

An efficient data-subspace inversion method for 2-D magnetotelluric data

Weerachai Siripunvaraporn* and Gary Egbert[‡]

ABSTRACT

There are currently three types of algorithms in use for regularized 2-D inversion of magnetotelluric (MT) data. All seek to minimize some functional which penalizes data misfit and model structure. With the most straightforward approach (exemplified by OCCAM), the minimization is accomplished using some variant on a linearized Gauss-Newton approach. A second approach is to use a descent method [e.g., nonlinear conjugate gradients (NLCG)] to avoid the expense of constructing large matrices (e.g., the sensitivity matrix). Finally, approximate methods [e.g., rapid relaxation inversion (RRI)] have been developed which use cheaply computed approximations to the sensitivity matrix to search for a minimum of the penalty functional. Approximate approaches can be very fast, but in practice often fail to converge without significant expert user intervention. On the other hand, the more straightforward methods can be prohibitively expensive to use for even moderate-size data sets. Here, we present a new and much more efficient variant on the OCCAM scheme. By expressing the solution as a linear combination of rows of the sensitivity matrix smoothed by the model covariance (the “representers”), we transform the linearized inverse

problem from the M -dimensional model space to the N -dimensional data space. This method is referred to as DASOCC, the data space OCCAM’s inversion. Since generally $N \ll M$, this transformation by itself can result in significant computational saving. More importantly the data space formulation suggests a simple approximate method for constructing the inverse solution. Since MT data are smooth and “redundant,” a subset of the representers is typically sufficient to form the model without significant loss of detail. Computations required for constructing sensitivities and the size of matrices to be inverted can be significantly reduced by this approximation. We refer to this inversion as REBOCC, the reduced basis OCCAM’s inversion. Numerical experiments on synthetic and real data sets with REBOCC, DASOCC, NLCG, RRI, and OCCAM show that REBOCC is faster than both DASOCC and NLCG, which are comparable in speed. All of these methods are significantly faster than OCCAM, but are not competitive with RRI. However, even with a simple synthetic data set, we could not always get RRI to converge to a reasonable solution. The basic idea behind REBOCC should be more broadly applicable, in particular to 3-D MT inversion.

INTRODUCTION

The magnetotelluric (MT) method for imaging crustal and upper mantle electrical conductivity has found increasing use in both geophysical exploration application (Orange, 1989; Vozoff, 1972) and in fundamental studies of large-scale tectonics (Jones, 1992; Wannamaker et al., 1994; Chen et al., 1996; Unsworth et al., 1999). Initial applications of MT were based on local 1-D interpretations, for which theories (Weidelt, 1972; Parker, 1980) and inversion methods (Jupp and Vozoff, 1975;

Constable et al., 1987; and Smith and Booker, 1988) are well developed. It is now clear that 2-D or even 3-D interpretation is essential for most real MT data sets. Over the past decade, very substantial progress has been made on the development of 2-D inversion methods. These have included straight-forward extensions of linearized search methods developed previously for 1-D regularized inversion (deGroot-Hedlin and Constable, 1990; Uchida, 1993), efficient approximate methods (Smith and Booker, 1991; Farquharson and Oldenburg, 1996), the subspace method (Oldenburg et al., 1993) and methods based on

Manuscript received by the Editor January 19, 1999; revised manuscript received September 17, 1999.

*Formerly Oregon State University, College of Oceanic and Atmospheric Science, Corvallis, Oregon 97331-5503; presently Dept. of Physics, Faculty of Science, Mahidd University, Rama VI Road, Rachathavek, Bangkok 10400, Thailand. E-mail: wsiripun@oce.orst.edu.

‡Oregon State University, College of Oceanic and Atmospheric Science, 104 Ocean Admin. Bldg., Corvallis, Oregon 97331-5503. E-mail: egbert@oce.orst.edu.

© 2000 Society of Exploration Geophysicists. All rights reserved.

direct iterative minimization of a regularized penalty functional (Rodi and Mackie, 2000). These programs are freely available and are widely used for interpretation of MT surveys. However, the available inversion algorithms are not without flaws. The fastest [e.g., rapid relaxation inversion (RRI)] intrinsically limit the model space search and can fail to converge (without substantial user intervention). The most general and flexible (e.g., OCCAM) run very slowly and require considerable computer memory.

We begin our discussion with a review of 2-D MT inversion methods, including a data space variant on OCCAM which we refer to as DASOCC, the data space Occam's inversion. We then describe the REBOCC (reduced basis Occam's inversion) algorithm, and demonstrate the stability and efficiency of the approach by inverting synthetic MT data. We use the synthetic data to compare the performance of the new scheme to some of the other methods, including RRI (Smith and Booker, 1991), the nonlinear conjugate gradient method (NLGG) (Rodi and Mackie, 2000), and the original Occam's inversion (OCCAM) (deGroot-Hedlin and Constable, 1990). Finally, we briefly consider an example application to field data from a dense MT profile across the San Andreas fault near Parkfield, California (Unsworth et al., 1997).

OVERVIEW OF INVERSION METHODS

We begin with a broad overview of previously developed 2-D MT inversion methods. To be explicit, we consider the earth as discretized into a series of M constant resistivity blocks, $\mathbf{m} = [m_1, m_2, \dots, m_M]$. There are N observed data $\mathbf{d} = [d_1, d_2, \dots, d_N]$, with estimated uncertainties $\mathbf{e} = [e_1, e_2, \dots, e_N]$. The fit of the theoretical model responses $\mathbf{F}[\mathbf{m}]$ to the observational data can be expressed as

$$X_d^2 = (\mathbf{d} - \mathbf{F}[\mathbf{m}])^T \mathbf{C}_d^{-1} (\mathbf{d} - \mathbf{F}[\mathbf{m}]), \quad (1)$$

where the superscript T represents matrix transpose, and \mathbf{C}_d is the data covariance matrix, which in practice is diagonal.

Because of the nonuniqueness of the inverse problem, an infinite number of models can produce the same misfit. Most modern MT inversion schemes resolve this nonuniqueness by seeking models that have minimum possible structure (in some sense) for a given level of misfit (Parker, 1994). This makes the inversion stable, with resulting models less likely to contain spurious features.

To quantify "model structure," we consider a model norm of the general form

$$X_m^2 = (\mathbf{m} - \mathbf{m}_0)^T \mathbf{C}_m^{-1} (\mathbf{m} - \mathbf{m}_0), \quad (2)$$

where \mathbf{m}_0 is a base (or prior) model, and \mathbf{C}_m is a model covariance matrix which characterizes the expected magnitude and smoothness of resistivity variations relative to \mathbf{m}_0 . Other approaches to minimum structure inversion are similar, though in some cases \mathbf{C}_m^{-1} is replaced by a model roughness operator. The minimum structure inverse problem is to minimize X_m^2 subject to $X_d^2 = X_*^2$, where X_*^2 is the desired level of misfit.

To solve this minimization problem, a Lagrange multiplier λ^{-1} can be introduced, resulting in an unconstrained func-

tional $U(\mathbf{m}, \lambda)$,

$$U(\mathbf{m}, \lambda) = (\mathbf{m} - \mathbf{m}_0)^T \mathbf{C}_m^{-1} (\mathbf{m} - \mathbf{m}_0) + \lambda^{-1} \{ (\mathbf{d} - \mathbf{F}[\mathbf{m}])^T \mathbf{C}_d^{-1} (\mathbf{d} - \mathbf{F}[\mathbf{m}]) - X_*^2 \}, \quad (3)$$

for which we seek stationary points (with respect to both \mathbf{m} and λ). Alternatively, we may consider the penalty functional $W_\lambda(\mathbf{m})$,

$$W_\lambda(\mathbf{m}) = (\mathbf{m} - \mathbf{m}_0)^T \mathbf{C}_m^{-1} (\mathbf{m} - \mathbf{m}_0) + \lambda^{-1} \{ (\mathbf{d} - \mathbf{F}[\mathbf{m}])^T \mathbf{C}_d^{-1} (\mathbf{d} - \mathbf{F}[\mathbf{m}]) \}. \quad (4)$$

In equation (4), λ acts to "trade off" between minimizing the norm of data misfit and the norm of the model (Tikhonov and Arsenin, 1977; Parker, 1994). When λ is large, the data misfit is de-emphasized, leading to a smoother model. In contrast, as $\lambda \rightarrow 0$, the inverse problem becomes closer to the ill conditioned least-squares inversion problem, resulting in an erratic model (see Parker, 1980).

Note that both U and W_λ have the same stationary points with respect to variations of the model, i.e., $\partial U / \partial \mathbf{m} = \partial W_\lambda / \partial \mathbf{m}$, where λ is fixed. Parker (1994) uses this to show that stationary points of equation (3) can be found by minimizing equation (4) for a series of λ values, and then choosing λ so that the misfit satisfies the constraint $X_d^2 = X_*^2$.

For linear $\mathbf{F}[\mathbf{m}]$, this is straightforward, since in this case (for fixed λ) $\partial U / \partial \mathbf{m} = 0$ is a linear system of equations which may be solved for \mathbf{m} . Because $\mathbf{F}[\mathbf{m}]$ is nonlinear for the MT inverse problem, iterative solution methods are required. We briefly consider some of the approaches which have been taken by previous workers and then outline our approach.

Rodi and Mackie (1999) provide a good review of several approaches, including a straightforward Gauss-Newton (GN) method. This approach is based on linearizing $\mathbf{F}[\mathbf{m}]$ with a Taylor series expansion,

$$\mathbf{F}[\mathbf{m}_{k+1}] = \mathbf{F}[\mathbf{m}_k + \delta \mathbf{m}] = \mathbf{F}[\mathbf{m}_k] + \mathbf{J}_k (\mathbf{m}_{k+1} - \mathbf{m}_k), \quad (5)$$

where k denotes iteration number, and $\mathbf{J}_k = (\partial \mathbf{F} / \partial \mathbf{m})|_{\mathbf{m}_k}$ is the $N \times M$ sensitivity matrix calculated at \mathbf{m}_k . Calculation of \mathbf{J}_k , which describes the perturbations in the data due to changes in the model, is described in detail by Rodi and Mackie (2000), Mackie and Madden (1993), and Rodi (1976). Substituting equation (5) in equation (4), we obtain

$$\begin{aligned} \tilde{W} &= (\mathbf{m}_{k+1} - \mathbf{m}_0)^T \mathbf{C}_m^{-1} (\mathbf{m}_{k+1} - \mathbf{m}_0) \\ &+ \lambda^{-1} \{ (\hat{\mathbf{d}}_k - \mathbf{J}_k (\mathbf{m}_{k+1} - \mathbf{m}_0))^T \\ &\times \mathbf{C}_d^{-1} (\hat{\mathbf{d}}_k - \mathbf{J}_k (\mathbf{m}_{k+1} - \mathbf{m}_0)) \}, \end{aligned} \quad (6)$$

where $\hat{\mathbf{d}}_k = \mathbf{d} - \mathbf{F}[\mathbf{m}_k] + \mathbf{J}_k (\mathbf{m}_k - \mathbf{m}_0)$. This \tilde{W} is then quadratic in \mathbf{m}_{k+1} and thus can be minimized exactly (for fixed λ). For numerical stability (Marquardt, 1963), damping is generally required to control step size for each iteration. The system of equations to be solved for each iteration then becomes

$$\begin{aligned} \mathbf{m}_{k+1} - \mathbf{m}_k &= [\lambda \mathbf{C}_m^{-1} + \mathbf{\Gamma}_k^m + \epsilon_k \mathbf{I}]^{-1} [\mathbf{J}_k^T \mathbf{C}_d^{-1} (\mathbf{d}_k - \mathbf{F}[\mathbf{m}_k]) \\ &- \lambda \mathbf{C}_m^{-1} (\mathbf{m}_k - \mathbf{m}_0)], \end{aligned} \quad (7)$$

where the "model space cross-product" matrix $\mathbf{\Gamma}_k^m = \mathbf{J}_k^T \mathbf{C}_d^{-1} \mathbf{J}_k$ is an $M \times M$ positive semidefinite symmetric matrix, \mathbf{I} is the

identity matrix, and ϵ_k is a damping parameter. Note that with the GN approach λ is fixed. Therefore, the algorithm will converge to a stationary point of equation (4), not equation (3). To achieve the stationary point of equation (3) (with respect to both λ and \mathbf{m}), the process would have to be repeated with different values of λ until the constraint $X_d^2 = X_*^2$ was satisfied.

The OCCAM approach, first proposed by Constable et al. (1987) (see also deGroot-Hedlin and Constable, 1993; Uchida, 1993; Parker, 1994), is also based on linearizing $\mathbf{F}[\mathbf{m}]$ and then solving for the stationary points of equation (6). Differentiating equation (6) with respect to \mathbf{m} and setting the result to zero leads to an iterative sequence of approximate solutions:

$$\mathbf{m}_{k+1}(\lambda) = [\lambda \mathbf{C}_m^{-1} + \mathbf{I}_m]^{-1} \mathbf{J}_k^T \mathbf{C}_d^{-1} \hat{\mathbf{d}}_k + \mathbf{m}_0. \quad (8)$$

The unique feature of the OCCAM approach is that the parameter λ is used in each iteration both as a step length control and a smoothing parameter. That is, equation (8) is solved for a series of trial values of λ and the misfit $X_d^2(\mathbf{m}_{k+1}(\lambda))$ for each λ is evaluated by solving the 2-D forward problem. As for the linear problem, λ should be chosen so that the condition $X_d^2 = X_*^2$ is met. Usually, in the early iterations, the true misfit X_d^2 is higher than the desired X_*^2 for all possible λ . The OCCAM process thus chooses the model with the minimum misfit as the basis for the next iteration. The process is then repeated until the misfit reaches the desired level. Parker (1994) called this process of bringing the misfit down to the target level phase I. Once the misfit reaches the desired level, phase II begins by keeping the misfit at the desired level, but varying λ to search for the model with smallest norm. Since the problem is nonlinear, the desired misfit may never be reached. However, in practice, improvement of the misfit from iteration to iteration can be expected, until a minimum is achieved.

Both GN and OCCAM share similar computational steps. For each iteration, \mathbf{J}_k must be calculated, and an $M \times M$ system of equations [equation (7) for GN and equation (8) for OCCAM] must be solved. These methods are thus very time-consuming (e.g. Smith and Booker, 1991; Rodi and Mackie, 2000). Furthermore, these methods require much memory to store the sensitivity and cross-product matrices. These computational inefficiencies are the result of strong dependence on the model space dimension M .

Several approaches have been proposed to avoid the heavy computational burden of the direct linearized search schemes. One approach is to use approximate sensitivities to eliminate calculation of the full sensitivity matrix. A good example in this category is the RRI introduced by Smith and Booker (1991). RRI turns the 2-D inverse problem into a series of 1-D inverse problems, by computing the approximate sensitivity of data at each site to variations of resistivity directly below the site. The model is updated by solving a series of 1-D inverse problems and interpolating horizontally to form the 2-D resistivity model. Fit to the data is tested with a full 2-D forward calculation, and step length is adjusted if necessary. The process is repeated until the misfit condition is met. Note that this approach eliminates both the 2-D sensitivity calculation and the need to solve a large $M \times M$ system of equations. Generally, RRI requires many iterations, but overall is very fast. RRI can handle very large 2-D data sets and has been applied to interpretation of many MT data sets (e.g., Unsworth et al., 1999). Oldenburg

and Ellis (1991) suggested a very similar approach based on using a series of 1-D inversions as an approximate inverse mapping (AIM) to map the data back to model space, followed by full calculation of the 2-D forward problem to assess model fit.

The efficiency of approximate inversion schemes comes at a price. Because of the incomplete search of the model space, schemes based on 1-D inversions can be insensitive to features that are not directly beneath the locations of measurement (Farquharson and Oldenburg, 1996). For example, if the data are dominated by a significant feature outside of the profile, the inversion may have difficulty finding any models which actually fit the data (e.g., Unsworth et al., 1999), or the inversion may insert a geologically unreasonable feature beneath the profile. In addition, inversion of vertical magnetic transfer functions is difficult with this approach, since vertical magnetic fields are sensitive to nearby structures, rather than features directly beneath the sites.

Farquharson and Oldenburg (1996) proposed a somewhat different approximate approach based on using 2-D sensitivities of a homogeneous or layered half-space, instead of the exact sensitivities. This scheme eliminates the need for calculation of \mathbf{J}_k but still requires inversion of large ($M \times M$) matrices. However, in many cases these simple approximations to the sensitivities are good enough to allow convergence of the inversion to an acceptable level of misfit. The effectiveness of this scheme depends on many factors: the complexity of structure of the true model, the closeness of the structure to the data sites, and the magnitude of the resistivity contrasts.

Another way to reduce the computational burden in 2-D MT inversion is provided by the subspace approach of Oldenburg et al. (1993): the inverse solution is sought in a low-dimensional subspace of the original M -dimensional model space. The success of this approach depends greatly on the choice of subspace basis vectors. Unfortunately, the proper choice is often not obvious, and a bad selection can lead to a poor solution. A combination of both approximate sensitivities and the subspace approach has been used by Oldenburg and Ellis (1993). Current implementation of OCCAM also allows for a simple sort of subspace approach, since the resistivity can be parameterized on a grid coarser than that used for numerical computations.

Rather than approximate the sensitivities, or impose prior restrictions on the model space, Mackie and Madden (1993) considered an approach based on a conjugate gradient (CG) relaxation solution of the linearized normal equations derived from equation (6). With the relaxation method, the actual computation of the sensitivity matrix can be avoided by using the fact that evaluating the gradient of the linearized penalty functional requires only one forward solution (per period) with a distributed set of sources either in the volume or on the surface. CG thus significantly reduces the computational requirements (both CPU time and memory) for each iteration, making attempts at even 3-D MT inversion feasible (Mackie and Madden, 1993). The model is updated for each iteration, and the CG relaxation solution process is repeated until the stationary point of equation (4) is reached. Note, however, that since CG is a descent method which does not directly use any information about curvature of the penalty functional, many iterations may be required compared to GN or OCCAM. Also, as for the GN approach, the entire process must be repeated for different values of λ to find a true minimum structure model which achieves a specified misfit X_*^2 .

Rodi and Mackie (2000) considered a variant on this CG approach. These authors applied NLCG directly to minimization of equation (4), with λ fixed. Similar to the CG method (Mackie and Madden, 1993), NLCG requires only a few forward solutions (per period) in each line minimization step. To improve the convergence rate, a simple preconditioner is used. Again, to obtain a norm minimizing solution with minimum structure at the desired misfit, one needs to minimize equation (4) for various values of λ . Rodi and Mackie (2000) show that the two descent methods (CG and NLCG) are comparable, and both are much more efficient than the GN method in terms of CPU time and memory requirements.

Here, we develop a variant on the OCCAM approach which is significantly more efficient than previously proposed methods. We begin by transforming the inverse problem from the model space into the data space, by expressing the solution as a linear combination of rows of the sensitivity matrix smoothed by the model covariance. This transformation reduces the size of the system of equations to be solved from $M \times M$ to $N \times N$. Since the number of model parameters M is often much larger than the number of data N , a significant decrease in both CPU time and memory can be achieved with this approach. More importantly, the data space formulation leads naturally to a simple approximation which can result in very significant computational savings in most cases.

Generally, MT data are smooth (in period, and for closely spaced sites, in space) and “redundant.” Therefore, in the data space approach, there is no need to use all of the sensitivities as basis functions. A subset is typically sufficient to construct the model without significantly loss of detail. With this approximation, it is unnecessary to compute all sensitivities, and the size of the system of equations that must be solved can be significantly reduced. We call this approach the REDuced Basis OCCam’s (REBOCC) inversion. Note that even though we construct the solution from subset of the smoothed sensitivities, the goal of the inversion remains to find the norm minimizing model subject to fitting *all* of the data well enough. As we shall discuss in more detail below, in the data space the choice of basis functions is very natural and is dictated by what features can be resolved by the available data. This is in contrast to the choice of a model subspace (Oldenburg et al., 1993), where the choice of subspace is rather arbitrary.

With careful implementation of forward modeling and sensitivity calculations, REBOCC runs in a fraction of the time required by methods such as GN or OCCAM and is also faster than DASOCC and NLCG. In addition, memory requirements are significantly reduced so that large data sets can be inverted with REBOCC on a standard workstation. The basic idea behind REBOCC generalizes readily to the 3-D case.

THE DATA SPACE OCCAM METHOD (DASOCC)

Parker (1994) shows that the minimizer of equation (6) for iteration k can be expressed as a linear combination of rows of the smoothed sensitivity matrix $\mathbf{C}_m \mathbf{J}_k^T$,

$$\mathbf{m}_{k+1} - \mathbf{m}_0 = \mathbf{C}_m \mathbf{J}_k^T \boldsymbol{\beta}_{k+1}, \quad (9)$$

where $\boldsymbol{\beta}_{k+1}$ is an unknown expansion coefficient vector of the basis functions $[\mathbf{C}_m \mathbf{J}_k^T]_j$; $j = 1, \dots, N$, which are sometimes referred to as the “representers” of the linearized data functionals for iteration k (e.g., Parker, 1994). Substituting equation (9)

into equation (6), we obtain

$$\begin{aligned} \tilde{W} = & \boldsymbol{\beta}_{k+1}^T \boldsymbol{\Gamma}_k^n \boldsymbol{\beta}_{k+1} + \lambda^{-1} \{ (\hat{\mathbf{d}}_k - \boldsymbol{\Gamma}_k^n \boldsymbol{\beta}_{k+1})^T \\ & \times \mathbf{C}_d^{-1} (\hat{\mathbf{d}}_k - \boldsymbol{\Gamma}_k^n \boldsymbol{\beta}_{k+1}) \}. \end{aligned} \quad (10)$$

Here $\boldsymbol{\Gamma}_k^n = \mathbf{J}_k \mathbf{C}_m \mathbf{J}_k^T$ is the $N \times N$ “data space cross-product” matrix, which is again symmetric and positive semidefinite. Differentiating equation (10) with respect to $\boldsymbol{\beta}$ and rearranging, the unknown expansion coefficients can be obtained as

$$\boldsymbol{\beta}_{k+1} = (\lambda \mathbf{C}_d + \boldsymbol{\Gamma}_k^n)^{-1} \hat{\mathbf{d}}_k. \quad (11)$$

The inverse problem thus becomes a search for the N real expansion coefficients $\boldsymbol{\beta}_{k+1}$, instead of the M -dimensional model, \mathbf{m}_{k+1} . Exactly as for the standard OCCAM, we can solve for $\boldsymbol{\beta}_{k+1}$, update the model, and then check the misfit for various values of λ . We again choose λ to achieve the minimum misfit if this exceeds the desired level X_*^2 (phase I) and use this model as the basis for the next iteration. Once the desired misfit is achieved, phase II begins to wipe out unnecessary features while keeping the misfit at the desired level.

We emphasize here that we have only transformed the inverse problem solution method from the model space to the data space. Solutions obtained in both spaces will be identical if we choose all parameters (i.e., λ and \mathbf{C}_m) the same. For brevity we refer to this variant on OCCAM as the Data Space OCCam (DASOCC) inversion. Note that a data space approach was also used by Smith and Booker (1988) in their treatment of the 1-D MT inverse problem.

The data space formulation offers several advantages. The most obvious is the reduction in the dimension of the system of equations which must be solved ($N \times N$ in the data space instead of $M \times M$ in the model space); Generally, $N \ll M$. This will be particularly true for the 3-D case.

Also, calculation in the model space requires \mathbf{C}_m^{-1} . Since it is not practical to specify a full $M \times M$ model covariance matrix \mathbf{C}_m and then compute the inverse, \mathbf{C}_m^{-1} is replaced by the first derivative roughness penalty in deGroot-Hedlin and Constable (1990). In the data space approach, \mathbf{C}_m is required, not its inverse. This offers some advantages since the model covariance can be readily used to include prior information such as an ocean or faults which should be fixed in the model. We discuss these issues in more detail in the model covariance section and in Appendix A.

REDUCED DATA SPACE OCCAM APPROACHES (REBOCC)

The data space formulation clearly shows that the solution is a linear combination of natural basis functions or representers. Each representer corresponds to a single data element (at a particular period, station, response, and mode). Just as the MT data should be smooth and redundant, the representers vary slowly with period and site location for a given response and mode. These basis functions are thus highly redundant, so that an excellent approximation to the solution can be found in a subspace of much lower dimension (Parker and Shure, 1982; see also Parker, 1994). This simple but critical concept of data redundancy can significantly speed up the inversion while also substantially decreasing memory requirements.

Prior to solving the inverse problem, we will select a subset of L (out of N) data for which representers will be calculated at each iteration. As a simple example, we could choose all data

for every other period (so $L = N/2$). Note that L can typically be considerably smaller than this, as we shall show later. For iteration $k + 1$, we seek solutions of the form

$$\mathbf{m}_{k+1} = \mathbf{C}_m \mathbf{G}_k^T \alpha_{k+1} + \mathbf{m}_0, \quad (12)$$

where α_{k+1} is the L -dimensional unknown coefficient vector for the reduced basis, and \mathbf{G}_k is the $L \times M$ subset sensitivity matrix.

To fit all of the data and to derive equations for α_{k+1} analogous to equations (10) and (11), we require the linearized relationship between $\delta \mathbf{m}$ and $\hat{\mathbf{d}}_k$. In fact, we do not strictly have this relationship, unless we calculate all of the sensitivities. However, the data vary smoothly, and so a data value would be well approximated by interpolation of “nearby” data (e.g., adjacent frequencies from the same site). In the same way, sensitivities vary smoothly with frequency and/or site location and can be interpolated from nearby sensitivities. We thus express the approximation to the full sensitivity matrix \mathbf{J}_k in terms of the subset sensitivity matrix \mathbf{G}_k using an interpolation matrix, \mathbf{B} , of size $N \times L$, i.e.,

$$\mathbf{J}_k \approx \mathbf{B} \mathbf{G}_k. \quad (13)$$

The interpolation matrix does not need to be very sophisticated. Recall that Farquharson and Oldenburg (1996) had success in 2-D MT inversion using a very simple sensitivity matrix generated from either a homogeneous or layered half-spaced, and RRI uses only 1-D approximate sensitivities. By comparison, the approximate sensitivity matrix generated by even a crude interpolation of a sparse subset of representers will actually be quite close to the exact sensitivity matrix \mathbf{J}_k (Siripunvaraporn, 1999).

Substituting equations (12) and (13) into equation (6), we find

$$\begin{aligned} \tilde{W} &= \alpha_{k+1}^T \mathbf{G}_k^T \alpha_{k+1} + \lambda^{-1} \{ (\hat{\mathbf{d}}_k - \mathbf{B} \mathbf{G}_k^T \alpha_{k+1})^T \\ &\quad \times \mathbf{C}_d^{-1} (\hat{\mathbf{d}}_k - \mathbf{B} \mathbf{G}_k^T \alpha_{k+1}) \}, \end{aligned} \quad (14)$$

where $\mathbf{G}_k^T = \mathbf{G}_k \mathbf{C}_m \mathbf{G}_k^T$ is the $L \times L$ “data subspace cross-product” matrix. To give the system of equation to be solved a form similar to equation (11), we follow Egbert et al. (1994) by decomposing $\mathbf{C}_d^{-1/2} \mathbf{B}$ into the $N \times N$ orthonormal matrix \mathbf{Q} , where $\mathbf{Q}^T = \mathbf{Q}^{-1}$, and the $N \times L$ matrix \mathbf{R} , i.e., $\mathbf{C}_d^{-1/2} \mathbf{B} = \mathbf{Q} \mathbf{R}$, where $\mathbf{Q} = [\bar{\mathbf{Q}} \mid \bar{\mathbf{Q}}_0]$, and $\mathbf{R}^T = [\bar{\mathbf{R}} \mid \mathbf{0}]$. Matrices $\bar{\mathbf{Q}}$ and $\bar{\mathbf{Q}}_0$ have dimensions $N \times L$ and $N \times N - L$ respectively. Matrix $\bar{\mathbf{R}}$ is the square $L \times L$ upper triangular matrix, and $\mathbf{0}$ is the $N - L \times L$ zero matrix, i.e., all elements are zeros.

Equation (14) becomes

$$\begin{aligned} \tilde{W} &= \alpha_{k+1}^T \bar{\mathbf{R}}^T \bar{\mathbf{Q}}_k \alpha_{k+1} + \lambda^{-1} \left\{ \left(\mathbf{C}_d^{-1/2} \hat{\mathbf{d}}_k - \mathbf{Q} \bar{\mathbf{R}}^T \alpha_{k+1} \right)^T \right. \\ &\quad \times \left. \left(\mathbf{C}_d^{-1/2} \hat{\mathbf{d}}_k - \mathbf{Q} \bar{\mathbf{R}}^T \alpha_{k+1} \right) \right\} \end{aligned} \quad (15)$$

where

$$\bar{\alpha}_{k+1} = (\bar{\mathbf{R}}^{-1})^T \alpha_{k+1} \quad (16)$$

and

$$\bar{\mathbf{R}}^T = \bar{\mathbf{R}} \mathbf{G}_k^T \mathbf{R}^T. \quad (17)$$

Inserting $\mathbf{Q} \mathbf{Q}^T = \mathbf{I}$ in between $\bar{\mathbf{R}}^T$ and $\bar{\alpha}_{k+1}$ on the right side of equation (15) and rearranging, this can be rewritten as

$$\begin{aligned} \tilde{W} &= \bar{\alpha}_{k+1}^T \bar{\mathbf{R}}^T \bar{\alpha}_{k+1} + \lambda^{-1} \{ X_{min}^2 \\ &\quad + (\hat{\mathbf{d}}_k - \bar{\mathbf{R}}^T \bar{\alpha}_{k+1})^T (\hat{\mathbf{d}}_k - \bar{\mathbf{R}}^T \bar{\alpha}_{k+1}) \}, \end{aligned} \quad (18)$$

where $\hat{\mathbf{d}}_k = \bar{\mathbf{Q}}^T \mathbf{C}_d^{-1/2} \hat{\mathbf{d}}_k$, and $X_{min}^2 = \|\bar{\mathbf{Q}}_0 \mathbf{C}_d^{-1/2} \hat{\mathbf{d}}_k\|^2 = \|\mathbf{C}_d^{-1/2} \hat{\mathbf{d}}_k\|^2 - \|\bar{\mathbf{Q}}^T \mathbf{C}_d^{-1/2} \hat{\mathbf{d}}_k\|^2$ is the approximate minimum achievable total square misfit for the selected basis. If we use all representers (i.e., $\mathbf{B} = \mathbf{I}$), then $X_{min}^2 = 0$. This corresponds to the fact that for a linear problem, we can fit the data exactly if we use all representers. This will not be true for the nonlinear MT problem. Thus X_{min}^2 only provides a very rough estimate of the magnitude of data misfit that might be achieved with the chosen reduced basis. In general, X_{min}^2 is high in the early iterations, and decreases to a constant in the later iterations.

Differentiating equation (18) with respect to $\bar{\alpha}$ and setting the result to zero, the unknown expansion coefficients can be obtained in a form similar to equation (11),

$$\bar{\alpha}_{k+1} = (\lambda \mathbf{I} + \bar{\mathbf{R}}^T)^{-1} \hat{\mathbf{d}}_k. \quad (19)$$

Again, just as in the model and the data space methods, after solving equation (19), we update the model using equations (12) and (16), then solve the forward problem to evaluate X_d^2 . The procedure is repeated to find the appropriate λ . The outer loop of the iterative minimization of equation (4) proceeds exactly as for OCCAM or for DASOCC.

Representer subsets for REBOCC

The success of the data subspace approach depends to some extent on the selection of data points which determine the representers used in the model expansion. Clearly it is necessary to select the basis to uniformly cover the full data set, so that the simple interpolation scheme used here is effective. Beyond this basic criteria, the optimal choice of data subsets is an issue that needs further study. Here, we offer some simple schemes based on our experience so far.

Two classes of data subsets which generally seem to work well are shown in Figure 1. In the first example, every p th period is chosen for all sites (“ p th-stripe” pattern). This pattern is safe to apply in almost all cases, since for physical consistency the

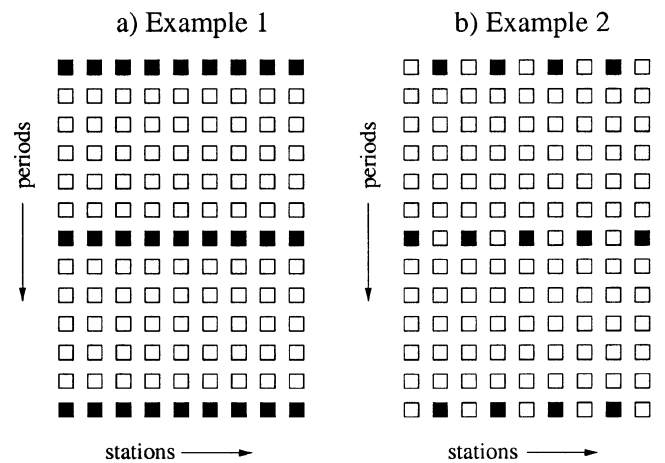


FIG. 1. Examples of subsets of data used to calculate the representers. Representers are calculated for data corresponding to the filled squares; the open squares represent the remaining data. (a) First example: “ p th-stripe” pattern where every p th period is selected for all sites (here $p = 6$). (b) Second example: “ p th:sth-checker” (with $p = 4$ and $s = 2$) pattern where the selected data form a checker pattern.

data must be a smooth function of period (e.g., Weidelt, 1972). Selecting at least one period per decade is our recommendation for this stripe pattern. Test runs with additional periods should be made (if feasible, given available computer resources) to verify that the data space has been adequately resolved.

This pattern offers a significant computation advantage: calculations for the sensitivity matrix need only be done for a reduced set of periods. For the 2-D problem, these sensitivities can be calculated by direct factorization of the coefficient matrix (see forward modeling section). With the p th-stripe pattern, this factorization need only be done for a small subset of periods. This idea is clearly directly transferable to the 3-D case.

In the second example every p th period and every s th station are selected in a staggered pattern to build the basis (“ p th- s th-checker” pattern). As noted above, because of the way we compute sensitivities, it is generally most efficient to use as many sites as possible (i.e. small s). However, this pattern can be used to reduce storage requirements for very large data sets.

Comparison of computational resource requirements

In this section, we summarize the computer resources required by each method (OCCAM, DASOCC, and REBOCC). In all three methods, most of the major computational costs lie in first calculating the sensitivity matrices. This cost should be equal for OCCAM and DASOCC (if the same method is used; see the forward modeling and sensitivity calculation section), and less for REBOCC, depending on the number of periods used. Second, computing the cross-product matrices requires about NM^2 , MN^2 , and ML^2 operations, respectively, for OCCAM, DASOCC, and REBOCC. Third, about $M^3/6$, $N^3/6$, and $L^3/6$ flops are required (for each λ) to solve equations (8), (11), and (19), respectively. In addition, extra calculations are required for REBOCC to factor \mathbf{B} (but this only has to be done once at the start of the inversion), and for computing the inner products of equation (17) about L^3 once every iteration.

Similarly, the memory required by each method is dominated by storage of two matrices: the sensitivity matrices (about NM , NM , and LM for OCCAM, DASOCC, REBOCC, respectively) and the cross-product matrices (about $M^2/2$, $N^2/2$, and $L^2/2$ for OCCAM, DASOCC, and REBOCC, respectively). For REBOCC, extra memory is required to store the $N \times L$ interpolation matrix \mathbf{B} .

As L approaches N , the extra calculations and memory required with the REBOCC method become significant. However, our experience shows that L need only be 10–30% of N to ensure convergence of REBOCC. Thus REBOCC can reduce memory requirements by at least 60% and CPU time by more than 80% compared to DASOCC.

REBOCC AND DASOCC: ALGORITHM DETAILS

The overall stability and efficiency of the DASOCC and REBOCC schemes depends on many details including forward modeling, the model covariance, the 1-D line search for the Lagrange multiplier λ , and static shift corrections. In this section we briefly describe our implementations of these parts of the inversion algorithms. Further details are provided in Appendix A and in Siripunvaraporn (1999).

Forward modeling and sensitivity calculation

Forward modeling is the heart of the inversion and thus must be reliable, fast, and accurate. It is used in two parts of the inversion—to compute the sensitivity matrix, and to compute responses for calculating the misfit.

Details on the second-order Maxwell’s equations that must be solved can be found in many previous publications on the MT method (e.g., Rodi, 1976). As in Smith and Booker (1991), we apply the finite difference (FD) method to these equations to form the discrete system $\mathbf{Ax} = \mathbf{b}$, where \mathbf{b} contains the terms associated with the known boundary values and the source fields, and \mathbf{x} represents the unknown electric or magnetic fields. Boundary conditions for the model domain are as in Smith and Booker (1991). The accuracy of the solution is controlled by the quality of the mesh. The reader is referred to Rodi (1976), deGroot-Hedlin and Constable (1990), and Smith and Booker (1991) for discussion of these issues.

Sensitivities for MT data (e.g., the apparent resistivity or phase) can be readily calculated in terms of sensitivities of the electric and magnetic field components at the surface (Mackie and Madden, 1993; Rodi, 1976). The surface field components in turn can always be expressed in the general form $\mathbf{a}^T \mathbf{x}$ (Rodi, 1976), where \mathbf{a}^T may depend upon \mathbf{m} , and \mathbf{x} is the discrete electric or magnetic field solution. These sensitivities may be calculated from

$$\frac{\partial(\mathbf{a}^T \mathbf{x})}{\partial m_j} = \frac{\partial \mathbf{a}^T}{\partial m_j} \mathbf{x} - \mathbf{a}^T \mathbf{A}^{-1} \left[\frac{\partial \mathbf{A}}{\partial m_j} \mathbf{x} \right] \quad j = 1, \dots, M. \quad (20)$$

The second term of the right side in equation (20) can be computed by solving the same system of equations required for the forward problem, but with a different right side, $(\partial \mathbf{A} / \partial m_j) \mathbf{x}$.

Rodi (1976) shows that there are two ways to calculate the field component sensitivities of equation (20). One requires solving the forward problem M times per period, once for each m_j . The other uses the reciprocity property of the forward problem and requires N_s (number of stations) forward solutions per period. Good reviews of the sensitivity calculations can be found in Rodi (1976), Mackie and Madden (1993), Rodi and Mackie (2000). Since $N_s \ll M$, the second approach is generally much more efficient. We use this approach for REBOCC and DASOCC.

The sparse system of linear equations $\mathbf{Ax} = \mathbf{b}$ can be solved in two general ways: by a direct method using the LU decomposition or with an iterative method (Press et al., 1992), such as preconditioned conjugate gradients (PCG). Both approaches have advantages. With a direct method, after \mathbf{A} is decomposed into lower (\mathbf{L}) and upper (\mathbf{U}) triangular matrices, solution (by forward and back substitution) is extremely fast. In REBOCC, the direct approach is thus used for constructing the sensitivity matrix where the same system of equations must be solved for multiple right sides. To solve for a single right side (e.g., when calculating the misfit during search for λ) an iterative method which takes advantage of sparseness is more efficient. We use an iterative approach in these circumstances. To use the classical PCG on the complex symmetric system (which is non-Hermitian), the conjugate transpose is not applied when computing the inner product, i.e., we use $\langle \mathbf{x}, \mathbf{y} \rangle = \mathbf{x}^T \mathbf{y}$ (Barrett et al., 1994). To speed up convergence, we use the incomplete LU decomposition level 3, ILU(3), (Kershaw, 1978; Smith and Booker, 1991) as a preconditioner.

Model covariance \mathbf{C}_m

In the penalty functional (4), the inverse of the model covariance is required to evaluate the norm $\mathbf{m}^T \mathbf{C}_m^{-1} \mathbf{m}$. Inversion of a nondiagonal $M \times M$ model covariance would not be computationally practical. Thus, for model space inversion approaches it is conventional to formulate the model norm in terms of a “roughness penalty” instead of the inverse of a covariance matrix [for example, deGroot-Hedlin and Constable (1990) used the first derivative roughness penalty in our equation (2)]. The relationship of some common roughness penalties to a generalized sort of model covariance is discussed by Wahba (1990).

With a data space approach, inversion of \mathbf{C}_m is not required [see equation (9)]. This gives us a great deal of freedom to define a model covariance which can include prior information of various sorts. However, it is not practical to even form an arbitrary full (nondiagonal) $M \times M$ model covariance matrix for typical 2-D problems. We thus develop a fairly general class of model covariance functions which allows reasonable flexibility for including prior information and, at the same time, allows for efficient computation of the matrix product $\mathbf{C}_m^T \mathbf{J}^T$ needed in equations (9) and (12) without actually computing \mathbf{C}_m explicitly. The approach is based on solving the 1-D diffusion equation, alternating between vertical and horizontal directions. See Appendix A for more technical details.

With this approach, it is not necessary to use a constant correlation length scale (or variances) throughout the model. Different length scales can be used in vertical and horizontal directions and in different parts of the model domain. The model covariance approach is thus very flexible, and allows rather general statistical specification of prior information. For example, geologic structures such as faults can be incorporated into the model by letting the smoothing length scale go to zero at the fault location, or part of the resistivity structure can be frozen by letting variances go to zero.

The choice of proper correlation length scale is important, but difficult to justify rigorously. Making length scales too large can result in difficulties in finding any models that fit the data. Choosing length scales too small can also result in problems (e.g., confusion of static shifts with deep lateral structure). Inevitably some experimentation with correlation length scales will be required. As a default strategy we make the vertical correlation length scale of each layer proportional to the depth of that layer, with the horizontal correlation length scale set to the maximum of the depth and the gap between stations. On the edges of the model domain, the length scale is set equal to the distance from sides of the model to the edge stations. This choice of correlation length scales coincides with the loss of the resolving power of the data at depth and near the boundaries of the model.

Line search for λ

At every iteration, we search for the λ that gives the model [defined by equations (12) and (19)] with the minimum misfit (phase I) or at the desired misfit (phase II). For each trial value of λ , the system of equations (19) must be solved, the model updated, and the misfit computed by solving the forward problem. Therefore, minimizing the number of λ s tried can help us to further reduce computational costs.

We use a relatively simple search method which takes advantage of several facts: (1) the misfit is a smooth function of λ ,

(2) the range of $\log_{10} \lambda$ is generally within the interval $[0, 6]$ (for our default model covariance; see Appendix A for details), and (3) the optimal choice of λ changes little between iterations.

For the first iteration, we start with three different equally spaced λ s covering one decade of λ . Misfits computed for these initial λ s will tell us whether we should go left or right, or stop if a minimum has already been bracketed. Usually, this scheme requires about 3–8 trial values of λ , about half of the 8–12 values of λ per iteration reported for the 2-D Occam’s inversion by deGroot-Hedlin and Constable (1990). With the prior knowledge of the previous iteration, the previous three bracketing points can be used as a starting point for the next iteration, and the process is repeated. Generally, the optimal λ does not change much between iterations, so only a small number of trial values of λ (but at least three) are required.

Once the desired misfit is achieved within the range of trial λ s bracketing the minimum, parabolic interpolation (Press et al., 1992) is used to locate λ providing the desired level of misfit. If two or more values of λ bracketing the minimum have the same (desired) misfit, we choose the larger λ , which usually corresponds to a smaller model norm.

In general, the desired misfit may never be reached, so that the smoothing process is not performed. The model with the minimum misfit (higher than the desired misfit) might contain some unnecessary features inserted by the inversion to make the fit better. In this case, we thus recommend an additional run with a higher target misfit to find the minimum norm model corresponding to a larger (and this time achievable) misfit.

Static shift correction

Shallow local inhomogeneities can distort the regional electric field, and cause a frequency-independent shift in the log apparent resistivity while leaving the phase unaffected. For REBOCC and DASOCC, static shifts can be incorporated as additional model parameters, which are automatically estimated by the inversion (if requested by the user). For each iteration static shifts for each site are estimated using the median residual [observed minus (undistorted) calculated] \log_{10} apparent resistivity. Then, the misfits are recalculated using the (distorted) calculated responses. More sophisticated approaches can be used to obtain the static shift factors (e.g., deGroot-Hedlin and Constable, 1993; Wu et al., 1993; Ogawa and Uchida, 1996), but tests with synthetic distorted data indicate that this simple scheme is effective.

NUMERICAL EXPERIMENTS

Synthetic data

To test DASOCC and REBOCC, we generated synthetic data from a 2-D model adapted from the COPROD2 inversion results of Wu et al. (1993). The model (Figure 2) consists of four layers, with three 1 ohm-m rectangular conductors embedded (see Figure 2 for details). The grid used to form the synthetic data was well discretized at 170 columns and 183 rows (with an additional 10 air layers for the TE mode) to ensure accuracy of the solution. Model responses, including apparent resistivity and phases for transverse magnetic (TM) and transverse electric (TE) modes, and also vertical magnetic field transfer functions (which we will refer to as tipper, or TP), were generated using the finite element forward modeling program of

Wannamaker et al. (1986). Data for 36 stations spaced at 3-km intervals and 31 periods increasing logarithmically from 1 s to 1000 s (about 10 periods per decade) were used for the inversion test. Two-percent Gaussian noise was added to the data prior to the inversion. Note that the three conductors are not clearly evident in either TE or TM model responses (Figure 2).

Data space and data subspace inversion

Our first experiment is to invert the full set of synthetic data with the data space (DASOCC) and data subspace (REBOCC) methods with different subsets of calculated representers for TM and combined TM+TE modes. For REBOCC, we consider three stripe patterns and one checker pattern. All inversions

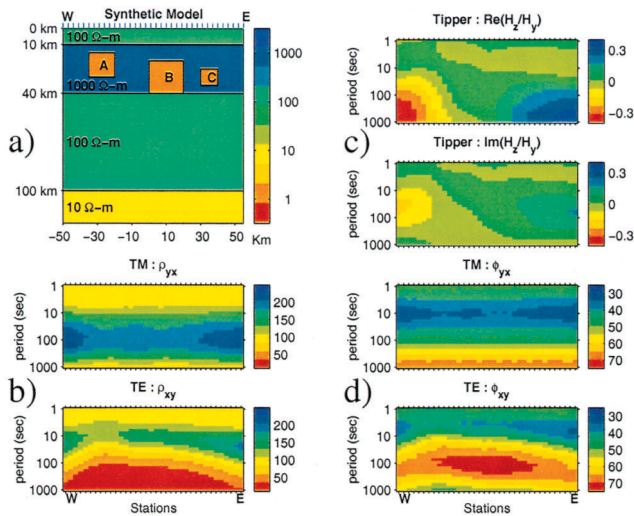


FIG. 2. (a) Model used to generate synthetic data responses. Inside the 1000 ohm-m resistive layer, there are three 1 ohm-m conductors: A, B, and C. Conductor A is 15 km \times 15 km, buried at 15 km depth. Conductor B is 20 km \times 20 km, buried at 20 km depth. Conductor C is 10 km \times 10 km, buried at 25 km depth. A and B are separated by 20 km, B and C are separated by only 10 km. Apparent resistivities and phases from TM and TE mode and tippers generated from the model are shown in (b)–(d), computed using the finite element forward modeling of Wannamaker et al. (1986) with 2% Gaussian noise added.

are run on a Sun UltraSparc I workstation. The desired rms is set to 1 (i.e., 2% misfit).

Due to limitations of computer resources, the model grid used for inverting the data is necessarily coarser. Using different forward modeling programs and grids for generating the data will help to reveal any systematic errors occurring in the inversion program. The model grid used for the inversion was thus discretized at 100 columns and 31 layers, plus 10 air layers for the TE mode. A 100 ohm-m half-space was used as a starting model for all inversions. The correlation length scales for all runs were set as described in the model covariance section. The horizontal smoothing length scale is set to 3 km (from the surface to about 3 km depth), whereas the vertical length scale is equal to the depth. At depths greater than 3 km, the horizontal and vertical length scales are both set equal to the depth.

Memory requirements for storing the cross-product and sensitivity matrices for all three methods (OCCAM, DASOCC, and REBOCC) and the interpolation matrix (for REBOCC only) are shown in Figure 3a. For TM+TE inversion memory requirements of the data space (DASOCC) approach exceeded those of the model space (OCCAM) approach because N is actually larger than M for this joint inversion case. However, using the reduced basis method, memory requirements can be significantly reduced (Figure 3a). For example, only 12 MBytes of memory are needed for a TM mode inversion using the 6th:2nd-checker pattern of calculated representers. Although the memory requirements for OCCAM are comparable to those of DASOCC, joint TM and TE inversion using this full data set with OCCAM is not practical. Performance results for OCCAM (deGroot-Hedlin and Constable, 1990) using a reduced data set will be discussed in the next section.

The convergence rates of DASOCC (dashed line) and REBOCC (solid lines) are plotted in Figure 3b and 3c. For both the TM and TM+TE inversions, the full basis inversion (DASOCC) requires much longer per iteration than any of the REBOCC inversions (by a factor of 6–35). Both the 10th-stripe and the 6th:2nd-checker data subset patterns have comparable numbers of representers (L), and require about the same amount of CPU time per iteration. However, the 6th:2nd-checker converges while the 10th-stripe does not. Clearly representers for more than one period per decade are required in this case. For the TM mode, all of the inversions converge

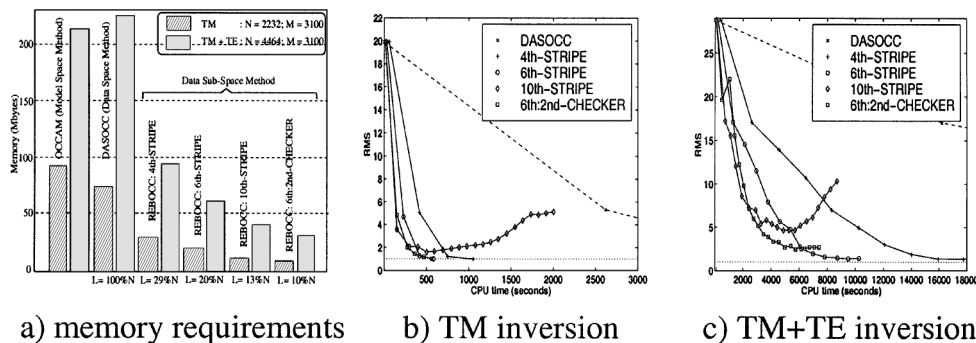


FIG. 3. (a) Memory requirements for OCCAM, DASOCC, and REBOCC for different representer subsets for TM and TM+TE modes. Note that double precision is used. (b) and (c) rms versus CPU time required for TM and TM+TE inversions with different representer subsets. Iteration numbers are indicated by the plotted symbols.

(except the 10th-stripe) to the desired misfit (5 iterations for the 6th:2nd-checker subset; 3 iterations for the rest). Models resulting from the different inversions are indistinguishable from one another. Only the models from the 6th-stripe subset are shown in Figure 4. For TM+TE inversion, none of the inversions are able to find a model with an rms of 1. Except for the 10th-stripe, all converge to some minimum level, which tends to get larger as L gets smaller. The minimum misfit is around 1.15 rms for the full basis, and 2.5 rms for the 6th:2nd-checker subset.

We used a 6th-stripe representer subset for additional single and joint mode inversions, using the same synthetic data with the same starting model. Figure 4 shows the results of these runs. The TM mode inversion reveals the layered host resistivity structure with little lateral variation. In contrast, the TE inversion reveals two conductors beneath the resistive layer. The second conductor on the east side of the model from $y = 0$ to 40 km corresponds to a combination of conductors B and C in the synthetic model. Interestingly, the tipper inversion displays the boundaries of the conductors more accurately than either single mode MT inversion. Even the small third conductor (C) can be distinguished. However, the layered host resistivity structure is poorly resolved with the tipper inversion. The joint inversion of TM and TE modes shows that the MT data is fit adequately with only two conductors inside the resistive layer (although the desired level of misfit is not quite achieved). Using tipper with the TE mode, and tipper with the TM mode, the inversions find all three conductors, which shows that tipper helps to resolve the smallest conductor in this synthetic case.

Finally, we inverted all data: TM, TE, and tipper ($N = 6696$, $M = 3100$, and $L = 1296$). The result is shown in Figure 4g. The inversion requires about 108 MBytes of memory and approximately 11 hours of run time (on a Sun UltraSparc 1) to obtain

a model with a rms of 1.05 at the 18th iteration. One can stop the inversion after the 8th iteration (about 5 hours) to obtain a model with an rms of about 2 (or 4% error), which is in fact indistinguishable from the fully converged run.

Comparison with other inversions

In the synthetic data example considered above, we used 10 periods per decade. This is comparable to the spacing of frequency bands obtained by most MT data processing schemes. However, it is a common practice for at least some inversion routines (e.g., NLCG and OCCAM) to use only 2–3 periods per decade.

To provide a fair comparison with these other inversions, we thus use a decimated data set with 3 periods per decade. This reduces the number of data to $N = 720$ for single mode inversions, and $N = 1440$ for joint inversion of TM+TE modes. A comparison using the full data set is given in Siripunvaraporn (1999). Here, we compare DASOCC and REBOCC with several other inversion programs: OCCAM of deGroot-Hedlin and Constable (1990), NLCG of Rodi and Mackie (2000), and RRI of Smith and Booker (1991). The same starting model (100 ohm-m half-space) and the same model mesh ($M = 100 \times 31 = 3100$) were used for all inversions. For timing comparisons, we run the inversions on the same machine, a Sun UltraSparc I. We use the 3rd-stripe subset for the REBOCC inversion (i.e., 4 out of 10 periods were used to calculate the representer).

Figure 5 shows the models produced by REBOCC, OCCAM, NLCG, and RRI for TM, TE, and TM+TE data sets (note that NLCG and RRI do not presently allow for inversion of tippers). The results from DASOCC are indistinguishable

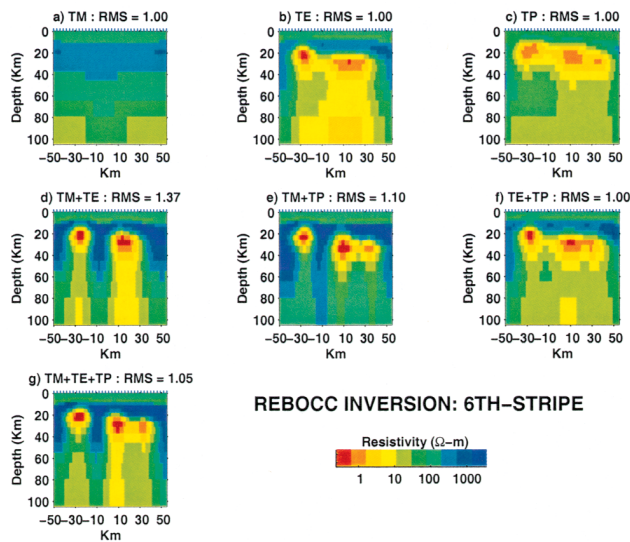


FIG. 4. Models produced by the 6th-stripe REBOCC inversion from various mode combination. The upper panels are the models from single mode inversion: (a) TM alone, (b) TE alone, (c) tipper alone. The lower panels are models from joint inversions: (d) TM and TE, (e) TM and tipper, (f) TE and tipper, and (g) TM, TE, and tipper.

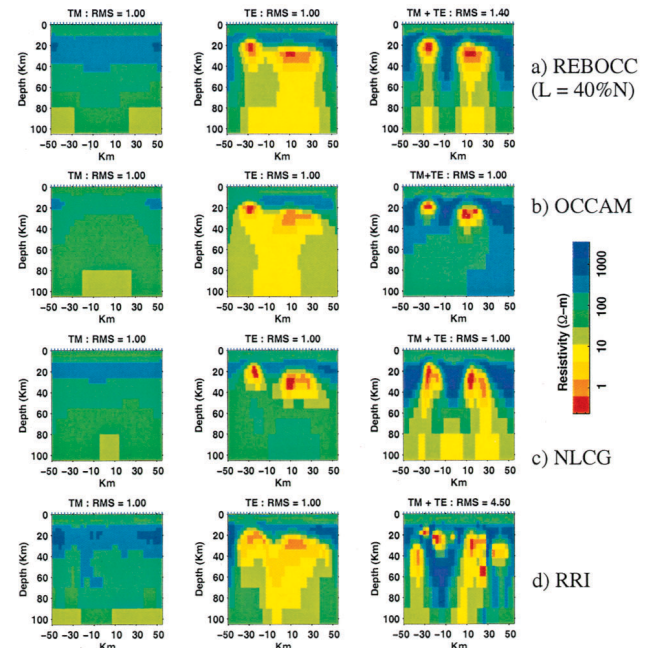


FIG. 5. Models obtained from the four inversions of TM, TE, and TM+TE mode discussed in the text. RRI fails to converge to a reasonable model for the TM+TE mode. The final result fits poorly and is very rough.

from REBOCC, and are thus not shown. In Figure 6, we plot rms misfit against CPU time (and iteration number) for RRI, REBOCC, DASOCC, NLCG, and OCCAM algorithms for TM inversion, TE inversion, and joint inversion of TM and TE modes. Note that the rms misfits in this section are calculated from the decimated data set. A slightly higher rms value can be expected for the full data set.

RRI, which is incredibly fast per iteration, requires 12 iterations to converge to the solution at the desired level of misfit for TM mode (27 iterations for TE mode). However, for the joint inversion of TM and TE, RRI fails to converge to a reasonable model. The minimum misfit achieved is quite large (4.53), and the final model is very rough. The small scale features between $y = -15$ km and 0 km (between conductor A and B in the synthetic model; Figure 2) are inserted in the early iterations and, once there, are very difficult to get rid of in the subsequent iterations (Smith and Booker, 1991). This confusion of the inversion is perhaps the result of the 1-D sensitivity approximation used by RRI. Beneath the sites from $y = -15$ to 0 km, the TM data (which is not sensitive to the isolated conductors) suggests a resistive structure; the TE data which is highly sensitive to the nearby conductors, tries to place conductive features beneath these sites.

OCCAM is the slowest method, as shown by the large misfit (the dashed curves at the top of the RMS plots of Figure 6) after all of other methods have converged to reasonable solutions. For single mode inversions, OCCAM requires about 20 hours per iteration for TM inversion (with a total of 5 iterations required to reach the desired rms) and about 22 hours per iteration (with a total of 9 iterations required) for TE inversion. For the joint TM+TE inversion, about 30 hours per iteration is required, with a total of 30 iterations to reach the desired RMS. The long run times per iteration of OCCAM result from larger matrices in the model space and from direct calculation of the sensitivity matrix elements, without making use of the reciprocity theorem.

For NLCG, the convergence rate depends on λ , and the optimal choice of λ will in general depend on the data set. In this synthetic data case, we found that $\lambda = 3$ resulted in the fastest convergence rate. NLCG converges to the desired misfit in 29 iterations (about 1100 seconds) for the TM mode, and 48 iterations (about 4400 seconds) for the TE mode (Figure 6). The same misfits of both single mode inversions can be accomplished at a lower CPU time with DASOCC, which requires about 650 seconds (3 iterations) for the TM inversion, and

about 1200 seconds (5 iterations) for the TE inversion. Convergence can be much more quickly accomplished with REBOCC 3rd-stripe, particularly for the TM mode inversion, where only about 320 seconds (3 iterations) is required. For the TE mode inversion, the CPU time of REBOCC is about 1000 seconds (8 iterations) which is comparable to that of DASOCC. Furthermore, recall that to find a true minimum structure solution (as REBOCC and DASOCC do), NLCG would have to be run with several trial values of λ (some of which might converge more slowly).

For the joint mode inversion, NLCG requires about 11 200 seconds (84 iterations) to reach an rms of 1.7 while REBOCC and DASOCC only need 3000 seconds (in 11 iterations) and 5000 seconds (in 7 iterations), respectively, to reach a similar level of misfit. However, in this example NLCG can reduce the rms to the desired rms of 1 after 157 iterations (21 000 seconds), whereas both DASOCC and REBOCC only result in a minimum rms of 1.1 (after 14 iterations, 10 500 seconds) and 1.4 (after 16 iterations, 4300 seconds), respectively. The failure to achieve the desired misfit for REBOCC and DASOCC probably is a result of the model covariance assumed. However in this case, a model with a rms of 2 and a model with a rms of 1 are indistinguishable.

Note that the convergence times plotted in Figure 6 for REBOCC reflect only phase I iterations where the goal is to bring the misfit down to the desired level (1.0 rms). Since for the joint (TM+TE) inversion the desired level cannot be reached, one should restart the inversion so that phase II can be completed with the desired misfit level set to a higher value (e.g., 1.5 rms starting with the model from the 14th iteration of phase I). phase II is necessary in order to obtain the minimum structure model (Parker, 1994). At the same time, for NLCG to obtain a model with minimum structure, several trial values of λ would be required.

It is difficult to compare the models produced by NLCG, RRI, OCCAM, DASOCC, and REBOCC (Figure 5) because of the difference of the smoothing parameters (model covariance). However, all of these inversions (except the joint mode inversion for RRI) reveal the main features, and produce responses that fit the data adequately.

Example with field data

Finally, we briefly consider application of REBOCC to a real data set—a high-density MT profile across the San Andreas

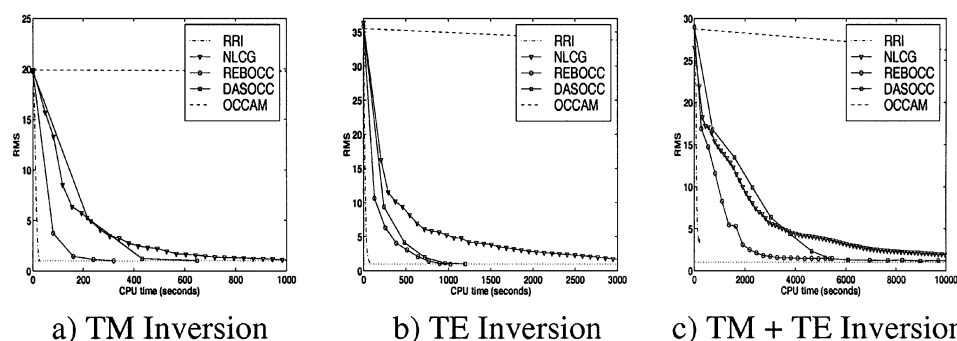


FIG. 6. Plots of rms versus CPU time of RRI, 6th-stripe REBOCC, and NLCG for TM, TE, and TM+TE inversions. Iteration numbers are indicated by the symbols.

fault (SAF) near Parkfield, California (Unsworth et al., 1997). This is a large data set with 55 TM mode responses, 37 TE mode responses, and 15 tippers each at 41 periods (from 0.01 s to 100 s). We discretized the model at 200 columns and 74 layers (plus 10 air layers for the TE mode), and used a 10 ohm-m half-space as a starting model. The static shift distortion parameters are set free so that the program can automatically adjust the values. An 8th-stripe subset ($L = 0.15N$) is used for joint inversion of TM, TE, and vertical magnetic transfer functions. The results of inversion are shown on Figure 7. After the 4th iteration (about 12 hours of CPU time on a Sun UltraSparc I), the inversion finds a model with a misfit of about 13% (an rms of about 2.7). A slightly better misfit can be obtained after a few more iterations. However, the norm of the model increases significantly when the data is fit better. Dimensionality analysis of the full impedance tensor suggests that the misfit of the model of Figure 7 is reasonable given the degree of 3-D complications in this data set (Siripunvaraporn et al., 1998). Thus, we prefer the model with the misfit at the 4th iteration. These CPU times are large, but this is a very large data set ($N = 8774$) and big model ($M = 14800$).

DISCUSSIONS AND CONCLUSIONS

The REBOCC inversion has been shown to be effective in practice. By using a relatively small subset of the representers, computational requirements (both memory and CPU time) can be substantially reduced. The choice of basis functions for REBOCC is very natural and is dictated by what features can

be resolved by the available data. Even though a small subset of representers is used to form the solution, we emphasize that the goal of the inversion is to find the norm minimizing model subject to fitting all of the data adequately. In particular with REBOCC, it is practical to use the full data set instead of a subset of periods, as is frequently done with OCCAM or NLCG.

In the numerical experiments with synthetic data with the full data set, REBOCC is significantly faster than most other inversion methods (NLCG, DASOCC, and OCCAM), but slower than RRI (Siripunvaraporn, 1999). However, we could not get RRI to converge for joint mode (TM+TE) inversion of our test data set (Siripunvaraporn, 1999). Possibly, a more experienced user of this program might be able to vary or adjust some parameters to make the inversion work successfully.

Experiments with the decimated data set (three periods per decade, comparable to common practice with NLCG or OCCAM) shows that REBOCC is still faster than DASOCC and NLCG. This is particularly true if we consider that the run times quoted for NLCG are for a single (and optimized) value of λ . To find a true minimum structure solution (as REBOCC and DASOCC do), we would have to run NLCG with several trial values of λ . Without a doubt all of these methods are significantly faster than OCCAM, but slower than RRI. However, as with the case of the full data set, we could not get RRI to converge for the joint mode inversion.

CG and NLCG are descent methods which make no use of the second derivative (Hessian) of the penalty functional. GN and OCCAM essentially calculate the full Hessian, while REBOCC make a very good (but not perfect) approximation. Our results indicate that this approximate calculation is worth the effort. This is particularly true for the 2-D case considered here, where direct LU factorization of the differential equation coefficient matrix is feasible, since in this case many sensitivities can be computed quickly once the factorization is complete.

Using the decimated data set might result in significantly less CPU time and memory. However, there is no guarantee that the model obtained in this way will fit the data that are omitted well enough. Models obtained by inverting a decimated data set might sometimes depend on which data are selected. This is much less likely to be the case for REBOCC, since we still require the model to fit all of the data. With only 2–3 points fit per decade, even a single bad estimate could result in significant model errors or difficulties in convergence. Requiring that all data be fit (as REBOCC does) should be expected to improve stability and robustness of the inversion.

The basic idea behind the REBOCC algorithm could also be applied to the 3-D inversion problem. We thus do not agree with the statement made by Rodi and Mackie (2000) that any sorts of inversion based on a sensitivity calculation will not be practical for realistic 3-D EM problems, even allowing for improvement of the computer hardware. A straightforward extension of REBOCC to 3-D inversion is readily apparent, although for calculating the sensitivity matrix the LU decomposition would probably be impractical (except on a supercomputer) and have to be replaced by a relaxation method.

The possibility of adapting the reduced basis method to other inversion approaches also deserves consideration. For example, one could try to solve equation (19) using conjugate gradients. As in the model space conjugate gradient approach used by Mackie and Madden (1993), the matrix multiplications

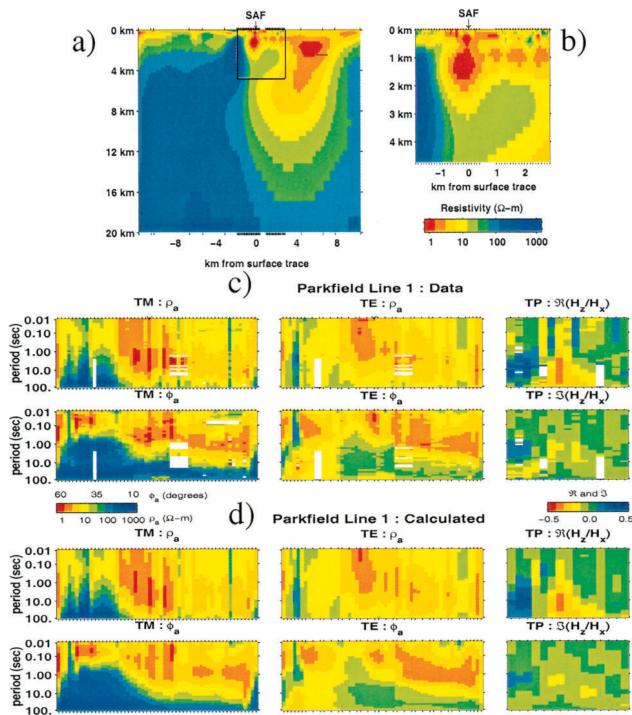


FIG. 7. Results from the 8th-stripe REBOCC inversion on a real data set from a high-density profile across the SAF: (a) the inverse model covering the whole profile; (b) a zoom of the rectangular region near the fault. Data sites are indicated by tick marks along the surface. (c) Measured data for TM, TE, and tippers; (d) corresponding calculated responses.

required can be reduced to solving the forward problem twice per inner-loop iteration. Using the reduced basis idea, forward calculations could be made for only a subset of frequencies. Sensitivity and cross-product matrices would not have to be constructed or stored. However, a full conjugate gradient solution would be required for each λ , so a true OCCAM-type approach would probably not be practical. A similar data space conjugate gradient scheme has been applied to oceanographic and meteorological inversion problems by Egbert (1997) and Bennett et al. (1996), respectively.

In summary, we believe that REBOCC is a significant advance in practical methods for solution of 2-D MT inverse problems, and that the basic ideas of our implementation should be highly relevant to the 3-D problem in some form.

ACKNOWLEDGMENTS

This research has been partially supported by the Royal Thai Embassy via the Development and Promotion of the Science and Technology (DPST) scholarship to W.S., and by grants from the US DOE (DE-FG03-96ER1495) and NSF (9614411-EAR). Comments from Randy Mackie, Chester Weiss, and two anonymous reviewers helped improve the manuscript.

REFERENCES

- Barret, R., Berry, M., Chan, T. F., Demmel, J., Donato, J., Dongarra, J., Eijkhout, V., Pozo, R., Romine, C., and Van der Vorst, H., 1994, Templates for the solution of linear systems: Building blocks for iterative methods: Soc. Ind. Appl. Math.
- Bennett, A. F., Chua, B. S., and Leslie, L. M., 1996, Generalized inversion of a global numerical weather prediction model: Meteorology and Atmospheric Physics, **60**, 165–178.
- Chen, L., Booker, J. R., Jones, A. G., Wu, N., Unsworth, M., Wei, W., and Tan, H., 1996, Electrically conductive crust in southern Tibet from INDEPTH magnetotelluric surveying: Science, **274**, 1694–1696.
- Constable, C. S., Parker, R. L., and Constable, C. G., 1987, Occam's inversion: A practical algorithm for generating smooth models from electromagnetic sounding data: Geophysics, **52**, 289–300.
- deGroot-Hedlin, C., and Constable, S., 1990, Occam's inversion to generate smooth, two-dimensional models from magnetotelluric data: Geophysics, **55**, 1613–1624.
- 1993, Occam's inversion and the North American Central Plains electrical anomaly: J. Geomag. Geoelectr., **45**, 985–999.
- Egbert, G. D., 1997, Tidal data inversion: interpolation and inference: Prog. Oceanog., **40**, 53–80.
- Egbert, G. D., Bennett, A. F., and Foreman, M. G., 1994, TOPEX/POSEIDON tides estimated using a global inverse model: J. Geophys. Res., **99**, 24 821–24 852.
- Farquharson, C. G., and Oldenburg, D. W., 1996, Approximate sensitivities for the electromagnetic inverse problem: Geophysics. J. Internat., **126**, 235–252.
- Jones, A. G., 1992, Electrical conductivity of the continental lower crust, in Fountain, D. M., Arculus, R. J., and Kay, R. W., Eds., Continental lower crust: Elsevier Science Publ. Co., Inc., 81–143.
- Jupp, D. L. B., and Vozoff, K., 1975, Stable iterative methods for the inversion of geophysical data: Geophys. J. Roy. Astr. Soc., **42**, 957–976.
- Kershaw, D. S., 1978, The Incomplete Cholesky-conjugate gradient method for the iterative solution of systems of linear equations: J. Comput. Phys., **26**, 43–65.
- Mackie, R. L., and Madden, T. R., 1993, Three-dimensional magnetotelluric inversion using conjugate gradients: Geophys. J. Internat., **115**, 215–229.
- Marquardt, D. W., 1963, An algorithm for least-squares estimation of nonlinear parameters: J. Soc. Indust. Appl. Math., **11**, 431–441.
- Oldenburg, D. W., and Ellis, R. G., 1991, Inversion of geophysical data using an approximate inverse mapping: Geophys. J. Internat., **105**, 325–353.
- 1993, Efficient inversion of magnetotelluric data in two dimensions: Phys. Earth Planet. Internat., **81**, 177–200.
- Oldenburg, D. W., McGillivray, P. R., and Ellis, R. G., 1993, Generalized subspace methods for large scale inverse problems: Geophys. J. Internat., **114**, 12–20.
- Ogawa, Y. and Uchida, T., 1996, A two-dimensional magnetotelluric inversion assuming Gaussian static shift: Geophysics, **126**, 69–76.
- Orange, A. S., 1989, Magnetotelluric Exploration for Hydrocarbons: Proc. IEEE, **77**, 287–317.
- Parker, R. L., 1980, The inverse problem of electromagnetic induction: existence and construction of solutions based upon incomplete data: J. Geophys. Res., **85**, 4421–4425.
- 1994, Geophysical inverse theory: Princeton Univ. Press.
- Parker, R. L., and Shure, L., 1982, Efficient modeling of the earth's magnetic field with harmonics splines: Geophys. Res. Lett., **9**, 812–815.
- Press, W. H., Teukolsky, S. A., Vetterling, W. T., and Flannery, B. P., 1992, Numerical Recipes in FORTRAN: The art of scientific computing, 2nd ed.: Cambridge Univ. Press.
- Rodi, W. L., 1976, A technique for improving the accuracy of finite element solutions for magnetotelluric data: Geophys. J. Roy. Astr. Soc., **44**, 483–506.
- Rodi, W. L., and Mackie, R. L., 2000, Nonlinear conjugate gradients algorithm for 2-D magnetotelluric inversion: Geophysics, in press.
- Siripunvaraporn, W., 1999, An Efficient data-subspace two-dimensional magnetotelluric inversion and its application to high resolution profile across the San Andreas faults at Parkfield, California: Ph.D. dissertation, Oregon State Univ.
- Siripunvaraporn, W., Egbert, G. D., Eisel, M., and Unsworth, M., 1998, A high resolution EM survey of the San Andreas fault (SAF): Local conductivity structure in a regional context: EOS, **79**, F227.
- Smith, J. T. and Booker, J. R., 1988, Magnetotelluric inversion for minimum structure: Geophysics, **53**, 1565–1576.
- 1991, Rapid inversion of two- and three-dimensional magnetotelluric data: J. Geophys. Res., **96**, 3905–3922.
- Tikhonov, A. N., and Arsenin, V. Y., 1977, Solutions of ill-posed problems: John Wiley and Sons, Inc.
- Uchida, T., 1993, Smooth 2-D inversion for magnetotelluric data based on statistical criterion ABIC: J. Geomag. Geoelectr., **45**, 841–858.
- Unsworth, M. J., Bedrosian, P., Egbert, G. D., and Eisel, M., 1997, 3-D electrical resistivity structure of the San Andreas fault at Parkfield, California: EOS, **78**, F696.
- Unsworth, M., Egbert, G., and Booker, J., 1999, High resolution electromagnetic imaging of the San Andreas fault in central California: J. Geophys. Res., **104**, 1131–1150.
- Vozoff, K., 1972, The magnetotelluric method in the exploration of sedimentary basins: Geophysics, **37**, 98–141.
- Wahba, G., 1990, Spline methods for observational data: Soc. Ind. Appl. Math.
- Wannamaker, P. E., Stodt, J. A., and Rijo, L., 1986, Two-dimensional topographic responses in magnetotelluric modeled using finite elements: Geophysics, **51**, 2131–2144.
- Wannamaker, P. E., Booker, J. R., Filloux, J. H., Jones, A. G., Jiracek, G. R., Chave, A. D., Tarits, P., Waff, H. S., Egbert, G. D., Young, C. T., Stodt, J. A., Martinez, M. G., Lwaw, L. K., Yukutake, T., Segawa, J. S., White, A., and Green, A. W., Jr, 1994, Magnetotelluric observations Across the Juan de Fuca subduction system in the EMSLAB Project: J. Geophys. Res., **94**, 14 111–14 125.
- Weidelt, P., 1972, The inverse problem of geomagnetic induction: Z. Geophys., **38**, 257–289.
- Wu, N., Booker, J. R., and Smith, J. T., 1993, Rapid two-dimensional inversion of COPROD2 data: J. Geomag. Geoelectr., **45**, 1073–1087.

APPENDIX A

MODEL COVARIANCE

Initially, we consider covariances of the general form

$$\mathcal{C}_{\mathcal{M}}(\mathbf{r}, \mathbf{r}') = \eta^2(\mathbf{r})\varrho(\mathbf{r} - \mathbf{r}'), \quad (\text{A-1})$$

where $\eta^2(\mathbf{r})$ is a prior model variance at \mathbf{r} in the model domain, and $\varrho(\mathbf{r} - \mathbf{r}') = \exp(-[(\mathbf{r} - \mathbf{r}')/r_e]^2)$ is the model correlation (with length scale r_e). For our initial discussion, we assume $\eta^2(\mathbf{r})$ is a constant. $\mathbf{C}_{\mathbf{m}}$ is the discrete representation of the covariance function $\mathcal{C}_{\mathcal{M}}(\mathbf{r}, \mathbf{r}')$. Multiplication of a vector \mathbf{u} in the discrete model parameter space (e.g., \mathbf{u} could be a column of \mathbf{J}^T) by $\mathbf{C}_{\mathbf{m}}$ is a discrete representation of the integral

$$\int \mathcal{C}_{\mathcal{M}}(\mathbf{r}, \mathbf{r}')\mathcal{U}(\mathbf{r}')d\mathbf{r}', \quad (\text{A-2})$$

in which $\mathcal{U}(\mathbf{r}')$ is smoothed by convolution with $\mathcal{C}_{\mathcal{M}}$, and \mathbf{u} is the discrete representation of \mathcal{U} .

Egbert et al. (1994) show that this integral can be computed (up to a scalar factor) by introducing a “pseudotime” t and solving the diffusion equation

$$\frac{\partial \mathcal{U}}{\partial t} = \gamma \nabla^2 \mathcal{U} \quad (\text{A-3})$$

with initial condition $\mathcal{U}(\mathbf{r}')$. Here ∇^2 is the 2-D Laplacian operator and γ is a diffusion parameter. If the pseudo diffusion time τ is chosen so that $r_e = \sqrt{4\gamma\tau}$, the integral (A-2) is given by $\mathcal{U}(\mathbf{r}, t = \tau)$. The matrix product $\mathbf{C}_{\mathbf{m}}\mathbf{u}$ can thus be computed by explicit time stepping of equation (A-3) on the model grid from $t = 0$ to τ . In general, allowing for a spatially varying variance $\eta^2(\mathbf{r})$, the scheme can be expressed as

$$\mathbf{C}_{\mathbf{m}} = \Sigma \mathbf{D}_0^{-\frac{1}{2}} \mathbf{D}^{\tau} \mathbf{D}_0^{-\frac{1}{2}} \Sigma, \quad (\text{A-4})$$

where Σ is the diagonal discrete model space variance matrix, \mathbf{D} is the sparse diffusion operator matrix, and \mathbf{D}_0 is a diagonal normalization matrix. The normalization factors on the diagonal of \mathbf{D}_0 are calculated from solving equation (A-3) by replacing \mathcal{U} with a delta function as an initial condition. Therefore, computation of this normalization matrix requires solving the

diffusion equation M times, but only once at the beginning of the process. Note that each step (i.e., multiplication by \mathbf{D}) represents a local smoothing of the field which requires minimal storage or computation time.

Egbert et al. (1994) applied this general approach to an oceanographic inverse problem. For REBOCC and DASOCC, we consider a slightly different approach based on a simple application of operator splitting methods (Press et al., 1992) which is more efficient and better suited to the MT inverse problem. With this approach, the 2-D diffusion equation is replaced by a series of 1-D problems alternating between vertical and horizontal directions:

$$\mathbf{C}_{\mathbf{m}} = \Sigma \mathbf{D}_0^{-\frac{1}{2}} \left[\mathbf{D}_{\mathbf{H}}^{\frac{1}{2}} \mathbf{D}_{\mathbf{V}} \mathbf{D}_{\mathbf{H}}^{\frac{1}{2}} \right]^{\tau} \mathbf{D}_0^{-\frac{1}{2}} \Sigma. \quad (\text{A-5})$$

Solutions of each 1-D diffusion equation (i.e., multiplication by $\mathbf{D}_{\mathbf{H}}$ and $\mathbf{D}_{\mathbf{V}}$) is rapid with a fully implicit method which is trivially implemented for a 1-D problem. This scheme is stable for arbitrarily long “time steps.” Since we do not need an accurate solution of the diffusion equation to define a reasonable model correlation function, we use only a small number of pseudotime steps, τ . Different length scales can be used in vertical and horizontal directions as described earlier.

The model covariance of equation (A-5) allows only for perturbations around the assumed background model, which may in fact be grossly wrong. We can make some allowance for this in the model covariance by adding a constant matrix \mathbf{K} (all elements are κ) so that the covariance matrix $\mathbf{C}_{\mathbf{m}}$ becomes

$$\mathbf{C}_{\mathbf{m}} = \Sigma \left(\mathbf{K} + \mathbf{D}_0^{-\frac{1}{2}} \left[\mathbf{D}_{\mathbf{H}}^{\frac{1}{2}} \mathbf{D}_{\mathbf{V}} \mathbf{D}_{\mathbf{H}}^{\frac{1}{2}} \right]^{\tau} \mathbf{D}_0^{-\frac{1}{2}} \right) \Sigma.$$

The addition of \mathbf{K} to the model covariance corresponds to allowing for uncertainties in the level of a constant background resistivity. In our experience, $\kappa = 1$ usually works well [with $\eta^2(\mathbf{r}) = 1$]; however, the limits of this choice need to be better understood.