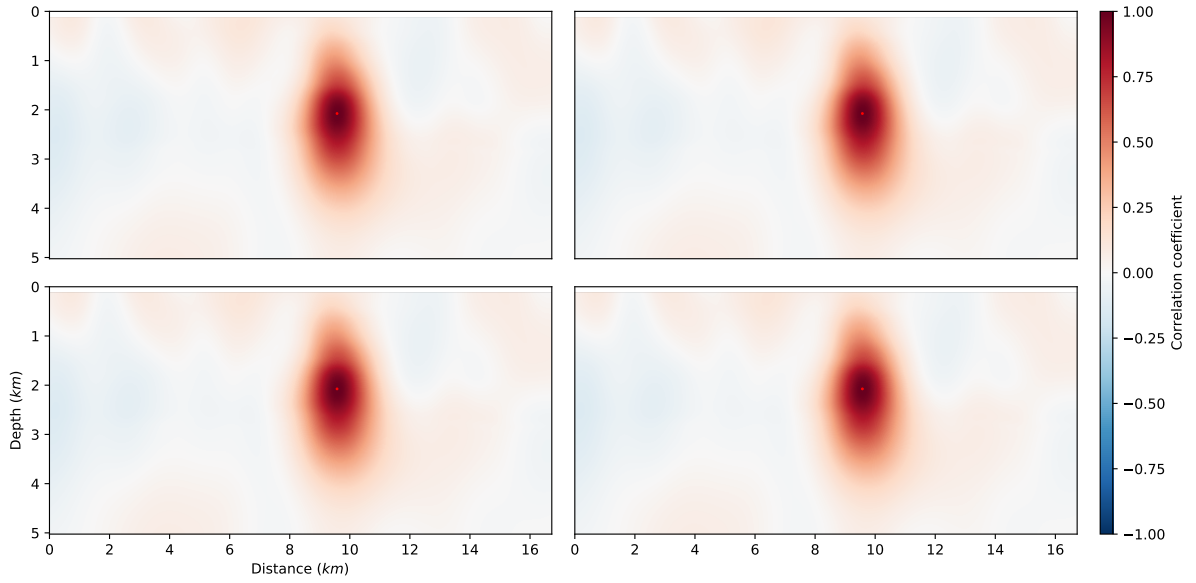


Figure 5.14: Off-diagonal elements of the initial joint-covariance matrix, plotted in the physical domain. The covariance matrix lines considered correspond to the parameter located at $z = 2.1$ km; $x = 9.6$ km. Top left: P-wave velocity correlation coefficient. Bottom right: density correlation coefficient. Bottom left and top right: correlation cross-talk terms.

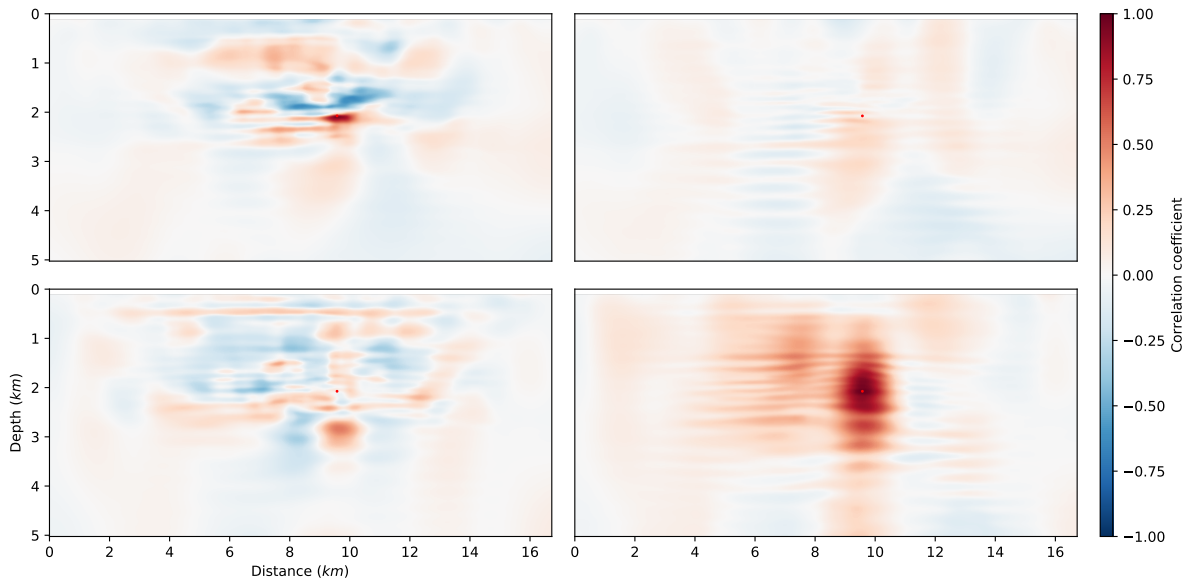


Figure 5.15: Off-diagonal elements of the posterior joint-covariance matrix, plotted in the physical domain. The covariance matrix lines considered correspond to the parameter located at $z = 2.1$ km; $x = 9.6$ km. Top left: P-wave velocity correlation coefficient. Bottom right: density correlation coefficient. Bottom left and top right: correlation cross-talk terms.

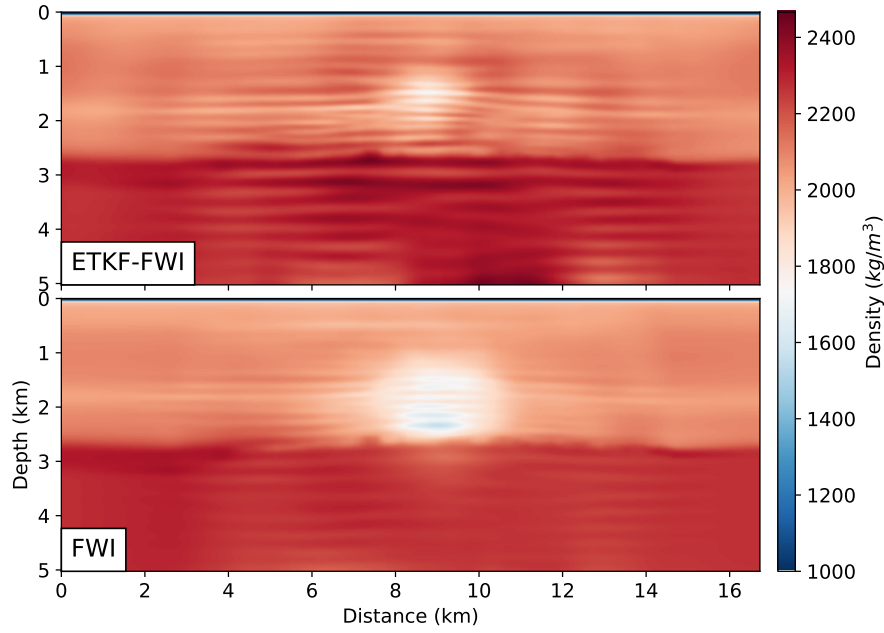


Figure 5.16: Comparison of ETKF-FWI (top) and FWI (bottom) density estimate with similar inversion setup (inversion parameters, regularization, acquisition geometry and data frequency groups).

Finally, we compare the estimated density, with an equivalent multiparameter FWI result, obtained with a similar inversion setup (data selection and processing, number of minimization steps, initial model) in Figure 5.16.

Contrarily to the velocity estimation, there are significant discrepancies between the density model recovered by the ETKF-FWI and its FWI equivalent. The density in the hydrocarbon layers is lower in the FWI estimate, while the ETKF-FWI result is characterized by a high wavenumber content and sharper density contrasts.

Motivated by these significant differences, we evaluate the quality of both the ETKF-FWI and the FWI solutions in the time-domain. To do so, we computed synthetic common-receiver gathers in both the ETKF-FWI and the FWI solutions, along with the initial models for reference. We then compare synthetics by evaluating their data-fit with the observed common receiver gather data in Figure 5.17. Time-domain synthetic common-receiver gathers are plotted in color over the black-and-white observed data after filtering with a 6 to 8 Hz band-pass filter.

On this visualization, synthetic blue arrivals should overlap white, observed arrivals, while red should be overlapped by black arrivals (and therefore not be visible). The blue color is hence indicative of good fit, while visible red is indicative of phases misalignment.

The FWI result (center) is significantly improving on the initial models (left), but the ETKF-FWI result (right) is exhibiting an overall better data fit. Late arrival diving waves, as well as near offset reflections, are improved (see red ellipses). It seems that the analysis step of the ETKF-FWI acts as a Hessian-like preconditioning term, allowing a better convergence, which might enhance parameter disambiguation.

While these preliminary results are a call for careful investigations, it seems that the analysis of the joint-space allows for better convergence of the ETKF-FWI scheme, compared to the classical

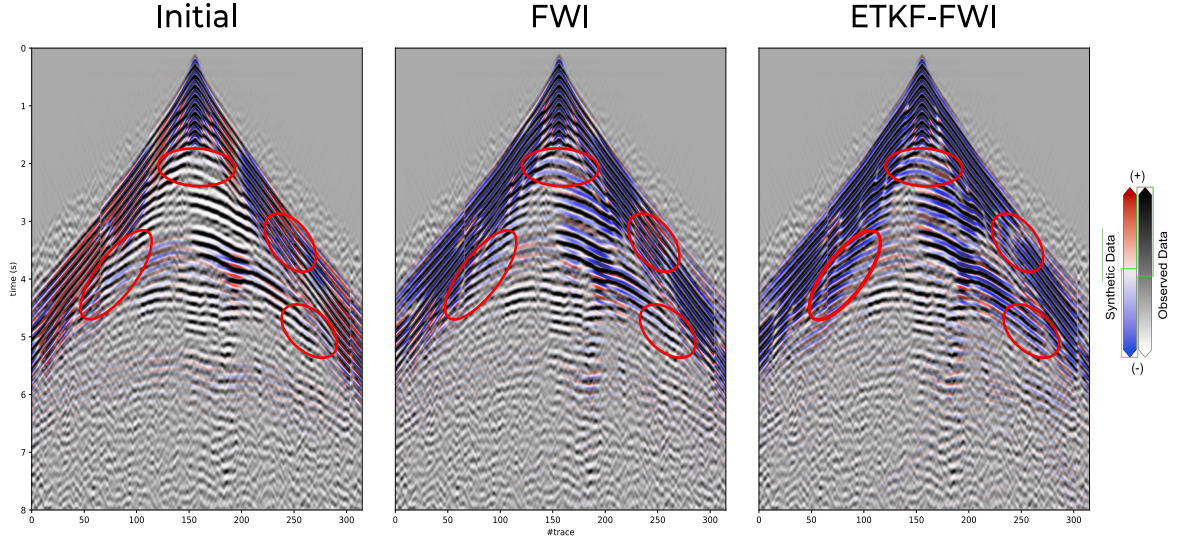


Figure 5.17: Data fit evaluated on a common receiver gather between (from left to right), the initial models, the FWI outcome, and the ETKF-FWI outcome. Blue arrivals denote a good data fit over corresponding white arrivals. Red arrivals overlapping white arrivals are indicative of misaligned phases. Major improvement areas granted by the ETKF-FWI results have been marked with red ellipses in all three common-receiver gathers.

FWI. These results prompt us to investigate the possibilities of extension of the methodology beyond mono-parameter inversion in future studies.

Conclusion

In this final chapter, we have presented a field-data application of the ETKF-FWI scheme. We were able to reproduce the observations made on the synthetic benchmark while evaluating the well-posedness of the filter when dealing with a more complex test-case. A multi-parameter application has also been proposed to complement the monoparameter inversion case.

We were able to evaluate the multi-parameter cross-talks, in both the variance maps and in the cross-correlation maps. We also have been able to underline the difference of uncertainty structure in the velocity and density reconstruction and proposed an explanation based on their respective sensitivity to different types of arrivals. Hence, body-waves, which are mostly responsible for the velocity reconstruction, might explain the geometrical spreading observed in the velocity variance map. Density reconstruction, however, is mostly governed by short-offset reflection arrivals, and hence, do not exhibit geometrical spreading, nor finite-frequency oscillations in the correlation maps.

Chapter 6

Conclusion and perspectives

Conclusion

The objective of this thesis has been to introduce a novel uncertainty estimation solution for FWI by combining the FWI and ensemble DA theories under a unified framework. A literature review on uncertainty estimation in FWI has been given in Chapter 1, along with a general theoretical introduction of the FWI concept. DA theory and its pairing with FWI have been introduced in Chapters 2 and 3, respectively, which allowed us to design our uncertainty estimation tool. In Chapters 4 and 5, we have demonstrated that the ETKF can be paired with a frequency-domain FWI quasi-Newton solver, and allows for uncertainty estimation of the tomographic solution. The resulting ETKF-FWI scheme can produce a robust state estimation while allowing us to recast our inversion problem in a local Bayesian framework. The results we have obtained so far are encouraging in several regards.

- Variance and correlation maps only require to store the ensemble to be computed, and they provide a straightforward way of evaluating the quality of convergence, the correlation links, and tradeoffs between parameters.
- The ETKF-FWI also allows integrating some form of data weighting terms in the whole tomographic process via the measurement noise matrix \mathbf{R} . If \mathbf{R} is appropriately set, the resulting uncertainty takes into account the physical properties of assimilated data.
- Finally, the extension perspectives offered by the DA framework and the full scalability of the method makes it a great candidate for uncertainty estimations.

This work also raised several questions regarding the application of the ETKF-FWI, that I want to re-emphasize as concluding matters:

How much of a problem is undersampling ? Regarding undersampling, its effects seem not too dramatic, as they do not affect the state estimate capabilities of the ETKF-FWI. We attribute this robustness to the inversion scheme that acts as our forecast, which is not expected to spread-out the ensemble members. The underestimation of variance and spurious correlations might be more of an issue, as they have a direct impact on our ability to use and interpret the quantitative covariance data. We have seen that variance underestimation could be solved with *covariance matrix inflation*, by artificially increasing the forecast covariance by a factor r to mitigate overconfidence in the forecast.

However, due to the necessity of evaluating an appropriate inflation parameter through trials and errors, its implementation in our case was limited.

We have also seen that our observation operator (wave equation modeling) is strongly non-local, which prevents us from applying covariance localization or local analysis to mitigate spurious correlation terms in the covariance matrix. Thus we propose to rely mostly on local covariance information, which seems to be preserved most of the time (as seen in correlation maps) and appears to be a reliable resolution proxy. Ultimately, the undersampling issues allowed us to address the validity of our low-rank approximation and evaluate its associated biases. We think this specific point should be investigated in any methodology proposal based on rank reduction or Hessian approximation, which is unfortunately not always discussed in current propositions among the uncertainty estimation literature.

How to characterize prior uncertainty, and define the initial ensemble ? Good practices, when it comes to initial ensemble building, may deserve entire research focus on its own. As it stands, we have adopted a pragmatical approach to generate initial perturbations, but defining "optimal" and how an optimal initial ensemble should be built, is an open question. One might advocate for producing higher variance initial ensembles, to allow further parameter exploration at the cost of stability and convergence. Another option would be to align with the tests we have set up by limiting the spread of the initial ensemble to ensure an optimal parameter estimation. To constrain a strict convergence, one might even choose to add perturbation in limited portions of the model only, to limit the chances of unphysical updates during the analysis. For instance, in our field-data test case, we could remove perturbations in the lower half of the domain, constrained by a small portion of data. This would prevent any unphysical updates driven by the data term during the analysis. With such questions, we think the initial model building deserves a careful investigation, as the options mentioned above might be logical choices depending on one's goals.

How is the quantitative uncertainty estimate reliable? It has to be reminded that uncertainty estimates are, at best, expressed both in terms of "local optimization" uncertainty and in the frame of finite-frequency wave propagation:

- Our method allows us to retrieve a local Bayesian solution to the tomographic problem. Therefore we do not explore the full extent of the possible solution, but rather the solutions fitting in a least-squares formalism, given the defined prior and the quality of observations.
- As the wave-propagation and the limited coverage act as a filter over the physical domain, it is not possible to estimate the absolute uncertainty values of the physical parameters. We instead think uncertainty should be expressed in terms of the optimal apparent macro-model as "seen" by the waves, in similar ways as Capdeville and Métivier (2018) suggestion for down-scaling and homogenization problems.

To get closer to absolute uncertainty estimation of the physical parameters, we think that calibrating our local uncertainty estimate over well-log information might be a solution. Recalling that we only obtain a filtered version of the model and its uncertainty estimation, it would require comparing our posterior variance, with filtered log data, at the adequate frequencies (as seen by the waves) and scale the covariance accordingly. This scaling solution could potentially mitigate the underestimation of variance values caused by undersampling, as well as mitigating the effects of our choice of prior ensemble. This solution would, however, be limited to crustal-scale exploration data, provided well-logs are available or not too distant from the target area.

Perspectives

While we have presented a comprehensive review of the ETKF-FWI capabilities, it is worth noting that there are several avenues for extending this method.

Better undersampling mitigation: In this thesis, we have opted for a pragmatical approach to mitigate variance underestimation, relying on multiplicative covariance inflation. This solution was, unfortunately, impractical due to the difficulty of tuning the inflation parameter. Assessing how the “ideal” r value changes depending on the case (subsurface complexity, domain size, model discretization, acquisition design, for instance) might be a first step into resolving variance underestimation in the ETKF-FWI. Hybrid inflation methods such as the one implemented by Whitaker and Hamill (2012) could potentially be an interesting take on inflation as they are based on re-evaluating the analysis ensemble rather than the forecast ensemble. In that case, one might be able to use well-log data to calibrate the analysis ensemble variance, to ensure that *in-situ* measurements fit in the ensemble spread.

Another solution to overcome inbreeding issues might come from methods such as the finite-size ensemble Kalman filter (Myrseth and Omre, 2010; Bocquet, 2011; Myrseth et al., 2013; Bocquet et al., 2015), as it has been designed to eliminate the need for inflation in ensemble DA. Alternatively, schemes relying on automatic/adaptive inflation such as Miyoshi (2011b) or Raanes et al. (2019) (the latter being a reformulation of the finite-size EnKF) should also be interesting takes on the problem, allowing on-line estimation of r . This would have the benefit of letting r evolve along with the dynamic axis, which might be beneficial for our frequency-dependent dynamics.

Beyond frequency domain FWI: The methodology would also benefit from moving beyond the frequency domain formalism into the time-domain FWI formalism, which is the current industry standard. While a time-domain extension would require to tackle a sizeable computational cost (mainly related to the cost of time-domain FWI), it could be the solution toward 3-D application, as FDFD modeling in 3-D remains expensive compared to FDTD approaches. The dynamic axis could be redesigned to mitigate the cost associated with time-domain FWI, based on data-managements strategies. Thus, instead of solving the FWI problem with the complete dataset for each of the ETKF-FWI forecasts, we could iterate through small subsets of data, according to a data-decimation strategy.

Apart from meeting industry standards, time-domain FWI would also allow us to consider time-based localization in the medium. As it is possible to predict the evolution of the wavefront at any point in time, we could exclude some parameters from the analysis based on the wavefront position. This should prevent distant unphysical updates in parts of the medium that have not been sampled by the propagating waves.

Multi-parameter inversions: Experimenting with other DA parameterization methods (such as the alternative propositions in Chapter 3) might also give new insights on FWI uncertainty estimation. We have seen, for instance, that propositions 2 and 3 could allow establishing links toward wavefield reconstruction inversion (van Leeuwen and Herrmann, 2013) or the ensemble Kalman inverse (Iglesias et al., 2013a). An analogous to WRI should result in a more convex misfit-function and reduce the sensitivity to the initial model, which in turn should allow relaxing the constraints for the initial ensemble generation. On the other hand, the adjoint-free EKI could potentially reduce the computational burden of the ensemble-based uncertainty estimation method, provided it can converge sufficiently fast.

With alternative parameterization also comes the possibility to investigate multi-parameter FWI further. We have seen in Chapter 5 that multi-parameter ETKF-FWI granted easy access to the cross-talk terms between inversion parameters, which is currently a challenging issue in multi-parameter FWI. It also seemed that extending the state space can lead to improvements in model recovery, as the forecast

covariance might play a preconditioning role during the analysis. While the multi-parameter field-data test leads us to this unexpected observation, investigations on this potential preconditioning effect should be taken from the ground-up, with simpler synthetic models.

Joint inversion, time-lapse: Our joint DA-FWI scheme also offers the possibility to combine geophysical exploration methodologies into a single inversion framework. By taking advantage of the natural sensor fusion capabilities of DA tools, one might design a filter that combines different types of geophysical data. In the context of georesources explorations, for example, one could integrate well-log observations along with seismic data in the analysis to potentially improve inversion results. Adding in-situ measurements in the analysis might also be the way to approach absolute quantitative uncertainty assessments, which has yet to be achieved in the field of uncertainty estimation in FWI.

Data assimilation also offers a natural way of tracking the evolution of systems with time. We can thus envision time-lapse or monitoring tomographic application, as it has already been proposed by Eikrem et al. (2019). Applied to reservoir monitoring, it might allow combining 4-D FWI, well-log time-series, and ground deformation into a complete monitoring solution.

Finally, note that the ETKF should be readily applicable to other tomographic problems, provided an adequate dynamic axis can be defined.

Cost and Applicability

Our literature review showed that there is a need for high-dimension uncertainty estimation in FWI. However, computational cost regarding uncertainty estimation tools is a real concern, as it adds-up with the already expensive computational cost of FWI. We have shown in this study that the ETKF-FWI fitted in this high-dimensionality paradigm, but it also takes advantage of an embarrassingly parallel problem to achieve full-scalability. As a closure, we compare the computational cost of our method with other uncertainty estimation methodologies in the literature.

First, we have not discussed how this methodology compares with global optimization approaches. We have shown in 1.2.1 that global optimization methods can accommodate non-convexities of the misfit function by sampling the entirety of the solution space, rather than sampling the cost function around the solution as we performed in the ETKF-FWI. While these methods seem very appealing, they have to rely on tricks to make this sampling possible and alleviate the curse of dimensionality problem they would face otherwise. These approaches are thus either limited to small problems (with a low number of unknown to sample) or rely on clever parameterizations (such as B-spline functions or Voronoi tessellation) to reduce the size of the search space. Nonetheless, most of these methodologies will require several thousands of samples (and thus as many partial-differential-equation (PDE) to solve), which makes them challenging to use as up to now. They also tend to produce very coarse solutions to the inverse problem (which nonetheless makes for great potential starting models for local uncertainty estimation, as shown in Sajeve et al. (2017b)). The philosophy of local and global approaches differs, as they propose to deal with very different but complementary aspects of uncertainty estimation.

We have seen in 1.2.2 that local approaches are based on rank-reduction methods. These approximations of the inverse Hessian operator in the vicinity of the solution, make sampling from the posterior covariance matrix affordable. Their low-rank approximation of the inverse Hessian operator, require to solve several forward and adjoint PDEs, typically several hundred to several thousand per frequencies (for example, Bui-Thanh et al. (2013) is evaluating 1400 PDE to estimate the first 700 eigenvalues of their global FWI application with hundreds of thousands of parameters). Fang et al. (2014) requires

to solve approximately 6000 forward modeling problems, with their MCMC sampling to produce an uncertainty estimate (with most of the cost coming from the sampling strategy). Zhu et al. (2016) is able to produce an uncertainty estimation along with the solution of the inverse problem at the minimal cost of 144 PDE resolution thanks to the assumption made on the structure of the Hessian operator. Though this cost is indeed reasonably low, it does not include the computational cost of the reverse time migration they are using to precondition their sampling. Finally, the number of PDE solved to sample the posterior covariance in Fang et al. (2018) proposition is the number of sources plus the number of receivers per frequencies (not including the number of PDE to solve the inverse problem). Besides, this method does seem to display challenging memories limitation as it requires to store the optimal wavefields in memory for each frequency bands, which may become challenging for large scale 3D application. The extension to uncertainty estimation of multi-parameter inversion also seems to be non-trivial in this extended domain FWI application, as only recent publications are addressing the multi-parameter aspect of wavefield reconstruction inversion (Aghamiry et al., 2019). Note also that the low-rank approximation methods of the propositions mentioned above (such as randomized Singular-Value-Decomposition, or Lanczos methods) are sequential by nature, which makes these uncertainty methods only as scalable as their PDE solver can get.

In comparison, the cost of ETKF-FWI in our applications ranges from 5000 to 18000 PDE solve (for the synthetic and field data cases, respectively), which might appear to be a daunting number (although convergence tests have shown we could potentially consider smaller ensemble size). However, unlike other methods, we are set to solve an embarrassingly parallel problem as all of our ensemble members are evolving independently during the bulk of the computational time (forecast step), which makes our problem not only scalable on the PDE solver but fully scalable on the ensemble size. Thanks to this advantage, and because of the development of hardware capacities towards the exascale and the current trend toward grid computation, we believe that the ETKF-FWI for uncertainty estimation can be a valuable approach even for large-scale FWI problems, as it is currently the case for DA applications.

Bibliography

- Aghamiry, H., Gholami, A., and Operto, S. (2019). Admm-based multi-parameter wavefield reconstruction inversion in VTI acoustic media with TV regularization. *Geophysical Journal International*, 219(2):1316–1333.
- Aki, K., Christoffersson, A., and Husebye, E. S. (1977). Determination of the three-dimensional seismic structure of the lithosphere. *Journal of Geophysical Research*, 82(2):277–296.
- Aleardi, M. and Mazzotti, A. (2016). 1d elastic full-waveform inversion and uncertainty estimation by means of a hybrid genetic algorithm–gibbs sampler approach. *Geophysical Prospecting*.
- Amestoy, P. R., Duff, I. S., and L’Excellent, J. Y. (2000). Multifrontal parallel distributed symmetric and unsymmetric solvers. *Computer Methods in Applied Mechanics and Engineering*, 184:501–520.
- Anderson, J. L. (2001). An ensemble adjustment kalman filter for data assimilation. *Monthly weather review*, 129(12):2884–2903.
- Anderson, J. L. (2003). A local least squares framework for ensemble filtering. *Monthly Weather Review*, 131(4):634–642.
- Anderson, J. L. (2007). An adaptive covariance inflation error correction algorithm for ensemble filters. *Tellus A: Dynamic Meteorology and Oceanography*, 59(2):210–224.
- Anderson, J. L. (2009). Spatially and temporally varying adaptive covariance inflation for ensemble filters. *Tellus A: Dynamic Meteorology and Oceanography*, 61(1):72–83.
- Anderson, J. L. and Anderson, S. L. (1999). A monte carlo implementation of the nonlinear filtering problem to produce ensemble assimilations and forecasts. *Monthly Weather Review*, 127(12):2741–2758.
- Asch, M., Bocquet, M., and Nodet, M. (2016). *Data assimilation: methods, algorithms, and applications*, volume 11. SIAM.
- Bardsley, J. M., Solonen, A., Haario, H., and Laine, M. (2014). Randomize-then-optimize: A method for sampling from posterior distributions in nonlinear inverse problems. *SIAM Journal on Numerical Analysis*, 36(4):A1895–A1910.
- Barkved, O., Buer, K., Halleland, K., Kjelstadli, R., Kleppan, T., and Kristiansen, T. (2003). 4d seismic response or primary production and waste injection at the valhall field. In *Extended Abstracts, 65th Annual EAGE Conference & Exhibition*.

- Barkved, O., Heavey, P., Kommedal, J. H., van Gestel, J.-P., ve, R. S., Pettersen, H., Kent, C., and Albertin, U. (2010). Business impact of full waveform inversion at valhall. *SEG Technical Program Expanded Abstracts*, 29(1):925–929.
- Bedle, H. and Lee, S. V. D. (2009). S velocity variations beneath north america. *Journal of Geophysical Research: Solid Earth*, 114(B7).
- Bellman, R. E. (2015). *Adaptive control processes: a guided tour*, volume 2045. Princeton university press.
- Bishop, C. H., Etherton, B. J., and Majumdar, S. J. (2001). Adaptive sampling with the ensemble transform kalman filter. part i: Theoretical aspects. *Monthly weather review*, 129(3):420–436.
- Bishop, C. H., Whitaker, J. S., and Lei, L. (2017). Gain form of the ensemble transform kalman filter and its relevance to satellite data assimilation with model space ensemble covariance localization. *Monthly Weather Review*, 145(11):4575–4592.
- Biswas, R. and Sen, M. (2017). 2d full-waveform inversion and uncertainty estimation using the reversible jump hamiltonian monte carlo. In *SEG Technical Program Expanded Abstracts 2017*, pages 1280–1285.
- Bocquet, M. (2011). Ensemble kalman filtering without the intrinsic need for inflation. *Nonlinear Processes in Geophysics*, 18(5):735–750.
- Bocquet, M. (2016). Localization and the iterative ensemble kalman smoother. *Quarterly Journal of the Royal Meteorological Society*, 142(695):1075–1089.
- Bocquet, M. and Farchi, A. (2019). On the consistency of the local ensemble square root kalman filter perturbation update. *Tellus A: Dynamic Meteorology and Oceanography*, 71(1):1–21.
- Bocquet, M., Raanes, P. N., and Hannart, A. (2015). Expanding the validity of the ensemble kalman filter without the intrinsic need for inflation. *Nonlinear Processes in Geophysics*, 22(6):645–662.
- Bocquet, M. and Sakov, P. (2012). Combining inflation-free and iterative ensemble kalman filters for strongly nonlinear systems. *Nonlinear Processes in Geophysics*, 19(3):383–399.
- Bocquet, M. and Sakov, P. (2014). An iterative ensemble kalman smoother. *Quarterly Journal of the Royal Meteorological Society*, 140(682):1521–1535.
- Bodin, T. and Sambridge, M. (2009). Seismic tomography with the reversible jump algorithm. *Geophysical Journal International*, 178(3):1411–1436.
- Bozdağ, E., Peter, D., Lefebvre, M., Komatitsch, D., Tromp, J., Hill, J., Podhorszki, N., and Pugmire, D. (2016). Global adjoint tomography: first-generation model. *Geophysical Journal International*, 207(3):1739–1766.
- Bozdağ, E., Trampert, J., and Tromp, J. (2011). Misfit functions for full waveform inversion based on instantaneous phase and envelope measurements. *Geophysical Journal International*, 185(2):845–870.
- Bui-Thanh, T., Burstedde, C., Ghattas, O., Martin, J., Stadler, G., and Wilcox, L. C. (2012). Extreme-scale uq for bayesian inverse problems governed by pdes. In *SC’12: Proceedings of the International Conference on High Performance Computing, Networking, Storage and Analysis*, pages 1–11. IEEE.

- Bui-Thanh, T. and Ghattas, O. (2012a). Analysis of the hessian for inverse scattering problems: I. inverse shape scattering of acoustic waves. *Inverse Problems*, 28(5):055001.
- Bui-Thanh, T. and Ghattas, O. (2012b). Analysis of the hessian for inverse scattering problems: II. inverse medium scattering of acoustic waves. *Inverse Problems*, 28(5):055002.
- Bui-Thanh, T., Ghattas, O., Martin, J., and Stadler, G. (2013). A computational framework for infinite-dimensional Bayesian inverse problems Part I: the linearized case with application to global seismic inversion. *SIAM Journal of Scientific Computing*, 35(6):A2494–A2523.
- Bunks, C., Salek, F. M., Zaleski, S., and Chavent, G. (1995). Multiscale seismic waveform inversion. *Geophysics*, 60(5):1457–1473.
- Burgers, G., van Leeuwen, P. J., and Evensen, G. (1998). Analysis scheme in the ensemble kalman filter. *Monthly weather review*, 126(6):1719–1724.
- Capdeville, Y. and Métivier, L. (2018). Elastic full waveform inversion based on the homogenization method: theoretical framework and 2-d numerical illustrations. *Geophysical Journal International*, 213(2):1093–1112.
- Carrassi, A., Vannitsem, S., Zupanski, D., and Zupanski, M. (2008). The maximum likelihood ensemble filter performances in chaotic systems. *Tellus A: Dynamic Meteorology and Oceanography*, 61(5):587–600.
- Cary, P. and Chapman, C. (1988). Automatic 1-D waveform inversion of marine seismic refraction data. *Geophysical Journal of the Royal Astronomical Society*, 93:527–546.
- Chada, N. K., Iglesias, M. A., Roininen, L., and Stuart, A. M. (2018). Parameterizations for ensemble kalman inversion. *Inverse Problems*, 34(5):055009.
- Chen, Y. and Oliver, D. S. (2012). Ensemble randomized maximum likelihood method as an iterative ensemble smoother. *Mathematical Geosciences*, 44(1):1–26.
- Chen, Z. (2003). Bayesian filtering: From kalman filters to particle filters, and beyond. *Statistics*, 182(1):1–69.
- Claerbout, J. (1971). Towards a unified theory of reflector mapping. *Geophysics*, 36:467–481.
- Cordua, K. S., Hansen, T. M., and Mosegaard, K. (2012). Monte Carlo full-waveform inversion of crosshole GPR data using multiple-point geostatistical a priori information. *Geophysics*, 77(2):H19–H31.
- Cosme, E., Brankart, J.-M., Verron, J., Brasseur, P., and Krysta, M. (2010). Implementation of a reduced rank square-root smoother for high resolution ocean data assimilation. *Ocean Modelling*, 33(1-2):87–100.
- Cox, H. (1964). On the estimation of state variables and parameters for noisy dynamic systems. *IEEE Transactions on automatic control*, 9(1):5–12.
- Cruse, E., Pica, A., Noble, M., McDonald, J., and Tarantola, A. (1990). Robust elastic non-linear waveform inversion: application to real data. *Geophysics*, 55:527–538.

- Crystalng, G.-H., Mclaughlin, D., Entekhabi, D., and Ahanin, A. (2011). The role of model dynamics in ensemble kalman filter performance for chaotic systems. *Tellus A: Dynamic Meteorology and Oceanography*, 63(5):958–977.
- Datta, D. and Sen, M. K. (2016). Estimating starting models for full waveform inversion using a global optimization method. *Geophysics*, 81(4):R211–R223.
- Desroziers, G., Berre, L., Chapnik, B., and Poli, P. (2005). Diagnosis of observation, background and analysis-error statistics in observation space. *Quarterly Journal of the Royal Meteorological Society: A journal of the atmospheric sciences, applied meteorology and physical oceanography*, 131(613):3385–3396.
- Devaney, A. (1984). Geophysical diffraction tomography. *Geoscience and Remote Sensing, IEEE Transactions on*, GE-22(1):3–13.
- Draper, N. R. and Smith, H. (1998). *Applied regression analysis*, volume 326. John Wiley & Sons.
- Du, Z., Querendez, E., and Jordan, M. (2012). Resolution and uncertainty in 3d stereotomographic inversion. In *Expanded Abstracts, 74th Annual EAGE Meeting (Copenhagen)*.
- Duane, S., Kennedy, A. D., Pendleton, B. J., and Roweth, D. (1987). Hybrid monte carlo. *Physics letters B*, 195(2):216–222.
- Duff, I. S. and Reid, J. K. (1983). The multifrontal solution of indefinite sparse symmetric linear systems. *ACM Transactions on Mathematical Software*, 9:302–325.
- Eikrem, K. S., Nævdal, G., and Jakobsen, M. (2019). Iterated extended kalman filter method for time-lapse seismic full-waveform inversion. *Geophysical Prospecting*, 67(2):379–394.
- Eliasson, P. and Romdhane, A. (2017). Uncertainty quantification in waveform-based imaging methods-a sleipner co2 monitoring study. *Energy procedia*, 114:3905–3915.
- Engquist, B. and Froese, B. D. (2014). Application of the Wasserstein metric to seismic signals. *Communications in Mathematical Science*, 12(5):979–988.
- Evensen, G. (1994). Sequential data assimilation with nonlinear quasi-geostrophic model using Monte Carlo methods to forecast error statistics. *Journal of Geophysical Research*, C5(99):143–162.
- Evensen, G. (2003). The ensemble kalman filter: theoretical formulation and practical implementation. *Ocean Dynamics*, 53:343–367.
- Evensen, G. (2009). *Data assimilation : The ensemble Kalman filter*. Springer.
- Evensen, G. and van Leeuwen, P. J. (2000). An ensemble kalman smoother for nonlinear dynamics. *Monthly Weather Review*, 128(6):1852–1867.
- Fang, Z., Herrmann, F. J., and Silva, C. D. (2014). Fast uncertainty quantification of 2D full-waveform inversion with randomized source subsampling. In *Expanded Abstracts, 76th Annual EAGE Conference & Exhibition, Amsterdam*. EAGE.
- Fang, Z., Silva, C. D., Kuske, R., and Herrmann, F. J. (2018). Uncertainty quantification for inverse problems with weak partial-differential-equation constraints. *Geophysics*, 83(6):R629–R647.

- Feller, W. (2008). *An introduction to probability theory and its applications*. John Wiley & Sons.
- Fichtner, A. (2011). *Full seismic waveform modelling and inversion*. Springer science & business media.
- Fichtner, A., Kennett, B. L. N., Igel, H., and Bunge, H. P. (2008). Theoretical background for continental- and global-scale full-waveform inversion in the time-frequency domain. *Geophysical Journal International*, 175:665–685.
- Fichtner, A., Kennett, B. L. N., Igel, H., and Bunge, H. P. (2009). Full seismic waveform tomography for upper-mantle structure in the Australasian region using adjoint methods. *Geophysical Journal International*, 179(3):1703–1725.
- Fichtner, A. and Trampert, J. (2011a). Hessian kernels of seismic data functionals based upon adjoint techniques. *Geophysical Journal International*, 185(2):775–798.
- Fichtner, A. and Trampert, J. (2011b). Resolution analysis in full waveform inversion. *Geophysical Journal International*, 187:1604–1624.
- Fichtner, A. and van Leeuwen, T. (2015). Resolution analysis by random probing. *Journal of Geophysical Research: Solid Earth*, pages 5549–5573. 2015JB012106.
- Fichtner, A., Zunino, A., and Gebraad, L. (2018a). Hamiltonian monte carlo solution of tomographic inverse problems. *Geophysical Journal International*, 216(2):1344–1363.
- Fichtner, A., Zunino, A., and Gebraad, L. (2018b). A tutorial introduction to the hamiltonian monte carlo solution of weakly nonlinear inverse problems.
- Fletcher, S. (2017). *Data Assimilation for the Geosciences*. Elsevier.
- Fletcher, S. J. and Zupanski, M. (2006). A hybrid multivariate normal and lognormal distribution for data assimilation. *Atmospheric Science Letters*, 7(2):43–46.
- French, S. W. and Romanowicz, B. A. (2015). Broad plumes rooted at the base of the earth’s mantle beneath major hotspots. *Nature*, 525(7567):95.
- Furrer, R. and Bengtsson, T. (2007). Estimation of high-dimensional prior and posterior covariance matrices in kalman filter variants. *Tellus A: Dynamic Meteorology and Oceanography*, 98(2):227–255.
- Gallagher, K., Charvin, K., Nielsen, S., Sambridge, M., and Stephenson, J. (2009). Markov chain monte carlo (mcmc) sampling methods to determine optimal models, model resolution and model choice for earth science problems. *Marine and Petroleum Geology*, 26(4):525–535.
- Gardner, G. H. F., Gardner, L. W., and Gregory, A. R. (1974). Formation velocity and density—the diagnostic basics for stratigraphic traps. *Geophysics*, 39:770–780.
- Gaspari, G. and Cohn, S. E. (1999). Construction of correlation functions in two and three dimensions. *Quarterly Journal of the Royal Meteorological Society*, 125(554):723–757.
- Gauthier, O., Virieux, J., and Tarantola, A. (1986). Two-dimensional nonlinear inversion of seismic waveforms: numerical results. *Geophysics*, 51(7):1387–1403.

- Gebraad, L. and Fichtner, D. A. (2018). Bayesian elastic full-waveform inversion using hamiltonian monte carlo. In *AGU Fall Meeting Abstracts*.
- Geman, S. and Geman, D. (1987). Stochastic relaxation, gibbs distributions, and the bayesian restoration of images. In *Readings in computer vision*, pages 564–584. Elsevier.
- Gholami, Y., Brossier, R., Operto, S., Prioux, V., Ribodetti, A., and Virieux, J. (2013). Which parametrization is suitable for acoustic VTI full waveform inversion? - Part 2: application to Valhall. *Geophysics*, 78(2):R107–R124.
- Gineste, M. and Eidsvik, J. (2017). Seismic waveform inversion using the ensemble kalman smoother. In *79th EAGE Conference and Exhibition 2017*.
- Gineste, M., Eidsvik, J., and Zheng, Y. (2019). Seismic waveform inversion using an iterative ensemble kalman smoother. In *Second EAGE/PESGB Workshop on Velocities*.
- Gineste, M., Eidsvik, J., and Zheng, Y. (2020). Ensemble-based seismic inversion for a stratified medium. *Geophysics*, 85(1):R29–R39.
- Green, P. J. (1995). Reversible jump markov chain monte carlo computation and bayesian model determination. *Biometrika*, 82(4):711–732.
- Grewal, M. S. and Andrews, A. P. (2001). *Kalman Filtering: Theory and Practice Using MATLAB*. John Wileys and Sons.
- Guzzi, R. (2015). *Data Assimilation: Mathematical Concepts and Instructive Examples*. Springer.
- Halko, N., Martinsson, P.-G., and Tropp, J. A. (2011). Finding structure with randomness: Probabilistic algorithms for constructing approximate matrix decompositions. *SIAM review*, 53(2):217–288.
- Hamill, T. M., Mullen, S. L., Snyder, C., Baumhefner, D. P., and Toth, Z. (2000). Ensemble forecasting in the short to medium range: Report from a workshop. *Bulletin of the American Meteorological Society*, 81(11):2653–2664.
- Hamill, T. M. and Whitaker, J. S. (2005). Accounting for the error due to unresolved scales in ensemble data assimilation: A comparison of different approaches. *Monthly Weather Review*, 133(11):3132–3147.
- Hamill, T. M. and Whitaker, J. S. (2011). What constrains spread growth in forecasts initialized from ensemble kalman filters? *Monthly Weather Review*, 139(1):117–131.
- Hamill, T. M., Whitaker, J. S., and Snyder, C. (2001). Distance-dependent filtering of background error covariance estimates in an ensemble kalman filter. *Monthly Weather Review*, 129(11):2776–2790.
- Harlim, J. and Hunt, B. (2007). A non-gaussian ensemble filter for assimilating infrequent noisy observations. *Tellus A*, 59(2):225–237.
- Harlim, J. and Hunt, B. R. (2005). Local ensemble transform kalman filter: An efficient scheme for assimilating atmospheric data. *preprint*.
- Hastings, W. K. (1970). Monte carlo sampling methods using markov chains and their applications.

- Hong, T. and Sen, M. (2009). A new mcmc algorithm for seismic waveform inversion and corresponding uncertainty analysis. *Geophysical Journal International*, 177:14–32.
- Horn, R. A. and Johnson, C. R. (2012). *Matrix analysis*. Cambridge university press.
- Houtekamer, P. and Mitchel, H. L. (1998). Data assimilation using an ensemble Kalman filter technique. *Monthly Weather Review*, 126(3):796–811.
- Houtekamer, P. L. and Mitchell, H. L. (2001). A sequential ensemble kalman filter for atmospheric data assimilation. *Monthly Weather Review*, 129(1):123–137.
- Houtekamer, P. L. and Mitchell, H. L. (2005). Ensemble kalman filtering. *Quarterly Journal of the Royal Meteorological Society: A journal of the atmospheric sciences, applied meteorology and physical oceanography*, 131(613):3269–3289.
- Houtekamer, P. L., Mitchell, H. L., Pellerin, G., Buehner, M., Charron, M., Spacek, L., and Hansen, B. (2005). Atmospheric data assimilation with an ensemble kalman filter: Results with real observations. *Monthly Weather Review*, 133(3):604–620.
- Houtekamer, P. L. and Zhang, F. (2016). Review of the ensemble kalman filter for atmospheric data assimilation. *Monthly Weather Review*, 144(12):4489–4532.
- Humpherys, J., Redd, P., and West, J. (2012). A fresh look at the kalman filter. *SIAM review*, 54(4):801–823.
- Hunt, B., Kostelich, E., and Szunyogh, I. (2007). Efficient data assimilation for spatiotemporal chaos: A local ensemble transform kalman filter. *Physica D: Nonlinear Phenomena*, 230(1):112–126.
- Hustedt, B., Operto, S., and Virieux, J. (2004). Mixed-grid and staggered-grid finite difference methods for frequency domain acoustic wave modelling. *Geophysical Journal International*, 157:1269–1296.
- Ide, K., Courtier, P., Ghil, M., and Lorenc, A. C. (1997). Unified notation for data assimilation: Operational, sequential and variational (gtspecial issue\data assimilation in meteorology and oceanography: Theory and practice). *Journal of the Meteorological Society of Japan. Ser. II*, 75(1B):181–189.
- Iglesias, M. A., Law, K. J., and Stuart, A. M. (2013a). Ensemble kalman methods for inverse problems. *Inverse Problems*, 29(4):045001.
- Iglesias, M. A., Law, K. J., and Stuart, A. M. (2013b). Evaluation of gaussian approximations for data assimilation in reservoir models. *Computational Geosciences*, 17(5):851–885.
- Jazwinski, A. H. (2007). *Stochastic processes and filtering theory*. Courier Corporation.
- Jin, L., Sen, M. K., and Stoffa, P. L. (2008). One-dimensional prestack seismic waveform inversion using ensemble kalman filter. In *SEG Technical Program Expanded Abstracts 2008*, pages 1920–1924. SEG.
- Jin, S., Madariaga, R., Virieux, J., and Lambaré, G. (1992). Two-dimensional asymptotic iterative elastic inversion. *Geophysical Journal International*, 108:575–588.
- Jordan, M. (2015). Estimation of spatial uncertainties in tomographic images. In *Expanded Abstracts, 77th Annual EAGE Meeting (Madrid)*.

- Kalman, R. (1960). A new approach to linear filtering and prediction problems. *Journal of basic Engineering*, 82(1):35–45.
- Kalman, R. E. and Bucy, R. S. (1961). New results in linear filtering and prediction theory. *Journal of basic engineering*, 83(1):95–108.
- Kalmikov, A. G. and Heimbach, P. (2014). A hessian-based method for uncertainty quantification in global ocean state estimation. *SIAM Journal on Scientific Computing*, 36(5):S267–S295.
- Kopp, R. E. and Orford, R. J. (1963). Linear regression applied to system identification for adaptive control systems. *Aiaa Journal*, 1(10):2300–2306.
- Labbe, R. R. (2016). Kalman and bayesian filters in python, [Online]. <https://github.com/rllabbe/Kalman-and-Bayesian-Filters-in-Python/> - (accessed 07/30/2019).
- Lailly, P. (1983). The seismic problem as a sequence of before-stack migrations. In Bednar, J., editor, *Conference on Inverse Scattering: Theory and Applications*. SIAM, Philadelphia.
- Lambaré, G., Virieux, J., Madariaga, R., and Jin, S. (1992). Iterative asymptotic inversion in the acoustic approximation. *Geophysics*, 57:1138–1154.
- Le Gland, F., Monbet, V., and Tran, V.-D. (2009). Large sample asymptotics for the ensemble kalman filter.
- Lee, E.-J. and Chen, P. (2013). Automating seismic waveform analysis for full 3-D waveform inversions. *Geophysical Journal International*, 194:572–589.
- Lee, K., Jung, S., and Choe, J. (2016). Ensemble smoother with clustered covariance for 3d channelized reservoirs with geological uncertainty. *Journal of Petroleum Science and Engineering*, 145:423–435.
- Leeuwenburgh, O., Evensen, G., and Bertino, L. (2005). The impact of ensemble filter definition on the assimilation of temperature profiles in the tropical pacific. *Quarterly Journal of the Royal Meteorological Society*, 131(631):3291–3300.
- Leonard, R. and Munns, J. (1987). *Valhall Field in Geology of Norwegian Oil and Gas Fields*. Graham and Trotman.
- Li, H., Kalnay, E., and Miyoshi, T. (2009). Simultaneous estimation of covariance inflation and observation errors within an ensemble kalman filter. *Quarterly Journal of the Royal Meteorological Society: A journal of the atmospheric sciences, applied meteorology and physical oceanography*, 135(639):523–533.
- Li, Y., Brossier, R., and Métivier, L. (2019). 3D frequency-domain elastic wave modeling using spectral element method with a parallel direct linear solver. In *Expanded Abstracts, 81th Annual EAGE Conference & Exhibition, London*, page Th R06 05. EAGE.
- Liberty, E., Woolfe, F., Martinsson, P.-G., Rokhlin, V., and Tygert, M. (2007). Randomized algorithms for the low-rank approximation of matrices. *Proceedings of the National Academy of Sciences*, 104(51):20167–20172.
- Liu, J. W. H. (1992). The multifrontal method for sparse matrix solution: theory and practice. *SIAM review*, 34(1):82–109.

- Liu, M. and Grana, D. (2018). Stochastic nonlinear inversion of seismic data for the estimation of petroelastic properties using the ensemble smoother and data reparameterization. *Geophysics*, 83(3):M25–M39.
- Liu, Q. and Peter, D. (2019). Square-root variable metric based elastic full-waveform inversion—part 2: uncertainty estimation. *Geophysical Journal International*, 218(2):1100–1120.
- Lorenc, A. C. (2003). The potential of the ensemble kalman filter for nwp—a comparison with 4d-var. *Quarterly Journal of the Royal Meteorological Society: A journal of the atmospheric sciences, applied meteorology and physical oceanography*, 129(595):3183–3203.
- Lorenz, E. N. (1963). Deterministic nonperiodic flow. *Journal of the atmospheric sciences*, 20(2):130–141.
- Luo, X. and Hoteit, I. (2013). Covariance inflation in the ensemble kalman filter: A residual nudging perspective and some implications. *Monthly Weather Review*, 141(10):3360–3368.
- Luo, Y. and Schuster, G. T. (1991). Wave-equation travelttime inversion. *Geophysics*, 56(5):645–653.
- Maggi, A., Tape, C., Chen, M., Chao, D., and Tromp, J. (2009). An automated time-window selection algorithm for seismic tomography. *Geophysical Journal International*, 178:257–281.
- Mandel, J. and Beezley, J. D. (2009). An ensemble kalman-particle predictor-corrector filter for non-gaussian data assimilation. In *International Conference on Computational Science*, pages 470–478. Springer.
- Mandel, J., Cobb, L., and Beezley, J. D. (2011). On the convergence of the ensemble kalman filter. *Applications of Mathematics*, 56(6):533–541.
- Marfurt, K. (1984). Accuracy of finite-difference and finite-element modeling of the scalar and elastic wave equations. *Geophysics*, 49:533–549.
- Martin, G. S., Wiley, R., and Marfurt, K. J. (2006). Marmousi2: An elastic upgrade for Marmousi. *The Leading Edge*, 25(2):156–166.
- Martin, J., Wilcox, L., Burstedde, C., and Ghattas, O. (2012). A stochastic Newton MCMC method for large-scale statistical inverse problems with application to seismic inversion. *SIAM Journal of Scientific Computing*, 34(3):A1460–A1487.
- Mazzotti, A., Bienati, N., Stucchi, E., Tognarelli, A., Aleardi, M., and Sajeva, A. (2016). Two-grid genetic algorithm full-waveform inversion. *The Leading Edge*, 35(12):1068–1075.
- Métivier, L. and Brossier, R. (2016). The SEISCOPE optimization toolbox: A large-scale nonlinear optimization library based on reverse communication. *Geophysics*, 81(2):F11–F25.
- Métivier, L., Brossier, R., Mérigot, Q., and Oudet, E. (2019). Graph space optimal transport for FWI: Auction algorithm, application to the 2d valhall case study. In *Expanded Abstracts, 81th Annual EAGE Conference & Exhibition, London*, page Tu R08 03. EAGE.
- Métivier, L., Brossier, R., Mérigot, Q., Oudet, E., and Virieux, J. (2016). An optimal transport approach for seismic tomography: Application to 3D full waveform inversion. *Inverse Problems*, 32(11):115008.

- Métivier, L., Brossier, R., Operto, S., and Virieux, J. (2017). Full waveform inversion and the truncated Newton method. *SIAM Review*, 59(1):153–195.
- Métivier, L., Brossier, R., Virieux, J., and Operto, S. (2013). Full Waveform Inversion and the truncated Newton method. *SIAM Journal On Scientific Computing*, 35(2):B401–B437.
- Metropolis, N., Rosenbluth, A. W., Rosenbluth, M. N., Teller, A. H., and Teller, E. (1953). Equation of state calculations by fast computing machines. *The journal of chemical physics*, 21(6):1087–1092.
- Mitchell, H. L. and Houtekamer, P. L. (2000). An adaptive ensemble kalman filter. *Monthly Weather Review*, 128(2):416–433.
- Mitchell, M. (1998). *An introduction to genetic algorithms*. MIT press.
- Miyoshi, T. (2011a). The gaussian approach to adaptive covariance inflation and its implementation with the local ensemble transform kalman filter. *Monthly Weather Review*, 139(9):1519–1535.
- Miyoshi, T. (2011b). The gaussian approach to adaptive covariance inflation and its implementation with the local ensemble transform kalman filter. *Monthly weather review*, 139(5):1519–1535.
- Mohorovičić, A. (1909). Das beben vom 8. x. 1909. *Jb. Met. Obs. Zagreb (Agram)*, 9:1–63.
- MUMPS team (2017). *MUMPS - Multifrontal Massively Parallel Solver users' guide - version 5.1.1 (March 21, 2017)*. ENSEEIHT-ENS Lyon, <http://mumps-solver.org>.
- Munns, J. W. (1985). The Valhall field: a geological overview. *Marine and Petroleum Geology*, 2:23–43.
- Myrseth, I. and Omre, H. (2010). Hierarchical ensemble kalman filter. *Spe Journal*, 15(02):569–580.
- Myrseth, I., Sætrom, J., and Omre, H. (2013). Resampling the ensemble kalman filter. *Computers & geosciences*, 55:44–53.
- Navon, I. M. (2009). Data assimilation for numerical weather prediction: a review. In *Data assimilation for atmospheric, oceanic and hydrologic applications*, pages 21–65. Springer.
- Nocedal, J. and Wright, S. J. (2006). *Numerical Optimization*. Springer, 2nd edition.
- Oke, P. R., Sakov, P., and Corney, S. P. (2007). Impacts of localisation in the enkf and enoi: experiments with a small model. *Ocean Dynamics*, 57(1):32–45.
- Operto, S., Brossier, R., Gholami, Y., Métivier, L., Prioux, V., Ribodetti, A., and Virieux, J. (2013). A guided tour of multiparameter full waveform inversion for multicomponent data: from theory to practice. *The Leading Edge*, Special section Full Waveform Inversion(September):1040–1054.
- Operto, S. and Miniussi, A. (2018). On the role of density and attenuation in 3D multi-parameter visco-acoustic VTI frequency-domain FWI: an OBC case study from the North Sea. *Geophysical Journal International*, 213:2037–2059.
- Operto, S., Miniussi, A., Brossier, R., Combe, L., Métivier, L., Monteiller, V., Ribodetti, A., and Virieux, J. (2015). Efficient 3-D frequency-domain mono-parameter full-waveform inversion of ocean-bottom cable data: application to Valhall in the visco-acoustic vertical transverse isotropic approximation. *Geophysical Journal International*, 202(2):1362–1391.

- Operto, S., Virieux, J., Ribodetti, A., and Anderson, J. E. (2009). Finite-difference frequency-domain modeling of visco-acoustic wave propagation in two-dimensional TTI media. *Geophysics*, 74 (5):T75–T95.
- Ortega, J. M. and Rheinboldt, W. C. (1970). *Iterative solution of nonlinear equations in several variables*, volume 30. Siam.
- Ott, E., Hunt, B. R., Szunyogh, I., Zimin, A. V., Kostelich, E. J., Corazza, M., Kalnay, E., Patil, D., and Yorke, J. A. (2004). A local Ensemble Kalman filter for atmospheric data assimilation. *Tellus A*, 56:415–428.
- Panning, M. P., Lekić, V., and Romanowicz, B. A. (2010). Importance of crustal corrections in the development of a new global model of radial anisotropy. *Journal of Geophysical Research: Solid Earth*, 115(B12).
- Petrie, R. E. (2008). Localization in the ensemble kalman filter. *Master's thesis, MSc. Atmosphere, Ocean and Climate, University of Reading, Reading, UK*.
- Petrie, R. E. and Dance, S. L. (2010). Ensemble-based data assimilation and the localisation problem. *Weather*, 65(3):65–69.
- Plessix, R. E. (2006). A review of the adjoint-state method for computing the gradient of a functional with geophysical applications. *Geophysical Journal International*, 167(2):495–503.
- Plessix, R. E. (2009). Three-dimensional frequency-domain full-waveform inversion with an iterative solver. *Geophysics*, 74(6):WCC53–WCC61.
- Plessix, R.-E., Baeten, G., de Maag, J. W., and ten Kroode, F. (2012). Full waveform inversion and distance separated simultaneous sweeping: a study with a land seismic data set. *Geophysical Prospecting*, 60:733 – 747.
- Pratt, R. G. (1990). Frequency-domain elastic modeling by finite differences: a tool for crosshole seismic imaging. *Geophysics*, 55(5):626–632.
- Pratt, R. G. (1999). Seismic waveform inversion in the frequency domain, part I: theory and verification in a physical scale model. *Geophysics*, 64:888–901.
- Pratt, R. G. and Goulty, N. R. (1991). Combining wave-equation imaging with traveltimes tomography to form high-resolution images from crosshole data. *Geophysics*, 56(2):204–224.
- Pratt, R. G., Shin, C., and Hicks, G. J. (1998). Gauss-Newton and full Newton methods in frequency-space seismic waveform inversion. *Geophysical Journal International*, 133:341–362.
- Pratt, R. G., Sirgue, L., Hornby, B., and Wolfe, J. (2008). Cross-well waveform tomography in fine-layered sediments - meeting the challenges of anisotropy. In *70th Annual EAGE Conference & Exhibition, Roma*, page F020.
- Pratt, R. G., Song, Z. M., Williamson, P. R., and Warner, M. (1996). Two-dimensional velocity models from wide-angle seismic data by wavefield inversion. *Geophysical Journal International*, 124:323–340.
- Pratt, R. G. and Worthington, M. H. (1990). Inverse theory applied to multi-source cross-hole tomography. Part I: acoustic wave-equation method. *Geophysical Prospecting*, 38:287–310.

- Press, W. H., Flannery, B. P., Teukolsky, S. A., and Vetterling, W. T. (1986). *Numerical recipes : the art of scientific computing*. Cambridge University Press.
- Prieux, V., Brossier, R., Operto, S., and Virieux, J. (2013a). Multiparameter full waveform inversion of multicomponent OBC data from Valhall. Part 1: imaging compressional wavespeed, density and attenuation. *Geophysical Journal International*, 194(3):1640–1664.
- Prieux, V., Brossier, R., Operto, S., and Virieux, J. (2013b). Multiparameter full waveform inversion of multicomponent OBC data from Valhall. Part 2: imaging compressional and shear-wave velocities. *Geophysical Journal International*, 194(3):1665–1681.
- Raanes, P. N., Bocquet, M., and Carrassi, A. (2019). Adaptive covariance inflation in the ensemble kalman filter by gaussian scale mixtures. *Quarterly Journal of the Royal Meteorological Society*, 145(718):53–75.
- Rao, V. and Sandu, A. (2015). A posteriori error estimates for the solution of variational inverse problems. *SIAM/ASA Journal on Uncertainty Quantification*, 3(1):737–761.
- Rawlinson, N. and Spakman, W. (2016). On the use of sensitivity tests in seismic tomography. *Geophysical Journal International*, 205(2):1221–1243.
- Ray, A., Sekar, A., Hoversten, G. M., and Albertin, U. (2016). Frequency domain full waveform elastic inversion of marine seismic data from the alba field using a bayesian trans-dimensional algorithm. *Geophysical Journal International*, 205(2):915–937.
- Rodell, M., Houser, P., Jambor, U., Gottschalck, J., Mitchell, K., Meng, C.-J., Arsenault, K., Cosgrove, B., Radakovich, J., Bosilovich, M., et al. (2004). The global land data assimilation system. *Bulletin of the American Meteorological Society*, 85(3):381–394.
- Sajeva, A., Aleardi, M., Galuzzi, B., Stucchi, E., Spadavecchia, E., and Mazzotti, A. (2017a). Comparing the performances of four stochastic optimisation methods using analytic objective functions, 1d elastic full-waveform inversion, and residual static computation. *Geophysical Prospecting*, 65:322–346.
- Sajeva, A., Aleardi, M., and Mazzotti, A. (2017b). Genetic algorithm full-waveform inversion: uncertainty estimation and validation of the results. *Bollettino di Geofisica Teorica ed Applicata*, 58(4).
- Sakov, P. and Bertino, L. (2011). Relation between two common localisation methods for the enkf. *Computational Geosciences*, 15(2):225–237.
- Sakov, P. and Oke, P. R. (2008). Implications of the form of the ensemble transformation in the ensemble square root filters. *Monthly Weather Review*, 136(3):1042–1053.
- Sakov, P., Oliver, D. S., and Bertino, L. (2012). An iterative enkf for strongly nonlinear systems. *Monthly Weather Review*, 140(6):1988–2004.
- Sen, M. and Stoffa, P. (1991). Non-linear one-dimensional seismic waveform inversion using simulated annealing. *Geophysics*, 56:1624–1638.
- Sen, M. K. and Biswas, R. (2017). Transdimensional seismic inversion using the reversible jump Hamiltonian Monte Carlo algorithm. *Geophysics*.

- Serfozo, R. (2009). *Basics of applied stochastic processes*. Springer Science & Business Media.
- Sirgue, L., Barkved, O. I., Dellinger, J., Etgen, J., Albertin, U., and Kommedal, J. H. (2010). Full waveform inversion: the next leap forward in imaging at Valhall. *First Break*, 28:65–70.
- Stewart, L. M., Dance, S. L., and LNichols, N. K. (2008). Correlated observation errors in data assimilation. *International journal for numerical methods in fluids*, 56(8):1521–1527.
- Stoffa, P. L. and Sen, M. K. (1991). Nonlinear multiparameter optimization using genetic algorithms: Inversion of plane-wave seismograms. *Geophysics*, 56(11):1794–1810.
- Tape, C., Liu, Q., Maggi, A., and Tromp, J. (2010). Seismic tomography of the southern California crust based on spectral-element and adjoint methods. *Geophysical Journal International*, 180:433–462.
- Tarantola, A. (1984). Linearized inversion of seismic reflection data. *Geophysical Prospecting*, 32:998–1015.
- Tarantola, A. (1987). *Inverse problem theory: methods for data fitting and model parameter estimation*. Elsevier, New York.
- Tarantola, A. (2005). *Inverse Problem Theory and Methods for Model Parameter Estimation*. Society for Industrial and Applied Mathematics, Philadelphia.
- Tejero, C. J., Dagnino, D., Sallarès, V., and Ranero, C. R. (2015). Comparative study of objective functions to overcome noise and bandwidth limitations in full waveform inversion. *Geophysical Journal International*, 203(1):632–645.
- Tikhonov, A. and Arsenin, V. (1977). *Solution of ill-posed problems*. Winston, Washington, DC.
- Tippett, M. K., Anderson, J. L., Bishop, C. H., Hamill, T. M., and Whitaker, J. S. (2003). Ensemble square root filters. *Monthly Weather Review*, 131(7):1485–1490.
- Toksöz, M. N. and Johnston, D. H. (1981). *Geophysics reprint series, No. 2: Seismic wave attenuation*. Society of exploration geophysicists, Tulsa, OK.
- Tran, K. and Hiltunen, D. (2011). Two-dimensional inversion of full waveforms using simulated annealing. *Journal of Geotechnical and Geoenvironmental Engineering*, 138(9):1075–1090.
- Trinh, P. T., Brossier, R., Métivier, L., Tavard, L., and Virieux, J. (2019). Efficient 3D time-domain elastic and viscoelastic Full Waveform Inversion using a spectral-element method on flexible Cartesian-based mesh. *Geophysics*, 84(1):R75–R97.
- Tromp, J., Tape, C., and Liu, Q. (2005). Seismic tomography, adjoint methods, time reversal and banana-doughnut kernels. *Geophysical Journal International*, 160:195–216.
- van Leeuwen, P. J. (1999). Comment on “data assimilation using an ensemble kalman filter technique”. *Monthly Weather Review*, 127(6):1374–1377.
- van Leeuwen, T. and Herrmann, F. J. (2013). Mitigating local minima in full-waveform inversion by expanding the search space. *Geophysical Journal International*, 195(1):661–667.
- van Leeuwen, T. and Mulder, W. A. (2010). A correlation-based misfit criterion for wave-equation traveltime tomography. *Geophysical Journal International*, 182(3):1383–1394.

- Vasco, D. W., Johnson, L. R., and Majer, E. L. (1993). Ensemble inference in geophysical inverse problems. *Geophysical Journal International*, 115(3):711–728.
- Virieux, J., Asnaashari, A., Brossier, R., Métivier, L., Ribodetti, A., and Zhou, W. (2017). An introduction to Full Waveform Inversion. In Grechka, V. and Wapenaar, K., editors, *Encyclopedia of Exploration Geophysics*, pages R1–1–R1–40. Society of Exploration Geophysics.
- Virieux, J. and Operto, S. (2009). An overview of full waveform inversion in exploration geophysics. *Geophysics*, 74(6):WCC1–WCC26.
- Virieux, J., Operto, S., Ben Hadj Ali, H., Brossier, R., Etienne, V., Sourbier, F., Giraud, L., and Haidar, A. (2009). Seismic wave modeling for seismic imaging. *The Leading Edge*, 28(5):538–544.
- Wang, X. and Bishop, C. H. (2003). A comparison of breeding and ensemble transform kalman filter ensemble forecast schemes. *Journal of the atmospheric sciences*, 60(9):1140–1158.
- Wang, X., Bishop, C. H., and Julier, S. J. (2004). Which is better, an ensemble of positive-negative pairs or a centered spherical simplex ensemble? *Monthly Weather Review*, 132(7):1590–1605.
- Warner, M. and Guasch, L. (2014). Adaptive waveform inversion - FWI without cycle skipping - theory. In *76th EAGE Conference and Exhibition 2014*, page We E106 13.
- Warner, M., Nangoo, T., Shah, N., Umpleby, A., and Morgan, J. (2013a). *Full-waveform inversion of cycle-skipped seismic data by frequency down-shifting*, chapter 176, pages 903–907. Society of Exploration Geophysics.
- Warner, M., Ratcliffe, A., Nangoo, T., Morgan, J., Umpleby, A., Shah, N., Vinje, V., Stekl, I., Guasch, L., Win, C., Conroy, G., and Bertrand, A. (2013b). Anisotropic 3D full-waveform inversion. *Geophysics*, 78(2):R59–R80.
- Weston, P. P., Bell, W., and Eyre, J. R. (2014). Accounting for correlated error in the assimilation of high-resolution sounder data. *Quarterly Journal of the Royal Meteorological Society*, 140(685):2420–2429.
- Whitaker, J. S. and Hamill, T. M. (2002). Ensemble data assimilation without perturbed observations. *Monthly Weather Review*, 130(7):1913–1924.
- Whitaker, J. S. and Hamill, T. M. (2012). Evaluating methods to account for system errors in ensemble data assimilation. *Monthly Weather Review*, 140(9):3078–3089.
- Wolfe, P. (1969). Convergence conditions for ascent methods. *SIAM Review*, 11.
- Wu, R. S. and Toksöz, M. N. (1987). Diffraction tomography and multisource holography applied to seismic imaging. *Geophysics*, 52:11–25.
- Yang, Y., Engquist, B., Sun, J., and Hamfeldt, B. F. (2018). Application of optimal transport and the quadratic Wasserstein metric to full-waveform inversion. *Geophysics*, 83(1):R43–R62.
- Yilmaz, Ö. (1993). *Seismic data processing*. Society of exploration geophysicists.
- Zhang, F., Snyder, C., and Sun, J. (2004). Impacts of initial estimate and observation availability on convective-scale data assimilation with an ensemble kalman filter. *Monthly Weather Review*, 132(5):1238–1253.

- Zhou, W. (2016). *Full Waveform Inversion of Early Arrivals and Reflections for Velocity Model Building and Case Study with Gas Cloud Effect*. PhD thesis, Univ. Grenoble Alpes.
- Zhou, W., Brossier, R., Operto, S., Virieux, J., and Yang, P. (2018). Velocity model building by waveform inversion of early arrivals and reflections: a 2d ocean-bottom-cable study with gas cloud effects. *Geophysics*, 83(2):R141–R157.
- Zhu, H., Bozdağ, E., and Tromp, J. (2015). Seismic structure of the European upper mantle based on adjoint tomography. *Geophysical Journal International*, 201(1):18–52.
- Zhu, H., Li, S., Fomel, S., Stadler, G., and Ghattas, O. (2016). A bayesian approach to estimate uncertainty for full-waveform inversion using a priori information from depth migration. *Geophysics*, 81(5):R307–R323.
- Zupanski, M. (2005). Maximum likelihood ensemble filter: Theoretical aspects. *Monthly weather review*, 133(6):1710–1726.
- Zupanski, M., Navon, I. M., and Zupanski, D. (2008). The maximum likelihood ensemble filter as a non-differentiable minimization algorithm. *Quarterly Journal of the Royal Meteorological Society: A journal of the atmospheric sciences, applied meteorology and physical oceanography*, 134(633):1039–1050.

