PONTIFÍCIA UNIVERSIDADE CATÓLICA
DO RIO DE JANEIRO

**Thiago de Menezes Duarte e Silva**

**Evaluating the impact of the inflation factors generation for the ensemble smoother with multiple data assimilation**

**Tese de Doutorado**

Thesis presented to the Programa de Pós–graduação em Matemática, do Departamento de Matemática da PUC-Rio in partial fulfillment of the requirements for the degree of Doutor em Matemática.

Advisor     :          Prof. Sinesio Pesco
Co-advisor: Prof. Abelardo Borges Barreto Jr.

Rio de Janeiro
August 2021

**Thiago de Menezes Duarte e Silva**

**Evaluating the impact of the inflation factors generation for the ensemble smoother with multiple data assimilation**

Thesis presented to the Programa de Pós–graduação em Matemática da PUC-Rio in partial fulfillment of the requirements for the degree of Doutor em Matemática. Approved by the Examination Committee:

**Prof. Sinesio Pesco**
Advisor
Departamento de Matemática – PUC-Rio

**Prof. Abelardo Borges Barreto Jr.**
Co-advisor
Departamento de Engenharia Mecânica – PUC-Rio

**Prof. Adolfo Puime Pires**
Universidade Estadual Norte Fluminense – UENF

**Prof. Alexandre A. Emerick**
Centro de Pesquisa e Desenvolvimento Leopoldo Américo
Miguez de Mello – CENPES

**Prof. Hélio Côrtes Vieira Lopes**
Departamento de Informática – PUC-Rio

**Prof. Márcio da Silveira Carvalho**
Departamento de Engenharia Mecânica – PUC-Rio

**Prof. Mustafa Onur**
The University of Tulsa – TU

Rio de Janeiro, August the 13th, 2021

**Thiago de Menezes Duarte e Silva**

Thiago is graduated in mathematics from the Fluminense Federal University in 2014 and a master of mathematics from Pontifical Catholic University in 2017.

## Acknowledgments

I would like to acknowledge my parents, Neinha and José, for everything they have done for me. I do not know if there exist better parents in the world. I am fortunate to have you as my parents.

To my advisors, Sinésio and Abelardo, for every support and motivation they gave me during these six years of post-graduation. Also, the moments of deconcentration certainly helped my understanding of the topics I was studying. Thank you for everything.

I want to extend my thanks to professor Mustafa Onur for accepting supervising me as a visiting scholar at the University of Tulsa and for all help he gave me during my stay there. Also, for assisting me to improve my understanding of petroleum engineering. I also would like to thank TUPREP and the TUPREP members for politely receiving me.

To Carlos for saving my life, turning my computer on every time I needed it during these dark times of lockdowns. I extend this special thanks to Creuza, Kátia, and Mariana.

To my friends from PUC-Rio, that stayed by my side from the master's degree until the end of the Ph.D. In particular, Renan V. Bela helped me many times with the IMEX simulator and suggested valuable comments to improve my work.

A special thanks to professor Simone Dantas for all her support since I was a young man until now, ten years later.

Finally, I want to dedicate a special thanks to Tahyz for sharing these student years with me, beginning with the bachelor, then master, and Ph.D. Thank you for tolerating all crises and the moments when we want to cease everything and run to the countryside. Also, for giving the idea of the study presented in the third chapter of this thesis, even by accident.

# Abstract

Silva, Thiago de Menezes Duarte; Pesco, Sinesio (Advisor); Barreto Jr., Abelardo Borges (Co-Advisor). **Evaluating the impact of the inflation factors generation for the ensemble smoother with multiple data assimilation**. Rio de Janeiro, 2021. 93p. Tese de Doutorado – Departamento de Matemática, Pontifícia Universidade Católica do Rio de Janeiro.

The ensemble smoother with multiple data assimilation (ES-MDA) gained much attention as a powerful parameter estimation method. The main idea of the ES-MDA is to assimilate the same data multiple times with an inflated data error covariance matrix. In the original ES-MDA implementation, these inflation factors, such as the number of assimilations, are selected *a priori*. The only requirement is that the sum of the inflation factors' inverses must be equal to one. Therefore, selecting them equal to the number of assimilations is a straightforward choice. Nevertheless, recent studies have shown a relationship between the ES-MDA update equation and the solution to a regularized inverse problem. Hence, the inflation factors play the role of the regularization parameter at each ES-MDA assimilation step. As a result, they have also suggested new procedures to generate these elements based on the discrepancy principle. Although several studies proposed efficient techniques to generate the ES-MDA inflation factors, an optimal procedure to generate them remains an open problem. Moreover, the studies diverge on which regularization scheme is sufficient to provide the best ES-MDA outcomes. Therefore, in this work, we address the problem of generating the ES-MDA inflation factors and their influence on the method's performance. We present a numerical analysis of the influence of such factors on the main parameters of the ES-MDA, such as the ensemble size, the number of assimilations, and the ES-MDA vector of model parameters update. With the conclusions presented in the aforementioned analysis, we propose a new procedure to generate ES-MDA inflation factors based on a regularizing scheme for Levenberg-Marquardt algorithms. It is shown through a synthetic two-dimensional waterflooding problem that the new method achieves better model parameters and data match compared to the other ES-MDA implementations available in the literature.

## Keywords

History matching; Uncertainty quantification; Ensemble smoother with multiple data assimilation; Reservoir characterization.

# Resumo

Silva, Thiago de Menezes Duarte; Pesco, Sinesio; Barreto Jr., Abelardo Borges. **Investigando o impacto da geração dos fatores de inflação para o *ensemble smoother* com múltipla assimilação de dados**. Rio de Janeiro, 2021. 93p. Tese de Doutorado – Departamento de Matemática, Pontifícia Universidade Católica do Rio de Janeiro.

O *ensemble smoother with multiple data assimilation* (ES-MDA) se tornou um poderoso estimador de parâmetros. A principal ideia do ES-MDA é assimilar os mesmos dados com a matriz de covariância dos erros dos dados inflada. Na implementação original do ES-MDA, os fatores de inflação e o número de assimilações são escolhidos *a priori*. O único requisito é que a soma dos inversos de tais fatores seja igual a um. Naturalmente, escolhendo-os iguais ao número de assimilações cumpre este requerimento. Contudo, estudos recentes mostraram uma relação entre a equação de atualização do ES-MDA com a solução para o problema inverso regularizado. Consequentemente, tais elementos agem como os parâmetros de regularização em cada assimilação. Assim, estudos propuseram técnicas para gerar tais fatores baseadas no princípio da discrepância. Embora estes estudos tenham propostos técnicas, um procedimento ótimo para gerar os fatores de inflação continua um problema em aberto. Mais ainda, tais estudos divergem em qual método de regularização é sufiente para produzir os melhores resultados para o ES-MDA. Portanto, nesta tese é abordado o problema de gerar os fatores de inflação para o ES-MDA e suas influências na performance do método. Apresentamos uma análise numérica do impacto de tais fatores nos parâmetros principais do ES-MDA: o tamanho do conjunto, o número de assimilações e o vetor de atualização dos parâmetros. Com a conclusão desta análise, nós propomos uma nova técnica para gerar os fatores de inflação para o ES-MDA baseada em um método de regularização para algorítmos do tipo Levenberg-Marquardt. Investigando os resultados de um problema de inundação de um reservatório 2D, o novo método obtém melhor estimativa tanto para os parâmetros do modelo tanto quanto para os dados observados.

## Palavras-chave

Ajuste de histórico; Quantificação de incertezas; Ensemble smoother com múltipla assimilação de dados; Caracterização de reservatórios.

# Table of contents

# List of figures

# List of tables

# List of Symbols

EnKF - Ensemble Kalman filter

ES - Ensemble smoother

ES-MDA - Ensemble smoother with multiple data assimilation

GN - Gauss-Newton

LM - Levenberg-Marquardt

# 1
# Introduction

Predicting the performance of an oil field plays a crucial role in reservoir engineering. The decisions that have to be made in developing and managing the reservoir depend on these pieces of information. The risks involved in this process are considerably high, demanding reasonable administration of uncertainty during an exploitation project. Likewise, reservoir simulation is an important tool when prognosticating an oil deposit's performance and administering uncertainty. As a result, constructing a robust reservoir model is, therefore, an important task. The process of characterizing the reservoir may be executed by incorporating dynamic observed data from a real field in a model, which is a popular technique called history matching.

The history matching approach consists of an inverse problem where the vector of reservoir model parameters is estimated given a vector of observations. This problem is intrinsically ill-posed because it is not possible to guarantee the uniqueness of a solution. A simple example illustrating this ill-posedness is a two-layered reservoir model where each layer has the same properties but permeability. The first layer has a permeability of 100 mD, and the second layer has a permeability of 200 mD. It is easy to show that the equivalent permeability of the system is equal to 150 mD. On the other hand, the same reservoir system with both layers with a permeability of 150 mD would provide the same equivalent permeability. The same result would occur if the permeabilities were exchanged within the layers. Following this simple example, it is straightforward to notice that there are infinite combinations of permeabilities that may result in the same equivalent permeability of 150 mD. Moreover, in the real field measurements, the data are regularly inaccurate and sometimes inconsistent. Therefore, the parameters estimated for a model in the history matching technique and their predictions are always uncertain.

This simple example exposes the difficulty of solving the inverse problem of history matching. A traditional approach to deal with the uncertainty of inverse problems is that one afforded by Bayesian statistics. The strategy consists of writing down the posterior conditional probability density function of the model parameters vector given a set of observed data. Therefore, the problem of estimating the model parameters reduces to sampling this

probability density function correctly. Moreover, these samples allow a good evaluation of uncertainty and managing predictions of the reservoir production.

Data assimilation methods applied to the history matching problem aim to minimize an objective function composed of two parts: one related to the data mismatch and another related to the model parameters mismatch. For practical history matching, the relation between the theoretical data and the model parameters is highly nonlinear. As a result, the objective function to be minimized to incorporate production data information into the reservoir model will be nonlinear. It transforms the problem of sampling the posterior probability density function of the vector of model parameters given the observations in solving a nonlinear minimization problem.

Many techniques are known to solve such minimization problems, and most of them are based on the computation of the gradient of the objective function. However, it is not often feasible to differentiate the nonlinear function that relates the model parameters to the theoretical data for reservoir simulation problems. One procedure to compute the gradient of the objective function is the adjoint method. However, such a method is not usually available for commercial simulators, which hinders the presentation of a general solution. Another way to differentiate might be the numerical perturbation methods, such as finite-difference. Nevertheless, the computational consumption of derivatives is linked to the number of parameters to be estimated. As the number of parameters in the history matching problem is often large, this technique might be costly.

Alternative methods such as ensemble-based are used to estimate the sensitivity matrix of the data from an ensemble. This approach considers the simulator a black box, which presents a general solution for the problem, enabling coupling the methods with any commercial or academic simulator. The sensitivity is easily computed by the simulator's outputs, and the computational cost is independent of the number of parameters. The only restriction to this approach is the size of the ensemble, which is problem-dependent. Estimating the sensitivity for a highly nonlinear system from an ensemble may require many members. In contrast, a moderate nonlinear system might demand a more modest ensemble.

Among the ensemble-based methods, the Ensemble Kalman Filter (EnKF) might be the most famous, with a large number of studies evaluating its performance in a large number of history matching problems. However, studies have noticed that the EnKF may result in non-physical values for the model parameters for highly nonlinear problems, e.g., providing negative pressure responses and saturations. As a result, many iterative forms of

the EnKF were proposed to deal with this inconsistency. Nevertheless, these iterative EnKF techniques demand increased computational time. Moreover, as the EnKF assimilates data sequentially in time, the simulator needs to be restarted at every assimilation step to update the state variables, raising an elevated computational expense.

The Ensemble Smoother (ES) appears as an alternative to the sequential data assimilation of the EnKF. Running the simulator from time zero until the end of the history time, the ES assimilates the data only once, providing a global and robust update of the model parameters vector. One of the main advantages of the ES is that restarting the simulator at every time step is not needed, which makes the ES implementation much easier than the EnKF. Moreover, as all data are assimilated in a unique update, the ES works as a traditional parameter estimation method, removing the inconsistency of update state vectors of the EnKF. Unfortunately, it has been proven that, at each time-step, implementing the EnKF update is similar to applying a Gauss-Newton iteration to the model state vector. On the other hand, as the ES assimilates the data with an individual update, applying the ES is similar to applying a single Gauss-Newton correction to the model parameters vector, proving ineffectiveness in affording reasonable estimates.

To overcome this issue of the ES, many iterative forms have been proposed in the literature. The most famous might be the ensemble smoother with multiple data assimilation (ES-MDA), which assimilates the same data multiple times with an inflated data error covariance matrix. The ES-MDA may be understood as multiple soft Gauss-Newton corrections to the model parameters instead of a single and potent update supplied by the ES. The ES-MDA has been proven to sample the posterior probability density function of the model parameters given a set of observations if the sum of the inverses of the data-error inflation factors is equal to one. In other words, the harmonic mean of the inflation factors must be equal to the number of assimilations. In the original ES-MDA implementation, these factors, such as the number of assimilations, must be selected a priori. Thus, a straightforward choice is setting all the inflation factors equal to the number of assimilations. However, this simple selection of the inflation factors might bring issues when applying the ES-MDA.

Different procedures to generate the ES-MDA inflation factors have been proposed in the last few years. However, it remains an open problem deciding which method achieves optimal performance and the best match of observed data and model parameters. Recently, the ES-MDA update equation has been compared to the solution to the regularized inverse problem. Thus, schemes

based on the discrepancy principle were largely used to generate the ES-MDA inflation factors. Therefore, this thesis is dedicated to studying the ES-MDA inflation factors following two paths: analyzing different techniques that have been proposed to generate these factors and proposing a new method to generate these factors. In the first approach, we analyze the generation of these elements using two popular procedures in generating the first inflation factor for the ES-MDA. Moreover, we present a numerical procedure that enables evaluate the quality of the inflation factor for the current problem and ensemble size. Following the results of the first part of this thesis, we present a new method to generate the ES-MDA inflation factors by setting the first and the last ones. The others are computed geometrically in decreasing order.

This thesis is segmented as follows: Chapter 2 supplies a summary of the Bayesian approach and how to generate the desired objective function to be minimized in the history matching problem; Chapter 3 displays the analysis of the inflation factors selection and their influences on the ES-MDA main parameters, such as the ensemble size and the number of assimilations, moreover, we discuss which method may be optimal to generate optimal ES-MDA outcomes; finally, Chapter 4 exhibits a new procedure to generate the ES-MDA inflation factors.

# 2
# Bayesian framework of the history matching problem

In Bayesian statistics, probability theory is used to estimate the uncertainty related to the data. This mechanism helps to deal with inverse problems due to the reduced amount of available data and the high number of model parameters to be estimated. Moreover, the data is almost inexact, containing a substantial quantity of measurement inaccuracies. Thus, a measure of determining the uncertainty is wishful. In reservoir engineering history matching problems, which is intrinsically an inverse problem, the data's inaccuracy is a relevant subject to reflect when estimating model parameters such that their production simulated responses match the observed data. As the reservoir models are often poor approximations of real ones, a probabilistic approach enables measuring the uncertainty related to the model and the data.

## 2.1
## Objective function

The starting spot of the Bayesian approach is considering $m \in \mathbb{R}^{N_m}$ the vector of model parameters and $d \in \mathbb{R}^{N_d}$ the vector of data. In this sense, $N_m$ is the number of model parameters, and $N_d$ is the number of available data. Consider the linear case, where the theoretical relation of the vector $m$ and the vector $d$ is given by:

$$d = Gm, \qquad (2\text{-}1)$$

where $G \in \mathbb{R}^{N_d \times N_m}$ is the sensitivity matrix of the data. Assign $d_{obs} \in \mathbb{R}^{N_d}$ to be the vector that contains a set of observations. Assuming that there are measurement errors $\epsilon \in \mathbb{R}^{N_d}$ in $d_{obs}$ and no model errors, we can write such vector as:

$$d_{obs} = Gm + \epsilon. \qquad (2\text{-}2)$$

The measurement errors $\epsilon$ are often assumed to have a Gaussian distribution with zero mean and covariance $C_d \in \mathbb{R}^{N_d \times N_d}$ [1]. Hence, the probability of

obtaining $d_{obs}$ from sampling the vector of model parameters $m$ can be deemed the probability of $\epsilon$, as:

$$f(d_{obs}|m) = f(\epsilon = d_{obs} - Gm). \tag{2-3}$$

As we assumed that $\epsilon$ have Gaussian distribution, we can describe the probability density function $f(\epsilon)$ as:

$$f(\epsilon) \propto \exp\left\{-\frac{1}{2}\left(Gm - d_{obs}\right)^T C_d^{-1}\left(Gm - d_{obs}\right)\right\}. \tag{2-4}$$

Additionally, if we consider that the model parameters $m$ are also uncertain and that they assume Gaussian distribution, it is possible to write their probability density function as:

$$f(m) \propto \exp\left\{-\frac{1}{2}\left(m - m_{pr}\right)^T C_m^{-1}\left(m - m_{pr}\right)\right\}, \tag{2-5}$$

where $m_{pr} \in \mathbb{R}^{N_m}$ is a prior estimate of the model parameters and $C_m \in \mathbb{R}^{N_m \times N_m}$ is the prior model covariance matrix. Finally, Bayes' rule permits one to write the probability density function of the model parameters $m$ conditioned by the data $d_{obs}$ as:

$$f(m|d_{obs}) = \frac{f(d_{obs}|m)f(m)}{f(d_{obs})} = \frac{f(d_{obs}|m)f(m)}{\int_\Omega f(d_{obs}|m)f(m)dm}. \tag{2-6}$$

The above equation can be reduced to:

$$f(m|d_{obs}) = aL(m|d_{obs})f(m), \tag{2-7}$$

where the function $L(m|d_{obs})$ is the likelihood function, which corresponds to $f(d_{obs}|m)$ and $a$ is a normalizing constant. From Equations (2-4) and (2-5), one can write down the probability density function depicted in Equation (2-7) in the following way:

$$f(m|d_{obs}) = a\exp\left\{-\mathcal{O}(m)\right\}, \tag{2-8}$$

where:

$$\mathcal{O}(m) = \mathcal{O}_m(m) + \mathcal{O}_d(m), \tag{2-9}$$

with:

$$\mathcal{O}_m(m) = (m - m_{pr})^T C_m^{-1} (m - m_{pr}).\tag{2-10}$$

and:

$$\mathcal{O}_d(m) = (Gm - d_{obs})^T C_d^{-1} (Gm - d_{obs}),\tag{2-11}$$

As the probability density function $f(m|d_{obs})$ in Equation (2-8) is assumed to have multivariate Gaussian distribution, it finds its maximum value when $\mathcal{O}(m)$ finds its minimum. The function $\mathcal{O}(m)$ measures the data and model parameters mismatch and, thus, it refers to the desired objective function to be minimized to achieve a vector of model parameters $m$ with the maximum probability of matching the observed data $d_{obs}$.

## 2.2
## Maximum *a posteriori* estimate

When the relation between the model parameters and the data is linear, one can prove that the constructed objective function $\mathcal{O}(m)$ is quadratic for any given vector $m \in \mathbb{R}^{N_m}$. Thus, to obtain a minimum of $\mathcal{O}(m)$, it is sufficient to calculate its gradient $\nabla\mathcal{O}(m)$ and set it to zero, as:

$$\nabla\mathcal{O}(m) = C_m^{-1} (m - m_{pr}) + G^T C_d^{-1} (Gm - d_{obs}) = 0.\tag{2-12}$$

Adding and subtracting $Gm_{pr}$ inside the data mismatch parenthesis, we obtain that:

$$\begin{aligned}
0 &= C_m^{-1} (m - m_{pr}) + G^T C_d^{-1} (Gm - Gm_{pr} + Gm_{pr} - d_{obs}) \\
&= C_m^{-1} (m - m_{pr}) + G^T C_d^{-1} (Gm - Gm_{pr}) + G^T C_d^{-1} (Gm_{pr} - d_{obs}) \\
&= C_m^{-1} (m - m_{pr}) + G^T C_d^{-1} G (m - m_{pr}) + G^T C_d^{-1} (Gm_{pr} - d_{obs}) \\
&= \left(C_m^{-1} + G^T C_d^{-1} G\right) (m - m_{pr}) + G^T C_d^{-1} (Gm_{pr} - d_{obs}).
\end{aligned}\tag{2-13}$$

From Equation (2-13), we obtain a value of $m$ such that $\nabla\mathcal{O}(m) = 0$. This vector corresponds to the maximum *a posteriori* estimate of the vector of model parameters and is denoted as $m_{map}$. This vector is computed in the following way:

$$m_{map} = m_{pr} + \left(C_m^{-1} + G^T C_d^{-1} G\right)^{-1} G^T C_d^{-1} (d_{obs} - Gm_{pr}).\tag{2-14}$$

Note that computing $m_{map}$ using Equation (2-14) requires solving an $N_m \times N_m$ matrix problem. However, for reservoir history matching, the number of model parameters $N_m$ is often significantly larger than the quantity of observed data $N_d$, i.e., $N_m \gg N_d$. Therefore, the practical application of Equation (2-14) may be unfeasible. On the other hand, note the following identity:

$$\left(C_m^{-1} + G^T C_d^{-1} G\right)^{-1} G^T C_d^{-1} = C_m G^T \left(GC_m G^T + C_d\right)^{-1}. \qquad (2\text{-}15)$$

Opposite to the left-hand side of Equation (2-15), the right-hand side of the same equation displays a matrix problem of size $N_d \times N_d$, which is more feasible when $N_d \ll N_m$, as explained before. Therefore, one may compute $m_{map}$ as follows:

$$m_{map} = m_{pr} + C_m G^T \left(GC_m G^T + C_d\right)^{-1} \left(d_{obs} - Gm_{pr}\right). \qquad (2\text{-}16)$$

To prove that the minimum of $\mathcal{O}(m)$, $m_{map}$, is unique, one must compute its Hessian function $\mathcal{H}(m)$, as:

$$\mathcal{H}(m) = \nabla\left(\nabla\mathcal{O}(m)\right) = C_m^{-1} + G^T C_d^{-1} G. \qquad (2\text{-}17)$$

As $C_m$ and $C_d$ are covariance matrices, and assumed to be positive-definite, one can attest that $\mathcal{H}(m)$ is also positive-definite [1]. Hence, the minimum of $\mathcal{O}(m)$ is unique. For the linear case, one may write the posterior probability density function of $m$ in terms of the $m_{map}$, as:

$$f(m) = \hat{a}\exp\left\{(m - m_{map})^T C_{map}^{-1} (m - m_{map})\right\}, \qquad (2\text{-}18)$$

where $C_{map}$ is called the posterior covariance matrix of the model parameters, and it is equal to the inverse of the Hessian function of $\mathcal{O}(m)$, as:

$$\begin{aligned} C_{map} &= \mathcal{H}^{-1} \\ &= \left(C_m^{-1} + G^T C_d^{-1} G\right)^{-1} \qquad (2\text{-}19) \\ &= C_m - C_m G^T \left(GC_m G^T + C_d\right)^{-1} GC_m, \end{aligned}$$

where the last equality is obtained by the Sherman-Morisson-Woodbury formula [2]. For the equations computed so far, there were assumed no model errors. However, it can be shown that such a theory is still valid for the case where there are model errors, and they have Gaussian distribution. In this case, the matrix $C_d$ in Equation (2-4) corresponds to the sum of the covariance matrices of measurement errors and model errors, i.e., $C_d = C_{d_1} + C_{d_2}$.

## 2.3
## Minimization for nonlinear problems

When a nonlinear function gives the relation between the model parameters and the theoretical data, i.e., $d = g(m)$, where $g : \mathbb{R}^{N_m} \to \mathbb{R}^{N_d}$, a minimization method is required to compute a minimum of $\mathcal{O}(m)$. Nevertheless, due to the nonlinearity of the function $g$, one cannot prove that such a minimum is unique. Moreover, it is expected that the objective function $\mathcal{O}(m)$ has multiple minima, depending on the initial guess given as input for the minimization method [1]. There are several procedures for minimizing a nonlinear function. Most of them are gradient-based, i.e., one must compute the nonlinear function gradient to obtain a minimizing vector. This section presents three popular techniques to solve the desired nonlinear least-squares minimization problem using line search and trust region algorithms [3].

### Characterizing a local minimum

Some concepts exposed in this section were used before to achieve a formula to compute $m_{map}$. However, for the nonlinear case, we chose to present a local minimum characterization to avoid misunderstanding. A straightforward strategy to determine whether a point $p \in \mathbb{R}^n$ is a local minimum of a specific function $f : \mathbb{R}^n \to \mathbb{R}$ might be investigating all the points in the immediate neighborhood of $p$. However, there are much more efficient methods to discover such property. One of them is given by the analysis of the gradient and the hessian of the function $f$. On the other hand, such research requires f to be smooth and twice continuously differentiable. Let us enunciate Taylor's theorem, which is a start to the minimization methods analysis.

**Teorema 2.1 *(Taylor's Theorem).***

*Suppose $f : \mathbb{R}^n \to \mathbb{R}$ continuously differentiable and $p \in \mathbb{R}^n$. Then, we have that:*

$$f(x + p) = f(x) + \nabla f(x + tp)^T p, \tag{2-20}$$

*for some $t \in (0, 1)$. Moreover, if $f$ is twice continuously differentiable, we have that:*

$$f(x + p) = f(x) + \nabla f(x)^T p + \frac{1}{2} p^T \nabla^2 f(x + tp)p, \qquad (2\text{-}21)$$

*for some $t \in (0, 1)$.*

Using Theorem 2.1, one may investigate conditions to characterize local minimizers of nonlinear functions by checking the gradient and the hessian values at such points. The next theorem exposes a first-order necessary condition regarding the gradient's behavior around a local minimizer.

**Teorema 2.2 *(First-order necessary condition).***

*If $x^*$ is a local minimum and $f$ is continuously differentiable in an open neighborhood of $x^*$, then $\nabla f(x^*) = 0$.*

The proof of Theorem 2.2 can be found in [3]. Also, to derive Equation (2-16), one must invoke Theorem 2.2, as mentioned in the text. The next theorem characterizes the hessian of a twice continuously differentiable function at a local minimizer.

**Teorema 2.3 *(Second-order necessary condition).***

*If $x^*$ is a local minimum and $\nabla^2 f$ is continuous in an open neighborhood of $x^*$, then $\nabla f(x^*) = 0$ and $\nabla^2 f(x^*)$ is positive-semidefinite.*

Again, one may find the proof for Theorem 2.3 in [3]. Just recall that a matrix $M \in \mathbb{R}^{n \times n}$ is said to be positive-definite if $p^T M p > 0$, for all $p \in \mathbb{R}^n$, such that $p \neq 0$, and $M$ is positive-semidefinite if $p^T M p \geq 0$. The next theorem presents sufficient conditions on the derivatives of $f$ that guarantee that $x^*$ is a local minimum.

**Teorema 2.4 *(Second-order sufficient condition).***

*Suppose that $\nabla^2 f$ is continuous in an open neighborhood of $x^*$, and that $\nabla f(x^*) = 0$ and $\nabla^2 f(x^*)$ is positive-definite. Then, $x^*$ is a strict local minimum of $f$.*

The previous theorems present interesting tools to characterize and recognize local minimizers for nonlinear functions. However, finding global minimizers is a difficult task for the nonlinear case. A good strategy for solving minimization problems is to define a convex objective function. In this case, it is possible to prove a useful result, shown in the next theorem and proved in [3].

**Teorema 2.5 *(Global minimizers for convex functions).***

*If $f$ is a convex function, any local minimizer $x^*$ is a global minimizer.*

**Objective function for the nonlinear case**

Section 2.1 displays the desired objective function's derivation to be minimized to solve the history matching problem for the linear case. In this part, we expose the objective function used for nonlinear history matching problems. In fact, there is no significant difference between the one used in the linear and the other used in the nonlinear case. A nonlinear function now gives the relation between the theoretical data and the model parameters. Therefore, we must change the second term in Equation (2-9) by replacing $Gm$ with $g(m)$. Therefore, we can rewrite the objective function $\mathcal{O}(m)$ in the following way:

$$\mathcal{O}(m) = \mathcal{O}_m(m) + \mathcal{O}_d(m), \tag{2-22}$$

with:

$$\mathcal{O}_m(m) = (m - m_{pr})^T C_m^{-1} (m - m_{pr}), \tag{2-23}$$

and:

$$\mathcal{O}_d(m) = (g(m) - d_{obs})^T C_d^{-1} (g(m) - d_{obs}). \tag{2-24}$$

In this case, we denote the maximum *a posteriori* estimate $(m_{map})$ as the following minimization problem:

$$m_{map} = \arg \min_{m \in \mathbb{R}^{N_m}} \mathcal{O}(m) \tag{2-25}$$

Because the objective function in Equations (2-22) to (2-24) is nonlinear, the minimization problem depicted in Equation (2-25) may have multiple solutions. When using gradient-based algorithms, the local minimum will be conditioned to the initial guess. It means that the history matching problem may have various $m_{map}$ estimates.

**Newton minimization method**

The well-known Newton direction to minimize a nonlinear function $f : \mathbb{R}^n \to \mathbb{R}$ can be derived by using Taylor's theorem (Theorem 2.1). Define the following function $m_k(p)$ as the second-order Taylor approximation of $f(x_k+p)$:

$$m_k(p) = f(x_k) + p^T \nabla f(x_k) + \frac{1}{2} p^T \nabla^2 f(x_k) p. \tag{2-26}$$

Let us assume that $\nabla^2 f(x_k)$ is positive-definite. Computing the derivative of $m_k$ and setting it to zero, one must achieve the following formula:

$$p_k^N = - \left( \nabla^2 f(x_k) \right)^{-1} \nabla f(x_k). \tag{2-27}$$

The search direction exposed in Equation (2-27) can be shown to be a descent direction [3]. Therefore, let us compute the Newton search direction to the nonlinear objective function $\mathcal{O}$ (Equation (2-22)). Let $m_k$ be the current guess:

$$\nabla_m \mathcal{O}(m_k) = \nabla_m \mathcal{O}_m(m_k) + \nabla_m \mathcal{O}_d(m_k). \tag{2-28}$$

From Equations (2-23) and (2-24), we obtain that:

$$\nabla_m \mathcal{O}_m(m_k) = C_m^{-1} \left( m_k - m_{pr} \right), \tag{2-29}$$

and

$$\nabla_m \mathcal{O}_d(m_k) = G_k^T C_d^{-1} \left( g(m_k) - d_{obs} \right), \tag{2-30}$$

where $G_k$ is the sensitivity matrix of $g$, evaluated at $m_k$, as follows:

$$G_k = \begin{bmatrix} \frac{\partial g_1(m_k)}{\partial m_1} & \cdots & \frac{\partial g_1(m_k)}{\partial m_{N_m}} \\ \vdots & \ddots & \vdots \\ \frac{\partial g_{N_d}(m_k)}{\partial m_1} & \cdots & \frac{\partial g_{N_d}(m_k)}{\partial m_{N_m}} \end{bmatrix}. \tag{2-31}$$

Now, let us compute the second derivative of the nonlinear objective function $\nabla_m^2 \mathcal{O}(m_k)$, as follows:

$$\begin{aligned} \nabla_m^2 \mathcal{O}(m_k) &= \nabla_m \left( \nabla_m \mathcal{O}_m(m_k) + \nabla_m \mathcal{O}_d(m_k) \right) \\ &= \nabla_m^2 \mathcal{O}_m(m_k) + \nabla_m^2 \mathcal{O}_d(m_k). \end{aligned} \tag{2-32}$$

From Equations (2-29) and (2-30), we obtain that:

$$\nabla_m^2 \mathcal{O}_m(m_k) = C_m^{-1}, \tag{2-33}$$

and

$$\nabla_m^2 \mathcal{O}_d(m_k) = G_k^T C_d^{-1} G + \nabla_m \left( G_k^T \right) C_d^{-1} \left( g(m) - d_{obs} \right). \tag{2-34}$$

Note that the term $\nabla_m \left( G_k^T \right)$ corresponds to the second derivative of the function $g$. Thus, the minimizer direction, given by the Newton procedure, is presented as follows:

$$\delta m_{k+1}^N = - \left( C_m^{-1} + G_k^T C_d^{-1} G + \nabla_m \left( G_k^T \right) C_d^{-1} \left( g(m_k) - d_{obs} \right) \right)^{-1}$$
$$\left( C_m^{-1} \left( m_k - m_{pr} \right) + G_k^T C_d^{-1} \left( g(m_k) - d_{obs} \right) \right). \tag{2-35}$$

Thus, the iterative process given by the Newton method is the following:

$$m_{k+1} = m_k + \mu_{k+1} \delta m_{k+1}^N. \tag{2-36}$$

**Gauss-Newton minimization method**

In the Gauss-Newton method, a simple change is made in the calculation of the second derivative of the objective function $\mathcal{O}$. The term involving a second-order derivative $\nabla_m \left( G_k^T \right)$ is neglected. Instead, the method uses an approximation as follows:

$$\nabla^2 \mathcal{O}(m_k) \approx C_m^{-1} + G_k^T C_d^{-1} G. \tag{2-37}$$

Neglecting the second-order derivative in the computation of the Hessian of $\mathcal{O}$ brings us some advantages. The first is that the time saved by calculating the second derivative of the function $g$ can be significant, hence, transforming the Gauss-Newton search direction faster to estimate. Another critical remark is that an explicit formula for the function $g$ is not available usually. Therefore, the first-order derivative must be computed numerically. Hence, a second-order derivative calculation is not feasible to be estimated. One must note that the Gauss-Newton and the Newton method behave similarly when the current guess $m_k$ is near a local minimum. There are other situations when the Gauss-Newton has advantages over the Newton method, which can be found in [3] and [1]. In reservoir engineering history matching problems, the covariance matrices $C_m$ and $C_d$ are often assumed to be symmetric and positive-definite. Consequently, their inverses are also symmetric and positive-definite. Therefore, one must observe that the Hessian approximation given by the Gauss-Newton algorithm is also positive-definite.

Using the approximation depicted in Equation (2-37), the minimizer direction given by the Gauss-Newton algorithm is given by:

$$\delta m_{k+1}^{GN} = -\left(C_m^{-1} + G_k^T C_d^{-1} G_k\right)^{-1}$$
$$\left(C_m^{-1}\left(m_k - m_{pr}\right)\right) + G_k^T C_d^{-1}\left(g(m_k) - d_{obs}\right)\right). \quad (2\text{-}38)$$

**Levenberg-Marquardt minimization method**

Equation (2-37) shows the Gauss-Newton approximation for the Hessian of the objective function $\mathcal{O}(m_k)$, omitting the second-order derivative. In the Levenberg-Marquardt algorithm, a scalar is added to the Gauss-Newton approximation in the following way:

$$\nabla^2\mathcal{O}(m_k) \approx C_m^{-1} + G_k^T C_d^{-1} G_k + \lambda_k I_m, \quad (2\text{-}39)$$

where $I_m \in \mathbb{R}^{N_m \times N_m}$ is a identity matrix, and $\lambda_k \in \mathbb{R}$, $\lambda_k > 0$ is the Levenberg-Marquardt parameter. Denoting the Gauss-Newton approximation of the Hessian by $H_k$, we can rewrite Equation (2-39) as follows:

$$\nabla^2\mathcal{O}(m_k) = H_k + \lambda_k I_m. \quad (2\text{-}40)$$

Thus, the iterative process of the Levenberg-Marquardt method becomes:

$$\delta m_{k+1}^{LM} = -\left(H_k + \lambda_k I\right)^{-1}$$
$$\left(C_m^{-1}\left(m_k - m_{pr}\right)\right) + G_k^T C_d^{-1}\left(g(m_k) - d_{obs}\right)\right). \quad (2\text{-}41)$$

The parameter controls both the search direction and the step size. If $\lambda$ is large, the Levenberg-Marquardt method will take a small step in the steepest descent direction. Whereas, if $\lambda$ is short, the process will behave similarly to a Gauss-Newton iteration. The choice of $\lambda$ depends on the problem and may be reduced iteration by iteration if the method is obtaining a minimum. Another significant improvement that the Levenberg-Marquardt method gives to the history matching problem is improving the condition number of the Hessian approximation. Recall that the condition number of a matrix is defined as follows:

**Definição 2.6** *Condition number of a matrix.*

*For square matrices A, the condition number of A, $\kappa(A)$, is given by:*

$$\kappa(A) = ||A|| ||A^{-1}||. \tag{2-42}$$

There is a convention that $\kappa(A) = \infty$ for all singular matrices. Note that underlying the regular euclidian norm, one can rewrite the condition number as:

$$\kappa(A) = ||A||_2 ||A^{-1}||_2 = \frac{\sigma_{max}(A)}{\sigma_{min}(A)}. \tag{2-43}$$

Where $\sigma_{max}(A)$ and $\sigma_{min}(A)$ corresponds to the maximum and the minimum singular values of $A$, respectively. Thus, one can write the condition number of the Levenberg-Marquardt Hessian approximation of $\mathcal{O}(m_k)$ as:

$$\kappa(H_k + \lambda_k I) = \frac{\sigma_{max}(H_k) + \lambda_k}{\sigma_{min}(H_k) + \lambda_k}. \tag{2-44}$$

For practical history matching problems, the derivatives are often computed numerically, resulting in errors in the approximation of the matrix $G_k$, hence, rising inaccuracies to the Hessian approximation. The condition number measures how the matrix errors can result in a lousy solution to linear systems. From Equation (2-44), one can note that the Levenberg-Marquardt parameter can fix the bad-conditioning problem when solving Equation (2-41). Moreover, suppose it is possible to compute the singular values of the matrix $H_k$. In that case, one can choose the Levenberg-Marquardt parameter to maintain the condition number fixed within the iterations.

## 2.4
## Ensemble Smoother

The ensemble-smoother (ES) is an ensemble-based method that works as a traditional parameters estimation. It uses an ensemble approximation of the sensitivity matrix to update the model parameters vector. Implementing the ensemble smoother is more comfortable than the previous gradient-based algorithms shown in the last subsection. Moreover, the ES does not require any derivatives, which eases coupling it with any simulator. This subsection will discuss the update equations of the ensemble smoother and the ensemble smoother with multiple data assimilation (ES-MDA), which is an iterative form of the ES.

Suppose that the relation between the model parameters $m \in \mathbb{R}^{N_m}$ and the theoretical data $d \in \mathbb{R}^{N_d}$ is given by a linear function of the form $d = g(m) = Gm$, where $G \in \mathbb{R}^{N_d \times N_m}$ is the sensitivity matrix of $g$. With this

assumption, Section 2.2 presented the maximum *a posteriori* estimate $m_{map}$ as follows:

$$m_{map} = m_{pr} + C_m G^T \left( G C_m G^T + C_d \right)^{-1} \left( d_{obs} - G m_{pr} \right). \tag{2-45}$$

In an ensemble-based method, we create an ensemble of model parameter vectors $\{m_j\}_{j=1}^{N_e}$, where $N_e \in \mathbb{N}^*$ is the ensemble size. Hence, we also produce an ensemble of theoretical data $\{d_j\}_{j=1}^{N_e}$, where $d_j = G m_j$. Thus, we can approximate the model parameters covariance matrix $C_m$ as follows:

$$C_m \approx \frac{1}{N_e - 1} \sum_{j=1}^{N_e} (m_j - \overline{m}) (m_j - \overline{m})^T, \tag{2-46}$$

where $\overline{m} = \frac{1}{N_e} \sum_{j=1}^{N_e} m_j$. Using this approximation, we can estimate the matrices $C_m G^T$ and $G C_m G^T$ using the ensemble members in the following way:

$$\begin{aligned}
C_m G^T &\approx \left( \frac{1}{N_e - 1} \sum_{j=1}^{N_e} (m_j - \overline{m}) (m_j - \overline{m})^T \right) G^T \\
&\approx \frac{1}{N_e - 1} \sum_{j=1}^{N_e} \left( (m_j - \overline{m}) (m_j - \overline{m})^T G^T \right) \\
&\approx \frac{1}{N_e - 1} \sum_{j=1}^{N_e} (m_j - \overline{m}) (G (m_j - \overline{m}))^T \\
&\approx \frac{1}{N_e - 1} \sum_{j=1}^{N_e} (m_j - \overline{m}) \left( d_j - \overline{d} \right)^T \\
&\approx C_{md}.
\end{aligned} \tag{2-47}$$

Note that the previous estimate for $C_m G^T$ requires assuming that $\overline{d} = G\overline{m}$. In a similar process, an estimate for $G C_m G^T$ is:

$$GC_m G^T \approx GC_{md}$$

$$\approx G \left( \frac{1}{N_e - 1} \sum_{j=1}^{N_e} (m_j - \overline{m}) \left( d_j - \overline{d} \right)^T \right)$$

$$\approx \frac{1}{N_e - 1} \sum_{j=1}^{N_e} G \left( (m_j - \overline{m}) \left( d_j - \overline{d} \right)^T \right)$$

$$\approx \frac{1}{N_e - 1} \sum_{j=1}^{N_e} \left( G (m_j - \overline{m}) \right) \left( d_j - \overline{d} \right)^T \tag{2-48}$$

$$\approx \frac{1}{N_e - 1} \sum_{j=1}^{N_e} \left( d_j - \overline{d} \right) \left( d_j - \overline{d} \right)^T$$

$$\approx C_{dd}.$$

Using the ensemble approximations for $C_m G^T$ and $GC_m G^T$, depicted in Equations (2-47) and (2-48), the $m_{map}$ can be rewritten as:

$$m_{map} = m_{pr} + C_{md} \left( C_{dd} + C_d \right)^{-1} \left( d_{obs} - G m_{pr} \right). \tag{2-49}$$

The primary objective of the ES is to update the vector of model parameters $m_j$ first created to build the ensemble approximations for the matrices $C_{md}$ and $C_{dd}$. Therefore, the update equation of the ES can be written as:

$$m_j^a = m_j^f + C_{md}^f \left( C_{dd}^f + C_d \right)^{-1} \left( d_{uc,j} - d_j^f \right), \quad j = 1, \cdots, N_e. \tag{2-50}$$

In the previous equation, the superscript $a$ refers to the analysis step an the superscript $f$ refers to the forward step. The vector $d_{uc,j} \in \mathbb{R}^{N_d}$ corresponds to an unconditional realization of the vector of observed data $d_{obs}$, i.e., $d_{uc,j} = d_{obs} + \epsilon_j$, where $\epsilon_j \sim \mathcal{N}(0, C_d)$. The updated ensemble corresponds to a sample of the posterior probability density function $f(m|d_{obs})$. Therefore, the ES tries to characterize $f(m|d_{obs})$ by estimating the first two statistical moments: the mean and variance.

The study of [4] showed that the ES update process is similar to applying a single Gauss-Newton iteration with full-step size and using the same ensemble approximation for the sensitivity matrix of the data. As most of the history matching problems are nonlinear, one single Gauss-Newton iteration is not expected to produce reasonable estimates for the model parameters. Thus, one cannot expect the ES to deliver trustable estimates for the vector of model

parameters. Many iterative forms of the ES have been proposed to improve its primary formulation. One of them is the ensemble smoother with multiple data assimilation (ES-MDA), which the idea is to assimilate the same data multiple times with an inflated data error covariance matrix.

**Ensemble smoother with multiple data assimilation**

The ES-MDA update equation is given by:

$$m_j^{k+1} = m_j^k + C_{md}^k \left( C_{dd}^k + \alpha_{k+1} C_d \right)^{-1} \left( d_{uc,j}^k - d_j^k \right), \tag{2-51}$$

where $j = 1, \cdots, N_e$, $k = 1, \cdots, N_a$. The variable $N_a \in \mathbb{N}^*$ refers to the number of assimilations. The set $\{\alpha_k\}_{k=1}^{N_a} \subset \mathbb{R}$ refers to the data error inflation factors, that must satisfy the following condition:

$$\sum_{k=1}^{N_a} \frac{1}{\alpha_k} = 1. \tag{2-52}$$

The requirement exhibited in Eq. (2-52) guarantees the ES-MDA to sample the posterior probability density function $f(m|d_{obs})$ correctly. A direct consequence of Equation (2-52) is that $\alpha_k \geq 1, \forall k = 1, \cdots, N_a$. Another useful procedure of writing the covariance matrices approximated from the ensemble is in the following way:

$$C_{md}^k = \Delta M^k \left( \Delta D^k \right)^T, \tag{2-53}$$

and

$$C_{dd}^k = \Delta D^k \left( \Delta D^k \right)^T, \tag{2-54}$$

where

$$\Delta M^k = \frac{1}{\sqrt{N_e - 1}} \left[ m_1^k - \overline{m}^k, \cdots, m_{N_e}^k - \overline{m}^k \right], \tag{2-55}$$

and

$$\Delta D^k = \frac{1}{\sqrt{N_e - 1}} \left[ d_1^k - \overline{d}^k, \cdots, d_{N_e}^k - \overline{d}^k \right]. \tag{2-56}$$

The matrix $\Delta M^k$ may be understood as the ensemble approximation of the square root of the model parameters' covariance matrix $C_m$ at the step $k$. Hence, we obtain that:

$$C_m \approx \Delta M^k \left(\Delta M^k\right)^T, \tag{2-57}$$

and

$$C_m^{1/2} \approx \Delta M^k. \tag{2-58}$$

Assuming the relationship between the model parameters and the data to be linear, we obtain that:

$$\Delta D^k = G \Delta M^k. \tag{2-59}$$

Because the variables contained in the vector of model parameters $m$ may have different dimensions and also the variance between them may be significant, [5] defined the dimensionless sensitivity matrix of the data pre multiplying the matrix $G$ by $C_d^{-1/2}$ and post multiplying the matrix $G$ by $C_m^{1/2}$, as follows:

$$G_D = C_d^{-1/2} G C_m^{1/2}. \tag{2-60}$$

As the matrix $\Delta M^k$ is an approximation of the square root of $C_m$, we obtain that:

$$G_D^k \approx C_d^{-1/2} G \Delta M^k. \tag{2-61}$$

From Equation (2-59), we obtain that:

$$G_D^k \approx C_d^{-1/2} \Delta D^k. \tag{2-62}$$

Note that the approximation of $G_D^k$, depicted in Equation (2-62), has dimensions $N_d \times N_e$. It is useful to rewrite Equation (2-51) in terms of the dimensionless sensitivity matrix $G_D^k$. Therefore, consider the matrix $C = C_{md}\left(C_{dd} + \alpha C_d\right)^{-1}$. We will hide the super and subscripts $k$ in this part for convenience. Using Equations (2-53) and (2-54), we can rewrite $C$ in the following way:

$$C = \Delta M \left(\Delta D\right)^T \left(\Delta D \left(\Delta D\right)^T + \alpha C_d\right)^{-1}. \tag{2-63}$$

Assuming that the data-error matrix $C_d$ is symmetric, and noting that $C_d = C_d^{1/2} I_{N_d} C_d^{1/2}$, where $I_{N_d} \in \mathbb{R}^{N_d \times N_d}$ is the identity matrix, we obtain that:

$$
\begin{aligned}
C &= \Delta M \left(\Delta D\right)^T \left(C_d^{1/2} \left(C_d^{-1/2} \Delta D \left(\Delta D\right)^T C_d^{-1/2} + \alpha I_{N_d}\right) C_d^{1/2}\right)^{-1} \\
&= \Delta M \left(\Delta D\right)^T C_d^{-1/2} \left(C_d^{-1/2} \Delta D \left(\Delta D\right)^T C_d^{-1/2} + \alpha I_{N_d}\right)^{-1} C_d^{-1/2}.
\end{aligned}
\tag{2-64}
$$

Noticing that $G_D^T = \Delta D^T C_d^{-1/2}$, directly from Equation (2-64), we achieve that:

$$C = \Delta M G_D^T \left(G_D G_D^T + \alpha I_{N_d}\right)^{-1} C_d^{-1/2}. \tag{2-65}$$

Replacing the matrix $C$ computed using Equation (2-65) in Equation (2-51), we obtain that:

$$\delta m_j^{k+1} = \left(G_D^k\right)^T \left(G_D^k \left(G_D^k\right)^T + \alpha_{k+1} I_{N_d}\right)^{-1} y_j^k, \tag{2-66}$$

where $j = 1, \cdots, N_e, k = 1, \cdots, N_a$, $\delta m_j^{k+1} = \left(\Delta M^k\right)^{\dagger} \left(m_j^{k+1} - m_j^k\right) \in \mathbb{R}^{N_e}$ is the dimensionless vector of update parameters, where the superscript $\dagger$ refers to the pseudo-inverse of the matrix $\Delta M^k$, and the vector $y_j^k = C_d^{-1/2} \left(d_{uc,j}^k - d_j^k\right) \in \mathbb{R}^{N_d}$ referring to the dimensionless vector of observed data.

# 3
# Influences of the inflation factors generation in the main parameters of the ES-MDA

The generation of the ES-MDA inflation factors has been the focus of several recent studies. Concurrently, recent researches have shown a relationship between inflation factors and the final ensemble estimates quality. They have also suggested techniques to generate these factors based on methods derived from the discrepancy principle. However, a procedure to efficiently generate ES-MDA inflation factors remains an open problem. Additionally, the studies diverge on what regularization method suffices to produce ES-MDA inflation factors that provide optimal final results. Therefore, this chapter presents an investigation of the generation of ES-MDA inflation factors. Two main paths will be investigated: selecting them constant, equal to the number of assimilations, and in a geometrically decreasing order. When selecting them geometrically, two techniques will be used to generate the first inflation factor: the regular discrepancy principle and a Levenberg-Marquardt regularization scheme. The main objective of this study is to examine the error propagation during the multiple data assimilation of the ES-MDA and the estimates' quality, considering only the generation of the inflation factors. Moreover, we numerically analyze their influence on the ES-MDA ensemble size and number of assimilations and how their choices affect the ES-MDA performance. The results presented in this chapter were published in Silva *et al.* (2021) [6].

## 3.1
## Introduction

The generation of ES-MDA inflation factors is the focus of several recent studies. Le *et al.* (2016) [7] noted that selecting constant inflation factors, equal to $N_a$, may lead to under and overcorrections in the final ensemble estimates when $N_a = 8$ and $N_a = 16$. Therefore, they proposed two automatic procedures to select such factors. The first one is based on the average objective function, in a way that the correction of the following step is not larger than a previously selected threshold. The second method uses a regularization scheme for Levenberg-Marquardt algorithms proposed by Hanke (1997) [8]. Similar to the first method of Le *et al.* (2016) [7], Emerick (2016) [9] also proposed a

method to generate the inflation factors for ES-MDA based on the average data mismatch function.

Rafiee and Reynolds (2017) [10] proposed to select the ES-MDA inflation factors geometrically in decreasing order, where the first factor was generated using a procedure derived from the Levenberg-Marquardt regularizing scheme of Hanke (1997) [8]. The following factors were selected geometrically in decreasing order. Their results suggest that the geometric generation of these factors enhances the ES-MDA final results. The studies of Evensen (2018) [11] and Emerick (2019) [12] supported this observation. Rafiee and Reynolds (2017) [10] also claim that it is possible to yield good ES-MDA final results even with a low number of assimilations, but with proper inflation factors selection. This claim was investigated in the work of Silva *et al.* (2021) [13]. Using the method of Rafiee and Reynolds (2017) [10], Silva *et al.* (2021) [13] showed that increasing $N_a$ from 4 to 8 did not produce significant improvements in the final results when estimating the reservoir skin properties. Emerick (2019) [12] proposed to select the last inflation factor *a priori*, and the previous ones are generated geometrically in increasing order. In contrast with the study of Rafiee and Reynolds (2017) [10], Emerick (2019) [12] tests the resulting first inflation factor using a standard discrepancy principle function [14]. Recently, Silva *et al.* (2021) [15] proposed a new method to compute the ES-MDA inflation factors by computing the first and the last factors using the formula of Rafiee and Reynolds (2017) [10]. The other inflation factors are selected geometrically in decreasing order.

The ES-MDA performance is intimately linked to its three primary parameters: the ensemble size, the inflation factors, and the number of assimilations. Additionally, the ES-MDA vector of model parameters update is a linear combination of the right singular vectors of the average dimensionless sensitivity matrix $G_D$ [5, 10]. Thus, changing any of these three main parameters results in a modified linear combination to compute the ES-MDA update. Hence, it is possible to analyze their influence on the ES-MDA performance. There are various researches regarding the ES-MDA achievements considering both the ensemble size [16, 13] and the number of assimilations [17, 10, 12]. Therefore, this study focus on examining the generation of the inflation factors.

The motivation of the present study comes from the research of Tavakoli and Reynolds (2010) [18]. They presented an in-depth investigation of the effects of the singular values of $G_D$ in minimizing the uncertainty. They concluded that the singular values are responsible for the uncertainty reduction in the vector of model parameters. Moreover, they claim that the largest singular values may produce an optimal parametrization in updating the vector

of model parameters when the singular values decay fast. The observations of Tavakoli and Reynolds (2010) [18] can also be joint with the study of Wang *et al.* (2010) [19], which claims that minimal singular values can result in notable errors when resolving matrix problems.

Considering the similarities between the ES-MDA and the Levenberg-Marquardt method, we can also mention the study of Shiranji and Emerick (2016) [20]. They present a comparison between the performance of the Gauss-Newton and the Levenberg-Marquardt methods applied to parameter estimation. They concluded that the search direction of the Levenberg-Marquardt algorithm has its main components in the order of the right singular vectors of $G_D$ that correspond to the largest singular values. In the Gauss-Newton algorithm, they showed that all right singular vectors have similar weights in the vector of the model update, even those corresponding to the smallest singular values. As a result, the final estimates of Gauss-Newton present more significant errors compared to the ones obtained by the Levenberg-Marquardt method. As noted by Rafiee and Reynolds (2017) [10], the update equation of ES-MDA has the same structure as the Levenberg-Marquardt algorithm.

Although several studies proposed different methods, efficiently generating the ES-MDA inflation factors remains an open problem. Thereby, this study provides a mathematical analysis of the inflation factors selection in the ES-MDA model parameters update and their influence in the ensemble size $N_e$ and the number of assimilations $N_a$. This analysis will consider the singular values and vectors of the average dimensionless sensitivity matrix $G_D$. The first inflation factor will be generated using the schemes suggested by Emerick (2019) [12], and Rafiee and Reynolds (2017) [10]: the discrepancy principle [14] and the regularizing system of Hanke (1997) [8]. The following inflation factors will be selected geometrically in decreasing order. It is essential to mention that only the first assimilation step will be analyzed.

This study's main contribution and novelty is presenting a method to numerically investigate the generation of the ES-MDA inflation factors before starting the data assimilation process. The suggested investigation approach is one of the main differences between the present study and the ones available in the literature investigating the inflation factors generation [10, 12]. Furthermore, we offer a numerical procedure to assess whether the inflation factor generation is suitable for the chosen ensemble size or the problem itself. The results and discussions presented in this study state an accurate connection of the inflation factors, the ES-MDA ensemble size $N_e$ and the number of assimilations $N_a$, and the singular values of the matrix $G_D$. The proposed method also allows the proper selection of inflation factors to

improve the ES-MDA final results and the ES-MDA performance.

This work shows that the vector of ES-MDA model parameters update is contained in the space spanned by the right singular vectors of the matrix $G_D$. Hence, if the singular values are decreasing ordered, it is desired that the first right singular vectors, corresponding to the most significant singular values, have more influence in the model parameters update. The numerical examples exhibited in this study demonstrate that the inflation factors have a massive impact on each coefficient's determination in the linear combination of the right singular vectors of $G_D$ that computes the ES-MDA model parameters update. Moreover, we show that selecting proper inflation factors may mitigate small singular values' effects on the multiple data assimilation. Thus, significantly diminishing error propagation within the ES-MDA iterations. Also, suitably selecting the inflation factors may also alleviate other ES-MDA current problems, such as ensemble collapse, spurious correlations due to small ensembles, and under or overcorrections in the final estimates. Finally, we intend to conclude which regularization method is recommended to attain optimal ES-MDA outcomes and performance.

This chapter is organized as follows: Section 3.2 presents the analytical formulation to compute the coefficient of each right singular vector of $G_D$ in the ES-MDA vector of model parameters update; Section 3.3 presents a brief theoretical background of regularization for nonlinear inverse problems, the methods used to generate the first inflation factor, and how to generate the following ones geometrically in decreasing order; finally, Section 3.4 presents the investigation of the influences of the inflation factors in the ES-MDA main parameters and performance.

## 3.2
## Singular Values Analysis

This section analyzes the effects of the average dimensionless sensitivity matrix $G_D$ in the ES-MDA update equation depicted in Equation (2-66). This analysis will be performed considering the singular values and the singular vectors of $G_D$. Furthermore, we aim to achieve a procedure to measure their impacts on the ES-MDA performance and final results. The assimilation process is not on focus in this section. Therefore, we will eliminate the super or subscripts $k$. Consider the singular value decomposition [2] of the matrix $G_D$ as follows:

$$G_D = U\Sigma V^T, \tag{3-1}$$

where the matrix $U \in \mathbb{R}^{N_d \times N_d}$ refers to the left singular vectors of the matrix $G_D$. The matrix $V \in \mathbb{R}^{N_e \times N_e}$ refers to the right singular vectors of the matrix $G_D$. The matrix $\Sigma \in \mathbb{R}^{N_d \times N_e}$ is a diagonal matrix where each diagonal entry refers to the singular values of $G_D$. Substituting Equation (3-1) in Equation (2-66), we obtain that:

$$\delta m_j = \left(U\Sigma V^T\right)^T \left(U\Sigma V^T \left(U\Sigma V^T\right)^T + \alpha I_{N_d}\right)^{-1} y_j. \tag{3-2}$$

Note that the matrices $U$ and $V$ are orthogonal, i.e., $UU^T = U^T U = I_{N_d}$ and $VV^T = V^T V = I_{N_e}$. Therefore, Equation (3-2) can be rewritten as:

$$\delta m_j = V\Sigma^T U^T \left(U\Sigma\Sigma^T U^T + \alpha I_{N_d}\right)^{-1} y_j. \tag{3-3}$$

Using the fact that matrix $U$ is orthogonal, we can rewrite Equation (3-3) as:

$$
\begin{aligned}
\delta m_j &= V\Sigma^T U^T \left(U\left(\Sigma\Sigma^T + \alpha I_{N_d}\right)U^T\right)^{-1} y_j \\
&= V\Sigma^T U^T \left(U^T\right)^{-1} \left(\Sigma\Sigma^T + \alpha I_{N_d}\right)^{-1} (U)^{-1} y_j \\
&= V\Sigma^T U^T U \left(\Sigma\Sigma^T + \alpha I_{N_d}\right)^{-1} U^T y_j \\
&= V\Sigma^T \left(\Sigma\Sigma^T + \alpha I_{N_d}\right)^{-1} U^T y_j.
\end{aligned}
\tag{3-4}
$$

From Equation (3-4), it follows straightforward that:

$$\delta m_j = \sum_{i=1}^{r} \left(\frac{\sigma_i}{\sigma_i^2 + \alpha} u_i^T y_j\right) v_i, \tag{3-5}$$

where $r \in \mathbb{N}^*$ is the rank of $G_D$ and $\{\sigma_i\}_{i=1}^{r} \subset \mathbb{R}$ refers to the set of singular values of $G_D$. Equation (3-5) states that the ES-MDA updated vector of model parameters $\delta m_j$ is contained in the space spanned by the singular vectors $v_i$, i.e., $\delta m_j$ can be described as a linear combination of $v_i$. Moreover, Equation (3-5) also presents an analytical formula to compute the individual coefficients of each singular vector $v_i$ for computing $\delta m_j$. Define the function $t_j(\sigma, u, \alpha) : \mathbb{R}^{N_d+2} \to \mathbb{R}$ that calculates each desired coefficient for the vector $v_i$ as:

$$t_j(\sigma, u, \alpha) = \frac{\sigma}{\sigma^2 + \alpha} u^T y_j. \tag{3-6}$$

Function $t_j$ is computed for each ensemble member $m_j$. Also, its variables are the singular values $\sigma_i$ of matrix $G_D$, the value of the inflation factor $\alpha$ at that assimilation step, and the left singular vectors $u_i$ of matrix $G_D$. Using the function $t_j$, exposed in Equation (3-6), we can rewrite Equation (3-5) as:

$$\delta m_j = \sum_{i=1}^{r} t_j \left( \sigma_i, u_i, \alpha \right) v_i. \tag{3-7}$$

All the variables of the function $t_j$ are selected in terms of the dimensionless sensitivity matrix $G_D$, except by the inflation factor $\alpha$, which the user selects. The only restriction is the one depicted in Equation (2-52). Although several studies proposed different methods to generate such factors efficiently, it remains an open problem deciding which one produces the best ES-MDA results. Equation (3-6) presents an analytical formula to compute each singular vector's coefficients in the vector of ES-MDA model parameters update. This formula is not entirely unfamiliar in the ensemble-based methods literature [10]. However, it has never been used to evaluate the quality of the ES-MDA inflation factors effects on the vector of model parameter updates, considering the matrix $G_D$. The study of Emerick (2019) [12] presents an analysis of the geometric generation of the ES-MDA inflation factors. However, in this study, we offer a different investigation of the production of such elements. Instead of comparing the ES-MDA results with varying inflation factors and deciding which one achieved the best outcome, we present a mathematical argument that numerically demonstrates our conclusion of which method may achieve the best result before starting data assimilation. Therefore, this strategy to assess the ES-MDA inflation factors generation considering the singular values and vectors of $G_D$ is one of this study's main contributions and novelty.

One can find a related strategy in the study of Shiranji and Emerick (2016) [20]. They present a procedure to compute each coefficient of right singular vectors of $G_D$ in the update equation of the Levenberg-Marquardt and Gauss-Newton algorithms. As a result, based on Tavakoli and Reynolds (2010) [18], and Wang *et al.* (2010) [19], we may generate inflation factors such that the singular vectors $v_i$ corresponding to the largest singular values $\sigma_i$ have more influence in the vector of model parameters update $\delta m_j$. In other words, we may select inflation factors $\{\alpha_k\}_{k=1}^{N_a}$ such that function $t_j$ computes the highest values for the greatest $\sigma_i$. This choice may lead to smaller error propagation in the multiple data assimilation of the ES-MDA. Moreover, it is desired that the singular vectors corresponding to the smallest singular values have minor importance in $\delta m_j$. Therefore, one can decide which choice for the inflation factor attends these requirements by applying Equation (3-6) to the

prior ensemble.

## 3.3
## Inflation factor generation from the discrepancy principle

This section presents a brief background of regularization for nonlinear inverse problems [21] and how we can link this theory to the ES-MDA update equation [10, 15]. The concept of regularization is crucial to yielding a significant analysis of the influence of the inflation factors selection in the ES-MDA performance from a mathematical perspective. Moreover, it was the base of the studies of Le *et al.* (2016) [7], Rafiee and Reynolds (2017) [10], Emerick (2019) [12], and Silva *et al.* (2021) [15] to present novel methods to generate ES-MDA inflation factors efficiently.

### Regularization for nonlinear inverse problems

Consider $x \in \mathbb{R}^{N_x}$ and $y \in \mathbb{R}^{N_y}$, where a nonlinear function $f : \mathbb{R}^{N_x} \to \mathbb{R}^{N_y}$ gives the relation between them, i.e., $y = f(x)$. To trace a parallel, reservoir history matching consists of solving the inverse problem of finding $x$ given $y$. However, there is no guarantee of the existence nor uniqueness of such a solution for a nonlinear function $f$. As a result, one may consider the following nonlinear minimization problem:

$$x^* = \min_{x \in \mathbb{R}^{N_x}} ||f(x) - y||_2^2. \tag{3-8}$$

The problem depicted in Equation (3-8) is widely studied in the mathematical literature and can be solved using different approaches. As an example, we can mention line search and trust-region techniques [3]. Nevertheless, in this study, we adopt a simple strategy of using a linearization of the function $f$ around the current guess $x^k$ obtained by the first-order Taylor series expansion, as follows:

$$f(x) \approx f(x^k) + A_k(x - x^k), \tag{3-9}$$

where

$$A_k = \begin{bmatrix} \frac{\partial f_1(x_k)}{\partial x_1} & \cdots & \frac{\partial f_1(x_k)}{\partial x_{N_x}} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_{N_y}(x_k)}{\partial x_1} & \cdots & \frac{\partial f_{N_y}(x_k)}{\partial x_{N_x}} \end{bmatrix}. \tag{3-10}$$

Considering $\hat{y}^k = f(x^k) - y$ and $\hat{x}^k = x - x^k$, we obtain that:

$$\hat{x}^{k+1} = \min_{\hat{x}} ||A_k\hat{x}^k - \hat{y}^k||^2. \tag{3-11}$$

Due to the nonlinearity of the function $f$ and the lack of a unique solution, one may search for the resolution of the Tikhonov regularized inverse problem [21], presented as follows:

$$\hat{x}^{k+1} = \min_{\hat{x}} \left\{ ||A_k\hat{x}^k - \hat{y}^k||^2 + \alpha_{k+1}||\hat{x}^k||^2 \right\}, \tag{3-12}$$

where $\alpha_{k+1} > 0$. Equation (3-12) has the following answer:

$$\hat{x}^{k+1} = A_k^T(A_kA_k^T + \alpha_{k+1}I)^{-1}\hat{y}^k. \tag{3-13}$$

The solution to the regularized problem depicted in Equation (3-13) is similar to the dimensionless ES-MDA update equation, presented in Equation (2-66), with $A_k = G_D^k$, $\hat{x}^{k+1} = \delta m_j^{k+1}$, and $\hat{y}^k = y_j^k$. Thus, note that Equation (2-66) displays a solution for the following regularized problem:

$$\delta m_j^{k+1} = arg\min_{\delta \overline{m}^k} \left\{ ||G_D^k\delta m_j^k - y_j^k||^2 + \alpha_{k+1}||\delta m_j^k||^2 \right\}, \tag{3-14}$$

**Inflation factor generation**

The inflation factors $\{\alpha_k\}_{k=1}^{N_a}$, will be generated geometrically in decreasing order, such that $\alpha_1$ will be computed using two popular procedures in the ensemble-based methods literature, basis for the studies of Le *et al.* (2016) [7], Rafiee and Reynolds (2017) [10] and Emerick (2019) [12]. Therefore, we consider that $k = 0$, referring to the first ES-MDA assimilation step. Consider the intricacy of updating the ensemble mean, $\overline{m}^0$, in the following minimization problem:

$$\delta\overline{m}^1 = arg\min_{\delta\overline{m}^0} \left\{ ||G_D^0\delta\overline{m}^0 - y^0||^2 + \alpha_1||\delta\overline{m}^0||^2 \right\}, \tag{3-15}$$

where $\delta\overline{m}^1 = (\Delta M^0)^{\dagger}(\overline{m}^1 - \overline{m}^0)$ and $y^0 = C_d^{-1/2}\left(d_{obs} - \overline{d}^0\right)$. As a result, we have that an answer to the problem depicted in Equation (3-15), given by Equation (2-66) with $k = 0$, can be written as:

$$\delta\overline{m}^1 = \left(G_D^0\right)^T \left(G_D^0\left(G_D^0\right)^T + \alpha_1 I_{N_d}\right)^{-1} y^0. \tag{3-16}$$

One can note that $\alpha_1$ in Equation (3-16) works as a regularization parameter to the regularized problem exposed in Equation (3-15) [14, 10]. This parameter may control the quality of the solution. If it is too little, the numerical answer may present high sensitivity and be excessively unstable. On the other hand, if $\alpha_1$ is exceedingly high, then the approximation error may be large [22]. Therefore, to guarantee that the regularization problem converges, one must produce $\alpha_1$ using the noise level $\eta$ [23, 10]. A usual assumption is that the data mismatch is higher than the noise level correlated with the data, i.e.:

$$||y^0|| > \eta. \tag{3-17}$$

In other words:

$$||C_d^{-1/2}\left(d_{obs} - \overline{d}^0\right)|| > \eta, \tag{3-18}$$

where the noise level correlated with the data can be described as:

$$\eta^2 \equiv ||C_d^{-1/2}\left(d_{obs} - d_{true}\right)||^2, \tag{3-19}$$

where $d_{true} = Gm_{true}$ and $m_{true}$ is the true vector of model parameters such that $d_{obs} = Gm_{true} + \epsilon$, with $\epsilon$ referring to the model error. Therefore, Equation (3-19) can be written as:

$$\eta^2 \equiv ||C_d^{-1/2}\epsilon||^2. \tag{3-20}$$

It is reasonable to consider that the noise vector $\epsilon$ assumes multivariate Gaussian distribution. Hence, $\eta^2$ attends a $\chi^2$ distribution with $N_d$ degrees of freedom. Accordingly, a rough approximation may be $\eta = \sqrt{N_d}$ [23]. Two rules will be applied to generate the regularization parameter to be used as an inflation factor for the first ES-MDA assimilation step. They are the following:

1) The discrepancy principle (DP):

$$||G_D^0 \delta\overline{m}^1 - y^0|| = \eta\tau, \tag{3-21}$$

where $\tau \geq 1$;

2) The Hanke condition (HC):

$$\rho^2||y^0||^2 \leq \alpha_1^2|| \left( G_D^0 \left( G_D^0 \right)^T + \alpha_1 I_{N_d} \right)^{-1} y^0||^2, \qquad (3\text{-}22)$$

where $\rho \in (0,1)$.

For Equation (3-21), we used $\tau = 1$, and for Equation (3-22), we used $\rho = 0.5$. For both equations, we used the traditional euclidean norm [2]. The regular discrepancy principle [14], exposed in Equation (3-21), was used in Emerick (2019) [12] to test the resulting value of $\alpha_1$. The rule exposed in Equation (3-22), proposed by Hanke (1997) [8], was used by Rafiee and Reynolds (2017) [10] to derive a formula to compute $\alpha_1$. Therefore, these rules will be used to compute $\alpha_1$ in the analysis presented by this study. The other inflation factors will be generated geometrically in decreasing order as follows:

$$\alpha_{k+1} = \gamma \alpha_k, \qquad (3\text{-}23)$$

with $\gamma \in (0,1)$. The number of assimilations $N_a$ will also be selected *a priori*. Thus, to guarantee that Equation (2-52) is satisfied, the value of $\gamma$ must be computed as the solution of the problem $f(\gamma) = 0$, where $\gamma \in (0,1)$ and $f$ is determined as follows:

$$f(\gamma) = \gamma^{-N_a} + \alpha_1 \left( 1 - \gamma^{-1} \right) - 1. \qquad (3\text{-}24)$$

In addition, the standard implementation of the ES-MDA with $\alpha_k = N_a$, $k = 1, \cdots, N_a$, will be tested. According to Rafiee and Reynolds (2017) [10], it is desired a low number of assimilations $N_a$ to maintain moderate the computational cost. They suggest using $N_a$ from 4 to 8. Thus, we compare the inflation factors generated using the conditions exhibited in Equations (3-21) and (3-22) with $\alpha_k = 4$, $k = 1, \cdots, 4$, and $\alpha_k = 8$, $k = 1, \cdots, 8$. Another significant remark is that both studies of Emerick (2019) [12] and Rafiee and Reynolds (2017) [10] derived formulas to compute $\alpha_1$ based on the procedures depicted in Equations (3-21) and (3-22). Instead, we will strictly apply such conditions to generate $\alpha_1$.

## 3.4
## Case study

In this section, we present an analysis of the inflation factor generation in the main parameters of the ES-MDA. We numerically examine their effects on the ensemble size $N_e$ and the number of assimilations $N_a$. We also present an investigation of the vector of model parameters update of ES-MDA $\delta\overline{m}^1$

for the different inflation factors selection. This analysis will be done in the history matching problem of water flooding a two-dimensional reservoir.

**Problem description**

The reservoir model is squared, with a total length of 1575 m in each direction. It is discretized in a 63×63×1 grid, where each gridblock has dimensions 25m×25m×15m. The reference log-permeability field was generated using a spherical covariance function of correlation length of 40 gridblocks. The prior estimate of the log-permeability field is equal to 5, with prior variance equals 1. Figure 3.1 displays the reference log-permeability field and the location of each well. These pieces of information were used to create ensembles of different sizes in this study. The vector of model parameters $m$ consists only of the gridblock log-permeability of each gridblock. The reservoir contains thirteen wells, nine of them are producing, and four are water injection wells. The producing wells operate at constant bottom-hole pressure (BHP) of 275 kgf/cm$^2$. The water injection wells run at a continuous BHP of 325 kgf/cm$^2$. The simulation time consists of 3600 days, with measurements every 150 days. The observed data consists of the water and oil rate (WR and OR) of the producing wells and the water injection rate of the water injection wells. For the synthetic measurements, it was added a Gaussian error of 3%. In this study, the matrix $C_d$ is diagonal, with every entry corresponding to the square of the standard deviation used to generate the errors in the observed measurements. We used a distance-based covariance localization during all data assimilation applied directly to the Kalman gain matrix [24]. We do not use any subspace inversion procedure if we do not intend to neglect singular vectors' impact corresponding to small singular values. In fact, the main objective is to analyze the inflation factors selection and their impact on the singular vectors components of $\delta \overline{m}^1$.

**Singular values**

With the information described in the previous subsection, we create four different ensembles of different sizes: $N_e = 25$, $N_e = 50$, $N_e = 100$, and $N_e = 500$. Figures 3.2 to 3.5 display the computed singular values of matrix $G_D$ estimated using the prior ensemble with different sizes. Note that the $G_D$ has dimension $N_d \times N_e$. Therefore, $rank(G_D) \leq min(N_d, N_e)$. One may note that the singular values present a fast decay rate for all cases. It indicates that only a few singular vectors $v_i$ (see Equation (3-5)) corresponding to the largest

Figure 3.1: Reference log-permeability field.



Figure 3.2: Computed singular values for $G_D$ estimated from the prior ensemble with $N_e = 25$.

singular values $\sigma_i$ may create an optimal parametrization for uncertainty quantification purposes [18]. Thus, as we do not use any parametrization or subspace inversion procedure, we can investigate the inflation factors selection and which available technique can neglect the effects of singular values corresponding to the smallest singular values.

**Ensemble size analyzis**

Knowing the singular values of the matrix $G_D$, it is possible to compute the coefficients $t_j$ (Equation (3-6)) of each singular vector $v_i$ in the vector of model parameters update calculated by the ES-MDA. Considering the problem of updating the mean log-permeability field in the first assimilation step of ES-MDA, $\overline{m}^0$, one may analyze the ensemble size's effects in each coefficient of the components of the vector $\delta\overline{m}^1$. We refer to the vector of coefficients of $v_i$

Figure 3.3: Computed singular values for $G_D$ estimated from the prior ensemble with $N_e = 50$.

Figure 3.4: Computed singular values for $G_D$ estimated from the prior ensemble with $N_e = 100$.

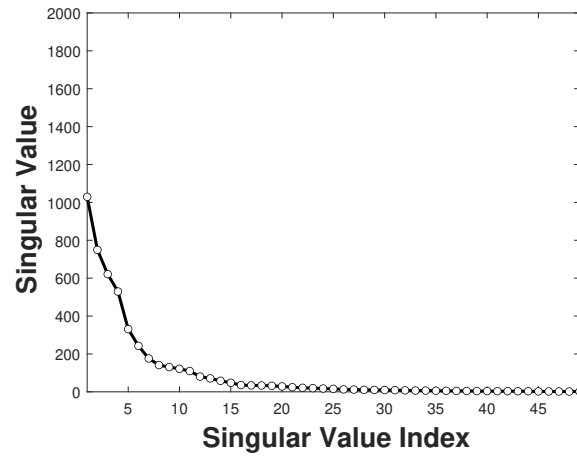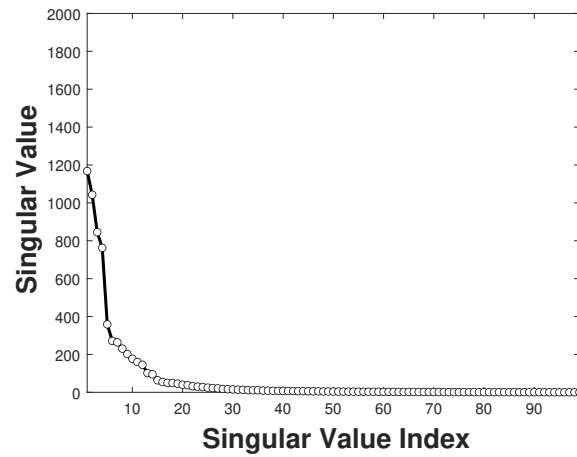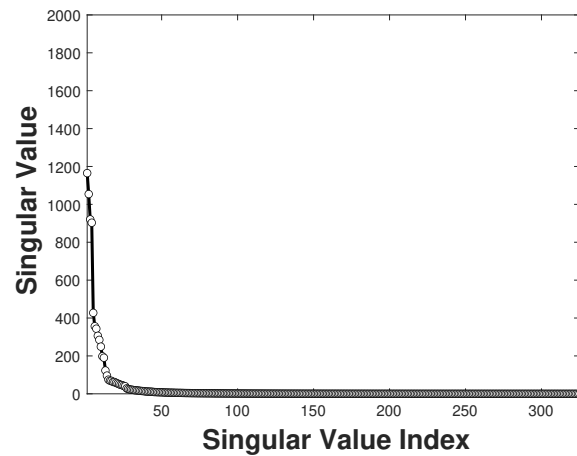Figure 3.5: Computed singular values for $G_D$ estimated from the prior ensemble with $N_e = 500$.

| $N_e$ | Discrepancy principle (DP) | Hanke condition (HC) |
|---|---|---|
| 25 | 59.36 | 341,294.00 |
| 50 | 799.02 | 417,960.60 |
| 100 | 839.61 | 726,932.02 |
| 500 | 1,916.34 | 997,308.71 |

Table 3.1: Computed $\alpha_1$ using Equations (3-21) and (3-22) for different ensembles.

in $\delta\overline{m}^1$ as $\bar{t}$.

The problem of investigating the ensemble size and their resulting estimates is well studied in the ensemble-based methods literature [16, 13]. Their primary approach is running the ES-MDA with varying ensembles sizes $N_e$ and comparing their results. This strategy is not wrong and provides useful insight into the ensemble size influences in data assimilation. One of the main reasons that a small ensemble cannot offer a fair characterization of the posterior probability density function may be the poor approximation of the covariance matrices $C_{md}$ and $C_{dd}$ in Equations (2-53) and (2-54). It causes spurious correlations of components of the vector $m$ that are far away from the observations. Moreover, the matrix $G_D$ estimated from a small ensemble is seriously rank deficient, limiting the degrees of freedom to assimilate data.

In this study, we provide additional information in the ensemble size investigation for the ES-MDA. Based on the work of Tavakoli and Reynolds (2010) [18] and Wang (2010) [19], observing the computed singular values in Figures 3.2 to 3.5, we may infer a supplementary conclusion on why small ensembles yield unsatisfactory final results for ES-MDA. For computing the coefficients $\bar{t}$, we use four different inflation factors. The first procedure is selecting them equal to $N_a$. Therefore, we select $\alpha_1 = 4$ and $\alpha_1 = 8$. The second approach is selecting $\alpha_1$ using the discrepancy principle (Equation (3-21)). The third and last method is using the scheme of Hanke (1997) [8] to estimate the first inflation factor (Equation (3-22)). Table 3.1 displays the computed values for $\alpha_1$ using each method for the different ensembles.

Figures 3.6 to 3.9 display the computed absolute values for the coefficients $\bar{t}$ in Equation (3-6). One may note that for the ensemble with $N_e = 25$, the right singular vectors $v_i$ corresponding to the smallest singular values have more influence on the model update $\delta\overline{m}^1$ when selecting the first inflation factor equal to $N_a$ or using the discrepancy principle. As noted by Wang (2010) [19], small singular values are responsible for significant errors. Thus, these three methods may provide inadequate estimates for the posterior ES-MDA ensemble. As the ensemble size grows, this issue is resolved. Shiranji and emerick (2016) [20] also checked this issue for Gauss-Newton algorithms.

Figure 3.6: Computed coefficients $\bar{t}$ (Equation (3-6)) for $v_i$ for $N_e = 25$ in linear and semi-log scales.



Figure 3.7: Computed coefficients $\bar{t}$ (Equation (3-6)) for $v_i$ for $N_e = 50$ in linear and semi-log scales.



Figure 3.8: Computed coefficients $\bar{t}$ (Equation (3-6)) for $v_i$ for $N_e = 100$ in linear and semi-log scales.

Figure 3.9: Computed coefficients $\bar{t}$ (Equation (3-6)) for $v_i$ for $N_e = 500$ in linear and semi-log scales.

Analyzing the coefficients of the components of $\delta\overline{m}^1$ estimated using the inflation factor computed utilizing the condition of Hanke (1997) [8], one may note that the vectors $v_i$ corresponding to the largest $\sigma_i$ are the ones that have more influence in the ES-MDA update, even when the ensemble size is small. It might reduce complications when updating model parameters. However, the issues regarding lousy estimates for the covariance matrices $C_{md}$ and $C_{dd}$ and the limited degrees of freedom of $G_D$ are still valid for small ensembles. Thus, we claim that small ensembles also forces the ES-MDA to compute an update vector that strongly considers right singular vectors corresponding to small singular values when badly selecting the inflation factor. This result supplies an interesting scheme. Using Equation (3-6), one can verify whether the inflation factor is proper for the ES-MDA for the chosen ensemble size before initiating data assimilation. Suppose the selected inflation factor computes an update vector that considers the smallest singular values majorly. In that case, it might end in poor ES-MDA outcomes, as explained previously. Therefore, the user can change the inflation factor selection before running the ES-MDA method. This process might be useful for applications problems when the ensemble size cannot be large due to computational or time obstacles.

**Singular values analysis**

The problem of generating the inflation factors for the ES-MDA has been on focus recently [7, 12, 10]. However, which method is the most adequate for data assimilation remains an open problem . Figures 3.6 to 3.9 also analyze which method may generate inflation factors for achieving optimal ES-MDA results. Following the same technique described in the previous subsection, one may investigate the effects of the right singular vectors $v_i$ in the ES-MDA model parameters update. It is expected that the vectors $v_i$ corresponding to

the largest $\sigma_i$ have more influence in the model parameters update vector for the ES-MDA to achieve the best results.

For all ensemble sizes, it is noteworthy that generating the first inflation factor using the method of Hanke (1997) [8] resulted in the weaker influence of singular vectors related to small singular values compared to the other inflation factors selection. On the other hand, this procedure produces considerably small coefficients $\bar{t}$ for the first singular vectors $v_i$. It might result in soft and limited update vectors $\delta\overline{m}^1$. Therefore, the first assimilation of ES-MDA using the procedure of Hanke (1997) [8] to generate $\alpha_1$ might not result in significant changes in the model parameters vector $m^0$. Differently, this aspect may be worthy for the multiple data assimilation of ES-MDA, where the prior ensemble contains a massive quantity of error approximation. Thus, the vector of model parameters is softly corrected during the ES-MDA assimilation steps.

When selecting inflation factors equal to the number of assimilations $N_a$, singular vectors corresponding to small singular values have substantial influence when the ensemble size is little. This issue is relatively resolved when the ensemble size increases. In the example shown in Figures 3.6 to 3.9, the mentioned problem is only resolved in the case where $N_e = 500$. On the other hand, larger ensembles may result in higher computational cost and time consumption, which hinders the assimilation process. Another solution to such a problem might be increasing the number of assimilations. In the example in Figures 3.6 to 3.9, increasing $N_a$ to 8 could not reach more reliable outcomes for any cases. Thus, one may select $N_a$ larger than 8. Still, it brings the same difficulty of computational cost and time. Therefore, we claim that equally selected inflation factors may result in poor approximations for model parameters in the ES-MDA with relatively small ensembles. This claim is supported by this study's observations, and the practical examples available in the ES-MDA literature, which equally selected inflation factors often presents under and overshooting in the model parameters [7, 10, 12].

The discrepancy principle appears to be a reasonable choice for generating inflation factors for ES-MDA. However, when the ensemble is small, as displayed in the case with $N_e = 25$, it presents the same issue as the one observed with equal inflation factors. An ensemble of size $N_e = 50$ might not resolve this issue thoroughly, but the problem is way smaller. The effects of singular vectors corresponding to smaller singular values might be neglected only with an ensemble of 100 members or higher. In the example of the ensemble with $N_e = 500$, singular vectors of index greater than 50 are almost entirely ignored.

One may note that a nearly large ensemble, e.g., $N_e = 500$, results in

smaller effects of singular vectors corresponding to small singular values, even when the inflation factors are selected equal to $N_a$. Moreover, also the singular vectors related to the largest singular values have moderate effects in the vector of model updates, e.g., smaller $\bar{t}$, which cannot be observed in other cases with smaller ensembles. This observation tells us that the inflation factor selection almost becomes pointless if the ensemble size is relatively big. This conclusion is in accordance with the definition of ES-MDA, which provides a correct sample of the posterior probability density function $f(m|d_{obs})$ if Equation (2-52) is satisfied and the ensemble size goes to infinity.

Overall, the inflation factor generated using the method of Hanke (1997) [8] seems to produce relatively small updates in the vector of model parameters due to the small coefficients of the first singular vectors. Using the discrepancy principle might escape this problem, but the ensemble size must not be small. When the inflation factors are equal to $N_a$, there are two possibilities that might yield good results: a large ensemble or a large number of assimilations. In both cases, computational cost and time consumption are increased. Now, let us analyze the update vector computed by each inflation factor selection and ensemble size.

**Model parameters update analysis**

Figures 3.10 to 3.13 display the vector of mean model parameters update $\delta\overline{m}^1$, computed for the different first inflation factor selection and different ensemble sizes. The figures display the vector $\delta\overline{m}^1$ computed using all available singular values and compare with the same vector calculated using only the singular values responsible for 90% of the sum of all singular values. In Figures 3.10 to 3.13, the top first column figure refers to the reference field, whereas the bottom first column figure corresponds to the prior mean. The following columns correspond to the different computations of the ES-MDA vector of update: $\alpha_1 = 4$, $\alpha_1 = 8$; $\alpha_1$ computed using the discrepancy principle, and $\alpha_1$ computed using the scheme of Hanke (1997) [8], respectively. The first line corresponds to computing the ES-MDA update vector by using all available singular values. In contrast, the second line corresponds to computing the ES-MDA update vector using the singular values responsible for 90% of the sum of all singular values.

As expected, when the ensemble is small, such as the example with $N_e = 25$, the vector $\delta\overline{m}^1$ might be unable to produce an appropriate update to the prior mean, especially when the inflation factor is equal to $N_a$ ($N_a = 4$ and $N_a = 8$). Although the update seems to be changing the prior mean

log-permeability towards the reference field, such an update's amplitude is excessive, which may cause under or overcorrections in the final ensemble mean. Moreover, the difference between $\delta\overline{m}^1$ computed adopting all singular values and the same vector utilizing just the largest ones is significant. This observation is explained by the effects of singular vectors related to small singular values in the update (see Figures 3.6 to 3.9), which may cause a considerable quantity of errors in the estimate. Although the parametrized update also seems unable to generate a reasonable estimation of the reference field, its amplitude is way smaller, which may alleviate under or overcorrections to the ES-MDA estimates.

When the ensemble grows to $N_e = 50$, the update becomes smaller for all inflation factors selection. Nevertheless, they still seem to be inefficient in producing reasonable estimates for the mean log-permeability field compared to the reference field. The difference between the full vector $\delta\overline{m}^1$ and the parametrized one is smaller in this case. However, it is still noticeable the effects of small singular values. This issue is relatively resolved when the ensemble size expands to $N_e = 100$, visually making an update vector that addresses the prior mean log-permeability towards the reference field. In this case, the effects of the inflation factors become very apparent. One may note that when the inflation factor is selected equal to the number of assimilations $N_a$, the update's amplitude is larger, computing the vector $\delta\overline{m}^1$ that clearly deviates from the parametrized one. This problem does not happen when we select $\alpha_1$ using the discrepancy principle or the Hanke condition, which states that the updates in these two cases are computed mostly in the direction of singular vectors corresponding to the largest singular values. This observation also holds for the ensemble with 500 members. However, the largest ensemble ($N_e = 500$) produces even smoother and smaller updates, which is useful for data assimilation. For all cases, the update made using the inflation factor selected using the method of Hanke (1997) [8] presented the smaller update, as expected by the previous discussion.

### ES-MDA results

Although the previous subsections aimed to analyze the effects of the inflation factors selection on the ensemble size $N_e$ and on the vector of update parameters $\delta\overline{m}^1$, a similar investigation to acquire information about their influence in the number of assimilations $N_a$ requires running the ES-MDA with different values for $N_a$. In this study, we use $N_a = 4$ and $N_a = 8$. These values are chosen because we wish to maintain a low number of

Figure 3.10: Update vector $\delta\overline{m}^1$ for each inflation factor selection for the case $N_e = 25$ and comparison with a low-order parametrization. The top-first column figure refers to the reference field, whereas the bottom-first column figure corresponds to the prior mean. The following columns correspond to the different computations of the ES-MDA vector of update: $\alpha_1 = 4$, $\alpha_1 = 8$; $\alpha_1$ computed using the discrepancy principle, and $\alpha_1$ computed using the scheme of Hanke (1997) [8], respectively. The first line corresponds to computing the ES-MDA update vector by using all available singular values. In contrast, the second line corresponds to computing the ES-MDA update vector using the singular values responsible for 90% of the sum of all singular values.



Figure 3.11: Update vector $\delta\overline{m}^1$ for each inflation factor selection for the case $N_e = 50$ and comparison with a low-order parametrization. Figures here have the same meaning as in Figure 3.10.



Figure 3.12: Update vector $\delta\overline{m}^1$ for each inflation factor selection for the case $N_e = 100$ and comparison with a low-order parametrization. Figures here have the same meaning as in Figure 3.10.
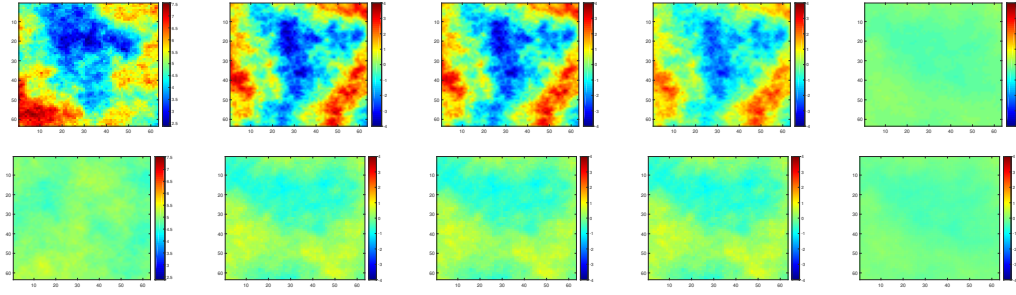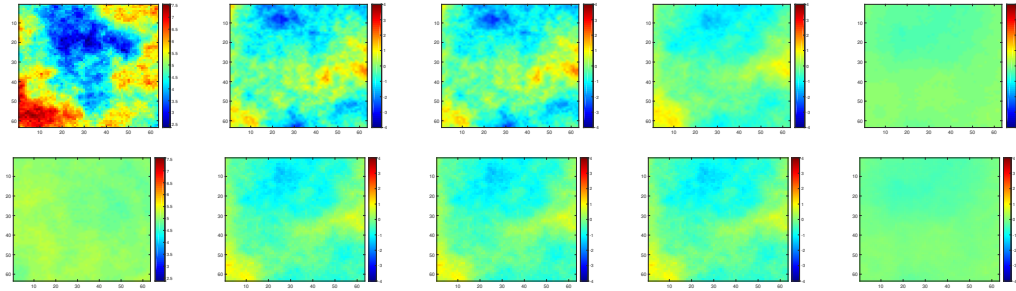
Figure 3.13: Update vector $\delta\overline{m}^1$ for each inflation factor selection for the case $N_e = 500$ and comparison with a low-order parametrization. Figures here have the same meaning as in Figure 3.10.

assimilations to avoid increased computational cost and time consumption [10]. The values of $\alpha_1$ computed for each ES-MDA run is exposed in Table 3.1. The other inflation factors will be computed geometrically in decreasing order, as displayed in Equation (3-23). As this study focuses exclusively on parameter approximation, we present only the model parameters' ES-MDA results.

When running the ES-MDA with the inflation factors equal to $N_a$, it will be referred to as 4x-EQL for $N_a = 4$ and 8x-EQL for $N_a = 8$. For the case where the first inflation factor is computed using the discrepancy principle, it will be referred to as 4x-DP for $N_a = 4$ and 8x-DP for $N_a = 8$. When using the scheme of Hanke (1997) [8] to generate the first inflation factor, it will be referred to as 4x-HC for $N_a = 4$ and 8x-HC for $N_a = 8$. In Figures 3.14 to 3.21, the first three columns correspond to the first three posterior ensemble members $m_1$, $m_2$, and $m_3$; the fourth column refers to the posterior ensembles means $\overline{m}$; the last column corresponds to the reference field. Additionally in those Figures, the first line corresponds to selecting $\alpha_1 = N_a$. The second line corresponds to selecting $\alpha_1$ using the discrepancy principle. The third line corresponds to selecting $\alpha_1$ using the scheme of Hanke (1997) [8]. We also display the root mean squared error (RMSE) for each ES-MDA final ensemble. The RMSE is computed using the following formula:

$$RMSE_j = \left( \frac{1}{N_m} \sum_{k=1}^{N_m} \left( m_{true,k} - m_{j,k} \right)^2 \right)^{1/2}, \tag{3-25}$$

Investigating the ES-MDA performance in the ensemble with $N_e = 25$ in Figure 3.14, one may note that 4x-EQL could not produce a reasonable estimate for the reference field. Moreover, the mean ensemble shows that the members do not vary, exposing the limited degrees of freedom in data assimilation due to the rank deficiency of $G_D$. This problem is well known in the ensemble-based data assimilation literature as *ensemble collapse* [25],

Figure 3.14: First three posterior ensemble members and mean for ES-MDA with $N_e = 25$ and $N_a = 4$. The first three columns correspond to the first three posterior ensemble members $m_1$, $m_2$, and $m_3$; the fourth column refers to the posterior ensembles means $\overline{m}$; the last column corresponds to the reference field. The first line corresponds to selecting $\alpha_1 = N_a$. The second line corresponds to selecting $\alpha_1$ using the discrepancy principle. The third line corresponds to selecting $\alpha_1$ using the scheme of Hanke (1997) [8].



Figure 3.15: First three posterior ensemble members and mean for ES-MDA with $N_e = 25$ and $N_a = 8$. Figures here have the same meaning as in Figure 3.14.

Figure 3.16: First three posterior ensemble members and mean for ES-MDA with $N_e = 50$ and $N_a = 4$. Figures here have the same meaning as in Figure 3.14.



Figure 3.17: First three posterior ensemble members and mean for ES-MDA with $N_e = 50$ and $N_a = 8$. Figures here have the same meaning as in Figure 3.14.



Figure 3.18: First three posterior ensemble members and mean for ES-MDA with $N_e = 100$ and $N_a = 4$. Figures here have the same meaning as in Figure 3.14.

Figure 3.19: First three posterior ensemble members and mean for ES-MDA with $N_e = 100$ and $N_a = 8$. Figures here have the same meaning as in Figure 3.14.



Figure 3.20: First three posterior ensemble members and mean for ES-MDA with $N_e = 500$ and $N_a = 4$. Figures here have the same meaning as in Figure 3.14.



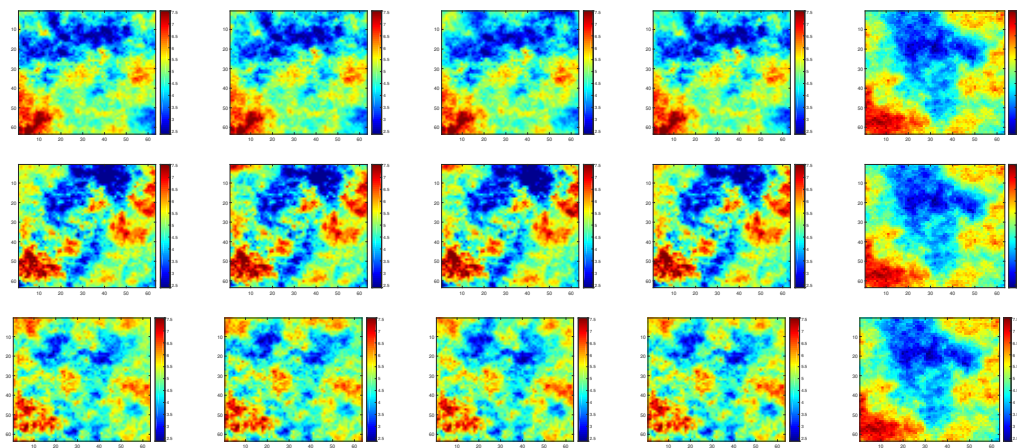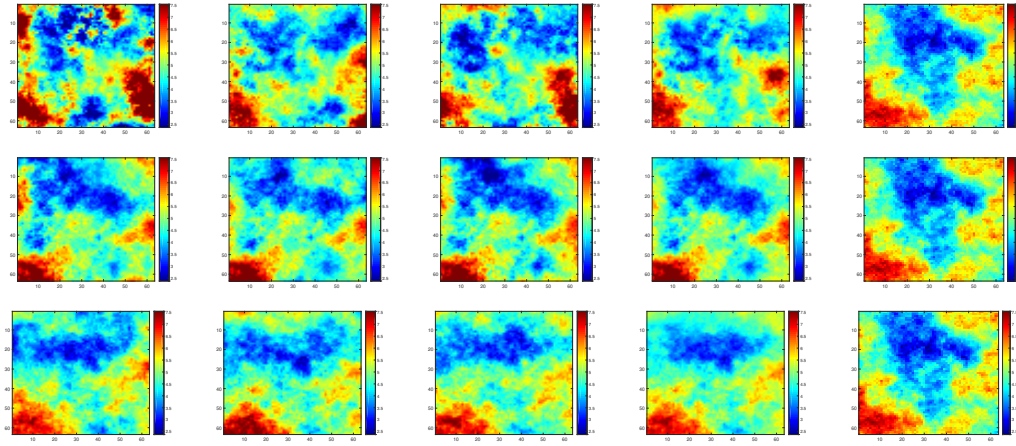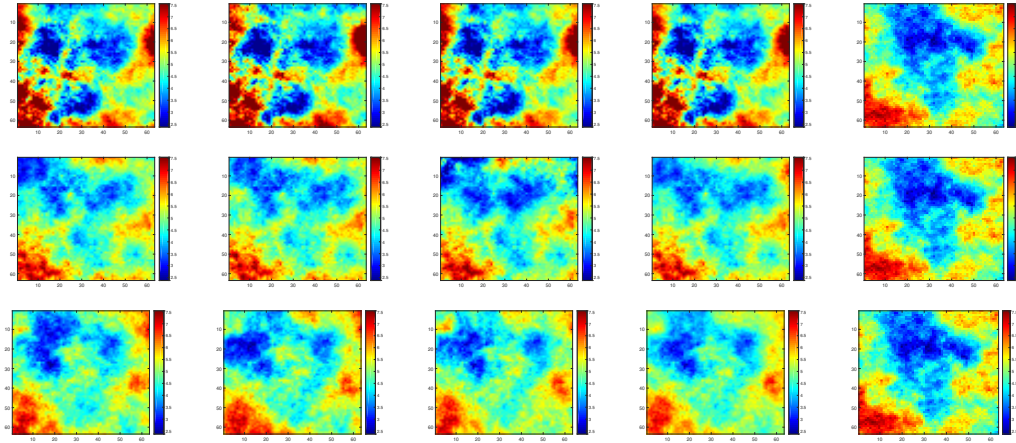Figure 3.21: First three posterior ensemble members and mean for ES-MDA with $N_e = 500$ and $N_a = 8$. Figures here have the same meaning as in Figure 3.14.
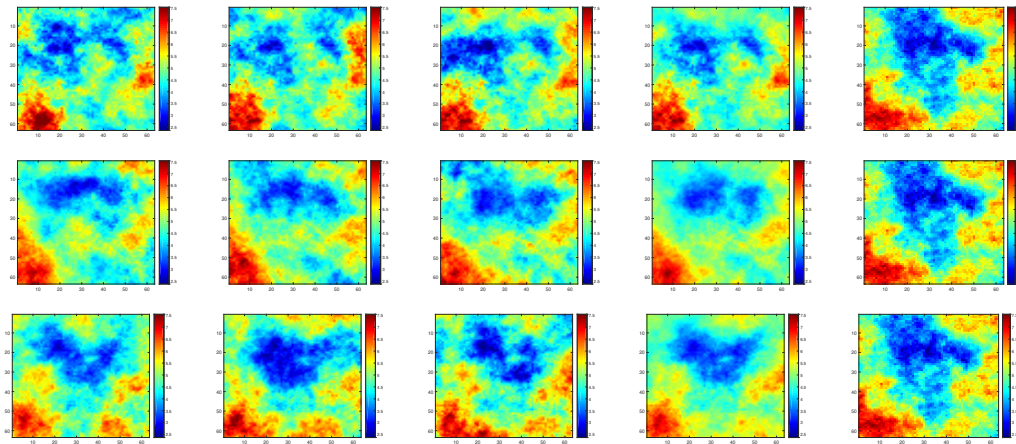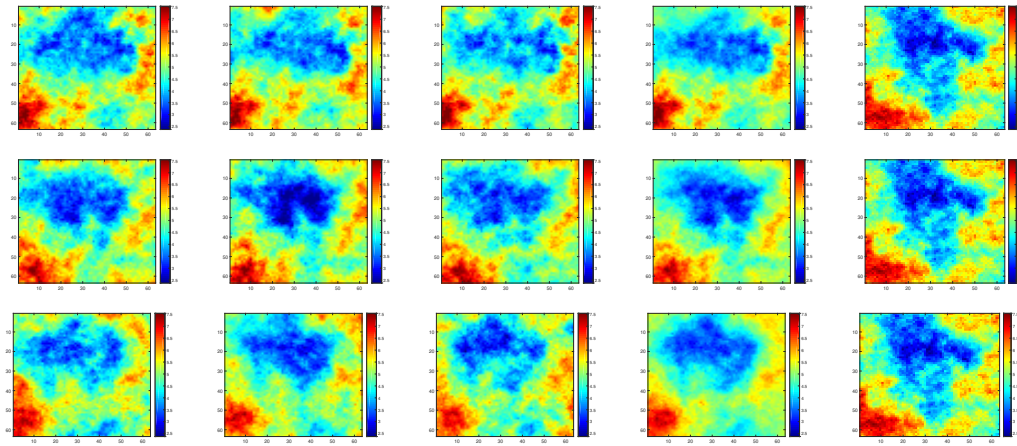
where the difference between the ensemble members is exceedingly minimal that they practically do not diversify. The inaccurate estimate is also because of the poor approximations of the covariance matrices $C_{md}$ and $C_{dd}$ and the influences of singular vectors corresponding to small singular values in the vector of model parameters update, as depicted in Figures 3.6 to 3.9. The use of a regularization method to determine the first inflation factor partially diminishes these problems. However, the ensemble means also indicate that the ensemble does not vary due to the pre-mentioned difficulty. The case where $N_a = 8$, in Figure 3.15, shows that a higher number of assimilations slightly alleviate the problems obtained with reduced iterations. However, the same troubles explained before are still valid for this case.

As the ensemble size increases to $N_e = 50$, the problem of limited degrees of freedom in data assimilation is marginally lessened due to the increased number of members in the first ensemble. This fact can be noted in the mean of each posterior ensemble in Figures 3.16 and 3.17. However, the quality of the estimates is still low. The use of the regularizing method to generate the inflation factors improved the ES-MDA estimates, particularly the technique of Hanke (1997) [8]. For both cases where $N_e = 25$ and $N_e = 50$, selecting the first inflation factor using Equation (3-22) provided a reasonable posterior ensemble even with all the small ensembles issues. Moreover, the results of 4x-HC and 8x-HC for these two cases were relatively similar, showing the method's consistency, which is demonstrated by the observations in Figures 3.6 to 3.9.

An ensemble of size $N_e = 100$ is commonly used in the ensemble-based methods literature [26, 7], capable of yielding reliable estimates for the history matching problem. An excessive augmented ensemble, such as in the example with $N_e = 500$, shows that the ES-MDA with all different inflation factors selection produced high-grade and comparable results. It can be explained by analyzing the effects of singular vectors corresponding to small singular values in the vector of model parameters update, strongly mitigated with a vast ensemble (see Figures 3.6 to 3.9). The good approximations of the matrices $C_{md}$ and $C_{dd}$, along with the nearly full-rankness of the matrix $G_D$ also improve the ES-MDA results for a relatively large ensemble. As a comparison, we present the computed RMSE for each posterior ensemble obtained by ES-MDA in Figures 3.22 to 3.25. In these figures, the black boxes correspond to the region between the percentiles 25 and 75; the solid red lines inside of the boxes correspond to the ensemble median; the solid red diamonds correspond to the ensemble mean; the red whiskers correspond to the ensemble outliers; the solid horizontal black lines correspond to the ensemble range.

Figure 3.22: Computed RMSE for prior and posterior ES-MDA ensembles for $N_e = 25$.



Figure 3.23: Computed RMSE for prior and posterior ES-MDA ensembles for $N_e = 50$.



Figure 3.24: Computed RMSE for prior and posterior ES-MDA ensembles for $N_e = 100$.

Figure 3.25: Computed RMSE for prior and posterior ES-MDA ensembles for $N_e = 500$.

The computed RMSE supports our claim that when the ensemble is small, using the method of Hanke (1997) [8] to generate the first inflation factor produces an excellent result with a reduced number of iterations. This observation is in accordance with the study of Hanke (2010) [27], which proves that the Hanke condition is of optimal order for Levenberg-Marquardt algorithms. Another critical remark is that selecting inflation factors equal to $N_a$ has proved unable to provide trustable approximations with little ensembles. One may note that when the ensemble is relatively large, e.g. $N_e = 500$, the inflation factors selection does not alter the ES-MDA final results significantly. Figures 3.22 to 3.25 also support the claim of Rafiee and Reynolds (2017) [10], stating that small number of assimilations with propper inflation factors selection yields good results for the ES-MDA. It can be noticed because for the cases where $N_e = 25$, $N_e = 50$, and $N_e = 100$, the 4x-EQL and 8x-EQL produced the poorest results compared to the other ES-MDA implementations.

One may note that with a relatively big ensemble, e.g., $N_e = 500$, the problem of ensemble collapse alleviates. It is because of the expanded degrees of freedom in data assimilation, where the matrices $C_{md}$ and $C_{dd}$ are almost full-rank for this problem. Moreover, the discrepancy in the results of 4x-EQL and 8x-EQL is not significant, exhibiting the high quality of approximating the average dimensionless sensitivity matrix $G_D$. Furthermore, as exposed in this study, large ensembles also avoid the effects of small singular values in the ES-MDA update equation. This fact brings an interesting discussion about a well-known conjecture in the ES-MDA literature that increasing the number of assimilations $N_a$ yields better results. This might very well be true, as exposed in the vast ES-MDA published studies [17, 7, 12, 28]. However, the theoretical and numerical results presented in this study indicate that larger ensembles are more desired than larger assimilation steps. It can be explained because large

ensembles naturally mitigate the effects of singular vectors corresponding to small singular values in the vector of model parameters update. Thus, providing optimal ES-MDA estimates. This fact can be observed by noting that doubling the number of assimilations $N_a$ for the case with $N_e = 500$ did not generate notable enhancements.

## 3.5
## Conclusions

This study investigates the influences of the inflation factors in the main parameters of the ES-MDA, such as the number of assimilations $N_a$ and the ensemble size $N_e$. Also, we investigate their influence on the vector of ES-MDA model parameters update, which is strictly associated with the method's performance. The primary objective is to conclude an optimal procedure to generate ES-MDA inflation factors considering only the model parameter estimations and the singular values and vectors of matrix $G_D$. Therefore, we examine the ES-MDA update equation in terms of the dimensionless sensitivity matrix's singular value decomposition to determine the effects of inflation factors in the model parameters.

The ES-MDA estimates can be composed as a linear combination of the dimensionless sensitivity matrix's right singular vectors. Moreover, we present an analytical formula to compute each right singular vector coefficient in the ES-MDA update vector. This procedure is the start point of this study. The proposed formula has all variables computed based on the matrix $G_D$, except by the inflation factor at that assimilation step, which is selected by the user. This technique explicitly computes each coefficient of the linear combination that determines the vector of model parameters update. Therefore, it is possible to assess the effects of each inflation factor generation on those coefficients.

Previous studies showed that minimal singular values could generate meaningful errors when resolving matrix problems. Moreover, if the singular values of $G_D$ have a fast decay rate, one can show that a parametrization containing just a few of the largest ones is optimal for uncertainty reducing purposes. Therefore, we demonstrate that the inflation factor generation substantially influences the weights of the right singular vectors of $G_D$ on the ES-MDA update equation. Consequently, we can determine which approach is optimal to generate inflation factors such that right singular vectors corresponding to small singular values have minimal effects on data assimilation.

Our numerical results show that if the ensemble size is relatively small, e.g., $N_e = 25$ or $N_e = 50$, the method of Hanke (1997) [8] provides optimal results even with a reduced number of assimilations, such as $N_a = 4$. It is

explained because such a method computes an ES-MDA update that partially neglects singular vectors corresponding to small singular values, resulting in diminishing the error propagation in data assimilation. It is also noticeable that selecting inflation factors equal to the number of assimilation $N_a$ produces an update that highly considers small singular values. Thus, enabling error propagation in the ES-MDA iterations with relatively small ensembles. When the ensemble is relatively large, e.g., $N_e = 500$, the inflation factors selection is not crucial to yield good approximations. It is explained because, as the ensemble increases, singular vectors corresponding to small singular values are neglected naturally. Finally, we also conclude that larger ensembles are more decisive than bigger assimilation steps in the ES-MDA. This fact holds not only because of the better ensemble approximations of covariance matrices but also because larger ensembles naturally push the ES-MDA almost to neglect small singular values in the update process. On the other hand, small ensembles force the ES-MDA to compute an update vector majorly in the direction of singular vectors corresponding to small singular values.

**4**

# A new procedure for generating data covariance inflation factors for ES-MDA

This study aims to introduce a new method for generating the data covariance inflation factors for ES-MDA. The main motivation of the presented study comes from the observation in Silva *et al.* (2021) [6] that applying the scheme of Hanke (1997) [8] to generate the first inflation factor resulted in the best ES-MDA performance. In the new method, the first inflation factor is generated using a Levenberg-Marquardt regularizing scheme. The last inflation factor is set by a parameter that limits its magnitude, computed using the singular values of the dimensionless sensitivity matrix estimated from the prior ensemble. As a result, the method computes the correct number of data assimilations that produces inflation factors such that the sum of their inverse is equal to one, as required by ES-MDA. It is shown through a synthetic two-dimensional water flooding history matching problem that the proposed methodology achieves both better model parameter match and data match with a smaller number of assimilations than the methods available in the literature. The results presented in this chapter were published in Silva *et al.* (2021) [15].

## 4.1
### Introduction

In the ES-MDA, the problem of selecting all the inflation factors equal to the number of data assimilations is that if it is not enough to provide reasonable results, one may need to restart the assimilation process with a higher number of assimilations, which increases computational time. Another problem is that it may cause overcorrections of the model parameters [12, 10, 7]. In a study presented by Le *et al.* (2016) [7], selecting all the inflation factors equal to the number of assimilations led to overshooting in the final permeability and porosity field when implementing the ES-MDA with 8 and 16 assimilation steps. Therefore, they proposed two adaptive procedures for selecting inflation factors. The first one, ES-MDA-RS, is based on the average data mismatch function, determining inflation factors that do not result in great changes in model parameters at each assimilation step. The second method, ES-MDA-

RLM, is based on a study proposed by Iglesias and Dawson (2013) [29], that uses a regularizing scheme proposed by Hanke (1997) [8] to compute the inflation factors at each assimilation step. Although the methods proposed by Le et al. (2016) [7] improved the performance of the ES-MDA, it often requires many iterations, which may hinder the application for large-scale problems [10]. Emerick (2016) [9] proposed a method to select the inflation factors adaptatively, based on the average data mismatch function multiplied by a factor smaller than one. The termination criterion is when the sum of the inverse of the inflation factors becomes equal to one. Emerick (2016) [9] tested the method in a real field case. However, the results were not significantly better than the standard ES-MDA.

Motivated by the studies of Le *et al.* (2016) [7] and Iglesias (2015) [30], Rafiee and Reynolds (2017) [10] presented an ES-MDA algorithm where the first inflation factor is computed based on the regularization condition for Levenberg-Marquardt algorithms of Hanke (1997) [8]. One of their objectives was to determine the number of assimilation *a priori* to maintain low computational cost and assimilation time. They state that selecting the number of assimilations from 4 to 8, with propper inflation factors, produce good results for the final ensemble. They also declared that selecting the inflation factors geometrically in decreasing order is good to improve ES-MDA results. In the study of Silva *et al.* (2021) [13], increasing the number of assimilation from 4 to 8 did not improve the ES-MDA outcomes significantly when using the method of Rafiee and Reynolds (2017) [10]. However, the first inflation factor computed by Rafiee and Reynolds (2017) [10] is often large. It implies that when the number of assimilations is close to 4, the last inflation factor is close to 1, i.e., the last assimilation is almost full-step.

Emerick (2019) [12] presented an ES-MDA algorithm that determines the inflation factor for the last iteration. Then, the algorithm computes the previous inflation factors geometrically in increasing order such that the sum of the inverse of them is equal to one. Computed value for the first inflation factor is tested against the Morozov Discrepancy Principle (MDP) [14]. If it is not satisfied, the method increases the number of assimilations in one and recomputes all the inflation factors again. However, recent studies [7, 10] observed that the ES-MDA update equation has a similar structure as a Levenberg-Marquardt algorithm. Therefore, the method of Hanke (1997) [8] can be applied to select the inflation factors Le *et al.* (2016) [7]. Moreover, as presented in Hanke (2010) [27], the scheme for nonlinear inverse problems of Hanke (1997) [8] is of optimal order, i.e., it provides optimal accuracy for such algorithms.

In this study, we present a novel method for generating the inflation factors for ES-MDA. We apply the analytical formula derived by Rafiee and Reynolds (2017) [10] to determine a lower bound for the first inflation factor and use the same equation to propose the inflation factor's computation for the last assimilation step using the singular values of the average sensitivity matrix estimated from the prior ensemble. In the calculation of the last inflation factor, we limit its magnitude to a previously selected threshold. Although we can mathematically prove such a procedure only for the linear-Gaussian case, the numerical examples presented in this study demonstrate that the proposed method is adequate for the nonlinear case. The other inflation factors are computed geometrically in decreasing order. The proposed algorithm then computes the correct number of data assimilations that produce inflation factors satisfying ES-MDA requirements. In addition, as the method of Emerick (2019) [12] has no efficient procedure to compute the last inflation factor, the proposed analytical formula introduced by this study can be used in the algorithm of Emerick (2019) [12].

The motivations for using the proposed method are the following: (i) if the number of assimilation runs is small (e.g., close to four), the last assimilation of the method of Rafiee and Reynolds (2017) [10] will be an approximately full-step update; (ii) there is no efficient way to compute the last inflation factor in the method proposed by Emerick (2019) [12]; (iii) as noted by Iglesias (2015) [30] and Rafiee and Reynolds (2017) [10], the regularization parameter of Hanke (1997) [8] usually decreases with the iterations; therefore, we may simulate this decreasing behavior in a geometric progression, computing the first and last inflation factor using the formula derived by Rafiee and Reynolds (2017) [10] *a priori*; (iv) the method of Hanke (1997) [8] has been proven to be of optimal order for Levenberg-Marquardt algorithms [27]. As it has a similar form as the ES-MDA update equation [7, 10], we conjecture that it may produce better outcomes than the discrepancy principle for generating ES-MDA inflation factors; v) finally, the study of Silva *et al.* (2021) [6] shows that generating the first inflation factor using the scheme of Hanke (1997) [8] results in optimal ES-MDA performance.

This chapter is organized as follows: Section 4.2 presents the methods of Rafiee and Reynolds (2017) [10] and Emerick (2019) [12]. In section 4.3, we describe the proposed method to generate the inflation factors for the ES-MDA and discuss its efficiency. Section 4.4 presents the results obtained by running the ES-MDA using the simplest implementation, the method of Rafiee and Reynolds (2017) [10], the method of Emerick (2019) [12], and the one proposed by this study in a synthetic two-dimensional waterflooding problem.

## 4.2
## Previous works

In this section, we review the methods of Rafiee and Reynolds (2017) [10] and Emerick (2019) [12] for generating the inflation factors for the ES-MDA.

### Method of *Rafiee and Reynolds (2017)*

Consider the problem of updating the ensemble mean in the first assimilation step $k = 0$, i.e., the vector $\overline{m}^0$. Rafiee and Reynolds (2017) [10] used the condition of Hanke (1997) [8] as a criterion to determine $\alpha_1$. Define the matrix $C = (G_D^0 (G_D^0)^T + \alpha_1 I)$, where $G_D^0 = C_d^{-1/2} \Delta D^0$, and consider $y^0 = C_d^{-1/2} \left( d_{obs} - \overline{d}^0 \right)$ in Equation (3-22). Thus, Equation (3-22) can be rewritten as:

$$\rho^2 \leq \alpha_1^2 \frac{||C^{-1} y^0||^2}{||y^0||^2}. \tag{4-1}$$

If $y^0$ is in the same direction of the $k$th singular vector of the matrix $C$, then Equation (4-1) reduces to:

$$\rho^2 \leq \frac{\alpha_1^2}{\left( \sigma_k^2 + \alpha_1 \right)^2}, \tag{4-2}$$

where $\sigma_k$ correspond to the $k$th singular value of $G_D^0$. The equality occurs when:

$$\alpha_1 = \frac{\rho}{1 - \rho} \sigma_k^2. \tag{4-3}$$

The largest value for $\alpha_1$ is obtained when $y^0$ is aligned with the first singular vector of the matrix $C$. Conversely, the smallest $\alpha_1$ is obtained when $y^0$ is aligned with the singular vector corresponding to the smallest singular value of $C$. The optimum value for $\alpha_1$ is between these two extremes. Therefore, Rafiee and Reynolds (2017) [10] computed $\alpha_1$ using the following equation:

$$\alpha_1 = \frac{\rho}{1 - \rho} \overline{\sigma}^2, \tag{4-4}$$

where

$$\overline{\sigma} = \frac{1}{N} \sum_{i=1}^{N} \sigma_i. \tag{4-5}$$

In Equation (4-5), $N$ is the number of non-zero singular values of $G_D^0$. Using $\alpha_1$ computed by Equation (4-4), the other inflation factors are selected geometrically in a decreasing order, as:

$$\alpha_k = \gamma^{k-1}\alpha_1 \quad \forall k = 1, \cdots, N_a, \qquad (4\text{-}6)$$

where $\gamma \in (0, 1]$. The computation of $\gamma$ can be done following Equation (2-52), by solving the following problem:

$$\sum_{k=1}^{N_a} \frac{1}{\gamma^{k-1}} = \alpha_1. \qquad (4\text{-}7)$$

As the left side of Equation (4-7) is the finite sum of the geometric progression with ratio $\gamma$, Equation (4-7) can be rewritten as:

$$\frac{1 - \gamma^{-N_a}}{1 - \gamma^{-1}} = \alpha_1. \qquad (4\text{-}8)$$

To solve this problem, define the function:

$$f_1(\gamma) = \frac{1 - \gamma^{-N_a}}{1 - \gamma^{-1}} - \alpha_1, \qquad (4\text{-}9)$$

and find $\gamma^*$ such that $f_1(\gamma^*) = 0$.

## Method of *Emerick (2019)*

Emerick (2019) [12] proposed to select the last inflation factor $\alpha_{N_a}$ and compute the previous ones geometrically in a increasing order, by using the following formula:

$$\alpha_k = \gamma^{k-N_a}\alpha_{N_a} \quad \forall k = 1, \cdots, N_a. \qquad (4\text{-}10)$$

The coefficient $\gamma \in (0, 1]$ can be computed by finding the root of $f_2$, defined as follows:

$$f_2(\gamma) = \frac{1 - \gamma^{N_a}}{1 - \gamma} - \alpha_{N_a}. \qquad (4\text{-}11)$$

The value of $\alpha_1$, computed from Equation (4-10) with proper $\gamma$, is tested against the Morozov's Discrepancy Principle (MDP) [14]. If $\alpha_1$ does not satisfy the MDP, the number of iterations $N_a$ is increased by one and Eqs. (4-10) and

(4-11) are applied again until a suitable value for $\alpha_1$ is obtained. To check the MDP, Emerick (2019) [12] defines a function as follows:

$$h(\alpha) = ||G_D^0 x_\alpha - y^0||^2 - (\tau\eta)^2, \tag{4-12}$$

where:

$$x_\alpha = \left(G_D^0\right)^T \left(G_D^0 \left(G_D^0\right)^T + \alpha I_d\right)^{-1} y^0. \tag{4-13}$$

He computes $\alpha^* > 0$ such that $h(\alpha^*) = 0$. Thus, it is sufficient to check if computed $\alpha_1$ is such that $\alpha_1 \geq \alpha^* > 0$. In Equation (4-12), Emerick (2019) [12] used $\alpha_{N_a} = 1.5$, $\tau = 1$ and $\eta = \sqrt{N_d}\sigma_y$, where $\sigma_y$ is the standard deviation of $y^0$.

## 4.3
## Proposed Methodology

In this section, we present a new formulation of generating the inflation factors for ES-MDA. Rafiee and Reynolds (2017) [10] used the scheme of Hane (1997) [8] to calculate the first inflation factor $\alpha_1$. However, the following ones are computed considering only Equation (2-52). Emerick (2019) [12] proposed a methodology to obtain these factors based on the selection of $\alpha_{N_a}$. However, Emerick (2019) [12] proposes no analytical procedure for efficiently computing $\alpha_{N_a}$. Considering the remarks about the decay rate of the method proposed by Rafiee and Reynolds (2017) [10], and the lack of a procedure to efficiently compute $\alpha_{N_a}$ in the method of Emerick (2019) [12], we attempt to present an analytical formula to compute the last inflation factor in advance. Besides, computing $\alpha_{N_a}$ as proposed by this study may be used in the algorithm proposed by Emerick (2019) [12].

### Method Formulation

Hereafter, we present the method to generate inflation factors where $\alpha_1$ and $\alpha_{N_a}$ are known. As a result, the method searches for the number of assimilation $N_a$ and the $\gamma$ ratio (Eqs. (4-6) and (4-10)) that produce inflation factors that satisfy Equation (2-52). In this study, we use $\rho = 0.5$, $\tau = 1$, and $\eta = \sqrt{N_d}$ in Eqs. (4-4) and (4-12) [23]. Suppose $\alpha_1$ and $\alpha_{N_a}$ known. Using Equation (4-6), we can compute the value of $\gamma$ from:

$$\alpha_{N_a} = \alpha_1 \gamma^{N_a - 1}, \tag{4-14}$$

as:

$$\gamma = \left( \frac{\alpha_{N_a}}{\alpha_1} \right)^{\frac{1}{N_a - 1}}. \tag{4-15}$$

However, if we generate inflation factors $\alpha_k$, $k = 1, \cdots, N_a$, using $\gamma$ computed by Equation (4-15) in Equation (4-6), we must not attend Equation (2-52). Therefore, we must compute $N_a$ that produces inflation factors that satisfy Equation (2-52). To do so, we use Equation (4-8) as follows:

$$1 - \frac{1}{\gamma^{N_a}} = \alpha_1 - \frac{\alpha_1}{\gamma}. \tag{4-16}$$

Taking all terms that contains $\gamma$ to the left hand side of Equation (4-16) and factoring out $\frac{1}{\gamma}$ in the left hand side:

$$\frac{1}{\gamma} \left( \alpha_1 - \frac{1}{\gamma^{N_a - 1}} \right) = \alpha_1 - 1. \tag{4-17}$$

Multiplying both sides of Equation (4-17) by $\frac{1}{\alpha_1}$:

$$\frac{1}{\gamma} \left( 1 - \frac{1}{\alpha_1 \gamma^{N_a - 1}} \right) = 1 - \frac{1}{\alpha_1}. \tag{4-18}$$

Using Equation (4-14) in Equation (4-18):

$$\frac{1}{\gamma} \left( 1 - \frac{1}{\alpha_{N_a}} \right) = 1 - \frac{1}{\alpha_1}. \tag{4-19}$$

Then, directly from Equation (4-19):

$$\gamma = \frac{1 - \frac{1}{\alpha_{N_a}}}{1 - \frac{1}{\alpha_1}}. \tag{4-20}$$

If $\alpha_1$ is large, Equation (4-20) provides an approximation of $\gamma$ as:

$$\gamma \approx 1 - \frac{1}{\alpha_{N_a}}. \tag{4-21}$$

Therefore, selecting a high value for $\alpha_{N_a}$ may lead to $\gamma \approx 1$, i.e., the algorithm may yield a high number of assimilations. This conclusion also holds for the method of Emerick (2019) [12]. For simplification, consider:

$$T_1 = \frac{1 - \frac{1}{\alpha_{N_a}}}{1 - \frac{1}{\alpha_1}}, \tag{4-22}$$

and

$$T_2 = \frac{\alpha_{N_a}}{\alpha_1}. \tag{4-23}$$

Finally, to compute $N_a$ it is needed to match Equation (4-15) and Equation (4-20), using Eqs. (4-22) and (4-23), as follows:

$$T_2^{\frac{1}{N_a-1}} = T_1. \tag{4-24}$$

Taking the logarithm of both sides of Equation (4-24):

$$\frac{1}{N_a - 1} \log T_2 = \log T_1. \tag{4-25}$$

Then, with a simple algebraic manipulation:

$$N_a = 1 + \frac{\log T_2}{\log T_1}. \tag{4-26}$$

The proposed method uses the values of $\gamma$ computed in Equation (4-20) and $N_a$ computed in Equation (4-26) for generating inflation factors satisfying Equation (2-52) with $\alpha_1$ and $\alpha_{N_a}$ previously determined. However, Equation (4-26) does not always assume integer values, as required by ES-MDA formulation [17]. To address this problem, observe that the scheme depicted in Equation (4-2), presents the following lower bound for $\alpha_1$:

$$\alpha_1 \geq \frac{\rho}{1 - \rho} \bar{\sigma}^2. \tag{4-27}$$

Rafiee and Reynolds (2017) [10] select $\alpha_1$ by searching for the equality of Equation (4-2). Thus, to alleviate the problem of computing non-integer values for $N_a$ in Equation (4-26), define the function $f_3$:

$$f_3(\alpha) = 1 + \frac{\log T_2(\alpha)}{\log T_1(\alpha)}, \tag{4-28}$$

where

$$T_1(\alpha) = \frac{1 - \frac{1}{\alpha_{N_a}}}{1 - \frac{1}{\alpha}}, \tag{4-29}$$

and

$$T_2(\alpha) = \frac{\alpha_{N_a}}{\alpha}, \tag{4-30}$$

and compute the minimum $\alpha^*$ such that:

$$\alpha^* \geq \frac{\rho}{1-\rho}\overline{\sigma}^2, \tag{4-31}$$

and $f_3(\alpha^*) \in \mathbb{N}$ and, finally, set $\alpha_1 = \alpha^*$.

In this study, we attempt to present a method to determine $\alpha_{N_a}$. The idea is to use the procedure of Rafiee and Reynolds (2017) [10], using the singular values of the dimensionless sensitivity matrix estimated from the prior ensemble $G_D^0$. They observed that $\alpha_1$ assumes its largest value if it is computed using the largest singular value of $G_D^0$. On the other hand, $\alpha_1$ assumes its smallest value if it is computed using the smallest singular value of $G_D^0$. Therefore, as we desire to compute inflation factors in decreasing order, we may assume that $\alpha_{N_a}$ is the smallest inflation factor. However, as a consequence of Equation (2-52), we need that $\alpha_{N_a} \geq 1$. Thus, we propose to select $\alpha_{N_a}$ such that:

$$\alpha_{N_a} = \frac{\rho}{1-\rho}\left(\frac{1}{p}\sum_{i=0}^{p}\sigma_{N-i}\right)^2, \tag{4-32}$$

where $p \in [1, N-1]$ is the minimum integer such that $\alpha_{N_a} > \mu_\alpha$, where $\mu_\alpha \in [1, N_a)$ is a threshold to select a minimum value for $\alpha_{N_a}$, and $N$ is the number of non-zeros singular values of $G_D^0$. The threshold $\mu_\alpha$ works as a lower bound for $\alpha_{N_a}$ and may control the quality of final results and the magnitude of $\alpha_{N_a}$. If one selects $\mu_\alpha = 1$, $\alpha_{N_a}$ might be close to one. Moreover, if $\alpha_{N_a} = 1$, we have the standard ES. On the other hand, if one selects high $\mu_\alpha$, as observed in Equation (4-21), it may lead to a large $N_a$, which increases computational cost. Another important remark is that $N_e$ must be sufficiently large for $G_D^0$ to have singular values close to one. Based on the examples examined in this study, selecting $N_e = 100$ was enough.

Generating $\alpha_{N_a}$ using Equation (4-32) is valid for data assimilation only in the linear-Gaussian case, where the dimensionless sensitivity matrix $G_D^0$ does not change during data assimilation, i.e., $y_j^k = G_D x_j^k$, with $G_D = G_D^k = G_D^0$, $\forall k = 1, \cdots, N_a$. The motivation for using Equation (4-32) to generate the last inflation factor in the linear-Gaussian case comes from the fact that the ES-MDA update equation has a similar structure as the Levenberg-Marquardt minimization algorithm [7, 10]. Therefore, the inflation factor $\alpha$ plays the role of the regularization parameter [12]. Thus, as the method of Hanke (1997) [8] has been proven of optimal order [27], it can be efficiently applied to compute the ES-MDA inflation factor. Nevertheless, this fact may not hold for the

nonlinear case. On the other hand, our numerical examples show that the proposed approach improves the ES-MDA outcomes with fewer assimilations than the others available in the literature. Moreover, the motivation to use such a procedure in the nonlinear case comes from the study of Silva *et al.* (2021) [6], which numerically shows that generating the first inflation factor using the scheme of Hanke (1997) [8] almost neglects the effects of minimal singular values, i.e., diminish error propagation on the multiple data assimilation processes. As a result, we conjecture that it may also happen for the last assimilation step, computing it by using the formula of Rafiee and Reynolds (2017) [10].

**Analysis of $f_3$**

Depending on the value of $\alpha_{N_a}$, it is not straightforward to predict what the values of $T_1(\alpha)$ (Equation (4-29)) and $T_2(\alpha)$ (Equation (4-30)) may assume. Therefore, we analyze the continuity and the limits of $f_3$ in the interval where it is continuous.

To verify if $f_3$ is a continuous function, we must check the values where $\log T_1(\alpha) \neq 0$. One can notice that $\log T_1(\alpha) = 0$ if and only if $T_1(\alpha) = 1$. Then, $T_1(\alpha) = 1$ implies that $\alpha = \alpha_{N_a}$. However, if $\alpha = \alpha_{N_a}$, we have that $\alpha_1 = \alpha_{N_a}$ and, using Equation (4-14), $\gamma = 1$ and the resulting method is the simplest formulation of the ES-MDA with $\alpha_k = N_a \ \forall \ k \in \{1, \cdots, N_a\}$. As we wish to compute inflation factors geometrically in decreasing order, we assume $\alpha > \alpha_{N_a}$ and guarantee that $f_3$ is a continuous function in $(\alpha_{N_a}, \infty)$. To analyze $f_3$ in $(\alpha_{N_a}, \infty)$ it is sufficient to check the limits of $f_3$ when $\alpha \to \alpha_{N_a}^+$ and $\alpha \to \infty$.

For the case when $\alpha \to \alpha_{N_a}^+$, we can see that:

$$\lim_{\alpha \to \alpha_{N_a}^+} T_1(\alpha) = \lim_{\alpha \to \alpha_{N_a}^+} T_2(\alpha) = 1. \tag{4-33}$$

Therefore, $\lim_{\alpha \to \alpha_{N_a}^+} f_3(\alpha)$ is not defined. Computing the following derivatives:

$$\frac{d}{d\alpha} \left( \log T_1(\alpha) \right) = -\frac{1}{\alpha \left( \alpha - 1 \right)}, \tag{4-34}$$

and

$$\frac{d}{d\alpha} \left( \log T_2(\alpha) \right) = -\frac{1}{\alpha}, \tag{4-35}$$

and applying L'Hôpital's rule:

$$\lim_{\alpha \to \alpha_{N_a}^+} f_3(\alpha) = 1 + \lim_{\alpha \to \alpha_{N_a}^+} \left( \frac{-\frac{1}{\alpha}}{-\frac{1}{\alpha(\alpha-1)}} \right) = 1 + \lim_{\alpha \to \alpha_{N_a}^+} \alpha - 1 = \alpha_{N_a}. \qquad (4\text{-}36)$$

The result shown in Equation (4-36) is consistent in the sense that if
$N_a = \alpha_{N_a}$, it refers to the simplest inflation factors selection for ES-MDA. For
the case $\alpha \to \infty$, note that:

$$\lim_{\alpha \to \infty} \log T_1(\alpha) = \lim_{\alpha \to \infty} \left( \log \left( 1 - \frac{1}{\alpha_{N_a}} \right) - \log \left( 1 - \frac{1}{\alpha} \right) \right)$$
$$= \log \left( 1 - \frac{1}{\alpha_{N_a}} \right). \qquad (4\text{-}37)$$

As $\alpha_{N_a} > 1$ (Equation (2-52)), we have that $1 - 1/\alpha_{N_a} \in (0, 1)$. Thus,
$\log \left( 1 - 1/\alpha_{N_a} \right) < 0$. One can notice that, if $\alpha \to \infty$, $T_2(\alpha) \to 0$ (Equation
(4-30)), implying that $\log T_2(\alpha) \to -\infty$. Thus, by Equation (4-37):

$$\lim_{\alpha \to \infty} f_3(\alpha) = \infty, \qquad (4\text{-}38)$$

what completes our analysis of $f_3$.

**Algorithm to compute $N_a$ and $\alpha^*$**

To compute $N_a$ and $\alpha^*$ attending Equation (4-31), we suggest a procedure
using a zero-finding method, e.g., the bisection method. First, define the
function:

$$f_{aux}(\alpha) = f_3(\alpha) - N_a, \qquad (4\text{-}39)$$

with previously defined $N_a$. From Equation (4-36):

$$\lim_{\alpha \to \alpha_{N_a}^+} f_{aux}(\alpha) = \alpha_{N_a} - N_a < 0, \qquad (4\text{-}40)$$

as $1 < \alpha_{N_a} < N_a$. From Equation (4-38):

$$\lim_{\alpha \to \infty} f_{aux}(\alpha) = \infty. \qquad (4\text{-}41)$$

Thus, we can use bisection to compute $\alpha^*$, such that $f_{aux}(\alpha^*) = 0$, in $(\alpha_{N_a}, \infty)$, and set $\alpha_1 = \alpha^*$. Note that we must obtain $\alpha^*$ respecting Equation (4-31). In other words, we need that $\alpha^* \in \left[ \frac{\rho}{1-\rho} \bar{\sigma}^2, \infty \right)$. Therefore, we must apply bisection in this interval. Hence, $N_a$ will be iteratively increased until we get $f_{aux}(\frac{\rho}{1-\rho} \bar{\sigma}^2) \leq 0$. Hereafter we present an algorithm for computing inflation factors and the number of data assimilation $N_a$. We denote the value of $\alpha_1$ computed using Equation (4-4) as $\alpha_1^0$ and the initial number of data assimilation as $N_a^0$.

---

**Algorithm 1** New procedure to generate ES-MDA inflation factors.

---

1. Set $N_a = N_a^0$, $\mu_\alpha$ and compute $\alpha_1^0$ using Equation (4-4);
2. Compute $\alpha_{N_a}$ using Equation (4-32);

3. While $f_{aux}(\alpha_1^0) \geq 0$ (Equation (4-39)), set $N_a = N_a + 1$;

4. Compute $\alpha^*$ using bisection in $f_{aux}$ and set $\alpha_1 = \alpha^*$;

5. Compute $\gamma$ using $\alpha_1$ and $\alpha_{N_a}$ in Equation (4-15) and generate $\{\alpha_i\}_{i=1}^{N_a}$ using Equation (4-6).

---

One can notice that if $\alpha^* = \alpha_1^0$, the method is similar to the method of Rafiee and Reynolds (2017) [10] with $N_a = N_a^0$. The choice of $\mu_\alpha$ affects the computed $N_a$. This fact will be demonstrated in the next section with the examples.

## 4.4
## Results and Discussion

In this section, we present a comparison of the results obtained when implementing ES-MDA with data covariance inflation factors generated using the methods of Rafiee and Reynolds (2017) [10], Emerick (2019) [12], the simplest choice with $\alpha_i = N_a \ \forall i = 1, \ldots, N_a$, and the method proposed in this study. For simplification, we follow the nomenclature defined by Emerick (2019) [12], where the method of Rafiee and Reynolds (2017) [10] is referred to as ES-MDA-GEO1; the method of Emerick (2019) [12] is referred to as ES-MDA-GEO2; selecting $\alpha_i = N_a \ \forall i = 1, \ldots, N_a$ is referred to as ES-MDA-EQL and the method of this study is referred to as ES-MDA-GEO3. The results are compared using the root mean square error (RMSE) and the data mismatch ($O_d$) of the final ensemble, exposed in Equation (4-42) and Equation (4-43):

$$RMSE = \left( \frac{1}{N_m} \sum_{k=1}^{N_m} (m_{true,k} - m_{j,k})^2 \right)^{1/2}, \qquad (4\text{-}42)$$

and

$$O_d(d_j^f) = \frac{1}{N_d}(d_j^f - d_{obs})^T C_d^{-1}(d_j^f - d_{obs}). \tag{4-43}$$

**Waterflooding example**

The methods will be evaluated using a synthetic two-dimensional water flooding history matching problem. The reservoir is squared with a total length of 1575 m in each direction. The reservoir domain is discretized in $63 \times 63 \times 1$ grid, where each grid block is 25m $\times$ 25m $\times$ 25m. The true permeability field (Figure 4.1) was generated using a spherical covariance function with a correlation length of 40 gridblocks. The prior model of the log-permeability is constant and equal to 5, with a prior variance of 1. These informations were used to create ensembles of size $N_e = 100$. We consider the vector of model parameters $m$ consisting of the gridblocks log-permeabilities. The reservoir model contains nine producing wells and four water injection wells. The water injection wells are operated at constant bottom-hole pressure (BHP) of 325 kgf/cm² and the producing wells are operated at constant BHP of 275 kgf/cm² for the whole simulation time, which is a period of 3600 days, with measurements every 150 days. The observed data consists of water and oil rate of the producing wells and water injection rate of water injection wells, where each well location is displayed in Figure 4.1, where circles corresponds to the oil producing wells and triangles corresponds to the water injection wells. For the synthetic measurements, it was added an error of 5% and we assume that the matrix $C_d$ is a diagonal matrix with each diagonal entry corresponding to the square of the standard deviation used to generate the errors in the observed measurements. During all data assimilation, the Schur-product-based covariance localization with isotropic correlation function of 40 gridblocks was applied directly to the Kalman gain matrix [24].

We created five different stochastic ensembles with the pieces of information described before, generated using distinct seeds each. The first method to be applied in the ensembles was ES-MDA-GEO3. Thus, we could discover the computed number of assimilations $N_a$ for each ensemble. Therefore, we can compare the proposed technique with the ES-MDA-GEO1 and the ES-MDA-EQL with the same number of assimilations. Because ES-MDA-GEO2 also computes the adequate number of assimilations to acquire good final results, it was not possible to analyze ES-MDA-GEO2 and ES-MDA-GEO3 with the same number of assimilations.

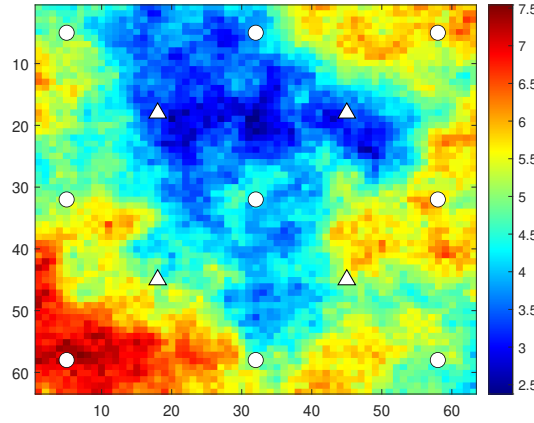For the five ensembles, we tested ES-MDA-GEO3 with two different

Figure 4.1: True log-permeability field (mD) for the two-dimensional test problem.

values of $\mu_\alpha$, $\mu_\alpha = 1.1$ and $\mu_\alpha = 1.2$. The inflation factors' values are exposed in Table 4.1. The main objective is to examine how the magnitude of $\alpha_{N_a}$ influences the quality of the final results. Selecting $\mu_\alpha > 1.2$ led to larger assimilation runs. As we wish to maintain $N_a$ between 4 and 8, we did not test any values for $\mu_\alpha$ greater than 1.2. This fact can be noticed in Table 4.1, where the higher values of $\alpha_{N_a}$ led to higher assimilation runs $N_a$. The values of the inflation factors and the $\gamma$ ratio of ES-MDA-GEO1 are exposed in Table 4.2.

For the ES-MDA-GEO2, setting $N_a = 4$ was sufficient to provide $\alpha_1$ that satisfies the MDP [14]. Figure 4.2 displays the root of the discrepancy function (Equation (4-12)) computed for each ensemble. The top-left Figure refers to Ensemble 1; the top-right Figure refers to Ensemble 2; the mid-left Figure refers to Ensemble 3; the mid-right Figure refers to Ensemble 4; the bottom Figure refers to Ensemble 5. As proposed by Emerick (2019) [12], we use $\alpha_{N_a} = 1.5$. Therefore, we obtained $\alpha_1 = 37.33$, $\alpha_2 = 12.78$, $\alpha_3 = 4.37$, and $\alpha_4 = 1.5$, with $\gamma = 0.3425$ for all ensembles.

**Case** $\mu_\alpha = 1.1$

In this case, all methods were tested using the same number of assimilations of ES-MDA-GEO3 for each ensemble (see Tables 4.1 and 4.2). The only exception is ES-MDA-GEO2, which was tested with $N_a = 4$ for all ensembles. Figure 4.3 exhibits the mean of the final ensemble members for each ES-MDA implementation. In Figure 4.3, the first line corresponds to the reference field; the first column corresponds to ES-MDA-EQL; the second column corresponds to ES-MDA-GEO1; the third column corresponds to ES-MDA-GEO2; the last column corresponds to ES-MDA-GEO3. The lines correspond

|  |  | Ensemble 1 | Ensemble 2 | Ensemble 3 | Ensemble 4 | Ensemble 5 |
|---|---|---|---|---|---|---|
| $\mu_\alpha = 1.1$ | $\alpha_1$ | 6,617.03 | 21,117.05 | 13,297.75 | 19,852.75 | 6,975.4 |
|  | $\alpha_2$ | 756.26 | 2,970.8 | 1,269.8 | 2,828.83 | 786.43 |
|  | $\alpha_3$ | 86.43 | 417.94 | 121.25 | 403.08 | 88.66 |
|  | $\alpha_4$ | 9.87 | 58.8 | 11.58 | 57.43 | 10 |
|  | $\alpha_5$ | 1.13 | 8.27 | 1.11 | 8.18 | 1.12 |
|  | $\alpha_6$ |  | 1.16 |  | 1.16 |  |
|  | $N_a$ | 5 | 6 | 5 | 6 | 5 |
|  | $\gamma$ | 0.1143 | 0.1407 | 0.0955 | 0.1425 | 0.1127 |
| $\mu_\alpha = 1.2$ | $\alpha_1$ | 12,276.25 | 23,439.67 | 13,754.12 | 20,560.37 | 8,659.24 |
|  | $\alpha_2$ | 2,662.23 | 4,541.66 | 2,924.1 | 4,075.51 | 1,465.85 |
|  | $\alpha_3$ | 577.33 | 878 | 621.66 | 807.85 | 248.14 |
|  | $\alpha_4$ | 125.2 | 170.5 | 132.16 | 160.13 | 42 |
|  | $\alpha_5$ | 27.15 | 33.03 | 28.09 | 31.74 | 7.11 |
|  | $\alpha_6$ | 5.88 | 6.4 | 5.97 | 6.29 | 1.21 |
|  | $\alpha_7$ | 1.27 | 1.24 | 1.27 | 1.24 |  |
|  | $N_a$ | 7 | 7 | 7 | 7 | 6 |
|  | $\gamma$ | 0.2169 | 0.1938 | 0.2126 | 0.1982 | 0.1693 |

Table 4.1: Computed inflation factors using Algorithm 1 for ES-MDA-GEO3 for the five different ensembles with $\mu_\alpha = 1.1$ and $\mu_\alpha = 1.2$.

to the five different ensembles. One may observe that the ES-MDA-EQL over-corrects the regions with high log-permeability. This fact is noticed in other studies about the ES-MDA implementation [7, 10]. It may be explained by the ill-conditioning of the ensemble approximations of covariance matrices. The overshooting become quite severe in the ensembles 1, 3, and 5, where $N_a = 5$, and slightly lightened in the ensembles 2 and 4, where $N_a = 6$. However, for all ensembles, this problem is still critical. The overshooting presented in the EQL method is mitigated when the inflation factors are selected using the discrepancy principle, which can be observed in the final results of ES-MDA-GEO1, GEO2, and GEO3. One may note that GEO1 and GEO3 obtained smoother solutions compared to the ones achieved through GEO2. This is because $\alpha_1$ is much greater for GEO1 and GEO3 than for GEO2. In fact, the results of GEO3 is lightly smoother than the ones of GEO1 for the same reason.

Figure 4.4 exposes the posterior standard deviation for the gridblocks log-permeabilities obtained from each ES-MDA implementation. In Figure 4.4, the first column corresponds to ES-MDA-EQL; the second column corresponds to ES-MDA-GEO1; the third column corresponds to ES-MDA-GEO2; the last column corresponds to ES-MDA-GEO3. The lines correspond to the five

|  |  | Ensemble 1 | Ensemble 2 | Ensemble 3 | Ensemble 4 | Ensemble 5 |
|---|---|---|---|---|---|---|
| $\mu_\alpha = 1.1$ | $\alpha_1$ | 4,696.23 | 4,802.48 | 4,816.39 | 4,589.93 | 4,888.64 |
|  | $\alpha_2$ | 586.54 | 919.67 | 597.61 | 887.36 | 604.24 |
|  | $\alpha_3$ | 73.25 | 176.11 | 74.15 | 171.55 | 74.68 |
|  | $\alpha_4$ | 9.15 | 33.72 | 9.2 | 33.16 | 9.23 |
|  | $\alpha_5$ | 1.14 | 6.45 | 1.14 | 6.41 | 1.14 |
|  | $\alpha_6$ |  | 1.23 |  | 1.24 |  |
|  | $N_a$ | 5 | 6 | 5 | 6 | 5 |
|  | $\gamma$ | 0.1249 | 0.1915 | 0.1241 | 0.1933 | 0.1236 |
| $\mu_\alpha = 1.2$ | $\alpha_1$ | 4,696.23 | 4,802.48 | 4,816.39 | 4,589.93 | 4,888.64 |
|  | $\alpha_2$ | 1,205.78 | 1,228.18 | 1,231.11 | 1,183.25 | 932.68 |
|  | $\alpha_3$ | 309.58 | 314.1 | 314.68 | 305.03 | 177.94 |
|  | $\alpha_4$ | 79.48 | 80.32 | 80.43 | 78.63 | 33.94 |
|  | $\alpha_5$ | 20.4 | 20.54 | 20.56 | 20.27 | 6.47 |
|  | $\alpha_6$ | 5.24 | 5.25 | 5.25 | 5.22 | 1.23 |
|  | $\alpha_7$ | 1.34 | 1.34 | 1.34 | 1.34 |  |
|  | $N_a$ | 7 | 7 | 7 | 7 | 6 |
|  | $\gamma$ | 0.2568 | 0.2557 | 0.2556 | 0.2578 | 0.1908 |

Table 4.2: Computed inflation factors and $\gamma$ ratio for ES-MDA-GEO1 for comparison with ES-MDA-GEO3 with $\mu_\alpha = 1.1$ and $\mu_\alpha = 1.2$.

different ensembles. Without a precise sample, it is not possible to estimate the standard deviation accurately. Nevertheless, we may obtain a trustable estimate using the ensemble. Overall, all methods presented standard deviation that does not vary too much, which means consistency. The ES-MDA-EQL obtained the lowest standard deviation values among all methods. It tells us that the log-permeability fields obtained by the EQL method do not vary from one ensemble to another. The ES-MDA-GEO1 and GEO3 presented higher values of the standard deviation among all of them. The methods provide different scenarios that match the observed data and the mean log-permeability, as exposed in Figure 4.3, allowing a more precise analysis when joint with seismic and geophysical data.

In Figure 4.5, the RMSE (Equation (4-42)) computed for each final ensemble is displayed in a boxplot. The top-left figure refers to ES-MDA-EQL; the top-right figure refers to ES-MDA-GEO1; the bottom-left figure refers to ES-MDA-GEO2; the bottom-right figure refers to ES-MDA-GEO3. The black boxes correspond to the percentiles 25 and 75. The solid red lines inside the boxes correspond to the median of the ensemble. The red diamonds correspond to the mean of the ensemble. The solid black lines outside the boxes correspond

to the extension of the ensemble, excluding the outliers. The red whiskers correspond to the outliers of the ensemble. The dashed black lines correspond to the median of the ensemble's medians. The dashed green lines correspond to the median of the ensemble's median obtained by ES-MDA-GEO3. We display the median of each ensemble's medians as a plan to visually compare the magnitude of the root mean squared error computed for each ensemble member and compare with the one obtained by ES-MDA-GEO3. One may see that the RMSE computed for ES-MDA-GEO3 with $\mu_\alpha = 1.1$ achieved the lowest values than the other ES-MDA implementations, indicating that the proposed method could estimate the reference log-permeability field more accurately.

Figure 4.6 exposes the posterior production data combining all the five ensembles obtained by the different ES-MDA implementations. The first column refers to well INJ-4; the second column refers to well PROD-6; the third column refers to well PROD-7. In contrast, the first line corresponds to ES-MDA-EQL; the second line corresponds to ES-MDA-GEO1; the third line corresponds to ES-MDA-GEO2; the fourth line corresponds to ES-MDA-GEO3. Blue solid lines correspond to the posterior ensemble data. The solid gray lines correspond to the prior ensemble data. The red circles correspond to the observed data. The approach of exposing the posterior production data by combining all ensembles is similar to the one presented in Le *et al.* (2016) [7]. Although it presents some outliers, one may note that the ES-MDA-EQL obtained the best match of observed data compared to the other techniques. The methods ES-MDA-GEO1 and GEO3 performed similarly. It is noticeable that, although it still presents a good match of the observed data, the ES-MDA-GEO2 obtained a considerable number of outliers. This fact can also be noted in Figure 4.7, which displays the computed data mismatch for all methods combining the five ensembles. Boxes, lines, and points here have the same meaning as in Figure 4.5. Figure 4.7 endorses the observation that ES-MDA-EQL obtained the best data match among all methods and that the GEO1 and GEO3 performed similarly, although the median of GEO3 is relatively lower than the one computed by GEO1.

**Case** $\mu_\alpha = 1.2$

This section compares the results obtained by the ES-MDA methods with the one obtained by ES-MDA-GEO3 with $\mu_\alpha = 1.2$. The results are not much different from the ones presented in the previous section. Therefore, we present only the statistical numerical measurements to avoid being repetitive.

Notice that the results for ES-MDA-GEO2 are the same as the ones exposed in the case where $\mu_\alpha = 1.1$. The other methods follow the same number of assimilations displayed in Tables 4.1 and 4.2.

Figure 4.8 displays the computed RMSE of the posterior ensemble obtained by the different ES-MDA implementations, combining all ensembles. The figures, boxes, lines, and points have the same meaning as in Figure 4.5. One may observe that the results were quite similar to the ones obtained when using $\mu_\alpha = 1.1$. However, the posterior ensemble of ES-MDA-GEO1 with higher assimilation runs resulted in the lowest RMSE among the other methods, including the one proposed in this study. On the other hand, no ensemble member of ES-MDA-GEO3 achieved RMSE higher than 1, including the outliers. This fact does not hold for the ES-MDA-GEO1.

The data mismatch function is displayed in Figure 4.9 in a boxplot. The figures, boxes, lines, and points here have the same meaning as in Figure 4.5. Again, ES-MDA-GEO1 and GEO3 presented similar mismatches. The posterior ensemble's median of ES-MDA-GEO1 presents a median slightly lower than the one obtained by ES-MDA-GEO3. The results obtained with $\mu_\alpha = 1.2$ led the ES-MDA-GEO3 to compute a larger value for $N_a$. However, we observed that the match of model parameters was not substantially increased. Therefore, we conclude that using $\mu_\alpha = 1.1$ provides a good trade-off between model parameters and observed data match with smaller assimilation runs for this problem.

**Analysis of final results**

In this section, we summarize the results obtained by all methods by combining the five ensembles. We show the median and the mean of the RMSE and the data mismatch function in Tables 4.3 and 4.4. Implementing ES-MDA-GEO3 with $\mu_\alpha = 1.1$ (Table 4.3) achieved the best model parameters match among the other ES-MDA implementations both for the mean and the median of the posterior ensembles. Concerning the data match, the ES-MDA-EQL produced the best match of the observed data. For the case where $\mu_\alpha = 1.2$ (Table 4.4), the ES-MDA-GEO3 again computed the lowest RMSE when combining the five ensembles. It can be inspected both in the mean and the median. However, the ES-MDA-GEO1 achieved a better match of observed data, which can be noted both in the mean and median of the ensemble. The ES-MDA-EQL achieved the best match of observed data again.

|  |  | EQL | GEO1 | GEO2 | GEO3 |
|---|---|---|---|---|---|
| RMSE | Mean | 0.7985 | 0.7569 | 0.777 | 0.7463 |
| | Median | 0.8013 | 0.7437 | 0.7663 | 0.7365 |
| $O_d$ | Mean | 0.216 | 0.5435 | 1.5066 | 0.4865 |
| | Median | 0.1364 | 0.4492 | 0.6107 | 0.3519 |

Table 4.3: Summary of the methods, comparing with ES-MDA-GEO3 with $\mu_\alpha = 1.1$.

|  |  | EQL | GEO1 | GEO2 | GEO3 |
|---|---|---|---|---|---|
| RMSE | Mean | 0.7907 | 0.7436 | 0.777 | 0.7316 |
| | Median | 0.7717 | 0.7402 | 0.7663 | 0.7314 |
| $O_d$ | Mean | 0.1078 | 0.2996 | 1.5066 | 0.3326 |
| | Median | 0.0972 | 0.247 | 0.6107 | 0.2605 |

Table 4.4: Summary of the methods, comparing with ES-MDA-GEO3 with $\mu_\alpha = 1.2$.

## 4.5
## Conclusions

In this study, we present a new procedure for generating inflations factors for the ES-MDA. The new method uses the analytical formula developed in the study of Rafiee and Reynolds (2017) [10] to compute both the first and the last ES-MDA inflation factors. The other inflation factors are generated geometrically in decreasing order. As a result, the method computes the correct number of assimilations that produces inflation factors such that the sum of their inverse is equal to one, as required by the ES-MDA.

The first and last inflation factors are computed based on the singular values of the average sensitivity matrix, computed using the prior ensemble. Although the proposed methodology is only valid for the linear-Gaussian case, we show with numerical examples that the proposed method is appropriate for the nonlinear case. Moreover, the numerical results suggest that selecting inflation factors as proposed leads to a better match of model parameters with a smaller number of iterations than the other ES-MDA implementations. In Equation (4-32), we determine the magnitude of the last inflation factor by setting the threshold $\mu_\alpha$. We tested two cases where $\mu_\alpha = 1.1$, and $\mu_\alpha = 1.2$. We

show that the choice of $\mu_\alpha$ is crucial for determining the number of assimilations $N_a$. This fact can be observed in Equation (4-21), where a high value of $\alpha_{N_a}$ may impose a high number of assimilations $N_a$.

The primary motivation for using inflation factors as proposed is based on the observations of Rafiee and Reynolds (2017) [10] and Iglesias (2015) [30] that the regularization parameter computed using the scheme of Hanke (1997) [8] usually decreases during the iterations of the ES. Therefore, we attempt to produce inflation factors based on the scheme of Hanke (1997) [8], but using the method of Rafiee and Reynolds (2017) [10] for the linear-Gaussian case. Furthermore, the ES-MDA update equation has the same arrangement as a Levenberg-Marquardt method to minimize the function $\mathcal{O}(m)$, depicted in Equation (2-22). Also, the scheme of Hanke (1997) [8] has been proven to be of optimal order, i.e., it provides optimal accuracy [27]. Hence, generating inflation factors as suggested by this work may yield better ES-MDA final results.

We applied the new method in a synthetic two-dimensional waterflooding test problem. Through the numerical results of five different stochastic ensembles, the ES-MDA-GEO3 yielded the best match of model parameters among the other ES-MDA implementations available in the literature when using $\mu_\alpha = 1.1$, with fewer assimilation steps. Excluding the ES-MDA-EQL, the proposed technique also achieved the best match of observed data. Comparing the results of $\mu_\alpha = 1.1$ and $\mu_\alpha = 1.2$, we observed that $\mu_\alpha = 1.1$ led to a smaller number of assimilations $N_a$ and provided good trade-off between final results and computational time for the tested problem.
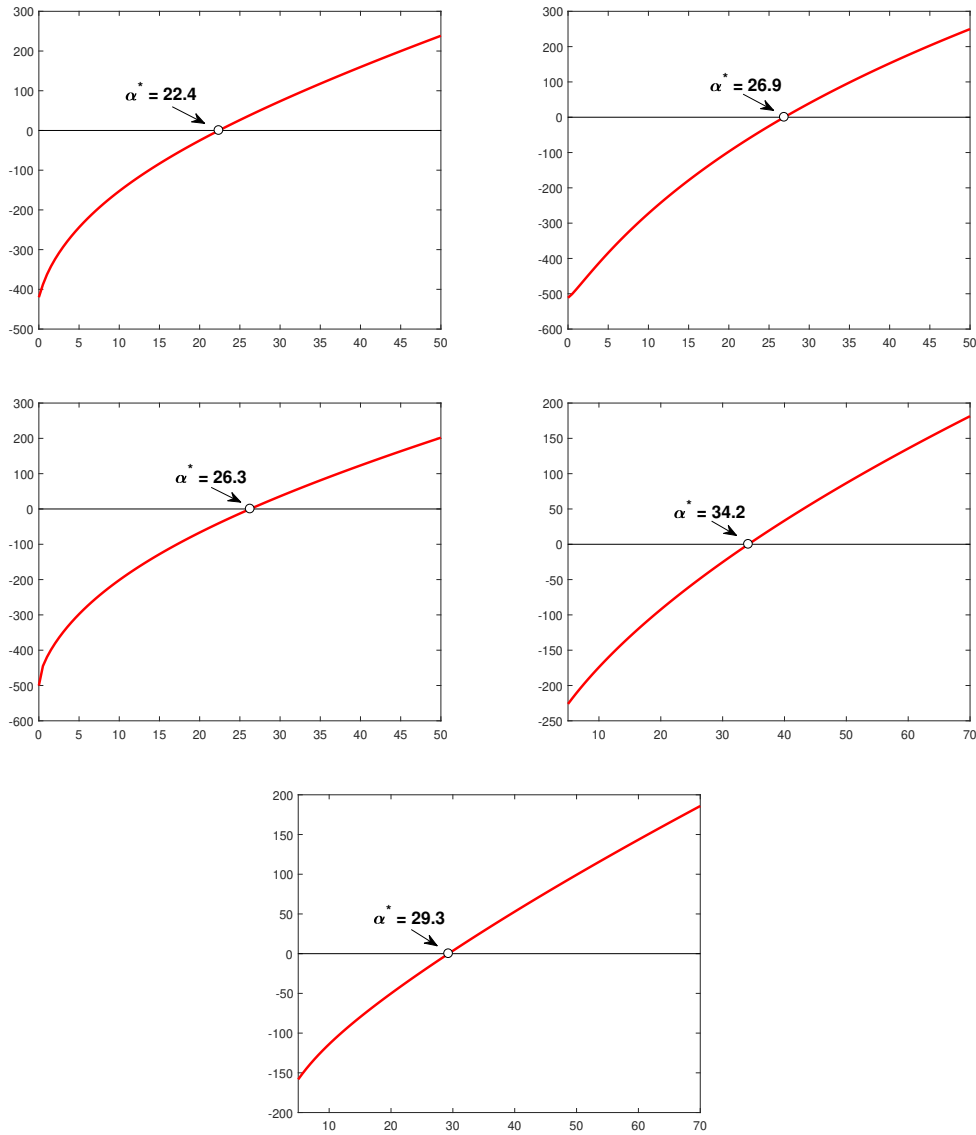
Figure 4.2: Root $\alpha^*$ of discrepancy function (Equation (4-12)) computed for ES-MDA-GEO2 for all five stochastic ensembles.
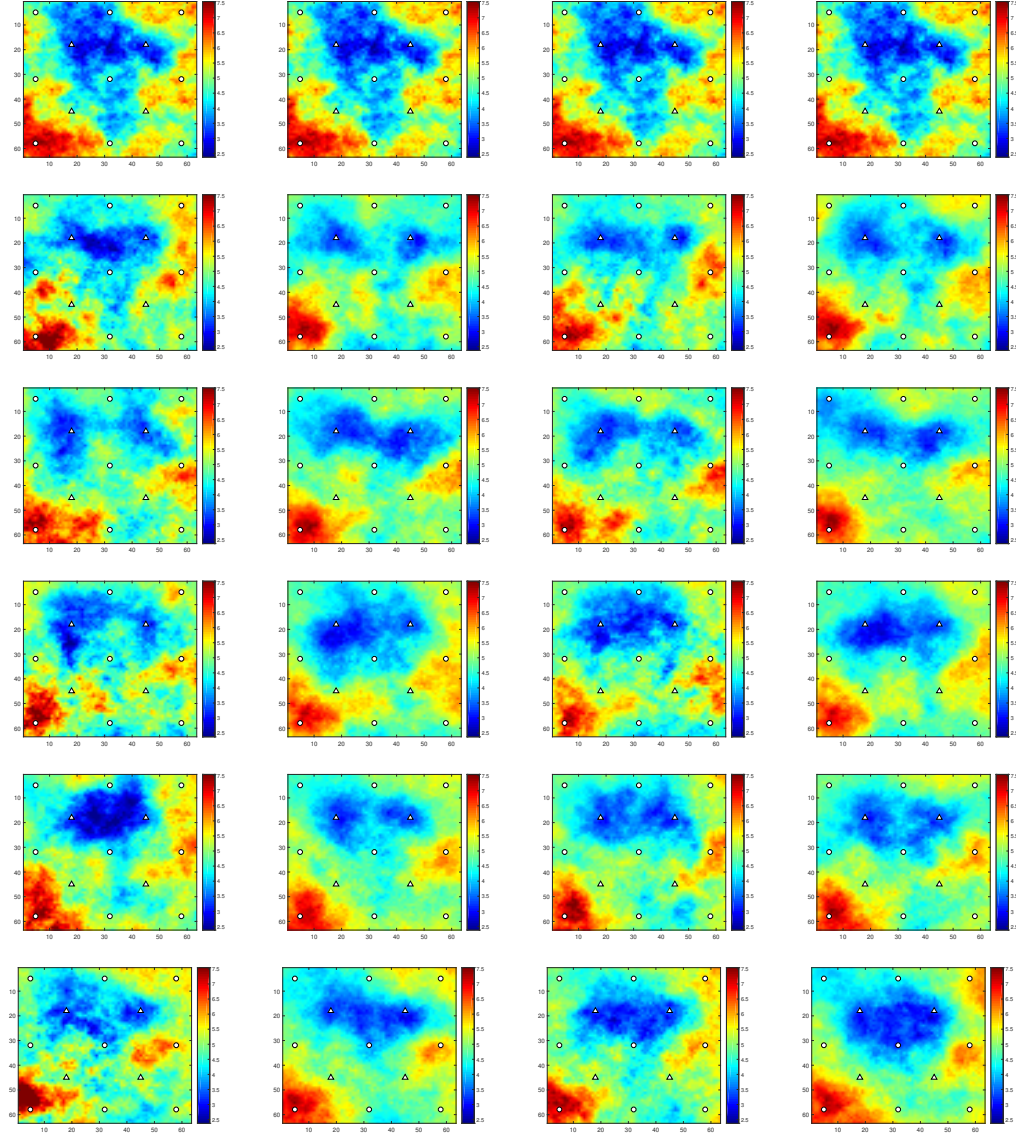
Figure 4.3: Ensemble means achieved with the methods, compared with ES-MDA-GEO3 using $\mu_\alpha = 1.1$. The first line corresponds to the reference field; the first column corresponds to ES-MDA-EQL; the second column corresponds to ES-MDA-GEO1; the third column corresponds to ES-MDA-GEO2; the last column corresponds to ES-MDA-GEO3. The lines correspond to the five different ensembles.
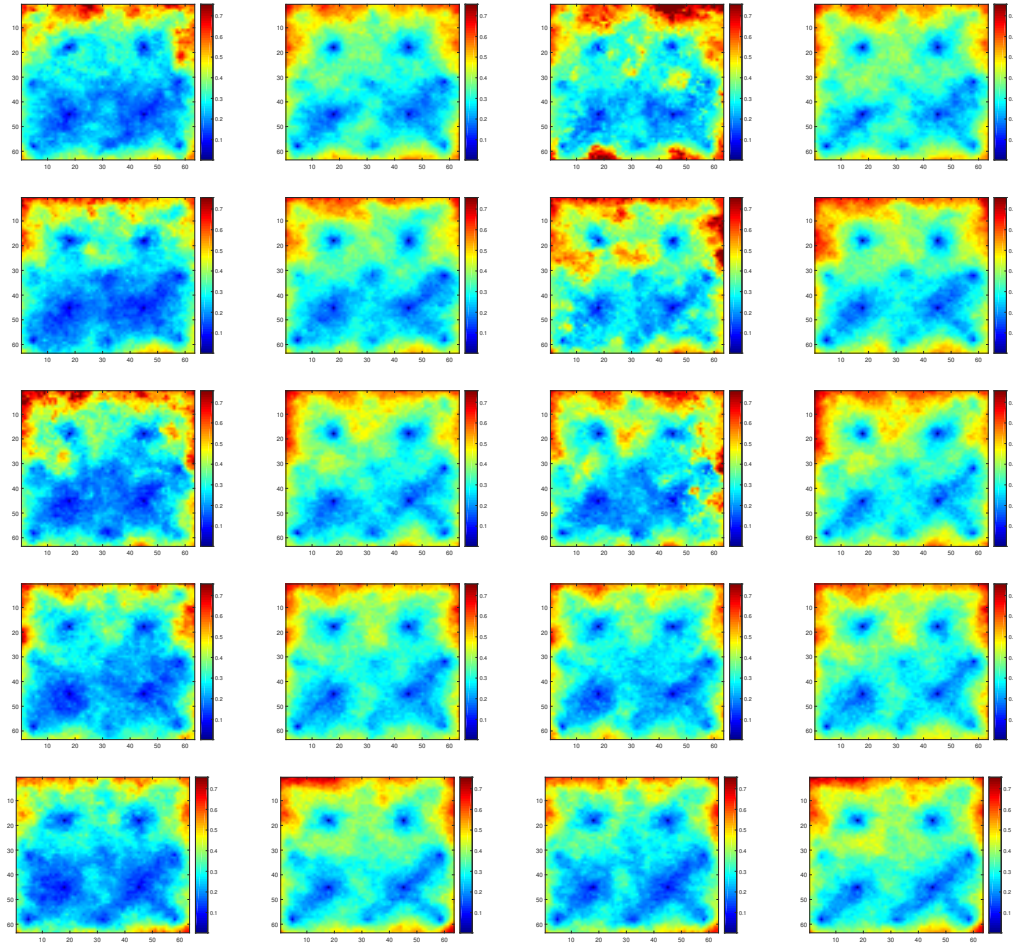
Figure 4.4: Posterior standard deviation for each ensemble obtained by different ES-MDA implementations. The first column corresponds to ES-MDA-EQL; the second column corresponds to ES-MDA-GEO1; the third column corresponds to ES-MDA-GEO2; the last column corresponds to ES-MDA-GEO3. The lines correspond to the five different ensembles.
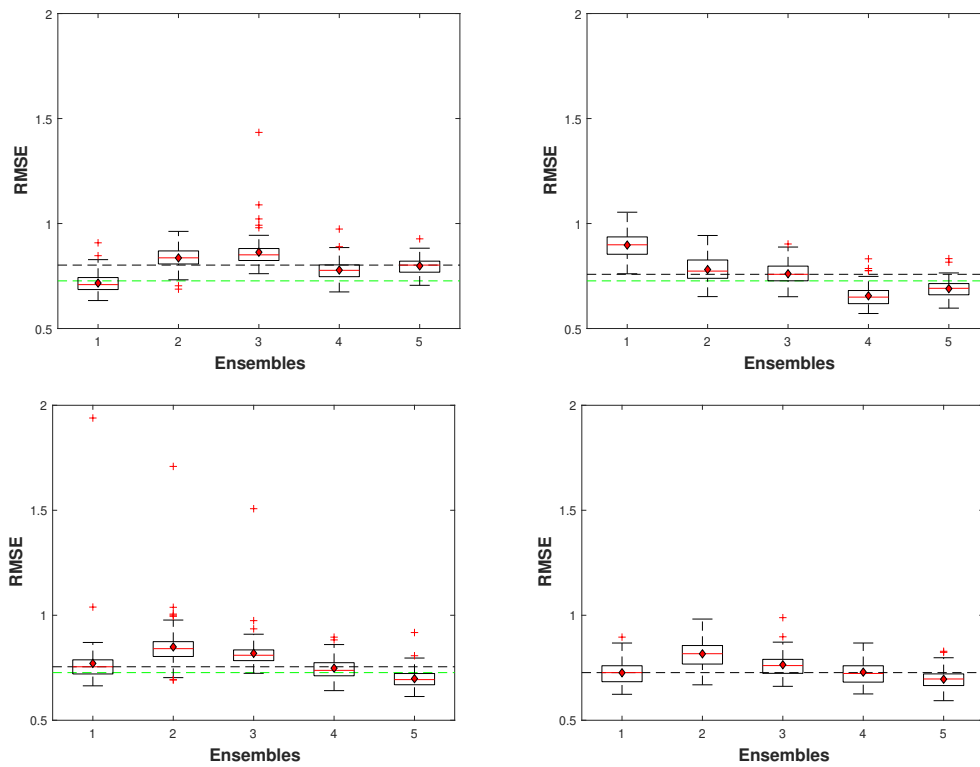
Figure 4.5: Computed RMSE comparing the ES-MDA implementations with ES-MDA-GEO3 with $\mu_\alpha = 1.1$. The top-left figure refers to ES-MDA-EQL; the top-right figure refers to ES-MDA-GEO1; the bottom-left figure refers to ES-MDA-GEO2; the bottom-right figure refers to ES-MDA-GEO3.
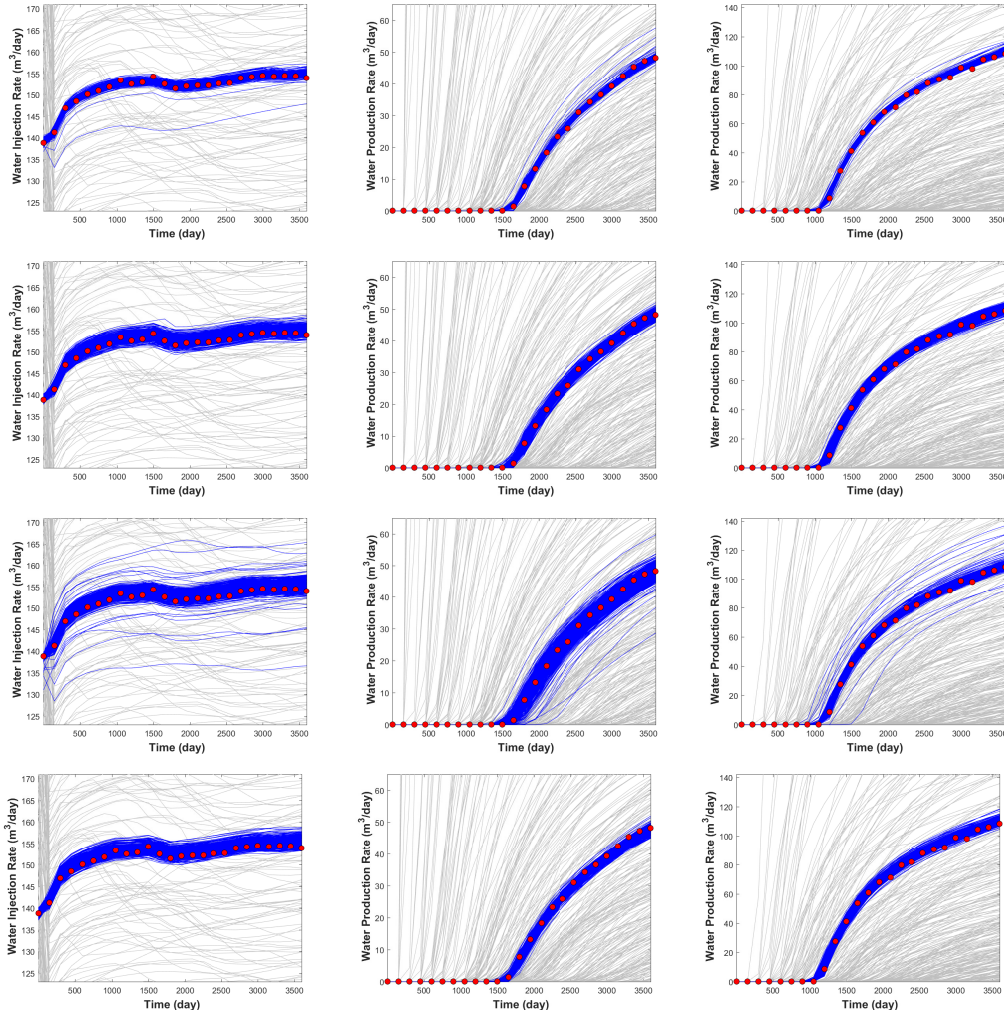
Figure 4.6: Production and injection water rates of wells INJ-4, PROD-6, and PROD-7 with different ES-MDA methods. The first column refers to well INJ-4; the second column refers to well PROD-6; the third column refers to well PROD-7. In contrast, the first line corresponds to ES-MDA-EQL; the second line corresponds to ES-MDA-GEO1; the third line corresponds to ES-MDA-GEO2; the fourth line corresponds to ES-MDA-GEO3.
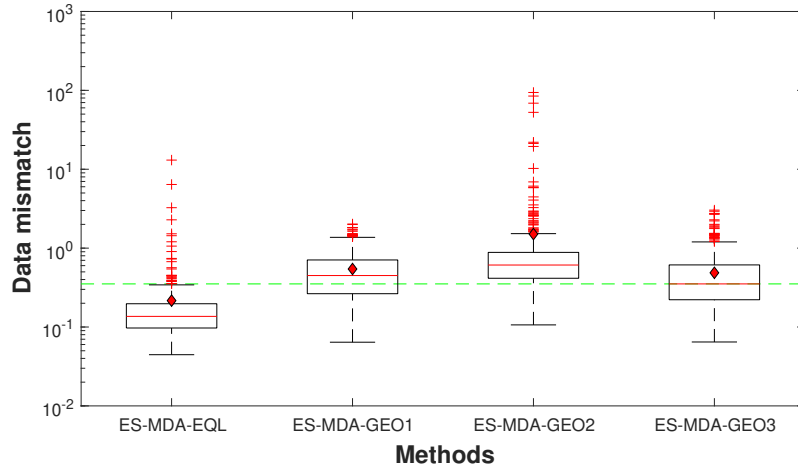
Figure 4.7: Posterior data mismatch for all methods, comparing with ES-MDA-GEO3 with $\mu_\alpha = 1.1$. Boxes, lines, and points here have the same meaning as in Figure 4.5.

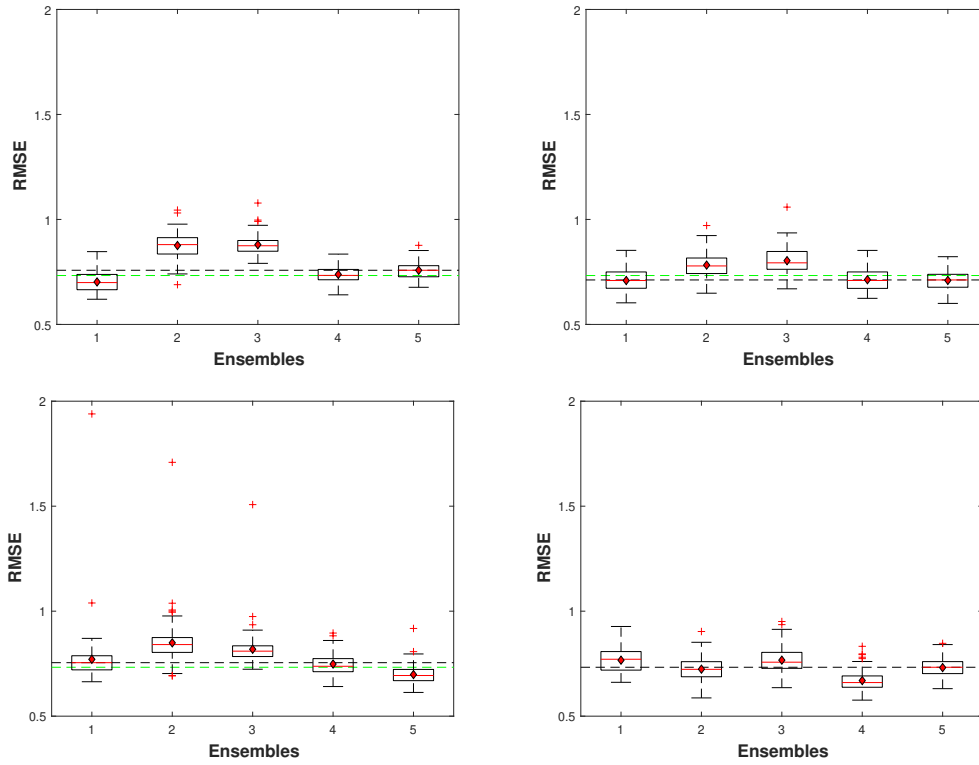Figure 4.8: Computed RMSE comparing the ES-MDA implementations with ES-MDA-GEO3 with $\mu_\alpha = 1.2$. Boxes, lines, and points here have the same meaning as in Figure 4.5.
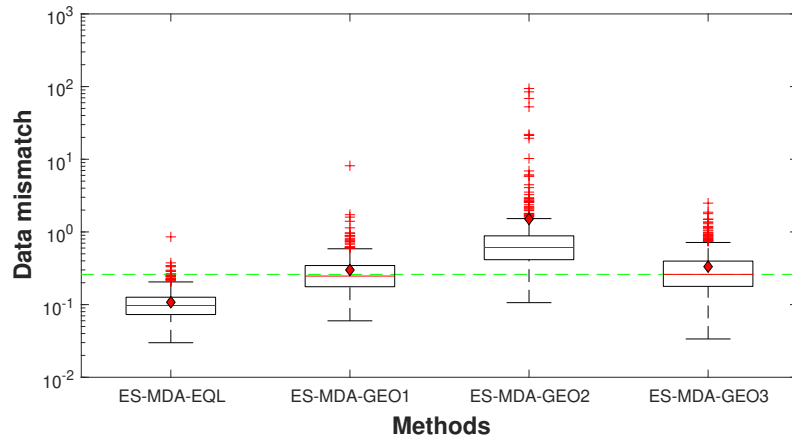
Figure 4.9: Data mismatch for all methods comparing with ES-MDA-GEO3 with $\mu_\alpha = 1.2$. The data mismatch was computed by combining all five ensembles. Boxes, lines, and points here have the same meaning as in Figure 4.5.

# 5
# Summary and Conclusions

In this thesis, we address the problem of generating the inflation factors for the ES-MDA. Recent studies have shown a relationship between the ES-MDA update equation and the solution to the regularized inverse problem. As a result, several numerical procedures that investigate the regularization parameter can be applied to examine the inflation factors and their impacts on the ES-MDA performance. In addition to this observation, the ES-MDA update equation has clear similarities with the Levenberg-Marquardt algorithm. Moreover, the first assimilation of the ES-MDA is similar to a Levenberg-Marquardt update. In this comparison, the inflation factor plays the role of the Levenberg-Marquardt parameter. Again, all mathematical procedures of this area can be used to evaluate the ES-MDA inflation factors.

The ES-MDA has three main parameters: the ensemble size, number of assimilations, and inflation factors. The relation between these parameters is quite unknown in the ensemble-based data assimilation literature. It is understood that a small ensemble may offer bad parameter approximations, and a small number of assimilations may result in poor estimates. However, the relation between the number of assimilations and ensemble size with the inflation factors was unexplored. This thesis's first study investigates this pre-mentioned relation between the inflation factors and the other ES-MDA main parameters. It is presented a numerical intimate connection between the ES-MDA main parameters. We also show that the inflation factors selection may infer the ES-MDA vector of model parameters. We conclude that the method of Hanke produces optimal ES-MDA outcomes, even when the ensemble is little with a small number of assimilations. Selecting the inflation factors equal to the number of assimilations enables error propagation in the multiple data assimilation process due to the effects of small singular vectors in the model update vectors. Finally, when the ensemble is nearly large, e.g., $N_e = 500$, the inflation factors selection becomes almost pointless, as long as Equation (2-52) is satisfied. The results presented in this study were published in [6].

As a consequence of the results of this thesis' first study, concluding that the method of [8] may provide the best ES-MDA performance, we attempt to produce the ES-MDA inflation factors based on that scheme. However, for

the nonlinear approach to the history matching problem, it is only possible to compute the inflation factor *a priori* for the first assimilation step. Using the results presented in [6], it would be desirous to computing all inflation factors based on the scheme of [8]. Assuming that the problem is linear, we propose a new procedure to compute the ES-MDA inflation factors by computing the first and the last inflation factors using the formula proposed by [10]. One explanation for using such a system is the observation that the regularizing parameter generated using the scheme of [8] often decreases within the iterations. Therefore, we attempt to simulate this decrease in a geometric progression by computing the first and the last inflation factors using the formula of [10]. Another explanation comes from the study of [27], which proves that the scheme of [8] is of optimal order. In other words, this scheme achieves optimal accuracy for model parameters. Thus, as the ES-MDA update equation has the same structure as the Levenberg-Marquardt algorithm, we showed with numerical examples that selecting inflation factors proposed by this thesis achieves optimal ES-MDA outcomes. Although the new method is only valid for the linear-Gaussian case, the numerical results show that the proposed methodology is proper for the nonlinear case. The results presented in this chapter were published in [15].

For future works, we intend to study an ES-MDA implementation where the number of assimilations and the inflation factors are selected adaptively. As a result, the only parameter to be determined *a priori* is the ensemble size. This procedure might be useful for history matching problems due to the difficulty of defining the ES-MDA main parameters before data assimilation. The first work of this thesis exposed a technique to assess the quality of the inflation factors before data assimilation starts. However, the other main parameters must be selected *a priori*. We also intend to investigate the effects of the inflation factors on the techniques to mitigate errors in ES-MDA, such as the subspace inversion procedure. This technique is used for parametrizing the sensitivity matrix, neglecting small singular values. Thus, diminishing error propagations and improving uncertainty quantification. As a result of the first study of this thesis, we showed that generating the first inflation factor using the scheme of [8] plays a similar role in data assimilation. Therefore, we believe that subspace inversion might be unnecessary if the inflation factors are selected using the scheme of [8]. However, we still need to investigate this claim thoroughly.

# Bibliography

[1] OLIVER, D. S.; REYNOLDS, A. C. ; LIU, N.. **Inverse theory for petroleum reservoir characterization and history matching**. Cambridge: Cambridge University Press, 2008.

[2] GOLUB, G. H.; LOAN, C. F. V.. **Matrix computations**. Johns Hopkins University Press, 2013.

[3] NOCEDAL, J.; WRIGHT, S.. **Numerical optimization**. Springer-Verlag New York, 2006.

[4] REYNOLDS, A. C.; ZAFARI, M. ; LI, G.. **Iterative forms of the ensemble Kalman filter**. Proceedings of the 10th European Conference on the Mathematics of Oil Recovery, Amsterdam, 4–7 September, 2006.

[5] ZHANG, F.; REYNOLDS, A. C. ; OLIVER, D. S.. **Evaluation of the reduction in uncertainty obtained by conditioning a 3D stochastic channel to multiwell pressure data**. Mathematical Geology, 34:715–742, 2002.

[6] SILVA, T. M. D.; PESCO, S. ; BARRETO JR., A. B.. **Influences of the inflation factors generation in the main parameters of the ensemble smoother with multiple data assimilation**. Journal of Petroleum Sciences and Engineering, 0:1–15, 2021.

[7] LE, D. H.; EMERICK, A. A. ; REYNOLDS, A. C.. **An adaptive ensemble smoother with multiple data assimilation for assisted history matching**. SPE Journal, 21(6):2195–2207, 2016.

[8] HANKE, M.. **A regularizing Levenberg–Marquardt scheme, with applications to inverse groundwater filtration problems**. Inverse Problems, 13(1):79–95, 1997.

[9] EMERICK, A. A.. **Analysis of the performance of ensemble-based assimilation of production and seismic data**. Journal of Petroleum Sciences and Engineering, 139:219–239, 2016.

[10] RAFIEE, J.; REYNOLDS, A. C.. **Theoretical and efficient practical procedures for the generation of inflation factors for ES-MDA**. Inverse Problems, 33(11):1–28, 2017.

[11] EVENSEN, G.. **Analysis of iterative ensemble smoothers for solving inverse problems**. Computational Geosciences, 22:885–908, 2018.

[12] EMERICK, A. A.. **Analysis of geometric selection of the data-error covariance inflation for ES-MDA**. Journal of Petroleum Sciences and Engineering, 182:1–17, 2019.

[13] SILVA, T. M. D.; BELA, R. V.; PESCO, S. ; BARRETO JR., A. B.. **ES-MDA applied to estimate skin zone properties from injectivity tests data in multilayer reservoirs**. Computers & Geosciences, 146:1–15, 2021.

[14] MOROZOV, V. A.. **Methods for solving incorrectly posed problems**. Springer-Verlag New York, 1984.

[15] SILVA, T. M. D.; PESCO, S.; BARRETO JR., A. B. ; ONUR, M.. **A new procedure for generating data covariance inflation factors for ensemble smoother with multiple data assimilation**. Computers & Geosciences, 150:1–14, 2021.

[16] RANAZZI, P. A.; SAMPAIO, M. A.. **Ensemble size investigation in adaptive ES-MDA reservoir history matching**. Journal of the Brazilian Society of Mechanical Sciences and Engineering, 41(413):1–11, 2019.

[17] EMERICK, A. A.; REYNOLDS, A. C.. **Ensemble smoother with multiple data assimilation**. Computers & Geosciences, 55:3–15, 2013.

[18] TAVAKOLI, R.; REYNOLDS, A. C.. **History matching with parameterization based on the singular value decomposition of a dimensionless sensitivity matrix**. SPE Journal, 15(2):495–508, 2010.

[19] WANG, Y.; LI, G. ; REYNOLDS, A. C.. **Estimation of depths of fluid contacts and relative permeability curves by history matching using iterative ensemble-Kalman smoothers**. SPE Journal, 15(2):509–525, 2010.

[20] SHIRANJI, M. G.; EMERICK, A. A.. **An improved TSVD-based Levenberg-Marquardt algorithm for history matching and comparison with Gauss-Newton**. Journal of Petroleum Science and Engineering, 143:258–271, 2016.

[21] KAIPIO, J. P.; SOMERSALO, E.. **Statistical and computational inverse problems**. Springer, 2005.

[22] HAMARIK, U.; PALM, R. ; RAUS, T.. **A family of rules for parameter choice in Tikhonov regularization of ill-posed problems with inexact noise level**. Journal of Computational and Applied Mathematics, 236(8):2146–2157, 2012.

[23] TARANTOLA, A.. **Inverse problem theory and methods for model parameter estimation**. Society for Industrial and Applied Mathematics, 2005.

[24] EMERICK, A. A.; REYNOLDS, A. C.. **Combining sensitivities and prior information for covariance localization in the ensemble Kalman filter for petroleum reservoir applications**. Computational Geosciences, 15:251–269, 2011.

[25] LUO, X.; BHAKTA, T.. **Automatic and adaptive localization for ensemble-based history matching**. Journal of Petroleum Sciences and Engineering, 184:1–18, 2020.

[26] EMERICK, A. A.; REYNOLDS, A. C.. **Investigation of the sampling performance of ensemble-based methods with a simple reservoir model**. Computational Geosciences, 17:325–350, 2013.

[27] HANKE, M.. **The regularizing Levenberg–Marquardt scheme is of optimal order**. Journal of Integral Equations and Applications, 22(2):259–283, 2010.

[28] CANCHUMUNI, S. W. A.; EMERICK, A. A. ; PACHECO, M. A.. **History matching geological facies models based on ensemble smoother and deep generative models**. Journal of Petroleum Sciences and Engineering, 177:941–958, 2019.

[29] IGLESIAS, M. A.; DAWSON, C.. **The regularizing Levenberg–Marquardt scheme for history matching of petroleum reservoirs**. Computational Geosciences, 17:1033–1053, 2013.

[30] IGLESIAS, M. A.. **Iterative regularization for ensemble data assimilation in reservoir models**. Computational Geosciences, 19:177–212, 2015.