

**Discussion on: Application of the variable projection scheme
for frequency-domain full-waveform inversion**

Tristan van Leeuwen¹, Aleksandr Y. Aravkin² and Felix J. Herrmann³

¹*Centrum Wiskunde & Informatica, Amsterdam, The Netherlands.*

²*IBM T.J. Watson Research, USA.*

³*University of British Columbia, Dept. of Earth, Ocean and Atmospheric Sciences,
Vancouver, Canada.*

(January 31, 2014)

Running head:

ABSTRACT

In the paper “Application of the variable projection scheme for frequency-domain full-waveform inversion”, Li et al. (2013) discuss a method for source estimation in the context of frequency-domain full-waveform inversion (FWI). This method is an extension of earlier work by Aravkin et al. (2012b) and Li et al. (2013) suggest similar improvements as concurrent work by Aravkin and van Leeuwen (2012). We have read the work by Li et al. (2013) with great pleasure but feel that their claims are not fully supported by numerical experiments. The goal of this discussion paper is twofold;

1. we would like to clarify some technical details of the extension as presented by Aravkin and van Leeuwen (2012) that are not discussed by Li et al. (2013), and
2. suggest additional numerical experiments to support the claims made by Li et al. (2013).

The remainder of this section contains the main discussion and refers to the appendices for more (technical) details and numerical results.

Li et al. (2013) discuss an automated procedure for estimating the source-weights in frequency domain FWI. The main idea is that these source-weights can be estimated on-the-fly during each iteration of FWI using a technique known as “variable projection”. A brief mathematical description of the problem is given in Appendix A (under “Problem statement”). The use of the variable projection method for source-estimation in FWI has a long history, as was noted by Li et al. (2013). As far as we know, it was first recognized by Aravkin et al. (2012b) that source-estimation is indeed an instance of variable projection, and this insight was used to generalize the concept of source-estimation to generic misfit functions.

Li et al. (2013) state that, when using the variable projection technique to compute the source-weights, the Jacobian¹ and Hessian contain correction terms which were not computed previously. They correctly point out that ignoring these correction terms does not affect the calculation of the gradient (i.e., multiplication of the adjoint of the Jacobian with the residual) so that gradient-descent methods can be applied without computing this correction term. This is exactly the regime discussed by Aravkin et al. (2012b), and therefore there are no gradient approximations used in this paper, contrary to what is suggested by Li et al. (2013). For the sake of completeness we have included a derivation in Appendix A (under “Computing the gradient”).

The Hessian does contain a correction term however, and explicit expressions of these terms are presented by Aravkin and van Leeuwen (2012). The variable projection method was also included in a frequency-domain waveform inversion method by van Leeuwen and Herrmann (2013), where an L-BFGS approximation of the *full* Hessian was used, which automatically includes the correction terms.

The expressions presented by Li et al. (2013) amount to a Gauss-Newton approximation of these correction terms. We include a derivation in Appendix A (under “Computing the Hessian”), along with an explicit Schur-complement expression for the reduced Hessian that was not presented in Aravkin and van Leeuwen (2012) nor by Li et al. (2013). From this expression, it is evident that ignoring correction terms amounts to *adding* a positive semidefinite correction term to the true Hessian, which is a common optimization technique (recall e.g. the Levenberg-Marquardt method). Hence, properly designed optimization strategies can ignore the correction terms, and still satisfy standard convergence guarantees (i.e. convergence to stationary points). The effect of including correction terms (that

¹which Li et al. (2013) refer to as “the gradient” in the introduction.

can possibly make the Hessian indefinite) is problem-specific, and in some cases, ignoring negative definite corrections in the context of Gauss-Newton methods can improve algorithm performance (Aravkin et al., 2012a).

The main contribution of Li et al. (2013) is the computation of a Gauss-Newton approximation of the correction terms in the Hessian. However, Li et al. (2013) do not compare the performance of their algorithm with a ‘standard’ approach that does not include the correction terms. Thus, the actual benefit of including the correction terms is yet to be investigated. To shed some light on this issue we present some numerical experiments on a toy-model depicted in 1, using both a reflection and a transmission setup. More details can be found in Appendix B. First, we compute the effect of the correction term by plotting the actual misfit as well a quadratic model based on the Hessian in figure 2. Secondly, we compare the use of an L-BFGS method to a Gauss-Newton method with and without correction terms. These experiments suggest that the correction terms can indeed affect the rate of convergence and the final results (see figures 3 and 4). However, the differences observed on two toy-problems are small and only occur when reducing the misfit by 3 orders of magnitude, a situation that rarely occurs in practice. The Matlab code used to conduct the numerical experiment can be found here: <https://github.com/tleeuwen/Variable-Projection-for-FWI.git>.

In conclusion, the correction terms suggested by Li et al. (2013) are a particular approximation of the full Hessian that follows from the variable projection approach (Aravkin and van Leeuwen, 2012). Numerical experiments on a toy problem suggest that including these terms may improve convergence but might not influence the final result dramatically. We expect the influence of the correction terms to be problem-dependent and we would like to encourage the authors to repeat this experiment with their code on a more realistic problem

and include the results in a reply.

APPENDIX A

THEORY

Problem statement

Li et al. (2013) consider the following optimization problem²

$$\min_{\mathbf{m}, \mathbf{c}} \Phi(\mathbf{m}, \mathbf{c}) = \sum_{k,j} \|\mathbf{d}_{k,j} - c_{k,j} \mathbf{s}_{k,j}(\mathbf{m})\|_2^2, \quad (\text{A-1})$$

where $\mathbf{d}_{k,j}$ is the observed data (organized in a vector) for the k^{th} frequency and j^{th} source, $\mathbf{s}_{k,j}(\mathbf{m})$ is the corresponding synthetic data (organized in a vector) for model \mathbf{m} and $c_{k,j}$ are the (possible complex) source weights, which are assembled in the vector \mathbf{c} . These source weights need to be estimated to calibrate the modelling code. Using the variable projection method entails projecting out the source-weights by solving

$$\min_{\mathbf{c}} \Phi(\mathbf{m}, \mathbf{c}), \quad (\text{A-2})$$

the solution of which we denote by $\bar{\mathbf{c}}(\mathbf{m})$. Note that the optimal source weight is implicitly defined by

$$\nabla_{\mathbf{c}} \Phi(\mathbf{m}, \bar{\mathbf{c}}) = 0. \quad (\text{A-3})$$

Substituting the optimal source-weight back into the objective we obtain a *reduced* objective in \mathbf{m} alone, written

$$\bar{\Phi}(\mathbf{m}) = \Phi(\mathbf{m}, \bar{\mathbf{c}}(\mathbf{m})). \quad (\text{A-4})$$

²we ignore the regularization and weighting introduced in Li et al. (2013) as this does not have an impact on the source estimation.

Next, we discuss the following three main algorithmic ingredients of the variable projection approach:

1. a method for obtaining $\bar{\mathbf{c}}(\mathbf{m})$,
2. an explicit expression for the gradient of $\bar{\Phi}(\mathbf{m})$, and
3. an explicit expression for the Hessian of $\bar{\Phi}(\mathbf{m})$.

In order to optimize the objective, one needs only the first two items, since this enables methods such as gradient-descent or L-BFGS. The last item (or one of several possible approximations) lets us use different (Gauss-) Newton variants.

Obtaining the source weights

The main driver behind the variable projection technique is that the inner optimization problem (A-2) is easy to solve. In this case, there is of course a closed-form solution:

$$c_{k,j} = \mathbf{s}_{k,j}^* \mathbf{d}_{k,j} / \|\mathbf{s}_{k,j}\|_2^2, \quad (\text{A-5})$$

where \cdot^* denotes the complex-conjugate transpose. See also equation (11) in Li et al. (2013). However, even if the problem does not have a closed-form solution (when using other misfit penalties, for example) one can devise an efficient procedure (e.g. by solving a number of independent scalar optimization problems in parallel) to obtain $c_{k,j}$. Some examples are discussed by Aravkin and van Leeuwen (2012) and Aravkin et al. (2012b).

Computing the gradient

Computation of the gradient can be done via the chainrule, yielding:

$$\nabla_{\mathbf{m}} \bar{\Phi}(\mathbf{m}) = \nabla_{\mathbf{m}} \Phi(\mathbf{m}, \bar{\mathbf{c}}) + \nabla_{\mathbf{c}} \Phi(\mathbf{m}, \bar{\mathbf{c}}) \nabla_{\mathbf{m}} \bar{\mathbf{c}}, \quad (\text{A-6})$$

from which it is immediately clear that the second term vanishes if we have really solved $\nabla_{\mathbf{c}}\Phi = 0$ to obtain $\bar{\mathbf{c}}$. In the quadratic case (i.e. with closed form solution given by Eq. (A-5)) this is certainly true and we have

$$\nabla_{\mathbf{m}}\bar{\Phi}(\mathbf{m}) = \nabla_{\mathbf{m}}\Phi(\mathbf{m}, \bar{\mathbf{c}}). \quad (\text{A-7})$$

Even if there is no closed-form solution for $\bar{\mathbf{c}}$, we can solve $\nabla_{\mathbf{c}}\Phi = 0$ up to arbitrary precision and safely ignore the correction term.

The same argument can be made by computing the Jacobian of the modelled data $\bar{c}_{k,j}\mathbf{s}_{k,j}$, which consists of two parts:

$$\frac{\partial \bar{c}_{k,j}(\mathbf{m})\mathbf{s}_{k,j}(\mathbf{m})}{\partial \mathbf{m}} = \bar{c}_{k,j}(\mathbf{m})\frac{\partial \mathbf{s}_{k,j}(\mathbf{m})}{\partial \mathbf{m}} + \mathbf{s}_{k,j}\frac{\partial \bar{c}_{k,j}}{\partial \mathbf{m}}. \quad (\text{A-8})$$

Compare this expressions with Li et al. (2013), equation (13). The main point of the original paper on variable projection by Golub and Pereyra (1973) is that the residual $\mathbf{d}_{k,j} - \bar{c}_{k,j}\mathbf{s}_{k,j}(\mathbf{m})$ is in the null-space of the adjoint of the second term in the Jacobian, so that this second term vanishes when computing the gradient. Indeed,

$$\left(\frac{\partial \bar{c}_{k,j}}{\partial \mathbf{m}}\right)^* \mathbf{s}_{k,j}^*(\mathbf{d}_{k,j} - \bar{c}_{k,j}\mathbf{s}_{k,j}) = 0, \quad (\text{A-9})$$

as can be verified by substituting the explicit expression for the optimal source-weight presented earlier. This was also noted by Rickett (2013).

Computing the Hessian

The Hessian can be computed by applying the chainrule to the gradient:

$$\nabla_{\mathbf{m}}^2\bar{\Phi}(\mathbf{m}) = \nabla_{\mathbf{m}}^2\Phi(\mathbf{m}, \bar{\mathbf{c}}) + \nabla_{\mathbf{m},\mathbf{c}}^2\Phi(\mathbf{m}, \bar{\mathbf{c}})\nabla_{\mathbf{m}}\bar{\mathbf{c}}, \quad (\text{A-10})$$

where $\nabla_{\mathbf{m},\mathbf{c}}^2$ denotes the mixed second derivative w.t.r. \mathbf{m} and \mathbf{c} . In this case none of the terms vanish, and we have to compute $\nabla_{\mathbf{m}}\bar{\mathbf{c}}$ in order to obtain the exact Hessian. These

expressions were presented in Aravkin and van Leeuwen (2012), and we suggested using only the first term as an approximation.

To better understand this approximation from an optimization point of view, we compute this term by considering the function $F(\mathbf{m}) = \nabla_{\mathbf{c}}\Phi(\mathbf{m}, \bar{\mathbf{c}})$, and computing its gradient:

$$\nabla_{\mathbf{m}}F(\mathbf{m}) = \nabla_{\mathbf{m},\mathbf{c}}^2\Phi(\mathbf{m}, \bar{\mathbf{c}}) + \nabla_{\mathbf{c}}^2\Phi(\mathbf{m}, \bar{\mathbf{c}})\nabla_{\mathbf{m}}\bar{\mathbf{c}}. \quad (\text{A-11})$$

Using the stationarity argument we find that $F(\mathbf{m}) = 0$, and hence $\nabla_{\mathbf{m}}F(\mathbf{m}) = 0$, yielding

$$\nabla_{\mathbf{m},\mathbf{c}}^2\Phi(\mathbf{m}, \bar{\mathbf{c}}) = -\nabla_{\mathbf{c}}^2\Phi(\mathbf{m}, \bar{\mathbf{c}})\nabla_{\mathbf{m}}\bar{\mathbf{c}}, \quad (\text{A-12})$$

which gives us

$$\nabla_{\mathbf{m}}\bar{\mathbf{c}} = -\left(\nabla_{\mathbf{c}}^2\Phi\right)^{-1}\nabla_{\mathbf{m},\mathbf{c}}^2\Phi. \quad (\text{A-13})$$

Finally, the Hessian of the reduced objective can be written as

$$\nabla_{\mathbf{m}}^2\bar{\Phi}(\mathbf{m}) = \nabla_{\mathbf{m}}^2\Phi - \nabla_{\mathbf{m},\mathbf{c}}^2\Phi\left(\nabla_{\mathbf{c}}^2\Phi\right)^{-1}\nabla_{\mathbf{m},\mathbf{c}}^2\Phi. \quad (\text{A-14})$$

A similar expression is also presented by Ruhe, A. and Wedin, P. A. (1980). Note also that this is precisely the Schur complement of $\bar{\mathbf{c}}$ in the full Hessian of $\Phi(\mathbf{m}, \mathbf{c})$, and equation (A-14) makes it very clear that ignoring the correction term amounts to a positive definite modification to $\nabla_{\mathbf{m}}^2\bar{\Phi}(\mathbf{m})$, or to its Gauss-Newton approximation, depending on the algorithm. Such a modification may have a positive rather than detrimental effect on the algorithm, as the correction term may actually render the Hessian indefinite, possibly causing the Gauss-Newton method to diverge. We therefore encourage its further investigation both from a theoretical and numerical perspective.

APPENDIX B

NUMERICAL EXPERIMENTS

In order to shed some preliminary light on the effect of the correction term, we conducted a few numerical experiments on a toy problem. The Matlab code used to perform the experiments can be found at: <https://github.com/tleeuwen/Variable-Projection-for-FWI>.
git.

The simulated data is given by

$$\mathbf{s}_{k,j}(\mathbf{m}) = PA_k(\mathbf{m})^{-1}\mathbf{q}_j, \quad (\text{B-1})$$

where $A_k = \omega_k^2 \text{diag}(\mathbf{m}) + \nabla^2$ is a 5-point finite-difference discretization of the Helmholtz operator (with absorbing boundary conditions) for the k^{th} frequency, \mathbf{q}_j is the j^{th} source function and P is the detection operator that samples the wavefield at the receiver locations.

The Jacobian is given by

$$J_{k,j} = \frac{\partial \mathbf{s}_{k,j}(\mathbf{m})}{\partial \mathbf{m}} = PA_k(\mathbf{m})^{-1}G_k(\mathbf{m})\text{diag}(\mathbf{u}_{k,j}), \quad (\text{B-2})$$

where $\mathbf{u}_{k,j} = A_k(\mathbf{m})^{-1}\mathbf{q}_j$ and $G_k(\mathbf{m})$ contains derivatives of A_k w.r.t. \mathbf{m} . For the experiments we ignore the contributions of the boundaries to the gradient, in which case $G_k(\mathbf{m}) = \omega_k^2 \text{diag}(\mathbf{w})$, where $\mathbf{w} = 1$ on the internal gridpoints and $\mathbf{w} = 0$ on the boundary.

The gradient of the source-weight is given by

$$\frac{\partial \bar{c}_{k,j}}{\partial \mathbf{m}} = \frac{1}{\|\mathbf{s}_{k,j}\|_2^2} J_{k,j}^* (\mathbf{d}_{k,j} - 2\bar{c}_{k,j}\mathbf{s}_{k,j}). \quad (\text{B-3})$$

Then, the Jacobian of the reduced objective is given by (cf. equation (A-8))

$$\bar{J}_{k,j} = \bar{c}_{k,j}J_{k,j} + \mathbf{s}_{k,j} \frac{\partial \bar{c}_{k,j}}{\partial \mathbf{m}}. \quad (\text{B-4})$$

The gradient and Hessian are now given by

$$\mathbf{g} = \sum_{k,j} \bar{J}_{k,j}^* (\mathbf{s}_{k,j} - \bar{c}_{k,j} \mathbf{s}_{k,j}), \quad (\text{B-5})$$

$$H = \sum_{k,j} \bar{J}_{k,j}^* \bar{J}_{k,j}. \quad (\text{B-6})$$

Toy problem

For the following experiments we use the velocity model depicted in figure 1. We use either a reflection or a transmission (cross-well) setup. For all the experiments we use a gridspacing of 20 m, three frequencies: $\{3, 5, 10\}$ Hz and 49 equispaced sources and receivers. The observed data is weighted with a random amplitude factor for each source and frequency. For all the experiments the initial model \mathbf{m}_0 is constant with a velocity of 2000 m/s.

Experiment 1

In order to asses how much the correction terms contribute to the Hessian we consider the quadratic model:

$$q(\mathbf{m}_0, \Delta\mathbf{m}) = \bar{\Phi}(\mathbf{m}_0) + \Delta\mathbf{m}^T \mathbf{g}(\mathbf{m}_0) + \frac{1}{2} \Delta\mathbf{m}^T H(\mathbf{m}_0) \Delta\mathbf{m} \quad (\text{B-7})$$

and compare this to the actual misfit $\bar{\Phi}(\mathbf{m}_0 + \Delta\mathbf{m})$. Figure 2 shows the quadratic models $q(\mathbf{m}_0, \alpha\Delta\mathbf{m})$ with and without correction terms as well as the actual misfit $\bar{\Phi}(\mathbf{m}_0 + \alpha\Delta\mathbf{m})$, with $\Delta\mathbf{m} = \nabla\bar{\Phi}(\mathbf{m}_0)/\|\nabla\bar{\Phi}(\mathbf{m}_0)\|_2$ as a function of α . For this particular setup, the correction terms do not contribute significantly. In the next experiment we investigate how the correction term affects the convergence when using a Gauss-Newton method.

Experiment 2

We solve the optimization problem by iteratively updating the model

$$\mathbf{m}_{i+1} = \mathbf{m}_i + \alpha_k \Delta \mathbf{m}_i. \tag{B-8}$$

The search direction $\Delta \mathbf{m}_i$ is either determined by applying the L-BFGS inverse Hessian to $-\mathbf{g}_i$ using a two-loop recursion (cf. Nocedal and Wright (2000), Algorithm 7.4), or by solving $H_i \Delta \mathbf{m}_i = -\mathbf{g}_i$ up to a prescribed relative tolerance (10^{-1} in this case) using CG, were, H_i can either be the conventional GN Hessian or contain the correction terms. We use a weak Wolfe linesearch to determine α_i and stop the iterations when the norm of the gradient drops below a preset relative tolerance (10^{-3} in this case). We compute the cost of the inversion in terms of the number of PDE solves; 2 PDE solves for a gradient computation and 3 for a Hessian-vector product.

Figure 3 shows the results for the transmission experiment. Here, the correction terms hardly affect the rate of convergence, although the results look slightly better when including the correction terms. For the reflection experiment, shown in figure 4, the correction terms significantly affect the rate of convergence but have little effect on the final reconstruction. In both cases the LBFGS results are nearly identical to the results obtained with the GN method with correction term, however, it is significantly cheaper (in terms of the number of PDE solves) as can be seen from Table 1.

ACKNOWLEDGEMENTS

This work was in part financially supported by the Natural Sciences and Engineering Research Council of Canada Discovery Grant (22R81254) and the Collaborative Research and Development Grant DNOISE II (375142-08). This research was carried out as part of

the SINBAD II project with support from the following organizations: BG Group, BGP, BP, CGG, Chevron, ConocoPhillips, ION, Petrobras, PGS, Total SA, WesternGeco, and Woodside.

REFERENCES

- Aravkin, A., J. Burke, and G. Pillonetto, 2012a, Robust and trend following kalman smoothers using student's t: Presented at the Proc. of SYSID.
- Aravkin, A., T. van Leeuwen, H. Calandra, and F. Hermann, 2012b, Source estimation for frequency-domain FWI with robust penalties: Presented at the EAGE expanded abstracts.
- Aravkin, A. Y., and T. van Leeuwen, 2012, Estimating nuisance parameters in inverse problems: *Inverse Problems*, **28**, 115016.
- Golub, G., and V. Pereyra, 1973, The differentiation of pseudo-inverses and nonlinear least squares problems whose variables separate: *SIAM Journal of Numerical Analysis*, **10**, 413–432.
- Li, M., J. Rickett, and A. Abubakar, 2013, Application of the variable projection scheme for frequency-domain full-waveform inversion: *Geophysics*, **78**.
- Nocedal, J., and S. Wright, 2000, *Numerical optimization*: Springer. Springer Series in Operations Research.
- Rickett, J., 2013, The variable projection method for waveform inversion with an unknown source function: *Geophysical Prospecting*, **61**, 874–881.
- Ruhe, A. and Wedin, P. A., 1980, Algorithms for Separable Nonlinear Least Squares Problems: *SIAM review*, **22**, 318–337.
- van Leeuwen, T., and F. J. Herrmann, 2013, Fast waveform inversion without source-encoding: *Geophysical Prospecting*, **61**, 10–19.

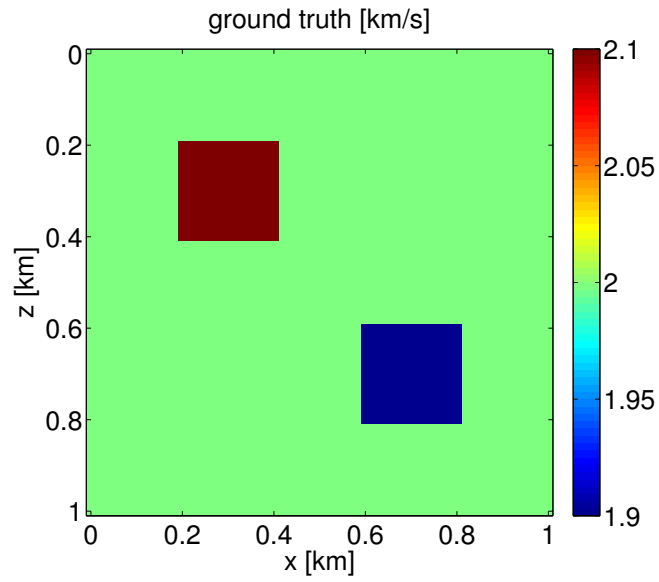


Figure 1: Velocity model used for experiments.

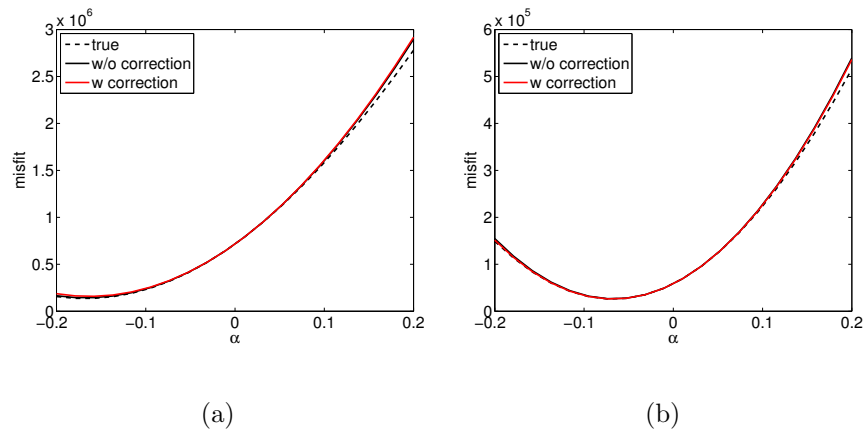


Figure 2: Actual misfit and quadratic model with and without correction terms for transmission (a) and reflection (b) configuration.

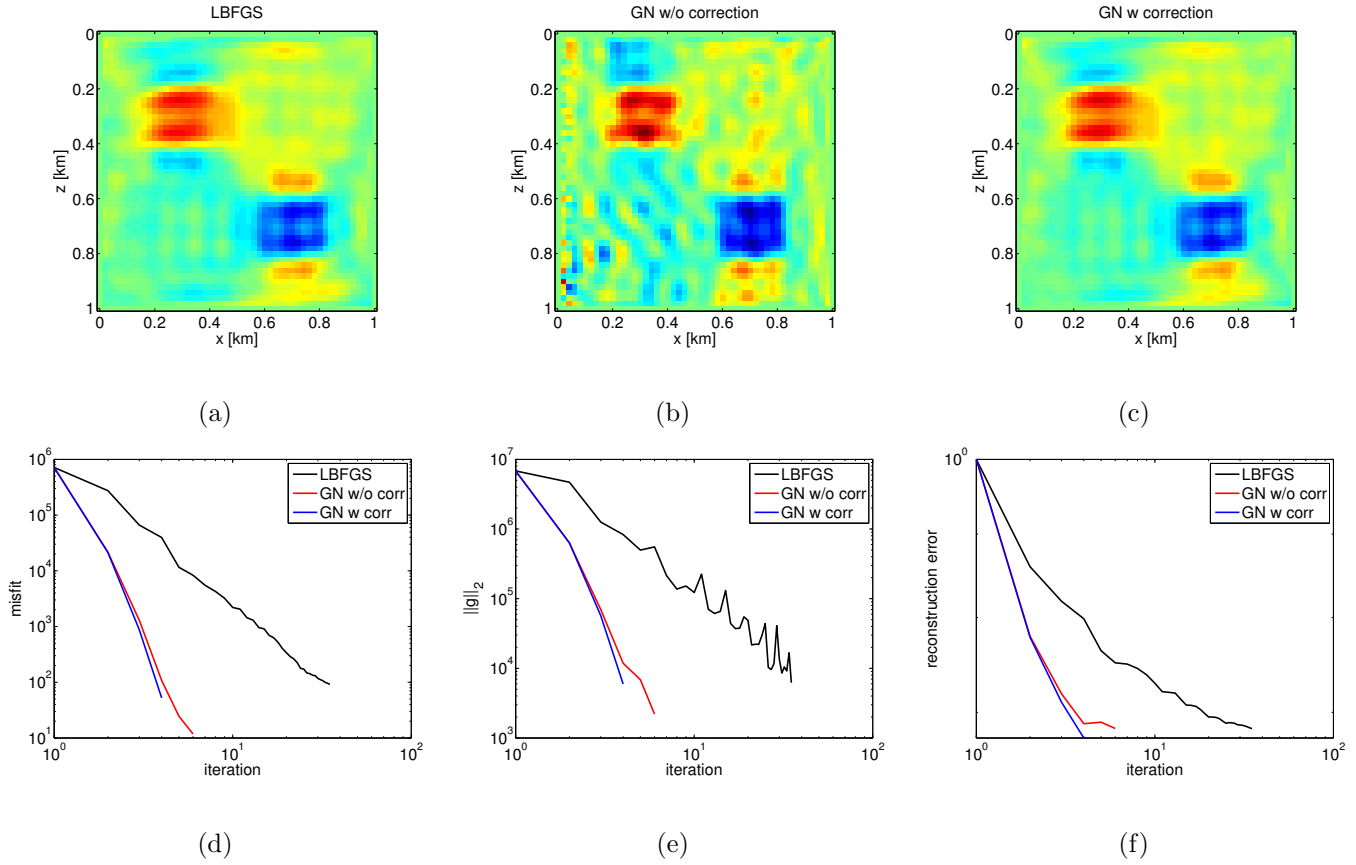


Figure 3: Inversion result for transmission configuration. The reconstructions with L-BFGS and GN with and without correction terms are shown in (a–c). The convergence histories in terms of the misfit, norm of the gradient and reconstruction error ($\|\mathbf{m}_k - \mathbf{m}_{\text{true}}\|_2 / \|\mathbf{m}_0 - \mathbf{m}_{\text{true}}\|$) are shown in (d–f).

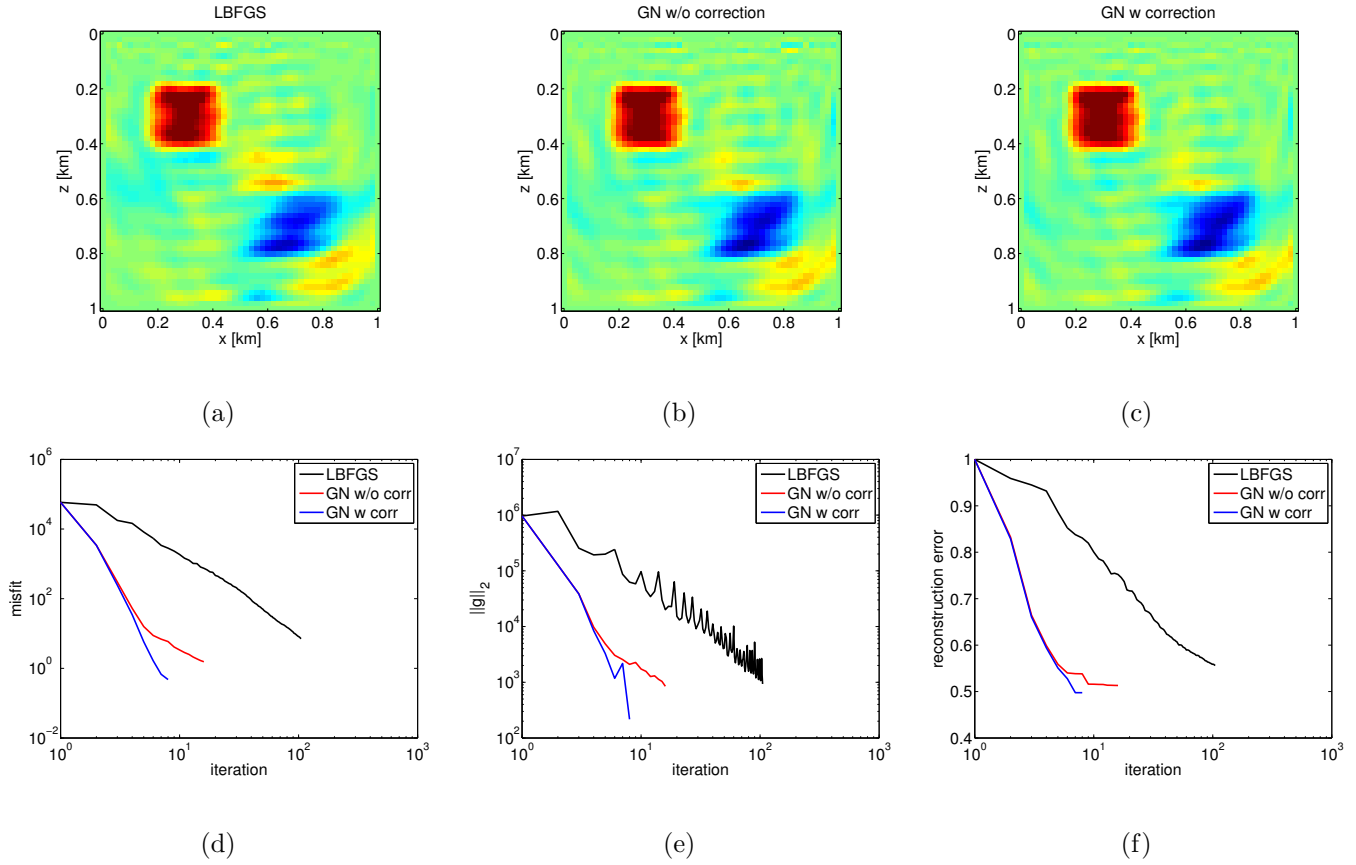


Figure 4: Inversion result for reflection configuration. The reconstructions with L-BFGS and GN with and without correction terms are shown in (a–c). The convergence histories in terms of the misfit, norm of the gradient and reconstruction error ($\|\mathbf{m}_k - \mathbf{m}_{\text{true}}\|_2 / \|\mathbf{m}_0 - \mathbf{m}_{\text{true}}\|$) are shown in (d–f).

Transmission	L-BFGS	GN w/o corr.	GN w corr.
Iterations	34	5	3
PDE solves	76	477	104

(a)

Reflection	L-BFGS	GN w/o corr.	GN w corr.
Iterations	104	15	7
PDE solves	222	1235	889

(b)

Table 1: Overview of iteration counts and computational cost in terms of the number of PDE-solves for the transmission (a) and reflection (b) experiments.