

# LINEAR AND NONLINEAR INVERSE PROBLEMS

ROEL SNIEDER AND JEANNOT TRAMPERT

*Dept. of Geophysics*

*Utrecht University*

*P.O. Box 80.021*

*3508 TA Utrecht*

*The Netherlands*

*email snieder@geo.uu.nl*

## 1. Introduction

An important aspect of the physical sciences is to make inferences about physical parameters from data. In general, the laws of physics provide the means for computing the data values given a model. This is called the “forward problem”, see figure 1. In the inverse problem, the aim is to reconstruct the model from a set of measurements. In the ideal case, an exact theory exists that prescribes how the data should be transformed in order to reproduce the model. For some selected examples such a theory exists assuming that the required infinite and noise-free data sets would be available. A quantum mechanical potential in one spatial dimension can be reconstructed when the reflection coefficient is known for all energies [Marchenko, 1955; Burridge, 1980]. This technique can be generalized for the reconstruction of a quantum mechanical potential in three dimensions [Newton, 1989], but in that case a redundant data set is required for reasons that are not well understood. The mass-density in a one-dimensional string can be constructed from the measurements of all eigenfrequencies of that string [Borg, 1946], but due to the symmetry of this problem only the even part of the mass-density can be determined. If the seismic velocity in the earth depends only on depth, the velocity can be constructed exactly from the measurement of the arrival time as a function of distance of seismic waves using an Abel transform [Herglotz, 1907; Wiechert, 1907]. Mathematically this problem is identical to the construction of a spherically symmetric quantum mechanical potential in three dimensions [Keller *et al.*, 1956]. However, the construction

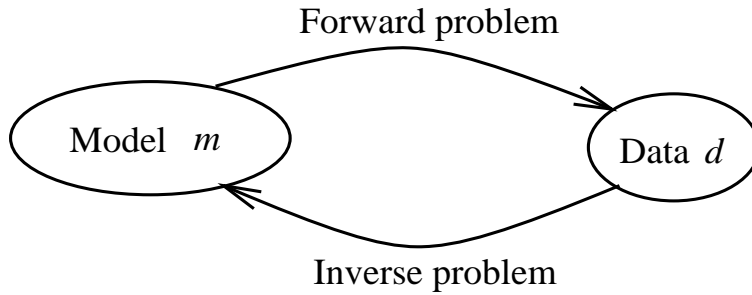


Figure 1. The traditional definition of the forward and inverse problems.

method of Herglotz-Wiechert only gives an unique result when the velocity increases monotonically with depth [Gerver and Markushevitch, 1966]. This situation is similar in quantum mechanics where a radially symmetric potential can only be constructed uniquely when the potential does not have local minima [Sabatier, 1973].

Despite the mathematical elegance of the exact nonlinear inversion schemes, they are of limited applicability. There are a number of reasons for this. First, the exact inversion techniques are usually only applicable for idealistic situations that may not hold in practice. For example, the Herglotz-Wiechert inversion presupposes that the velocity in the earth depends only on depth and that the velocity increases monotonically with depth. Seismic tomography has shown that both requirements are not met in the earth's mantle [Nolet *et al.*, 1994]. Second, the exact inversion techniques often are very unstable. The presence of this instability in the solution of the Marchenko equation has been shown explicitly by Dorren *et al.* [1994]. However, the third reason is the most fundamental. In many inverse problems the model that one aims to determine is a continuous function of the space variables. This means that the model has infinitely many degrees of freedom. However, in a realistic experiment the amount of data that can be used for the determination of the model is usually finite. A simple count of variables shows that the data cannot carry sufficient information to determine the model uniquely. In the context of linear inverse problems this point has been raised by Backus and Gilbert [1967, 1968] and more recently by Parker [1994]. This issue is equally relevant for nonlinear inverse problems.

The fact that in realistic experiments a finite amount of data is available to reconstruct a model with infinitely many degrees of freedom necessarily means that the inverse problem is not unique in the sense that there are many models that explain the data equally well. The model obtained from the inversion of the data is therefore not necessarily equal

to the true model that one seeks. This implies that the view of inverse problems as shown in figure 1 is too simplistic. For realistic problems, inversion really consists of two steps. Let the true model be denoted by  $m$  and the data by  $d$ . From the data  $d$  one reconstructs an estimated model  $\tilde{m}$ , this is called the *estimation problem*, see figure 2. Apart from estimating a model  $\tilde{m}$  that is consistent with the data, one also needs to investigate what relation the estimated model  $\tilde{m}$  bears to the true model  $m$ . In the *appraisal problem* one determines what properties of the true model are recovered by the estimated model and what errors are attached to it. The essence of this discussion is that *inversion* = *estimation* + *appraisal*. It does not make much sense to make a physical interpretation of a model without acknowledging the fact of errors and limited resolution in the model [Trampert, 1998].

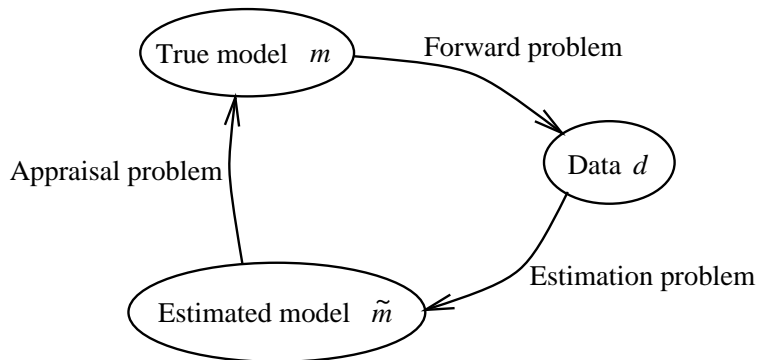


Figure 2. The inverse problem viewed as a combination of an estimation problem plus an appraisal problem.

In general there are two reasons why the estimated model differs from the true model. The first reason is the non-uniqueness of the inverse problem that causes several (usually infinitely many) models to fit the data. Technically, this model null-space exists due to inadequate sampling of the model space. The second reason is that real data (and physical theories more often than we would like) are always contaminated with errors and the estimated model is therefore affected by these errors as well. Therefore model appraisal has two aspects, non-uniqueness and error propagation.

Model estimation and model appraisal are fundamentally different for discrete models with a finite number of degrees of freedom and for continuous models with infinitely many degrees of freedom. Also, the problem of model appraisal is only well-solved for linear inverse problems. For this reason the inversion of discrete models and continuous models is treated separately, and the case of linear inversion and nonlinear inversion is also

treated independently. In section 2 linear inversion for a finite number of model parameters is discussed. This is generalized in section 3 to deal with linear inverse problems for continuous models with infinitely many degrees of freedom. In reality many inverse problems are not really linear, but often these problems can be linearized by making a suitable approximation. In section 4 the single-scattering approximation is derived. This technique forms the basis of imaging tools used in reflection seismology. Rayleigh's principle, as treated in section 5, is the linearization that forms the basis for the inversion for the Earth's structure using normal-mode frequencies. The linearization technique of seismic travel time tomography is based on Fermat's principle, which is treated in section 6. Non-linear inverse problems are significantly more difficult than linear inverse problems. It is shown in section 7 that non-linearity can be a source of ill-posedness. Presently, there is no satisfactory theory for the appraisal problem for nonlinear inverse problems. In section 8 three methods are presented that can be used for the nonlinear appraisal problem. However, neither of these methods is quite satisfactory, which indicates that nonlinear inverse problem theory is a field with important research challenges.

## 2. Solving finite linear systems of equations

As argued in the previous section, the inverse problem maps a finite number of data onto a model. In most practical applications in geophysics the model is a continuous function of the space coordinates and therefore has infinitely many degrees of freedom. For the moment we will ignore this and will assume that the model can be characterized by a finite number of parameters. We will return to the important case of models that are infinitely dimensional in section 3.

### 2.1. LINEAR MODEL ESTIMATION

For a finite-dimensional model, the model parameters can be ordered in a vector  $\mathbf{m}$ , and similarly the data can be ordered in a vector  $\mathbf{d}$ . The matrix  $\mathbf{A}$  relates the data to the model through the product  $\mathbf{A}\mathbf{m}$ . This matrix is often referred to as the theory operator. Indeed, it contains all the information on physics and mathematics we have chosen to model in the given problem. In practice, the data are contaminated with errors  $\mathbf{e}$ , so that the recorded data and the model are related by:

$$\mathbf{d} = \mathbf{A}\mathbf{m} + \mathbf{e} \quad (1)$$

It should be noted there often there is an certain arbitrariness in the choice of the model parameters that are contained in the model vector  $\mathbf{m}$ . For



example, if one wants to describe the density in the earth one could choose a model where the Earth's mantle and the core have a uniform density, in that case there are two model parameters. Alternatively, one could expand the density in the Earth in a large amount of eigenfunctions defined on the sphere such as spherical harmonics for lateral variations and polynomials for depth variations, in that case there are much more model parameters. These two different parameterizations of the same model correspond to different model parameters  $\mathbf{m}$  and to a different matrix  $\mathbf{A}$ . This example illustrates that the model  $\mathbf{m}$  is not necessarily the true model,<sup>1</sup> but that the choice of the model parameters usually contains a restriction on the class of models that can be constructed. Below we will refer to  $\mathbf{m}$  as the true model regardless of the difficulties in its definition.

From the recorded data one makes an estimate of the model. Since this estimate in practice will be different from the true model the estimated model is denoted by  $\tilde{\mathbf{m}}$ . There are many ways for designing an inverse operator that maps the data on the estimated model [e.g. *Menke, 1984; Tarantola, 1987; Parker, 1994*]. Whatever estimator one may choose, the most general linear mapping from data to the estimated model can be written as:

$$\tilde{\mathbf{m}} = \mathbf{A}^{-g} \mathbf{d} \quad (2)$$

The operator  $\mathbf{A}^{-g}$  is called the *generalized inverse* of the matrix  $\mathbf{A}$ . In general, the number of data is different from the number of model parameters. For this reason  $\mathbf{A}$  is usually a non-square matrix, and hence its formal inverse does not exist. Later we will show how the generalized inverse  $\mathbf{A}^{-g}$  may be chosen, but for the moment  $\mathbf{A}^{-g}$  does not need to be specified. The relation between the estimated model  $\tilde{\mathbf{m}}$  and the true model  $\mathbf{m}$  follows by inserting (1) in expression (2):

$$\tilde{\mathbf{m}} = \mathbf{A}^{-g} \mathbf{A} \mathbf{m} + \mathbf{A}^{-g} \mathbf{e} \quad (3)$$

The matrix  $\mathbf{A}^{-g} \mathbf{A}$  is called the *resolution kernel*, this operator is given by:

$$\mathbf{R} \equiv \mathbf{A}^{-g} \mathbf{A} \quad (4)$$

Expression (3) can be interpreted by rewriting it in the following form:

<sup>1</sup>We urge the reader to formulate a definition of the concept "true model." It is not so difficult to formulate a vague definition such as "the true model is the model that corresponds to reality and which is only known to the gods." However, we are not aware of any definition that is operational in the sense that it provides us with a set of actions that could potentially tell us what the true model really is.

$$\tilde{\mathbf{m}} = \mathbf{m} + \underbrace{(\mathbf{A}^{-g} \mathbf{A} - \mathbf{I}) \mathbf{m}}_{\text{Limited Resolution}} + \underbrace{\mathbf{A}^{-g} \mathbf{e}}_{\text{Error propagation}} \quad (5)$$

In the ideal case, the estimated model equals the true model vector:  $\tilde{\mathbf{m}} = \mathbf{m}$  meaning that our chosen parameters, ordered in vector  $\mathbf{m}$ , may be estimated independently from each other. The last two terms in equation (5) account for *blurring* and *artifacts* in the estimated model. The term  $(\mathbf{A}^{-g} \mathbf{A} - \mathbf{I}) \mathbf{m}$  describes the fact that components of the estimated model vector are linear combinations of different components of the true model vector. We only retrieve averages of our parameters and “blurring” occurs in the model estimation as we are not able to map out the finest details. In the ideal case this term vanishes; this happens when  $\mathbf{A}^{-g} \mathbf{A}$  is equal to the identity matrix. With (4) this means that for perfectly resolved model parameters the resolution matrix is the identity matrix:

$$\text{Perfect resolution: } \mathbf{R} = \mathbf{I} \quad (6)$$

As noted earlier, usually there is a certain ambiguity in the definition of the model parameters that define the vector  $\mathbf{m}$ . The resolution operator tells us to what extent we can retrieve the model parameters independently from the estimation process. However, the resolution matrix does not tell us completely what the relation between the estimated model and the real underlying physical model is, because it does not take into account to what extent the choice of the model parameters has restricted the model that can be obtained from the estimation process.

The last term in (5) describes how the errors  $\mathbf{e}$  are mapped onto the estimated model.<sup>2</sup> These errors are not known deterministically, otherwise they could be subtracted from the data. A statistical analysis is needed to describe the errors in the estimated model due to the errors in the data. When the data  $d_j$  are uncorrelated and have standard deviation  $\sigma_{d_j}$ , the standard deviation  $\sigma_{m_i}$  in the model estimate  $\tilde{m}_i$ , resulting from the propagation of data errors only, is given by:

$$\sigma_{m_i}^2 = \sum_j \left( A_{ij}^{-g} \sigma_{d_j} \right)^2 \quad (7)$$

Ideally, one would like to obtain both: a perfect resolution and no errors in the estimated model. Unfortunately this cannot be achieved in practice. The error propagation is, for instance, completely suppressed by using the generalized inverse  $\mathbf{A}^{-g} = 0$ . This leads to the (absurd) estimated model  $\tilde{\mathbf{m}} = 0$  which is indeed not affected by errors. However, this

<sup>2</sup>As shown by *Scales and Snieder* [1998] the concept of errors in inverse problems is not as simple as it appears.

particular generalized inverse has a resolution matrix given by  $\mathbf{R} = 0$ , which is far from the ideal resolution matrix given in (6). Hence in practice, one has to find an acceptable trade-off between error-propagation and limitations in the resolution.

## 2.2. LEAST-SQUARES ESTIMATION

Let us for the moment consider the case where the number of independent data is larger than the number of unknowns. In that case, the system  $\mathbf{d} = \mathbf{A}\mathbf{m}$  cannot always be satisfied for any given model  $\mathbf{m}$  because of possible errors contained in the data vector making the equations inconsistent. As an example, let us consider the following problem. We have two masses with weight  $m_1$  and  $m_2$ . The weighing of the first mass yields a weight of 1 (kilo). Then one measures the second mass to find a weight of 2. Next, one weighs the masses together to find a combined weight of 2. The system of equations that corresponds to these measurements is given by:

$$\begin{aligned} m_1 &= d_1 = 1 \\ m_2 &= d_2 = 2 \\ m_1 + m_2 &= d_3 = 2 \end{aligned} \tag{8}$$

The matrix  $\mathbf{A}$  for this problem is given by:

$$\mathbf{A} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \\ 1 & 1 \end{pmatrix} \tag{9}$$

It is clear that this system of equations cannot be satisfied. It cannot be true that the first mass has a weight  $m_1 = 1$ , and the second mass has a weight  $m_2 = 2$  while the combined mass is equal to  $m_1 + m_2 = 2$ . Clearly errors have been made during the measurements, but there is no reason to discard one of the three equations in favor of the other two. This problem is illustrated graphically in figure 3. The three equations (8) correspond to the three solid lines in the  $(m_1, m_2)$ -plane. The fact that the three lines do not intersect in a single point signifies that the linear equations are inconsistent. The inverse problem of determining the two masses thus consists in reconciling these equations in a meaningful way.

A common way to estimate a model is to seek the model  $\tilde{\mathbf{m}}$  that gives the best fit to the data in the sense that the difference, measured by the  $L_2$ -norm, between the data vector  $\mathbf{d}$  and the recalculated data  $\mathbf{A}\tilde{\mathbf{m}}$  is made as small as possible. This means that the least-squares solution is given by the model that minimizes the following cost function:

$$S = \|\mathbf{d} - \mathbf{A}\mathbf{m}\|^2 \tag{10}$$

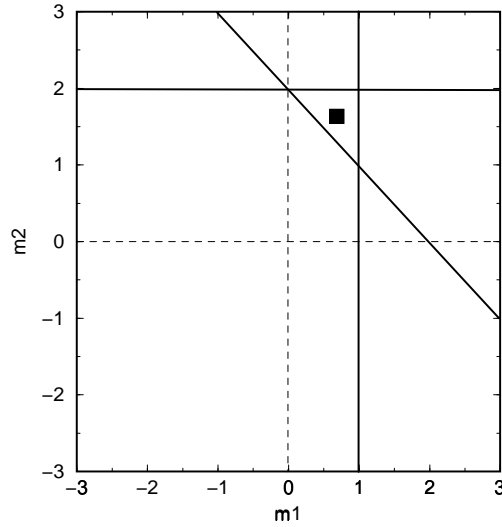


Figure 3. Geometrical interpretation of the linear equations (8).

As shown in detail by *Strang* [1988] this quantity is minimized by the following model estimate:

$$\tilde{\mathbf{m}} = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{d} \quad (11)$$

In the example of figure 3 the least-squares solution is the point in the  $(m_1, m_2)$ -plane that has the smallest distance to the three lines in that figure, this point is indicated by a black square. Using the matrix (9) one readily finds that the least-squares estimator of the problem (8) is given by:

$$\tilde{\mathbf{m}} = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{d} = \frac{1}{3} \begin{pmatrix} 2 & -1 & 1 \\ -1 & 2 & 1 \end{pmatrix} \mathbf{d} \quad (12)$$

For the used data vector this means that the estimated model is given by:

$$\begin{aligned} \tilde{m}_1 &= 2/3 \\ \tilde{m}_2 &= 5/3 \end{aligned} \quad (13)$$

### 2.3. MINIMUM NORM ESTIMATION

In some problems the number of unknowns is less than the number of parameters. Consider for example the situation where there are two masses  $m_1$  and  $m_2$  and one has measured only the combined weight of these

masses:

$$m_1 + m_2 = d = 2 \quad (14)$$

The matrix that corresponds to this system of one equation is given by:

$$\mathbf{A} = \begin{pmatrix} 1 & 1 \end{pmatrix} \quad (15)$$

Graphically this problem is depicted in figure 4. Clearly any model vector lying on the solid line fits the equation (14) exactly. There are thus infinitely many solutions, provided the masses are positive, that exactly fit the data. A model estimate can be defined by choosing a model that fits the data exactly and that has the smallest  $L_2$ -norm, this model is indicated by in figure 4 by the black square.

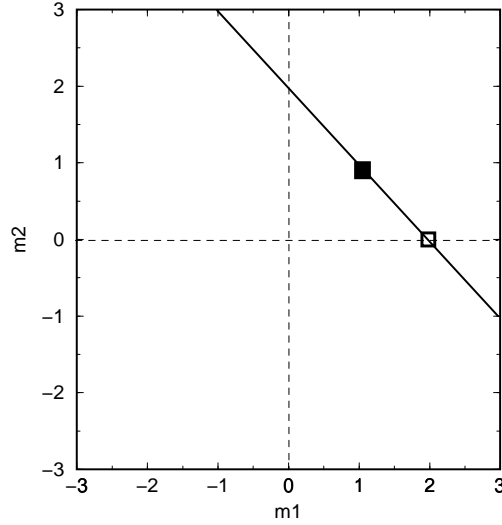


Figure 4. Geometrical interpretation of the linear equation (14) with two unknowns.

For a general underdetermined system of equations the minimum norm solution is defined as the model that fits the data exactly,  $\mathbf{A}\mathbf{m} = \mathbf{d}$ , and that minimizes  $\|\mathbf{m}\|^2$ . Using Lagrange multipliers one can show that the minimum-norm solution is given by:

$$\tilde{\mathbf{m}} = \mathbf{A}^T (\mathbf{A}\mathbf{A}^T)^{-1} \mathbf{d} , \quad (16)$$

a detailed derivation is given by *Menke* [1984]. One readily finds that the minimum norm solution of “system” (14) is given by

$$m_1 = m_2 = 1 . \quad (17)$$

## 2.4. MIXED DETERMINED PROBLEMS

In the least-squares estimation, we assumed that we had enough information to evaluate all model parameters, even though contradictions occurred due to measurement errors. The problem is then purely over-determined and as a consequence  $\mathbf{A}^T \mathbf{A}$  is regular. In the minimum norm solution, we assumed no contradictions in the available information, but we don't have enough equations to evaluate all model parameters. This is the case of a purely under-determined problem and here  $\mathbf{A} \mathbf{A}^T$  is regular. The most common case, however, is that we have contradictory information on some model parameters, while others cannot be assessed due to a lack of information. Then neither  $\mathbf{A}^T \mathbf{A}$  nor  $\mathbf{A} \mathbf{A}^T$  can be inverted and the problem is ill-posed. Even if the inverse matrices formally exist, they are often ill-conditioned meaning that small changes in the data vector lead to large changes in the model estimation. This means that errors in the data will be magnified in the model estimation. Clearly a trick is needed to find a model that is not too sensitive on small changes in the data. To this effect, *Levenberg* [1944] introduced a damped least-squares solution. From a mathematical point of view, ill-posedness and ill-conditioning result from zero or close to zero singular values of  $\mathbf{A}$ .

Suppose one has a matrix  $\mathbf{M}$  with eigenvalues  $\lambda_n$  and eigenvectors  $\hat{\mathbf{v}}_n$ :

$$\mathbf{M} \hat{\mathbf{v}}_n = \lambda_n \hat{\mathbf{v}}_n \quad (18)$$

One readily finds that the matrix  $(\mathbf{M} + \gamma \mathbf{I})$  has eigenvalues  $(\lambda_n + \gamma)$ :

$$(\mathbf{M} + \gamma \mathbf{I}) \hat{\mathbf{v}}_n = (\lambda_n + \gamma) \hat{\mathbf{v}}_n \quad (19)$$

This means that the eigenvalues of a matrix can be raised by adding the scaled identity matrix to the original matrix. This property can be used to define the *damped least-squares solution*:

$$\tilde{\mathbf{m}} = \left( \mathbf{A}^T \mathbf{A} + \gamma \mathbf{I} \right)^{-1} \mathbf{A}^T \mathbf{d} \quad (20)$$

Since the matrix  $\mathbf{A}^T \mathbf{A}$  has positive eigenvalues<sup>3</sup> its eigenvalues are moved away from zero when the constant  $\gamma$  is positive. Alternatively, the solution (20) can be found by minimizing the following cost function:

<sup>3</sup>That the eigenvalues of  $\mathbf{A}^T \mathbf{A}$  are positive follows from the following identity:  $(\mathbf{x} \cdot \mathbf{A}^T \mathbf{A} \mathbf{x}) = (\mathbf{A}^T \mathbf{x} \cdot \mathbf{A} \mathbf{x}) = (\mathbf{A} \mathbf{x} \cdot \mathbf{A} \mathbf{x}) = \|\mathbf{A} \mathbf{x}\|^2 \geq 0$ . When  $\mathbf{x}$  is the eigenvector  $\hat{\mathbf{u}}^{(n)}$  of  $\mathbf{A}^T \mathbf{A}$  with eigenvalue  $\mu_n$ , this expression can be used to show that  $\mu_n \|\hat{\mathbf{u}}^{(n)}\|^2 = \mu_n (\hat{\mathbf{u}}^{(n)} \cdot \hat{\mathbf{u}}^{(n)}) = (\hat{\mathbf{u}}^{(n)} \cdot \mathbf{A}^T \mathbf{A} \hat{\mathbf{u}}^{(n)}) \geq 0$ , hence the eigenvalues  $\mu_n$  are positive.

$$S = \|\mathbf{d} - \mathbf{A}\mathbf{m}\|^2 + \gamma \|\mathbf{m}\|^2 \quad (21)$$

This expression clearly shows what the effect of the damping is. Minimizing the first term of (21) amounts to finding the model that gives the best fit to the data. Minimizing the last term of (21) amounts to finding the model with the smallest norm. In general we cannot minimize both terms simultaneously, but in minimizing (21) we comprise in finding a model that both fits the data reasonably well and whose model size is not too large. The parameter  $\gamma$  controls the emphasis we put on these conflicting requirements and for this reason it is called the *trade-off parameter*.

For a number of applications the following matrix identity is extremely useful:<sup>4</sup>

$$\boxed{\left(\mathbf{A}^T \mathbf{B}^{-1} \mathbf{A} + \mathbf{D}^{-1}\right)^{-1} \mathbf{A}^T \mathbf{B}^{-1} = \mathbf{D} \mathbf{A}^T \left(\mathbf{B} + \mathbf{A} \mathbf{D} \mathbf{A}^T\right)^{-1}} \quad (22)$$

In this expression  $\mathbf{B}$  and  $\mathbf{D}$  are regular square matrices, whereas  $\mathbf{A}$  needs not to be square. This expression can be used to show that when damping or regularization is used, the least-squares solution and the minimum-norm solution (both supplied with a damping term) are identical. To see this use (22) with  $\mathbf{B}^{-1} = \mathbf{I}$  and  $\mathbf{D}^{-1} = \gamma \mathbf{I}$ . It then follows that

$$\left(\mathbf{A}^T \mathbf{A} + \gamma \mathbf{I}\right)^{-1} \mathbf{A}^T \mathbf{d} = \frac{1}{\gamma} \mathbf{A}^T \left(\mathbf{I} + \mathbf{A} \frac{1}{\gamma} \mathbf{A}^T\right)^{-1} \mathbf{d} = \mathbf{A}^T \left(\mathbf{A} \mathbf{A}^T + \gamma \mathbf{I}\right)^{-1} \mathbf{d}. \quad (23)$$

The left hand side corresponds to the damped least-squares solution (20) while the right hand side is the damped version of the minimum-norm solution (16). This implies that when damping is applied the least-squares solution and the minimum-norm solution are identical.

## 2.5. THE CONSISTENCY PROBLEM FOR THE LEAST-SQUARES SOLUTION

The least-squares solution appears to provide an objective method for finding solutions of overdetermined problems. However, there is trouble ahead. To see this, let us consider the overdetermined system of equations (8). Mathematically, this system of equations does not change when we

<sup>4</sup>This identity follows from the identity  $\mathbf{A}^T + \mathbf{A}^T \mathbf{B}^{-1} \mathbf{A} \mathbf{D} \mathbf{A}^T = \mathbf{A}^T \mathbf{B}^{-1} \mathbf{A} \mathbf{D} \mathbf{A}^T + \mathbf{A}^T$ . Write the first term on the left hand side as  $\mathbf{A}^T \mathbf{B}^{-1} \mathbf{B}$  and the last term on the right hand side as  $\mathbf{D}^{-1} \mathbf{D} \mathbf{A}^T$ . The resulting expression can then be written as  $\mathbf{A}^T \mathbf{B}^{-1} (\mathbf{B} + \mathbf{A} \mathbf{D} \mathbf{A}^T) = (\mathbf{A}^T \mathbf{B}^{-1} \mathbf{A} + \mathbf{D}^{-1}) \mathbf{D} \mathbf{A}^T$ . The expression (22) then follows by multiplying on the left with  $(\mathbf{B} + \mathbf{A} \mathbf{D} \mathbf{A}^T)^{-1}$  and by multiplying on the right with  $(\mathbf{A}^T \mathbf{B}^{-1} \mathbf{A} + \mathbf{D}^{-1})^{-1}$ .

multiply the last equation with a factor two. The following two systems of equations thus are completely equivalent:

$$\left. \begin{array}{rcl} m_1 = d_1 & = & 1 \\ m_2 = d_2 & = & 2 \\ m_1 + m_2 = d_3 & = & 2 \end{array} \right\} \Leftrightarrow \left\{ \begin{array}{rcl} m_1 = d_1 = d'_1 & = & 1 \\ m_2 = d_2 = d'_2 & = & 2 \\ 2m_1 + 2m_2 = 2d_3 = d'_3 & = & 4 \end{array} \right. \quad (24)$$

The matrices of the original system and the new equivalent system are given by

$$\mathbf{A} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \\ 1 & 1 \end{pmatrix} \quad \text{and} \quad \mathbf{A}' = \begin{pmatrix} 1 & 0 \\ 0 & 1 \\ 2 & 2 \end{pmatrix} \quad (25)$$

In this section the unprimed quantities denote the original system of equations while the primed quantities refer to the transformed system of equations. One readily finds that the least-squares solution (11) of the original system and the transformed system are given by

$$\tilde{\mathbf{m}} = \frac{1}{3} \begin{pmatrix} 2 & -1 & 1 \\ -1 & 2 & 1 \end{pmatrix} \mathbf{d} \quad \text{and} \quad \tilde{\mathbf{m}} = \frac{1}{9} \begin{pmatrix} 5 & -4 & 2 \\ -4 & 5 & 2 \end{pmatrix} \mathbf{d}', \quad (26)$$

Using the numerical values of the original data vector  $\mathbf{d}$  and the transformed data vector  $\mathbf{d}'$  this leads to the following estimates of the model:

$$\tilde{\mathbf{m}} = \begin{pmatrix} 2/3 \\ 5/3 \end{pmatrix} \quad \text{and} \quad \tilde{\mathbf{m}} = \begin{pmatrix} 5/9 \\ 14/9 \end{pmatrix} \quad (27)$$

The problem is that these two estimators of the same model are *different*! This is surprising because the original system of equations and the transformed system of equations in (24) are mathematically equivalent. The reason that the two solutions are different is that the metric in the original data space and in the transformed data space has been changed by the transformation. This is a different way of saying that distances are measured in different ways in the least-squares criteria for solving the two systems of equations. Since the least-squares solution minimizes distances it makes sense that the least-squares solution changes when the metric (or measuring unit) of the data space is changed. This implies that the least-squares solution is not as objective as it appeared at first sight, because arbitrary transformations of the system of equations lead to different least-squares solutions!

For the least-squares solution the generalized inverse is given by  $\mathbf{A}^{-g} = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T$ . One readily deduces that both for the original system and



the transformed system the resolution kernel is equal to the identity matrix:  $\mathbf{R} = \mathbf{A}^{-g} \mathbf{A} = \mathbf{I}$  and  $\mathbf{R}' = \mathbf{A}^{-g'} \mathbf{A}' = \mathbf{I}$ . Hence both systems have perfect resolution! The reader may be interested to pause and explain how this can be reconciled with the fact that the estimated models in (27) are different.

The reason for this discrepancy lies in the error propagation term  $\mathbf{A}^{-g} \mathbf{e}$  in (5). We know that errors must be present in the data used in the systems defined in expression (24) because the equations are inconsistent. After scaling the equations, data and errors are reconciled in different ways in the two systems of equations, so that different model estimators are obtained. *It is thus the presence of inconsistencies in the system of equations caused by errors that creates a dependence of the least-squares solution to arbitrary scaling operations.*

Let us now consider the properties of the least-squares solution under transformations of the data vector and the model vector in a more general way. The initial system of equations is given by

$$\mathbf{A} \mathbf{m} = \mathbf{d} \quad (28)$$

This expression is not quite correct because we ignored the errors  $\mathbf{e}$  which will always be present. This is the reason why the above expression can not exactly be satisfied, and we will have to seek the least-squares solution to this system of equations. Let us consider a transformation of the model parameters through a transformation matrix  $\mathbf{S}$ :

$$\mathbf{m}' = \mathbf{S} \mathbf{m} , \quad (29)$$

and a transformation of the data vector with a transformation matrix  $\mathbf{Q}$ :

$$\mathbf{d}' = \mathbf{Q} \mathbf{d} . \quad (30)$$

Assume that  $\mathbf{S}$  has an inverse, the transformed system of equations then is given by

$$\mathbf{Q} \mathbf{A} \mathbf{S}^{-1} \mathbf{m}' = \mathbf{Q} \mathbf{d} = \mathbf{d}' . \quad (31)$$

The original system of equations (28) has the least-squares solution

$$\boxed{\tilde{\mathbf{m}}^{(1)} = \left( \mathbf{A}^T \mathbf{A} \right)^{-1} \mathbf{A}^T \mathbf{d}} \quad (32)$$

The solution of the transformed system of equations (31) follows from the same expression, setting  $\mathbf{A}' = \mathbf{Q} \mathbf{A} \mathbf{S}^{-1}$  and replacing  $\mathbf{d}$  by  $\mathbf{Q} \mathbf{d}$ . This, however, gives the solution to the transformed model vector  $\mathbf{m}'$ . In order to compare this solution with the model estimate (32) we need to transform back to the original model space, using the relation  $\mathbf{m} = \mathbf{S}^{-1} \mathbf{m}'$ .

The least-squares solution  $\tilde{\mathbf{m}}^{(2)}$  that follows from the transformed system is then given by:

$$\tilde{\mathbf{m}}^{(2)} = \mathbf{S}^{-1} \left( \mathbf{S}^{T-1} \mathbf{A}^T \mathbf{Q}^T \mathbf{Q} \mathbf{A} \mathbf{S}^{-1} \right)^{-1} \mathbf{S}^{T-1} \mathbf{A}^T \mathbf{Q}^T \mathbf{Q} \mathbf{d} \quad (33)$$

Assuming again that the appropriate inverses exist, this expression can be simplified by repeatedly applying the matrix identity  $(\mathbf{NM})^{-1} = \mathbf{M}^{-1} \mathbf{N}^{-1}$  to the term  $\left( \mathbf{S}^{T-1} \mathbf{A}^T \mathbf{Q}^T \mathbf{Q} \mathbf{A} \mathbf{S}^{-1} \right)^{-1}$ , giving  $\left( \mathbf{S}^{T-1} \mathbf{A}^T \mathbf{Q}^T \mathbf{Q} \mathbf{A} \mathbf{S}^{-1} \right)^{-1} = \left( \mathbf{S}^{-1} \right)^{-1} \left( \mathbf{A}^T \mathbf{Q}^T \mathbf{Q} \mathbf{A} \right)^{-1} \left( \mathbf{S}^{T-1} \right)^{-1} = \mathbf{S} \left( \mathbf{A}^T \mathbf{Q}^T \mathbf{Q} \mathbf{A} \right)^{-1} \mathbf{S}^T$ . The least-squares solution of the transformed system can then be written as:

$$\boxed{\tilde{\mathbf{m}}^{(2)} = \left( \mathbf{A}^T \mathbf{Q}^T \mathbf{Q} \mathbf{A} \right)^{-1} \mathbf{A}^T \mathbf{Q}^T \mathbf{Q} \mathbf{d}} \quad (34)$$

Comparing this with the least-squares solution  $\tilde{\mathbf{m}}^{(1)}$  in expression (32) of the original system one finds that the least-squares solution is invariant for:

- Transformations of the model vector altogether.
- Transformations of the data vector if  $\mathbf{Q}^T \mathbf{Q} = \mathbf{I}$ .

The first property can be understood if we recall that the cost function in the least-squares problem does not minimize the model length but only the data misfit. The last property can be understood by comparing the quantities that are minimized in the original system and for the transformed system. For the original system one minimizes:

$$S = \|\mathbf{d} - \mathbf{A}\mathbf{m}\|^2 = ((\mathbf{d} - \mathbf{A}\mathbf{m}) \cdot (\mathbf{d} - \mathbf{A}\mathbf{m})) , \quad (35)$$

while in the transformed system one minimizes

$$\begin{aligned} S' = \|\mathbf{Q}\mathbf{d} - \mathbf{Q}\mathbf{A}\mathbf{m}\|^2 &= (\mathbf{Q}(\mathbf{d} - \mathbf{A}\mathbf{m}) \cdot \mathbf{Q}(\mathbf{d} - \mathbf{A}\mathbf{m})) \\ &= (\mathbf{Q}^T \mathbf{Q}(\mathbf{d} - \mathbf{A}\mathbf{m}) \cdot (\mathbf{d} - \mathbf{A}\mathbf{m})) \end{aligned} \quad (36)$$

These two quantities are identical when the transformation  $\mathbf{Q}$  is unitary, i.e. when  $\mathbf{Q}^T \mathbf{Q} = \mathbf{I}$ . This result stems from the property that unitary matrices do not affect the norm of a vector.

## 2.6. THE CONSISTENCY PROBLEM FOR THE MINIMUM-NORM SOLUTION

Consistency problems do not only arise for the least-squares solution, the minimum-norm solution suffers from the same problem. As an example

let us return to the underdetermined “system” of equations (14). The minimum-norm solution of this problem is given by

$$\boxed{\tilde{m}_1 = 1 \quad , \quad \tilde{m}_2 = 1} \quad (37)$$

Now carry out a transformation from the model vector  $\mathbf{m}$  to a new model vector  $\mathbf{m}'$ :

$$\begin{aligned} m'_1 &= m_1 + m_2 \\ m'_2 &= m_2 \end{aligned} \quad (38)$$

For this new model vector the “system” of equations is given by:

$$\tilde{m}'_1 = d = 2 \quad (39)$$

Note that this transformed model vector brings out the fact that the system is undetermined much more clearly than the original system (14) because the new system imposes no constraint whatsoever on the model parameter  $m'_2$ . The minimum-norm solution of the transformed equation (39) is given by  $\tilde{m}'_1 = 2$ ,  $\tilde{m}'_2 = 0$ . With the transformation (38) this solution in the transformed model space corresponds to the following solution in the original model space:

$$\boxed{\tilde{m}_1 = 2 \quad , \quad \tilde{m}_2 = 0} \quad (40)$$

This solution is shown by the open square in figure 4. Note that this solution differs from the minimum-norm solution (37) for the original system of equations. The reason for this discrepancy is similar to the consistency problem for the least-squares problem in section 2.5; the transformation (38) has changed the metric of model space, so that distances in the original model space and the transformed model space are measured in different ways. For this reason the minimum norm solution of the original problem and transformed problem are different.

We could carry out a similar general analysis for the transformation properties of the minimum norm solution under general transformations of the model vector and data vector as we carried out for the least-squares solution in section 2.5. However, in practice one applies regularization to the equations. As shown in equation (23) the damped least-squares solution and the damped minimum-norm solution are identical. For this reason the general transformation properties are treated in the next section for the damped least-squares solution.

## 2.7. THE NEED FOR A MORE GENERAL REGULARIZATION

The analysis of the transformation properties of the damped least-squares is completely analogous to the analysis of the undamped least-squares solution of section 2.5. Ignoring errors for the moment, the linear system of equations is given by (28):  $\mathbf{A}\mathbf{m} = \mathbf{d}$ , and the transformation of the model vector and data vector is given by (29) and (30) respectively:  $\mathbf{m}' = \mathbf{S}\mathbf{m}$  and  $\mathbf{d}' = \mathbf{Q}\mathbf{d}$ . Assuming again that  $\mathbf{S}^{-1}$  exists, the transformed system of equations is given by (31):  $\mathbf{QAS}^{-1}\mathbf{m}' = \mathbf{Qd}$ .

The damped least squares solution of the original system is given by:

$$\tilde{\mathbf{m}}^{(1)} = \left( \mathbf{A}^T \mathbf{A} + \gamma \mathbf{I} \right)^{-1} \mathbf{A}^T \mathbf{d} \quad (41)$$

Analogously to (34) the damped least-squares solution of the transformed equations is given by:

$$\tilde{\mathbf{m}}^{(2)} = \mathbf{S}^{-1} \left( \mathbf{S}^{T-1} \mathbf{A}^T \mathbf{Q}^T \mathbf{Q} \mathbf{A} \mathbf{S}^{-1} + \gamma \mathbf{I} \right)^{-1} \mathbf{S}^{T-1} \mathbf{A}^T \mathbf{Q}^T \mathbf{Q} \mathbf{d} \quad (42)$$

The damping parameter here is not necessarily equal to the damping parameter in the original damped least-squares solution, but for our purpose we do not need to make the distinction. Expression (42) can be simplified using the same steps as in the derivation of (34). Writing the term  $\gamma \mathbf{I}$  as  $\gamma \mathbf{I} = \gamma \mathbf{S}^{T-1} \mathbf{S}^T \mathbf{S} \mathbf{S}^{-1}$ , it follows that

$$\tilde{\mathbf{m}}^{(2)} = \left( \mathbf{A}^T \mathbf{Q}^T \mathbf{Q} \mathbf{A} + \gamma \mathbf{S}^T \mathbf{S} \right)^{-1} \mathbf{A}^T \mathbf{Q}^T \mathbf{Q} \mathbf{d} \quad (43)$$

This expression points to a fundamental problem: the damping term in the model space  $\mathbf{m}'$  is given by the identity matrix  $\gamma \mathbf{I}$  (see (42)) and the damping term is  $\gamma \mathbf{S}^T \mathbf{S}$  when expressed in terms of the original model vector  $\mathbf{m}$  (see (43)). This implies that the damping  $\gamma \mathbf{I}$  is not invariant for transformations of the model parameters. The terms  $\mathbf{Q}^T \mathbf{Q}$  appear when a transformation of the data vector is carried out. This implies that the damped least-squares solution is in general not invariant under transformations of the data vector or model vector.

There is therefore a need for a more general regularization which allows to change model and data space in a consistent manner so that the solution is coordinate independent. Such a general regularization can be found from (43) by setting  $\mathbf{Q}^T \mathbf{Q} = \mathbf{W}_d$  and by defining  $\mathbf{S}^T \mathbf{S} = \mathbf{W}_m$ . The general least-squares solution is then given by:

$$\tilde{\mathbf{m}} = \left( \mathbf{A}^T \mathbf{W}_d \mathbf{A} + \gamma \mathbf{W}_m \right)^{-1} \mathbf{A}^T \mathbf{W}_d \mathbf{d} . \quad (44)$$

This solution minimizes the following cost function:

$$S = (\mathbf{d} - \mathbf{A}\mathbf{m})^T \mathbf{W}_d (\mathbf{d} - \mathbf{A}\mathbf{m}) + \gamma \mathbf{m}^T \mathbf{W}_m \mathbf{m} \quad (45)$$

This expression shows that in general the weight matrices  $\mathbf{W}_d$  and  $\mathbf{W}_m$  can be anything (as long as they are positive definite to ensure that  $S$  has minima). Written in this way,  $\gamma$  may be seen as a trade-off parameter which compromises between two characteristics of the model: its size and its disagreement with the data. Both independent properties of the model cannot be arbitrary small together, hence there is a need for a balance. The choice of an optimum  $\gamma$ , however, is not an easy question. We have shown explicitly that when you start with a simple damped least-squares solution you can transform that problem into a more generally regularized least-squares solution in a different coordinate system and vice versa.

This implies that there is no reason to favor the damped least-squares solution over the more general least-squares solution (44). In fact, most inverse problems are ill-posed (partly underdetermined and partly overdetermined) and ill-conditioned (small errors in the data causes large variations in the model) which goes hand in hand with large null-spaces and hence non-unique solutions. Regularization is thus needed, but there is a large ambiguity in its choice [*Scales and Snieder, 1997*]. This reflects the fundamental difficulty that one faces in solving inverse problems: solving the system of equations is a minor problem, compared to choosing the regularization.

One approach is to use Bayesian statistics where one treats the inverse problem from a statistical point of view combining a-priori information about the data and the model with the data that are actually measured [*Tarantola and Valette, 1982a; Tarantola and Valette, 1982b*]. The weight matrices reflect true physical a-priori information (in a statistical sense) that one has of the data and the model, *independent* of the measured data. This includes for example the statistical noise characteristics of the instrument that produced the data, as well as information of the model and data that follow from other arguments. (For example, the mass-density in the Earth must be positive.) In such a Bayesian approach the weight matrices are given by

$$\mathbf{W}_d = \mathbf{C}_d^{-1} \quad , \quad \gamma \mathbf{W}_m = \mathbf{C}_m^{-1} \quad , \quad (46)$$

where  $\mathbf{C}_d^{-1}$  and  $\mathbf{C}_m^{-1}$  are the a-priori covariance matrices for the data and model respectively:

$$\mathbf{C}_d = \left\langle (\mathbf{d} - \langle \mathbf{d} \rangle) (\mathbf{d} - \langle \mathbf{d} \rangle)^T \right\rangle \quad , \quad (47)$$

$$\mathbf{C}_m = \left\langle (\mathbf{m} - \langle \mathbf{m} \rangle) (\mathbf{m} - \langle \mathbf{m} \rangle)^T \right\rangle \quad . \quad (48)$$

In these expressions the brackets  $\langle \dots \rangle$  denote the expectation value. In this interpretation the estimator (44) corresponds to the most likely a-posteriori model when the error distribution is Gaussian. The statistical basis of Bayesian inversion leads to an objective solution if one respects the rule that the a-priori information has a true physical meaning. In practice however, one should realize that the choice of the a-priori distribution of the data and model is very often subjective as well. The reader can find further details in the column “To Bayes or not to Bayes” of *Scales and Snieder* [1997].

A different approach is to define the misfit function in such a way that it favours models with given properties (small, smooth, ...) [Parker, 1994]. Choosing a-priori information then amounts to defining an appropriate norm in which the data misfit and any given property of the model are measured. In our case, the weight matrices would then define a particular metric for the  $L_2$ -norm. As an example of choosing the weight matrices  $\mathbf{W}_m$  the use of Occam’s inversion is quite common [Constable *et al.*, 1987] where one seeks the smoothest model that is consistent with the data. Instead of putting a constraint on the model length, one seeks the square of its gradient to be as small as possible, i.e. the last term in (45) is a discretization of  $\|\nabla m\|^2 = \int (\nabla m \cdot \nabla m) dV = - \int m \nabla^2 m dV$ ,<sup>5</sup> and hence  $\mathbf{W}_m$  corresponds to a discretized form of the Laplacian  $-\nabla^2$ .

## 2.8. THE TRANSFORMATION RULES FOR THE WEIGHT MATRICES

One of the fundamental requirements of an inverse solution should be that the results of the inversion are independent of arbitrary scalings applied to the model vector or data vector. Alas, this requirement is often ignored which can render comparisons of different models quite meaningless. For practical implications see *Trampert and L  v  que* [1990] and *Trampert et al.* [1992]. Here we derive how the weight matrices  $\mathbf{W}_m$  and  $\mathbf{W}_d$  should scale under such transformations for the least-squares solution to remain invariant.

Let us first consider the scaling (29) of the model vector:  $\mathbf{m}' = \mathbf{S}\mathbf{m}$ . Under this transformation the model term in the least-squares quantity (45) transforms as

$$\mathbf{m}^T \mathbf{W}_m \mathbf{m} = \mathbf{m}'^T \mathbf{S}^{T-1} \mathbf{W}_m \mathbf{S}^{-1} \mathbf{m}' = \mathbf{m}'^T \mathbf{W}_m' \mathbf{m}' , \quad (49)$$

with

<sup>5</sup>Note that we tacitly assumed in the last identity that there are no nonzero terms arising from the boundary conditions. A formal treatment based on Green’s theorem allows for the incorporation of nonzero boundary terms.

$$\mathbf{W}'_m = \mathbf{S}^{T-1} \mathbf{W}_m \mathbf{S}^{-1} . \quad (50)$$

Under this transformation rule for the model weight matrix the least-squares criterion is unchanged, hence the least-squares solution is not changed when the model weight matrix  $\mathbf{W}_m$  is transformed. It is of interest to note that this rule implies that for Bayesian inversions, where the weight matrix is the inverse of the a-priori model covariance matrix ( $\gamma \mathbf{W}_m = \mathbf{C}_m^{-1}$ ), the covariance matrix should transform as

$$\mathbf{C}'_m = \mathbf{S} \mathbf{C}_m \mathbf{S}^T \quad (51)$$

One easily verifies from definition (48) that this is indeed the transformation rule for covariance operators.

Next let us consider how the transformation (30) for the data vector  $\mathbf{d}' = \mathbf{Q}\mathbf{d}$  affects the transformation of the data weight matrix  $\mathbf{W}_d$ . The matrix  $\mathbf{A}$  scales under this transformation as  $\mathbf{A}' = \mathbf{Q}\mathbf{A}$ . Under this transformation the data term in the least-squares quantity (45) transforms as

$$\begin{aligned} & (\mathbf{d} - \mathbf{A}\mathbf{m})^T \mathbf{W}_d (\mathbf{d} - \mathbf{A}\mathbf{m}) \\ &= (\mathbf{d}' - \mathbf{A}'\mathbf{m})^T \mathbf{Q}^{T-1} \mathbf{W}_d \mathbf{Q}^{-1} (\mathbf{d}' - \mathbf{A}'\mathbf{m}) \\ &= (\mathbf{d}' - \mathbf{A}'\mathbf{m})^T \mathbf{W}'_d (\mathbf{d}' - \mathbf{A}'\mathbf{m}) \end{aligned} \quad (52)$$

with

$$\mathbf{W}'_d = \mathbf{Q}^{T-1} \mathbf{W}_d \mathbf{Q}^{-1} . \quad (53)$$

For a Bayesian inversion the data weight matrix is the inverse of the data covariance matrix ( $\mathbf{W}_d = \mathbf{C}_d^{-1}$ ), so that for a Bayesian inversion  $\mathbf{C}_d$  should transform as

$$\mathbf{C}'_d = \mathbf{Q} \mathbf{C}_d \mathbf{Q}^T . \quad (54)$$

Note again that this is the correct transformation rule for a covariance matrix defined in (47). This implies that the Bayesian viewpoint, where  $\mathbf{W}_m$  and  $\mathbf{W}_d$  are the inverses of the model and data covariance matrices, ensures that the solution is invariant under transformations of the model vector and/or the data vector.

Although we have derived in which way the weight matrices  $\mathbf{W}_m$  and  $\mathbf{W}_d$  should transform under transformations of model and data vectors, this does by no means imply that these matrices can be defined in a unambiguous way. An ill-posed and/or ill-conditioned inverse problem can only be solved if one is willing to control the solution by imposing a regularization term. In general, there is no unique recipe for choosing the

weight matrices  $\mathbf{W}_m$  and  $\mathbf{W}_d$ . It is the subjective input of the user that determines the choice of these matrices.

## 2.9. SOLVING THE SYSTEM OF LINEAR EQUATIONS

It should be noted that the least-squares solution always requires solving a set of linear algebraic equations. For instance, equation (44) may be written as

$$\left(\mathbf{A}^T \mathbf{W}_d \mathbf{A} + \gamma \mathbf{W}_m\right) \tilde{\mathbf{m}} = \mathbf{A}^T \mathbf{W}_d \mathbf{d}. \quad (55)$$

This represents a square system of linear equations, the so-called normal equations, of the form  $\mathbf{B}\mathbf{x} = \mathbf{y}$ . If we are merely interested in the estimation part of the problem,  $\mathbf{B}$  doesn't need to be inverted. If we are also interested in the appraisal part of the problem (and we always should), it must be realized that  $\mathbf{B}$  needs to be inverted at the cost of additional computer time. Many standard subroutine packages are available and *Press et al.* [1989] give a good and practical introduction to the subject. The reader should realize, however, that the system  $\mathbf{B}\mathbf{x} = \mathbf{y}$  may become quite large for realistic geophysical problems which makes it worthwhile to consider a specialized routine best suited to the nature of  $\mathbf{B}$  (symmetric, banded, sparse, ...). The dimension of the set of normal equations is also worth considering. Remember that the matrix  $\mathbf{A}$  has the dimension  $(N \times M)$ , where  $N$  is the number of data and  $M$  the number of model parameters. System (55) has the dimension of the model space, but using (22) we may obtain a strictly equivalent system of the dimension of the data space. Choosing the smallest dimension to write the normal equation can save quite some computer time. Most techniques for solving the set of algebraic equations directly work with the matrix  $\mathbf{B}$  as a whole requiring sufficient memory space to hold the matrix. But in global travel time tomography, for instance, these dimensions may become extremely large ( $N > 10^6$  and  $M > 10^5$ ) so that iterative methods need to be employed which only work on parts of  $\mathbf{B}$  at a time. Another problem which frequently occurs is that even though regularization is included in  $\mathbf{B}$ , it is singular or numerically very close to singular. A powerful technique, called Singular Value Decomposition (SVD), can diagnose precisely what the problem is and will give a useful numerical answer. SVD is the most effective tool in inverse theory to understand why a certain result has been obtained.

Iterative methods or SVD need not to work on square systems and may thus directly use the matrix  $\mathbf{A}$ . In this context it is useful to realize that the generalized least squares solution (44) is equivalent to the simple least squares solution of the system



$$\begin{pmatrix} \mathbf{W}_d^{1/2} \mathbf{A} \\ \dots \\ \sqrt{\gamma} \mathbf{W}_m^{1/2} \end{pmatrix} \mathbf{m} = \begin{pmatrix} \mathbf{W}_d^{1/2} \mathbf{d} \\ \dots \\ 0 \end{pmatrix}. \quad (56)$$

For a discussion on the meaning of a square root of a positive definite matrix the reader is referred to *Tarantola* [1987]. Keeping also in mind a certain freedom in choosing weighting matrices (see 2.7), the user might want to define directly  $\mathbf{W}^{1/2}$  rather than  $\mathbf{W}$ . Expression (56) shows that regularization has the effect of adding extra rows to the system of linear equations, but the enlarged system is still of the form  $\mathbf{B}\mathbf{x} = \mathbf{y}$ , where the matrix  $\mathbf{A}$  and the data vector  $\mathbf{d}$  are augmented by extra rows that account for the regularization.  $\mathbf{B}$  is not square anymore as in the case of normal equations. We will now illustrate in more detail the essence of singular value decomposition and iterative techniques as applied to the system  $\mathbf{B}\mathbf{x} = \mathbf{y}$ .

### 2.9.1. Singular value decomposition

Singular value decomposition was developed by *Lanczos* [1961], this technique is a generalization of the eigenvector decomposition of matrices to non-square matrices. Let us first consider a real symmetric  $N \times N$ -square matrix  $\mathbf{B}$  with eigenvectors  $\hat{\mathbf{v}}^{(n)}$  and eigenvalues  $\lambda_n$ . For such a matrix the eigenvectors form an orthonormal set, hence any vector  $\mathbf{x}$  can be projected on these eigenvectors:  $\mathbf{x} = \sum_{n=1}^N \hat{\mathbf{v}}^{(n)} (\hat{\mathbf{v}}^{(n)} \cdot \mathbf{x})$ . When  $\mathbf{B}$  acts on this expression the result can be written as:

$$\mathbf{B}\mathbf{x} = \sum_{n=1}^N \lambda_n \hat{\mathbf{v}}^{(n)} (\hat{\mathbf{v}}^{(n)} \cdot \mathbf{x}) = \mathbf{y}. \quad (57)$$

Decomposing the vector  $\mathbf{y}$  using the same eigenvectors  $\hat{\mathbf{v}}^{(n)}$  gives  $\mathbf{y} = \sum_{n=1}^N \hat{\mathbf{v}}^{(n)} (\hat{\mathbf{v}}^{(n)} \cdot \mathbf{y})$ , and inserting this into expression (57) yields the following expansion for the solution vector  $\mathbf{x}$ :

$$\mathbf{x} = \sum_{n=1}^N \frac{1}{\lambda_n} \hat{\mathbf{v}}^{(n)} (\hat{\mathbf{v}}^{(n)} \cdot \mathbf{y}). \quad (58)$$

It can be seen that small eigenvectors can lead to instabilities in the solution  $\mathbf{x}$ . Singular value decomposition generalizes this expansion to non-square matrices. Details of this technique are given by *Lanczos* [1961] and by *Aki and Richards* [1980].

Now consider the following non-square system of equations:

$$\underbrace{\mathbf{B}}_{\substack{N \times M \\ \text{matrix}}} \underbrace{\mathbf{x}}_{\substack{M \\ \text{rows}}} = \underbrace{\mathbf{y}}_{\substack{N \\ \text{rows}}} \quad (59)$$

Singular value decomposition is based on an expansion of  $\mathbf{x}$  in an orthonormal set of eigenvectors  $\hat{\mathbf{v}}^{(n)}$  and of  $\mathbf{y}$  in an orthonormal set  $\hat{\mathbf{u}}^{(n)}$ . These vectors cannot be the eigenvectors of  $\mathbf{B}$  because this matrix is not square, hence it does not have any eigenvectors. Instead, these vectors are related by the following relation:

$$\mathbf{B}\hat{\mathbf{v}}^{(n)} = \lambda_n \hat{\mathbf{u}}^{(n)} \quad , \quad \mathbf{B}^T \hat{\mathbf{u}}^{(n)} = \lambda_n \hat{\mathbf{v}}^{(n)} \quad (60)$$

It can easily be seen that the vectors  $\hat{\mathbf{v}}^{(n)}$  are the eigenvectors of  $\mathbf{B}^T \mathbf{B}$  while the vectors  $\hat{\mathbf{u}}^{(n)}$  are the eigenvectors of  $\mathbf{B} \mathbf{B}^T$ , hence these vectors can readily be determined.  $\mathbf{B}^T \mathbf{B}$  and  $\mathbf{B} \mathbf{B}^T$  share the same nonzero eigenvectors  $\lambda_n^2$ . The  $\lambda_n$  are called the singular values of  $\mathbf{B}$ . When  $\mathbf{B}$  acts on the vector  $\mathbf{x}$  the result can be written as

$$\mathbf{B}\mathbf{x} = \sum_{n=1}^P \lambda_n \hat{\mathbf{u}}^{(n)} \left( \hat{\mathbf{v}}^{(n)} \cdot \mathbf{x} \right) . \quad (61)$$

The upper limit of the summation is determined by the number of eigenvalues that are nonzero because the vanishing eigenvalues do not contribute to the sum. This number  $P$  can be significantly less than the dimension of the problem:  $P \leq N$  and  $P \leq M$ .

It is convenient to arrange the vectors  $\hat{\mathbf{u}}^{(n)}$  and  $\hat{\mathbf{v}}^{(n)}$  as the columns of matrices  $\mathbf{U}$  and  $\mathbf{V}$ , the eigenvectors from index  $P$  onwards correspond to zero eigenvalues and need to be included to make  $\mathbf{U}$  and  $\mathbf{V}$  complete:

$$\mathbf{U} = \left( \begin{array}{cccccc} \vdots & \vdots & & \vdots & \vdots & \vdots \\ \hat{\mathbf{u}}^{(1)} & \hat{\mathbf{u}}^{(2)} & \dots & \hat{\mathbf{u}}^{(P)} & \hat{\mathbf{u}}^{(P+1)} & \dots & \hat{\mathbf{u}}^{(N)} \\ \vdots & \vdots & & \vdots & \vdots & & \vdots \end{array} \right) , \quad (62)$$

$\underbrace{\hspace{10em}}_{\mathbf{U}_p} \qquad \underbrace{\hspace{10em}}_{\mathbf{U}_0}$

$$\mathbf{V} = \left( \begin{array}{cccccc} \vdots & \vdots & & \vdots & \vdots & \vdots \\ \hat{\mathbf{v}}^{(1)} & \hat{\mathbf{v}}^{(2)} & \dots & \hat{\mathbf{v}}^{(P)} & \hat{\mathbf{v}}^{(P+1)} & \dots & \hat{\mathbf{v}}^{(M)} \\ \vdots & \vdots & & \vdots & \vdots & & \vdots \\ \underbrace{\hspace{10em}}_{\mathbf{V}_p} & \underbrace{\hspace{10em}}_{\mathbf{V}_0} \end{array} \right), \quad (63)$$

The *orthogonality* of the eigenvectors implies that  $\mathbf{U}^T \mathbf{U} = \mathbf{I}$  and  $\mathbf{V}^T \mathbf{V} = \mathbf{I}$ . The *completeness* of the eigenvectors implies that  $\mathbf{U} \mathbf{U}^T = \mathbf{I}$  and  $\mathbf{V} \mathbf{V}^T = \mathbf{I}$ . Since the orthogonality of the eigenvectors also holds in the subspaces spanned by  $\mathbf{U}_p$  and  $\mathbf{V}_p$  we have  $\mathbf{U}_p^T \mathbf{U}_p = \mathbf{I}$  and  $\mathbf{V}_p^T \mathbf{V}_p = \mathbf{I}$ . However, the vectors in these subspaces do in general not form a complete set, so that in general  $\mathbf{U} \mathbf{U}_p^T \neq \mathbf{I}$  and  $\mathbf{V} \mathbf{V}_p^T \neq \mathbf{I}$ .

The generalization of (61) to non-square systems can be written as

$$\mathbf{B} = \left( \begin{array}{cc} \mathbf{U}_p & \mathbf{U}_0 \end{array} \right) \left( \begin{array}{cc} \boldsymbol{\Sigma} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{array} \right) \left( \begin{array}{c} \mathbf{V}_p^T \\ \mathbf{V}_0^T \end{array} \right), \quad (64)$$

where the matrix  $\boldsymbol{\Sigma}$  is given by:

$$\boldsymbol{\Sigma} = \left( \begin{array}{cccc} \lambda_1 & 0 & \dots & 0 \\ 0 & \lambda_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & & \lambda_p \end{array} \right) \quad (65)$$

It follows from (61) that the eigenvectors  $\hat{\mathbf{u}}^{(n)}$  that correspond to a vanishing eigenvalue do not contribute when  $\mathbf{B}$  acts on a vector. These eigenvectors are ordered in the sub-matrix  $\mathbf{U}_0$ . This is equivalent to the statement that according to the representation (64) the matrix  $\mathbf{B}$  can be constructed from  $\mathbf{U}_p$  and  $\mathbf{V}_p$  alone.  $\mathbf{U}_0$  and  $\mathbf{V}_0$  are dark spots of the space not illuminated by operator  $\mathbf{B}$ . Since  $\mathbf{U}_0^T \mathbf{B} \mathbf{x} = \mathbf{0}$  the predicted data  $\mathbf{B} \mathbf{x}$  are orthogonal to the subspace spanned by  $\mathbf{U}_0$ , see figure 5. This means that any components in the data vector that lie in  $\mathbf{U}_0$  cannot be explained by *any* model. These components of the data vector necessarily correspond to errors in the data or errors in the operator  $\mathbf{B}$  as a description of the physical problem. It is for this reason that  $\mathbf{U}_0$  is called the *data-null-space*. In a least squares inversion one aims at minimizing the data misfit. Minimizing the data misfit then amounts to finding a model that produces a data vector in the subspace  $\mathbf{U}_p$  that is closest to the true data. It follows from figure 5 that this is achieved by simply projecting the components of  $\mathbf{U}_0$  contained in the data out of the problem. This is

exactly what is done by limiting the sum over eigenvalues in (64) to the nonzero eigenvalues only. Of course, when  $\mathbf{U}_0$  is empty, one can always find  $\mathbf{x}$  which explains the data  $\mathbf{y}$  exactly because  $\mathbf{U}_p$  spans the complete data space.

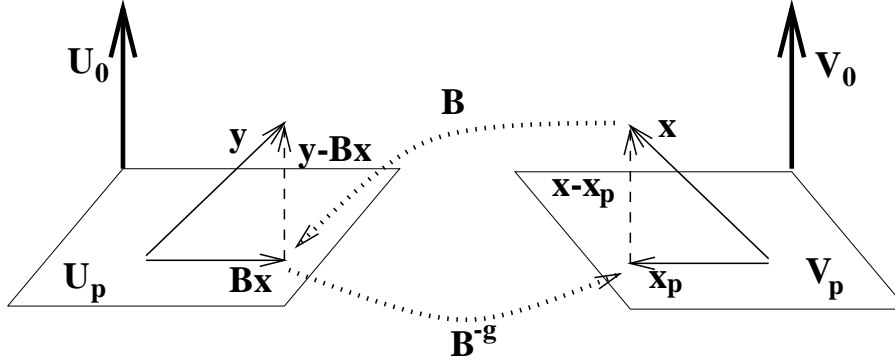


Figure 5. Geometrical interpretation of singular value decomposition. Note that starting from a model vector  $\mathbf{x}$  and inverting the corresponding data  $\mathbf{B}\mathbf{x}$  one obtains an estimated model vector  $\mathbf{x}_p$  that in general differs from  $\mathbf{x}$  because the resolution operator  $\mathbf{R} = \mathbf{B}^{-g}\mathbf{B}$  in general is not equal to the identity matrix.

In a similar way, the restriction of the summation over eigenvalues to the nonzero eigenvalues has the effect that the model estimate lies in the subspace spanned by  $\mathbf{V}_p$ , but that the estimated model has no component in  $\mathbf{V}_0$ . Any component of the model in  $\mathbf{V}_0$  does not affect the data because  $\mathbf{B}\mathbf{V}_0 = \mathbf{0}$ . This means that  $\mathbf{V}_0$  defines the *model-null-space*. The data have no bearing on the components of the model vector that lie in  $\mathbf{V}_0$ . Setting the component of the model vector in the model-null-space equal to zero then implies that in the model estimation one only takes the nonzero eigenvalues into account. Expanding  $\mathbf{x}$  in the vectors  $\hat{\mathbf{v}}^{(n)}$  and  $\mathbf{y}$  in the vectors  $\hat{\mathbf{u}}^{(n)}$  and taking only the nonzero eigenvalues into account one can thus generalize the solution (58) in the following way to non-square systems:

$$\mathbf{x} = \sum_{n=1}^P \frac{1}{\lambda_n} \hat{\mathbf{v}}^{(n)} \left( \hat{\mathbf{u}}^{(n)} \cdot \mathbf{y} \right) . \quad (66)$$

Using the matrices  $\mathbf{U}_p$  and  $\mathbf{V}_p$  this result can also be written as:

$$\mathbf{x} = \mathbf{V}_p \mathbf{\Sigma}^{-1} \mathbf{U}_p^T \mathbf{y} , \quad (67)$$

with  $\mathbf{\Sigma}^{-1}$  given by

$$\Sigma^{-1} = \begin{pmatrix} 1/\lambda_1 & 0 & \cdots & 0 \\ 0 & 1/\lambda_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 1/\lambda_p \end{pmatrix} \quad (68)$$

Similar to the forward problem, the inverse problem is not a function of  $\mathbf{U}_0$  and  $\mathbf{V}_0$ . If both of these subspaces are zero, the operator  $\mathbf{B}$  has an exact inverse. If  $\mathbf{U}_0$  exists, one can show that the residual  $\mathbf{y} - \mathbf{B}\mathbf{x}$  is perpendicular to  $\mathbf{B}\mathbf{x}$  and hence the residual is minimum as for the least-squares solution. If  $\mathbf{V}_0$  exists, solution (67) has no component in  $\mathbf{V}_0$  and is therefore of minimum norm.

Clearly, small errors in  $\mathbf{y}$  can lead to large errors in  $\mathbf{x}$  when multiplied with  $1/\lambda_n$  and the singular value is small. This process of error magnification can be controlled by limiting the summation in (66) to eigenvalues that differ significantly from zero. Alternatively, one can replace  $1/\lambda_n$  by  $\lambda_n / (\lambda_n^2 + \gamma)$  with  $\gamma$  a positive constant. One can show that this is equivalent to the damped least-squares solution (20). See for example *Matsu'ura and Hirata* [1982] for a discussion on these different strategies. It should be noted that cutting off or damping small eigenvalues leads to different results. This makes it virtually impossible to quantitatively compare solutions of the same problem obtained with these fundamentally different strategies. One of the main reasons for the popularity of singular value decomposition is the control that one can exert over the error propagation in the solution. The drawback is that one needs to determine the eigenvectors of a matrix. For realistic large-scale problems ( $P > 10^4$ ) this may require a prohibitive amount of CPU time. On the other hand, once the eigenvectors are calculated, resolution and error propagation are obtained at virtually no cost, since they merely involve the matrix multiplication  $\mathbf{V}_p \mathbf{V}_p^T$ .

### 2.9.2. Iterative least-squares

The least-squares solution of the system  $\mathbf{B}\mathbf{x} = \mathbf{y}$  is, as shown in (11), given by  $\mathbf{x} = (\mathbf{B}^T \mathbf{B})^{-1} \mathbf{B}^T \mathbf{y}$ . Note that the matrix  $\mathbf{B}$  may contain some form of regularization as seen in expression (56). In practice, given the large-scale of many inverse problems, the matrix  $\mathbf{B}^T \mathbf{B}$  may not fit in the computer memory. It is for this reason that one has developed iterative techniques which improve existing estimates of the solution.

Suppose one has in the  $n$ -th iteration of an iterative process a model estimate  $\mathbf{x}_n$  and that one seeks an update  $\delta\mathbf{x}_n$  such that the new model estimate  $\mathbf{x}_{n+1} = \mathbf{x}_n + \delta\mathbf{x}_n$  is a better estimate of the model. Inserting

this expression into the relation  $\mathbf{B}\mathbf{x} = \mathbf{y}$  gives an expression for the model update:

$$\mathbf{B}\delta\mathbf{x}_n = (\mathbf{y} - \mathbf{B}\mathbf{x}_n) \quad (69)$$

Note that the right hand side of this expression is the residual of the model estimate  $\mathbf{x}_n$ , i.e. the difference  $\mathbf{y} - \mathbf{B}\mathbf{x}_n$ , is a measure to what extend the model estimate  $\mathbf{x}_n$  does not explain the data. Expression (69) prescribes how the model should be updated in order to reduce the data residual. The least squares-solution of this expression is given by

$$\delta\mathbf{x}_n = \left(\mathbf{B}^T\mathbf{B}\right)^{-1} \mathbf{B}^T (\mathbf{y} - \mathbf{B}\mathbf{x}_n) \quad (70)$$

However, we have not gained anything yet because this expression is just as difficult to solve as the original equation and we still need to deal with the inverse of  $\mathbf{B}^T\mathbf{B}$ .

The advantage of solving the problem iteratively is that in such an approach one can replace the inverse  $\left(\mathbf{B}^T\mathbf{B}\right)^{-1}$  by a suitably chosen estimate  $\mathbf{P}$  of this inverse, i.e. one computes a model update using the following expression:

$$\mathbf{x}_{n+1} = \mathbf{x}_n + \mathbf{P}\mathbf{B}^T (\mathbf{y} - \mathbf{B}\mathbf{x}_n) \quad (71)$$

The operator  $\mathbf{P}$  is called a preconditioning operator. If one sets  $\mathbf{P} = \left(\mathbf{B}^T\mathbf{B}\right)^{-1}$  one retrieves the full solution in one step, but in that case one needs to compute  $\left(\mathbf{B}^T\mathbf{B}\right)^{-1}$ , which is what we wanted to avoid. Recognizing that  $\mathbf{B}^T (\mathbf{y} - \mathbf{B}\mathbf{x}_n)$  is the direction of descent of the cost function at  $\mathbf{x}_n$ , one may, on the other end of the spectrum, choose  $\mathbf{P} = c\mathbf{I}$ , where  $c$  is a constant derived from a least-squares criterion to ensure the steepest possible descent [e.g. *Tarantola*, 1984]. In practice one has to find a balance between using an advanced preconditioning operator (that may be difficult to compute but that leads to a solution with a few iterations), and a simple preconditioning operator (that is easy to compute but that may require many iterations). The most commonly used algorithms in geophysics are SIRT (Simultaneous Iterative Reconstruction Technique) and LSQR (Least-Squares Conjugate Gradient). SIRT has the drawback of introducing implicit regularization into the solution [*van der Sluis and van der Vorst*, 1987] and if corrected for significantly decreases the convergence speed [*Trampert and L  v  que*, 1990]. A more successful balance is achieved by the LSQR algorithm [*Paige and Saunders*, 1982a, 1982b; *van der Sluis and van der Vorst*, 1987]. An iterative scheme that mimics the properties of SVD is given by *Nolet and Snieder* [1990].

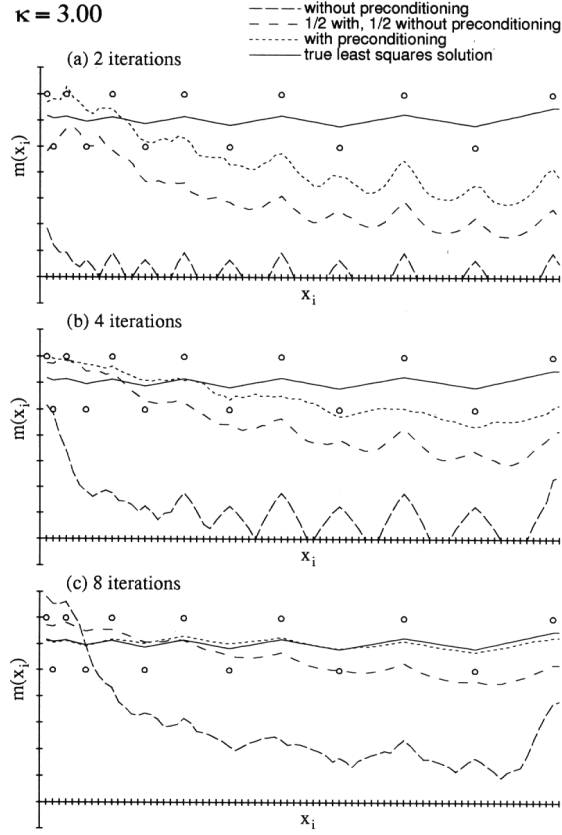


Figure 6. One-dimensional example of the convergence of a conjugate-gradient inversion at (a) iteration 2, (b) iteration 4 and (c) iteration 8. The solid line represents the true least-squares solution. The long-dashed lines are the interactive solutions without preconditioning while the dotted solutions are obtained using preconditioning.

Note that in the iterative least-squares algorithm (71) there is no need to invert a matrix, one only needs to multiply with the matrices  $\mathbf{P}$  and  $\mathbf{B}^T$ , which can be done row by row. In many practical problems such as seismic tomography, the matrix  $\mathbf{B}$  is very sparse which means that most of the matrix elements are zero, see section 6.1. For such a matrix the iterative least-squares algorithm is particularly efficient because the matrix multiplications only works on the nonzero matrix elements. The implementation of this idea in seismic tomography was developed by Clayton and Comer, [1983] and by Nolet [1985].

Despite the efficiency of the iterative least-squares problems one should

be aware that the convergence of this process may be very slow. An example from *VanDecar and Snieder* [1994] is shown in figure 6. In this example one fits a function whose values are given by the dots in figure 6 by a model defined as the samples of that function at the  $x$ -values indicated by the tick marks along the horizontal axis. This problem obviously is ill-posed because the function is not defined between the data points. In order to define a solution, Occam's method is used where  $\mathbf{W}_m$  in (56) is a discretized version of the gradient operator. The true regularized least-squares solution of this problem is shown by the solid line. The iterative least-squares solution based on the conjugate-gradient algorithm after 2, 4 and 8 iterations is shown by the long-dashed line in the panels of figure 6.

It is clear that the convergence of the iterative least-squares algorithm is slow. The reason for this is easy to see. At every iteration, every point in the model interacts only with its neighbouring points through the damping term  $\mathbf{W}_m$  in the matrix. This means that a model element can only interact with a model element  $n$  positions further down the line after  $n$  iterations. In other words, the damping term imposes a global smoothness constraint on the solution, but in each iteration the model elements interact only locally. For this reason an inordinate number of iterations are required. *VanDecar and Snieder* [1994] developed a preconditioning operator that leads to a very rapid convergence, this is indicated by the dotted lines in figure 6. Note that this iterative solution has superior convergence properties.

A general drawback of all iterative techniques is that most of its advantages are lost if one is interested in the appraisal part of the problem since this involves the knowledge of  $(\mathbf{B}^T \mathbf{B})^{-1}$ , which an iterative technique doesn't explicitly compute. A way around this is to solve the system  $M$  times, each time replacing the data by a column of  $\mathbf{B}$  for the resolution matrix for instance [*Trampert and L  v  que*, 1990].

### 3. Linear inverse problems with continuous models

Up to this point, the model vector was finite dimensional. In many inverse problems the model is a continuous function of the space coordinates; it therefore has infinitely many degrees of freedom. For example, in inverse problems using gravity data one wants to determine the mass density  $\rho(\mathbf{r})$  in the earth. This is in general a continuous function of the space coordinates. In realistic inverse problems one has a finite amount of data. A simple variable count shows that it is impossible to determine a continuous model with infinitely many degrees of freedom from a finite amount of data in a unique way.



In order to make the problem manageable we will restrict ourselves in this section to linear inverse problems. For such problems the data and the model are related by

$$d_i = \int G_i(x)m(x)dx + e_i . \quad (72)$$

The notation in this section is one-dimensional, but the theory can be used without modification in higher dimensions as well. The model is a continuous model but the data vector is discrete and in practice has a finite dimension. The kernel  $G_i(x)$  plays the same role as the matrix  $\mathbf{A}$  in (1). Note that the data are contaminated with errors  $e_i$ .

Since the forward problem is linear, the estimated model is obtained by making a linear combination of the data (this is the most general description of a linear estimator):

$$\tilde{m}(x) = \sum_i a_i(x)d_i . \quad (73)$$

The coefficients  $a_i(x)$  completely specify the linear inverse operator. By inserting (72) in (73) one arrives again at a relation between the true model  $m(x)$  and the estimated model  $\tilde{m}(x)$ :

$$\tilde{m}(x) = \underbrace{\int R(x, x')m(x')dx'}_{\text{Finite resolution}} + \underbrace{\sum_i a_i(x)e_i}_{\text{Error propagation}} , \quad (74)$$

with the resolution kernel  $R$  given by

$$R(x, x') = \sum_i a_i(x)G_i(x') . \quad (75)$$

The first term in (74) accounts for averaging that takes places in the mapping from the true model to the estimated model. It specifies through the resolution kernel  $R(x, x')$  what spatial resolution can be attained. In the ideal case the resolution or averaging kernel is a delta function:  $R(x, x') = \delta(x - x')$ . The resolution kernel, however, is a superposition of a finite amount of data kernels  $G_i(x')$ . These data kernels are in general continuous functions, and since a delta function cannot be constructed from the superposition of a finite number of continuous functions, the resolution kernel will differ from a delta function. In Backus-Gilbert theory [Backus and Gilbert, 1967; 1968] one seeks the coefficients  $a_i(x)$  in such a way that the resolution kernel resembles a delta function as well as possible given a certain criterion which measures this resemblance.

The second term in (74) accounts for error propagation. A proper treatment of this term needs to be based on statistics. It is shown by *Backus and Gilbert* [1970] that one cannot simultaneously optimize the resolution and suppress the error propagation and that one has to seek a trade-off between finite resolution and error propagation. The above mentioned work tries to explain the data exactly, even if they contain errors. *Gilbert* [1971] extended the theory to explain data within their error bars only.

As mentioned above, the Backus-Gilbert strategy finds the coefficients  $a_i(x)$  in equation (73) by imposing a condition on the averaging kernel. *Tarantola and Valette* [1982] solve the problem in a Bayesian framework. They introduce prior information on the data (Gaussian error distribution) and prior assumptions (also Gaussian) on the unknown function  $m(x)$  which also yields the coefficients  $a_i(x)$ . In this approach, the resolution or averaging kernel is a consequence of the *a priori* information, but generally different from a delta function. They further show that the Backus-Gilbert approach is contained in Tarantola-Valette solution.

In summary, in both strategies the infinite dimensional problem is transformed into a finite dimensional problem by seeking local averages of the true model.

### 3.1. CONTINUOUS MODELS AND BASIS FUNCTIONS

Another approach is to evoke basis functions which amounts to changing the parameterization of the model. In general, any continuous model can be written as a superposition of a complete set of basis functions:

$$m(x) = \sum_{j=1}^{\infty} m_j B_j(x) . \quad (76)$$

In many global geophysical applications spherical harmonics are used to represent the seismic velocity or the density inside the earth, since they form a natural basis to describe a function on the sphere. In that case the  $B_j(x)$  are the spherical harmonics and the sum over  $j$  stands for a sum over degree  $l$  and angular order  $m$ . The advantage of such an expansion is that one now deals with a discrete vector  $m_j$  of expansion coefficients rather than with a continuous function. However, the basis functions  $B_j(x)$  only form a complete set when the sums is over infinitely many functions, and hence the problem has shifted from dealing with a continuous function to dealing with a vector with infinitely many components.

Inserting the expansion (76) into the forward problem (72) we may write

$$d_i = \sum_{j=1}^{\infty} A_{ij} m_j + e_i , \quad (77)$$

where the matrix elements  $A_{ij}$  are the projection of the data kernels onto the basis functions:

$$A_{ij} = \int G_i(x) B_j(x) dx \quad (78)$$

In practice, one cannot deal with infinitely many coefficients, and it is customary to ignore the fact that a model vector has infinitely many dimensions by taking only the first  $L$  basis functions into account. The resulting  $L$ -dimensional model vector will be denoted by  $\mathbf{m}_L$ , and a similar notation  $\mathbf{A}_L$  is used for the first  $L$  rows of the matrix  $\mathbf{A}$ . The solution of the resulting finite dimensional model vector can be found using any technique shown in section 2. Recall that a general least-squares solution may be found by minimizing

$$S_L = (\mathbf{d} - \mathbf{A}_L \mathbf{m}_L)^T \mathbf{C}_d^{-1} (\mathbf{d} - \mathbf{A}_L \mathbf{m}_L) + \mathbf{m}_L^T \mathbf{C}_{mL}^{-1} \mathbf{m}_L \quad (79)$$

as a function of the truncated model vector  $\mathbf{m}_L$ . The weighting operators  $\mathbf{C}_d$  and  $\mathbf{C}_m$  may or may not be given the interpretation of covariance operators (see section 2.7). The resulting model estimate is then given by

$$\tilde{\mathbf{m}}_L = \left( \mathbf{A}_L^T \mathbf{C}_d^{-1} \mathbf{A}_L + \mathbf{C}_{mL}^{-1} \right)^{-1} \mathbf{A}_L^T \mathbf{C}_d^{-1} \mathbf{d} \quad (80)$$

### 3.2. SPECTRAL LEAKAGE, THE PROBLEM

Truncating the model vector after the first  $L$  elements may appear to be a convenient way to reduce the problem to a finite dimensional one. However, there are problems associated with this truncation. Analogously to (2) any linear estimator of the first  $L$  coefficients of the infinite vector  $\mathbf{m}$  can be written as:

$$\tilde{\mathbf{m}}_L = \mathbf{A}_L^{-g} \mathbf{d} . \quad (81)$$

Let us now divide the sum over model elements in (77) as a sum over the first  $L$  elements that we are interested in and the remaining model elements:

$$\mathbf{d} = \mathbf{A}_L \mathbf{m}_L + \mathbf{A}_{\infty} \mathbf{m}_{\infty} + \mathbf{e} . \quad (82)$$

In this expression  $\mathbf{m}_\infty$  denotes the infinitely dimensional vector with elements  $(m_{L+1}, m_{L+2}, \dots)$ . Inserting this expression in (81) yields the relation between the estimated  $L$  model coefficients  $\tilde{\mathbf{m}}_L$  and the true model  $\mathbf{m}$  (for a discussion on the objectiveness of the concept of true model, see section 2.5):

$$\tilde{\mathbf{m}}_L = \mathbf{m}_L + \underbrace{(\mathbf{A}_L^{-g} \mathbf{A}_L - \mathbf{I}) \mathbf{m}_L}_{\text{Limited Resolution}} + \underbrace{\mathbf{A}_L^{-g} \mathbf{A}_\infty \mathbf{m}_\infty}_{\text{Spectral leakage}} + \underbrace{\mathbf{A}_L^{-g} \mathbf{e}}_{\text{Error propagation}} \quad (83)$$

The last three terms in this expression account for deviations in the estimated model from the true model. The second term and the last term are identical to the corresponding terms in expression (5) for finite-dimensional problems, accounting for finite resolution within the components of the vector  $\mathbf{m}_L$  and error propagation respectively. The term  $\mathbf{A}_L^{-g} \mathbf{A}_\infty \mathbf{m}_\infty$  has no counterpart in (5) and is responsible for a spurious mapping of the model coefficients  $(m_{L+1}, m_{L+2}, \dots)$  onto the estimated model coefficients  $(\tilde{m}_1, \dots, \tilde{m}_L)$ . Since the basis function expansion in many problems is that of a spectral expansion, the mapping from the coefficients  $\mathbf{m}_\infty$  onto the first  $L$  coefficients  $\tilde{\mathbf{m}}_L$  will be referred to as *spectral leakage* [Snieder *et al.*, 1991].

An example of spectral leakage is shown in the figures 7-10. Suppose a function on a line, defined in the interval  $0 \leq x < 1$ , is expanded in normalized sines and cosines which have at most three wavelengths in the interval. This means that in this example  $L = 7$ . The function is sampled at given locations and the sampling is twice as dense on the subintervals  $0 \leq x < 0.25$  and  $0.5 \leq x < 0.75$  than on the remaining part of the line. The inverse problem consists in determining the expansion coefficients in the finite set of basis functions on the line given the sampled data points. Details of this example may be found in Snieder [1993].

In figure 7 the sampled function is a sine wave with exactly three wavelengths in the given interval. The sampling points are indicated by circles. The reconstructed function is given by the solid line and is indistinguishable from the true function. In figure 8 the input function is a sine wave with six wavelengths on the interval, this input function is indicated by the dashed line. Because of the orthogonality of the trigonometric functions, one would expect this function to have no projection onto the seven basis functions that are used in the inversion. Nevertheless, the reconstructed function shown by the solid line in figure 8 differs significantly from zero; it has about 50% of the magnitude of the input functions. As shown in Snieder [1993] the expansion coefficients have errors of about 100%! This difference between the estimated model and zero is entirely due to spectral leakage because the first  $L$  model components

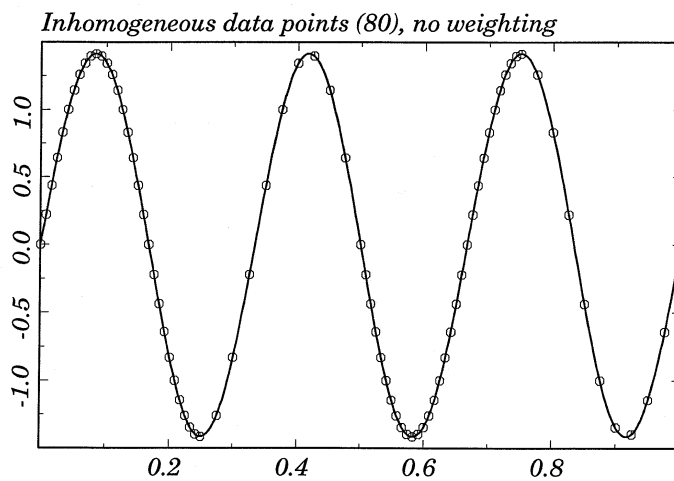


Figure 7. Unweighted least-squares fit of 80 data points (circles) sampling a sine wave with three wavelength along the interval. The estimated model is shown by the thick solid line, and is undistinguishable from the true projection on the basis functions.

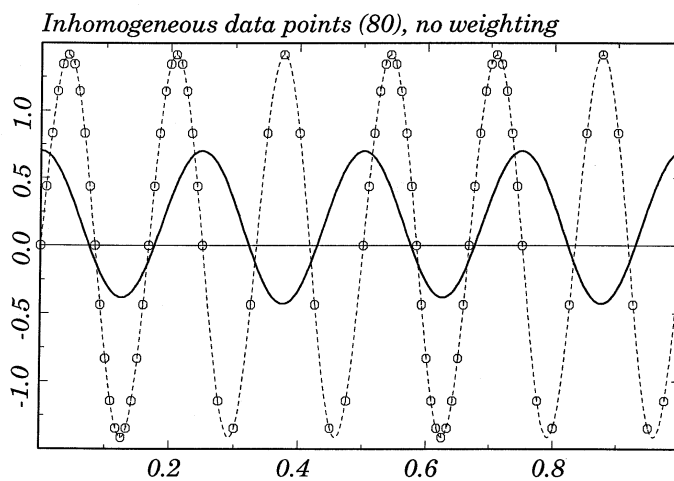


Figure 8. Unweighted least-squares fit of 80 data points (circles) sampling a sine wave with six periods (dashed line). The estimated model is shown by the thick solid line, the true projection on the basis functions is shown by the thin solid line.

$m_L$  are equal to zero and there are no data errors in this example so that

only the term  $\mathbf{A}_L^{-g} \mathbf{A}_\infty \mathbf{m}_\infty$  in (83) can give rise to errors in the model reconstruction.

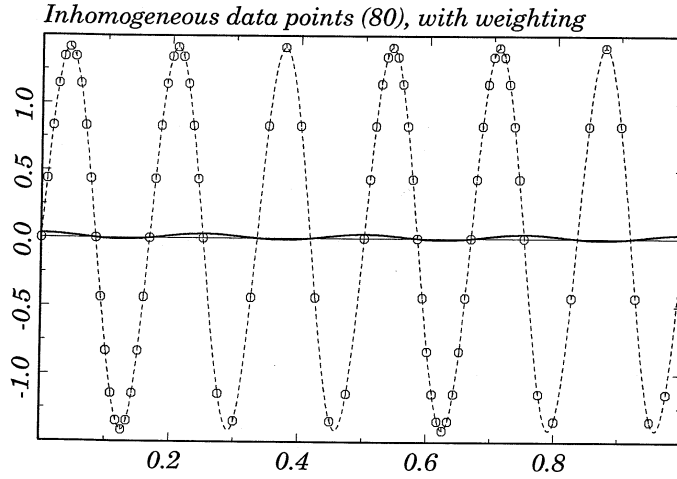


Figure 9. As the previous figure but for an inversion where each data point is with the inverse of the distance to neighbouring points.

An interesting observation is that in the interval  $0 \leq x < 0.25$  where the sampling is dense the input model and the reconstructed model are in phase whereas on the interval  $0.25 \leq x < 0.5$  where the sampling is sparser the input model and reconstructed model are out of phase. This means that the least-squares criterion has selected a good fit in the densely sampled parts of the interval at the expense of a poorer fit in the sparsely sampled parts of the interval. This suggests a cure by down-weighting the contribution of the points in the densely sampled parts of the interval. Figure 9 shows the reconstructed function when the data points are weighted with weights that are inversely proportional to the sampling distance [Snieder, 1993], it can be seen that the reconstructed model indicated by the solid line is close to its true value (which is equal to zero). Weighting the data is the key to suppressing spectral leakage, we will return to this issue in section 3.3. The reason for the spectral leakage is that the orthogonality relation of the basis functions is weighted by the data kernels (i.e.  $\sum_i \int G_i(x) B_j(x) dx \int G_i(y) B_k(y) dy$ ) and that the basis functions are not orthogonal for this inner product, i.e. this quantity is not equal to  $\delta_{jk}$ .

Let us momentarily return to the unweighted inversion of figure 8. In that example 80 data points have been used. Note that the input function

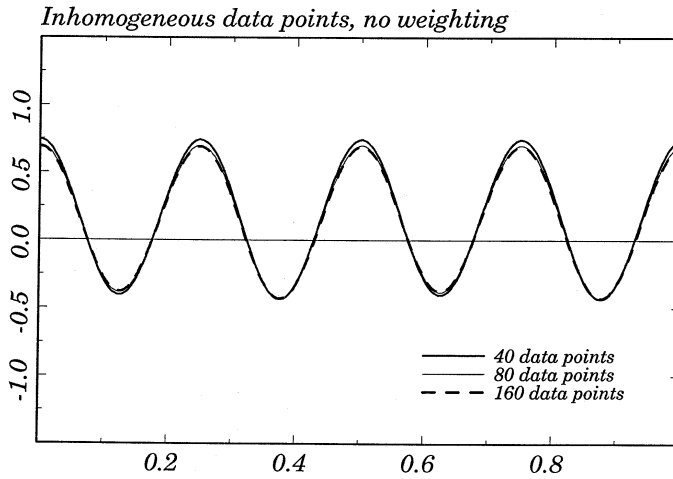


Figure 10. The estimated model when the true model is a sine wave with six wavelengths along the interval for inversions with different numbers of sampling points with the same sampling density along the line.

is not undersampled in any part of the interval, in other words there is no aliasing [Claerbout, 1976] occurring in this problem. In figure 10 the reconstructed function is shown when 40, 80 or 160 data points have been used in the inversion. In all examples the density of data points was the same as in figure 9. Spectral leakage is not a problem due to the number of data points. Adding more data points while keeping the sampling density constant does not help to suppress spectral leakage.

*Spectral leakage is a fundamentally different issue than aliasing!*

In some studies, the reliability of estimated models has been studied by dividing the data set randomly in two parts, keeping thus the sampling density constant, and comparing the models reconstructed from the two halved data sets [e.g. Woodhouse and Dziewonski, 1984]. The fact that spectral leakage does not reduce when more data points are used implies that this test does not detect artifacts due to the variability in the data coverage, and therefore may give an overly optimistic impression of the reliability of the estimated model.

### 3.3. SPECTRAL LEAKAGE, THE CURE

The method presented in this section for suppressing spectral leakage was developed by Trampert and Snieder [1996]. A key element in the analysis

is to acknowledge that more than the first  $L$  basisfunctions are needed to describe the model. This is explicitly done by minimizing

$$S = (\mathbf{d} - \mathbf{A}_L \mathbf{m}_L - \mathbf{A}_\infty \mathbf{m}_\infty)^T \mathbf{C}_d^{-1} (\mathbf{d} - \mathbf{A}_L \mathbf{m}_L - \mathbf{A}_\infty \mathbf{m}_\infty) + \mathbf{m}_L^T \mathbf{C}_{mL}^{-1} \mathbf{m}_L + \mathbf{m}_\infty^T \mathbf{C}_{m\infty}^{-1} \mathbf{m}_\infty \quad (84)$$

The first term accounts for the minimization of the data misfit. Both, the first  $L$  basis functions contribute to the data misfit with  $\mathbf{A}_L \mathbf{m}_L$ , as well as the remaining basis functions through the term  $\mathbf{A}_\infty \mathbf{m}_\infty$ . The regularization terms  $\mathbf{m}_L^T \mathbf{C}_{mL}^{-1} \mathbf{m}_L$  and  $\mathbf{m}_\infty^T \mathbf{C}_{m\infty}^{-1} \mathbf{m}_\infty$  play a crucial role because they control to what extent the data misfit is distributed over the first  $L$  basisfunctions and over the remaining basis functions.

One has to minimize expression (84) both with respect to  $\mathbf{m}_L$  and  $\mathbf{m}_\infty$ . Setting the derivatives  $\partial S / \partial \mathbf{m}_L$  and  $\partial S / \partial \mathbf{m}_\infty$  both equal to zero and eliminating  $\mathbf{m}_\infty$  from these equations leads to the following estimate of the vector  $\mathbf{m}_L$ :

$$\tilde{\mathbf{m}}_L = \left( \mathbf{A}_L^T \mathbf{C}_{dL}^{-1} \mathbf{A}_L + \mathbf{C}_{mL}^{-1} \right)^{-1} \mathbf{A}_L^T \mathbf{C}_{dL}^{-1} \mathbf{d}, \quad (85)$$

where the new operator  $\mathbf{C}_{dL}$  is given by

$$\mathbf{C}_{dL} = \mathbf{C}_d + \mathbf{A}_\infty \mathbf{C}_{m\infty} \mathbf{A}_\infty^T. \quad (86)$$

Expression (85) is very similar to the model estimate (80) obtained by simply ignoring all basis functions beyond degree  $L$ . The only difference is that the data weighting operator is given by (86) rather than by  $\mathbf{C}_d$ , but this difference is essential. The data weighting operator is positive definite. This implies that the new data weighting operator  $\mathbf{C}_{dL}$  is always larger than the old operator  $\mathbf{C}_d$ . The new operator depends on the kernels  $\mathbf{A}_\infty$  of the basisfunctions that are not inverted for, as well as the model weighting operator  $\mathbf{C}_{m\infty}$  of these basis functions. From a Bayesian point of view these facts are all related by the observation that in this approach the data are partly explained by the infinite model vector  $\mathbf{m}_\infty$ . In the inversion for  $\mathbf{m}_L$  only, the parts of the data which may be explained by  $\mathbf{m}_\infty$  should be considered as noise, because they allow a variation in the data that is not cause by  $\mathbf{m}_L$ . For this reason this approach leads to an increase of the data variance  $\mathbf{C}_{dL}$ .

Although (86) looks simple, it is not trivial to evaluate because the product  $\mathbf{A}_\infty \mathbf{C}_{m\infty} \mathbf{A}_\infty^T$  contains matrices of infinite dimensions. When the model weighting matrix  $\mathbf{C}_{m\infty}$  is chosen to be proportional to the identity matrix ( $\mathbf{C}_{m\infty} = \gamma \mathbf{I}$ ) there is a simple alternative because one only needs to evaluate the product  $\mathbf{A}_\infty \mathbf{A}_\infty^T$ . The  $ij$ -element of  $\mathbf{A}_\infty \mathbf{A}_\infty^T$  can be written as  $(\mathbf{A}_\infty \mathbf{A}_\infty^T)_{ij} = \sum_{k=L+1}^{\infty} A_{ik} A_{jk}$ . Using expression (78) for the matrix



elements and interchanging the summation and the integration one can show that

$$\left(\mathbf{A}_\infty \mathbf{A}_\infty^T\right)_{ij} = \int dx \int dy G_i(x) G_j(y) \sum_{k=L+1}^{\infty} B_k(x) B_k(y) . \quad (87)$$

The sum over  $k$  can be written as a sum over all  $k$ -values minus the sum over the first  $L$  values:

$$\sum_{k=L+1}^{\infty} B_k(x) B_k(y) = \sum_{k=1}^{\infty} B_k(x) B_k(y) - \sum_{k=1}^L B_k(x) B_k(y) . \quad (88)$$

The basis functions form a complete set. Because of the closure relation  $\sum_{k=1}^{\infty} B_k(x) B_k(y) = \delta(x - y)$  [e.g. p. 157 of *Merzbacher, 1970*] this can be written as:

$$\sum_{k=L+1}^{\infty} B_k(x) B_k(y) = \delta(x - y) - \sum_{k=1}^L B_k(x) B_k(y) . \quad (89)$$

Inserting this in (87) leaves after carrying out the  $y$ -integration in the first term and using (78) for the second term:

$$\left(\mathbf{A}_\infty \mathbf{A}_\infty^T\right)_{ij} = \Gamma_{ij} - \left(\mathbf{A}_L \mathbf{A}_L^T\right)_{ij} , \quad (90)$$

where the Gram matrix  $\Gamma$  of our inverse problem is given by

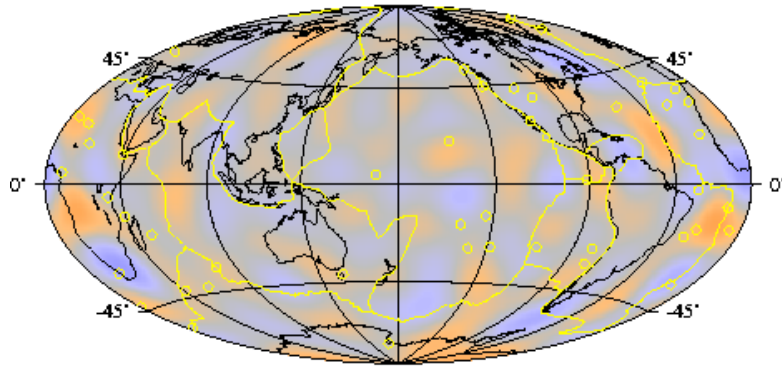
$$\Gamma_{ij} = \int G_i(x) G_j(x) dx . \quad (91)$$

The first term in (90) can be computed by direct integration, whilst the second term now only involves multiplication of matrices of *finite* dimension. It should, however, be pointed out that not all problems in geophysics possess a well defined Gram matrix.

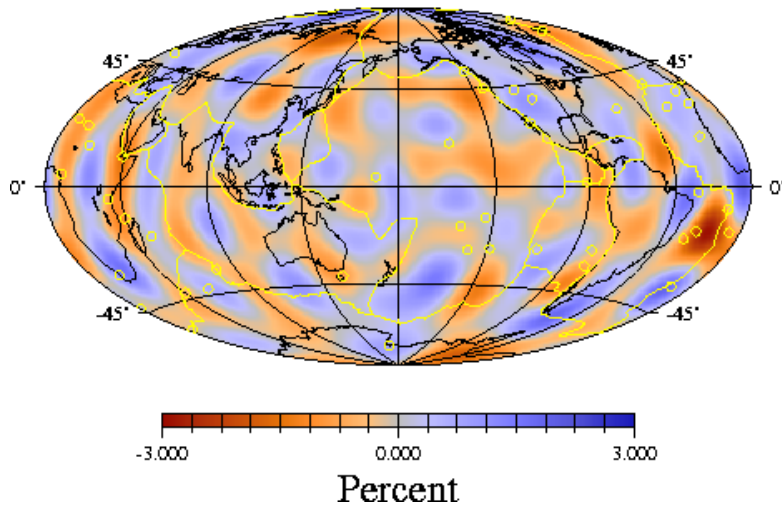
### 3.4. SPECTRAL LEAKAGE AND GLOBAL TOMOGRAPHY

The burning question is of course if spectral leakage is a serious problem in geophysics. We believe that the answer to this question is yes and illustrate this with a surface wave tomography problem. A surface wave of a given frequency sees an average structure of the outer-most layer of the Earth. The penetration of the wave depends on its frequency. The mapping problem concerns only a two-dimensional quantity (phase velocity as function of latitude and longitude) and is thus easier to represent than

true model - anti leakage solution



true model - least squares solution



*Figure 11.* Example of phase velocity mapping with inhomogeneous sampling. The true model represents up to 6% perturbations with respect to a spherically average reference model. Plotted is the difference between true and estimated model: simple least-squares inversion (bottom) and anti-leakage inversion (top).

a three-dimensional problem. To quantify the effect of leakage, we need to know the true answer to the problem.

We designed a synthetic experiment where we used a station-event distribution based on the global seismicity from 1989 together with the available global seismic stations. We have chosen arbitrarily 3000 paths giving a dense but uneven coverage (for details see *Trampert and Snieder, [1996]*). Next we took a Rayleigh wave phase velocity model [*Trampert and Woodhouse, 1995*] for a period of 80 seconds expressed in terms of spherical harmonics up to degree and order 40 and computed synthetic data for our chosen ray geometry. We added 10% random noise to the data and performed a classical least squares inversion up to degree and order 12. In figure 11 (bottom panel) we show the difference between the true degree 12 input model and the recovered model. This signal is far from zero and reaches in most places on the globe half the amplitude of the true model. If we applied the anti-leakage weighting defined in expression(86) the amplitudes of the artefacts are significantly reduced, see the top panel of figure 11. The anti-leakage operator does not completely prevent the bias because of the errors which we introduced in the synthetic data. To understand this, recall from equation (86) that the description of data errors and the anti-leakage operator add up linearly and thus any imperfect description of data errors will influence the anti-leakage. This suggests that our current tomographic models are likely to carry a bias from small-scale structures that are not accounted for by our current parameterizations.

#### 4. The single scattering approximation and linearized waveform inversion

In the previous sections, the inverse problem was described when the forward problem was linear. Unfortunately, most wave propagation problems are nonlinear in the sense that the relation between the wave field and the medium is nonlinear. However, in many practical situations this relation can be linearized, notably with the single-scattering approximation (this section), with Fermat's theorem (section 6) and Rayleigh's principle (section 5).

##### 4.1. THE BORN APPROXIMATION

Many wave propagation problems can symbolically be written in the form

$$Lu = F \tag{92}$$

In this expression  $u$  is the wave field,  $F$  accounts for the source that excites the waves and  $L$  is a differential operator that describes the wave propagation. For example, for acoustic waves in a medium with constant density (92) is given by

$$\left( \nabla^2 + \frac{\omega^2}{c^2(\mathbf{r})} \right) u(\mathbf{r}) = F(\mathbf{r}) , \quad (93)$$

so that  $L = \nabla^2 + \omega^2/c^2(\mathbf{r})$ . Often, the operator  $L$  can be decomposed into an operator  $L_0$  that correspond to a reference medium for which we can solve the problem easily plus a small perturbation  $\varepsilon L_1$  that accounts for a perturbation of the medium:

$$L = L_0 + \varepsilon L_1 . \quad (94)$$

The small parameter  $\varepsilon$  has been introduced to facilitate a systematic perturbation approach. As an example, the operator  $L$  in (93) can be decomposed using

$$\frac{1}{c^2(\mathbf{r})} = \frac{1}{c_0^2} (1 + \varepsilon n(\mathbf{r})) , \quad (95)$$

where  $c_0$  is a constant velocity that described the mean properties of the medium well while  $n(\mathbf{r})$  accounts for the velocity perturbations. Using this decomposition one obtains  $L_1 = n(\mathbf{r})\omega^2/c_0^2$ .

The wavefield is perturbed by the perturbation of the medium, this perturbation can be written as a regular perturbation series:

$$u = u_0 + \varepsilon u_1 + \varepsilon^2 u_2 + \dots \quad (96)$$

A systematic perturbation approach is obtained by inserting (94) and (96) in (92) and by collection the terms that are of equal power in the perturbation strength  $\varepsilon$ . The terms proportional to  $\varepsilon^0$  and  $\varepsilon^n$  respectively lead to the following equations:

$$L_0 u_0 = F , \quad (97)$$

$$L_0 u_n = -L_1 u_{n-1} \quad \text{for} \quad n \geq 1 . \quad (98)$$

Equations (97) and (98) are of the form  $L_0 u = \text{forcing term}$  and hence can be solved using the Green's function  $G$  of the unperturbed problem, i.e. by using the Green's function that satisfies

$$L_0 G(\mathbf{r}, \mathbf{r}') = \delta(\mathbf{r} - \mathbf{r}') . \quad (99)$$

For vector waves the Green's function should be replaced by a second order tensor, but the principles are the same. Using this Green's function the solution of (97) and (98) is given by

$$u_0 = GF \quad \text{and} \quad u_n = -GL_1 u_{n-1} \quad (n \geq 1). \quad (100)$$

The solution can now be constructed by solving the above equation recursively and by inserting the result in (96):

$$u = \underbrace{u_0}_{\text{Unperturbed wave}} - \underbrace{\varepsilon GL_1 u_0}_{\text{Single scattered wave}} + \underbrace{\varepsilon^2 GL_1 GL_1 u_0}_{\text{Double scattered wave}} + \cdots \quad (101)$$

In this formulation the total wavefield is written as a sum over the unperturbed wave, the waves that are scattered once by the inhomogeneity  $L_1$ , the waves that are scattered twice by the perturbation  $L_1$  and all the higher order scattered waves. The series (101) is called the *Neumann series* in scattering theory.

The Born approximation  $u^{Born}$  consists in truncating this multiple scattering series after the single scattered wave:

$$u^{Born} = u_0 - \varepsilon GL_1 u_0. \quad (102)$$

The great advantage of this approximation is that the scattered waves are given by  $-\varepsilon GL_1 u_0$ , hence the scattered waves now depend linearly on the perturbation of the medium that is contained in the operator  $L_1$ . This means that the scattered waves are related in this approximation linearly to the perturbation of the medium. This makes it possible to use the theory of linear inverse problems as shown in the section 2 for the solution of this problem. In doing so, one must keep in mind that the Born approximation ignores multiple scattering effects. When such effects are present in the data one should be extremely cautious in using the Born approximation.

#### 4.2. INVERSION AND MIGRATION

As one of the simplest examples let us return to the acoustic wave equation (93) with the perturbed velocity given in (95). The unperturbed problem is then given by  $(\nabla^2 + k^2)u = F$  where the constant wavenumber is given by  $k = \omega/c_0$ . The Green's function for this differential equation is given by

$$G(\mathbf{r}, \mathbf{r}') = -\frac{e^{ik|\mathbf{r} - \mathbf{r}'|}}{4\pi |\mathbf{r} - \mathbf{r}'|}. \quad (103)$$

When a point source with spectrum  $F(\omega)$  is located in  $\mathbf{r}_s$  the unperturbed wave is given by

$$u_0(\mathbf{r}) = - \frac{e^{ik|\mathbf{r} - \mathbf{r}_s|}}{4\pi |\mathbf{r} - \mathbf{r}_s|} F(\omega) . \quad (104)$$

In the Born approximation, the scattered waves are scattered only once, it follows from (102) that the single-scattered waves at a receiver position  $\mathbf{r}_r$  are given by:

$$u_s(\mathbf{r}_r) = - \frac{1}{(4\pi c_0)^2} \int \frac{e^{ik|\mathbf{r}_r - \mathbf{r}|}}{|\mathbf{r}_r - \mathbf{r}|} n(\mathbf{r}) \frac{e^{ik|\mathbf{r} - \mathbf{r}_s|}}{|\mathbf{r} - \mathbf{r}_s|} dV F(\omega) . \quad (105)$$

Suppose one measures these single-scattered waves, as is done in a seismic reflection experiment. The inverse problem then consists in finding the model perturbation  $n(\mathbf{r})$  given the recorded data  $u_s(\mathbf{r}_r)$ . The theory of section 2 can be used for this when one discretizes the volume integral in (105). After discretization the scattering integral by dividing the volume in small cells, this integral can be written in the form

$$u_i = \sum_j A_{ij} n_j . \quad (106)$$

Since in a realistic imaging experiment one uses many source positions and records the wavefield at many receivers,  $u_i$  stands for the reflected wave for source-receiver pair and frequency component  $\#i$ . Because of the discretized volume integral,  $n_j$  denotes the perturbation of  $1/c^2(\mathbf{r})$  in cell  $\#j$ . The matrix elements  $A_{ij}$  are given by

$$A_{ij} = - \frac{1}{(4\pi c_0)^2} \int_{cell\ j} \frac{e^{ik_i|\mathbf{r}_{ri} - \mathbf{r}|}}{|\mathbf{r}_{ri} - \mathbf{r}|} \frac{e^{ik_i|\mathbf{r} - \mathbf{r}_{si}|}}{|\mathbf{r} - \mathbf{r}_{si}|} dV F(\omega_i) . \quad (107)$$

In principle, one can solve the linear system (106) by brute force. However, in many realistic problems the size of the system of equations is so large that this is practically speaking impossible. This is notably the case in seismic exploration where the number data and the number of model parameters are exceedingly large. For problems of this scale, iterative solutions of the linear system of equations seems the only realistic way of obtaining a solution. In fact, it is in practice only possible to carry out the first step of such an iterative process. Let us find an estimated model by using the first step of the iterative solution (71) of section 2.9.2 using the preconditioning operator  $\mathbf{P} = \text{const.} \cdot \mathbf{I}$ . In this approach the

estimated model is given by:

$$\tilde{\mathbf{n}}(\mathbf{r}) \sim \mathbf{A}^\dagger \mathbf{d} \sim \sum_{\substack{\text{sources} \\ \text{receivers} \\ \text{frequencies}}} \frac{e^{-i\omega(|\mathbf{r}_r - \mathbf{r}| + |\mathbf{r} - \mathbf{r}_s|)/c}}{|\mathbf{r}_r - \mathbf{r}| |\mathbf{r} - \mathbf{r}_s|} d_{rs}(\omega), \quad (108)$$

where  $d_{rs}(\omega)$  denotes the scattered wave for source  $s$ , recorded at receiver  $r$  with frequency  $\omega$ . For simplicity it is assumed here that the source signal is a delta-pulse in the time domain, so that  $F(\omega) = 1$ , and that all cells have equal volume:  $V_j = \text{const}$ . Note that the transpose  $\mathbf{A}^T$  has been replaced by the Hermitian conjugate  $\mathbf{A}^\dagger$ .<sup>6</sup> The reason for this is that the analysis of section 2 was for real matrices. A similar analysis for complex matrices shows that the results are identical provided the transpose  $\mathbf{A}^T$  is replaced by its complex conjugate. The summation in the matrix product effectively leads to a summation over all sources, receivers and frequency components because all these were labelled by the index  $i$  in (106).

It is instructive to consider the summation over frequencies in (108). At each frequency the data  $d_{rs}(\omega)$  are multiplied with  $\exp(-i\omega\tau)$  with  $\tau = (|\mathbf{r}_r - \mathbf{r}| + |\mathbf{r} - \mathbf{r}_s|)/c$ . This means that the frequency summation has the form of a Fourier transform so that up to a constant, the frequency summation gives the data in the time domain at time  $t = (|\mathbf{r}_r - \mathbf{r}| + |\mathbf{r} - \mathbf{r}_s|)/c$ :

$$\tilde{\mathbf{n}}(\mathbf{r}) \sim \sum_{\substack{\text{sources} \\ \text{receivers}}} \frac{1}{|\mathbf{r}_r - \mathbf{r}| |\mathbf{r} - \mathbf{r}_s|} d_{rs}(t = \frac{(|\mathbf{r}_r - \mathbf{r}| + |\mathbf{r} - \mathbf{r}_s|)}{c}) \quad (109)$$

This expression implies that the image can be constructed by summing the data over all source-receiver pairs, and by considering for each target point the data at a time needed to travel from the source location  $\mathbf{r}_s$  to the target point  $\mathbf{r}$  to the receiver location  $\mathbf{r}_r$ .

This is the procedure followed in the imaging of seismic reflection data called ‘‘Kirchhoff migration’’ [Claerbout, 1985; Yilmaz, 1987]. Effectively one sums in such an approach over all the available data that have a travel time consistent with a scatterer at the target point  $\mathbf{r}$ . The only difference with the classical Kirchhoff migration is the presence of the geometrical spreading terms  $1/(|\mathbf{r}_r - \mathbf{r}| |\mathbf{r} - \mathbf{r}_s|)$  that are not included in Kirchhoff migration. In non-destructive testing this approach is known as Synthetic Aperture Focussing Technique (SAFT) [Mayer *et al.*, 1990].

<sup>6</sup>The complex conjugate of the transpose is by definition the Hermitian conjugate:  $A_{ij}^\dagger = A_{ji}^*$ .

The main conclusion is that Kirchhoff migration (up to some constants) corresponds to the first step of an iterative solution of the linearized scattering equation. The derivation of this section is a simplified version of the derivation given by *Tarantola* [1984] who incorporates other source signals than a delta pulse.

#### 4.3. THE BORN APPROXIMATION FOR TRANSMISSION DATA

There is a widespread belief that the Born approximation can only be used for truly scattered waves. However, the Born approximation can also be used to account for effects of medium perturbations on transmitted waves, but with a domain of applicability to transmission problems that is smaller than for reflection or true scattering problems. To see this, assume that the velocity has a small constant perturbation  $\delta c$  and that a wave propagates over a distance  $L$ . The wave will then experience a phase shift  $\exp i\varphi$ , given by  $\varphi = -(\omega/c^2) L\delta c$ . In the Born approximation this perturbation is replaced by  $\exp i\varphi \approx 1 + i\varphi$ . This is only a good approximation when the phase shift  $\varphi$  is much less than a cycle. In practice, this limits the use of the Born approximation for the inversion of transmission data. Note that even when the velocity perturbation is small, the requirement that the phase shift is small breaks down for sufficiently large propagation distances  $L$ .

This does not imply that the Born approximation cannot be used for the inversion of transmitted waves. For surface waves propagating in the earth, the wavelength can be very long. (A Rayleigh wave at a period of 75s has a wavelength of about 300km.) This means that these waves do not travel over very many wavelengths when they propagate for a few thousand kilometers. In addition, the heterogeneity in the Earth's mantle is not very large. This makes it possible to carry out inversions of transmitted surface wave data using the Born approximation.

The Born approximation for surface waves with the resulting scattering and mode-conversion coefficients is derived for a flat geometry by *Snieder* [1986a, 1986b], the generalization to a spherical geometry can be found in *Snieder and Nolet* [1987]. Although the Green's function and the scattering and mode-conversion coefficients are different for elastic surface waves than for acoustic waves, the Born approximation leads to a linearized relation between the perturbation  $\delta u$  of the waveform and the perturbation  $m$  of the Earth model:

$$\delta u_{rs}(\omega) = \iiint K_{rs}(\mathbf{r}, \omega) m(\mathbf{r}) dV , \quad (110)$$

where  $K_{rs}(\mathbf{r}, \omega)$  contains the surface wave Green's function of the incoming and outgoing wave at the scatterer location  $\mathbf{r}$  for source-receiver pair



“ $rs$ ” and frequency  $\omega$ . The perturbation  $m(\mathbf{r})$  stands for the perturbation in the elasticity tensor and/or the density. After discretization the volume integral can be written as a linear system of equations:

$$\delta u_i = \sum_j K_{ij} m_j, \quad (111)$$

where  $\delta u_i$  stands for the perturbation of the surface wave at source-receiver pair  $\#i$  and frequency  $\omega_i$ , while  $m_j$  stands for the perturbation of the Earth model in cell  $\#j$ . This linear system can be solved numerically, and the specific implementation for the inversion of surface wave data is shown by *Snieder* [1988a].

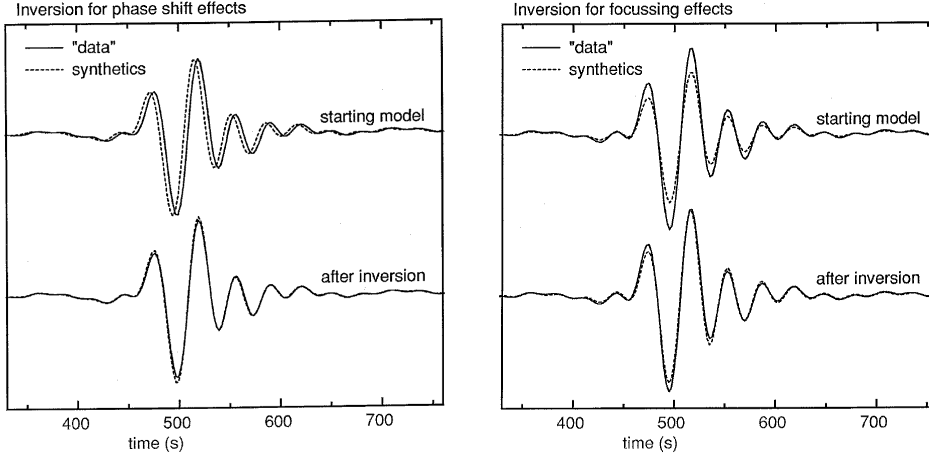
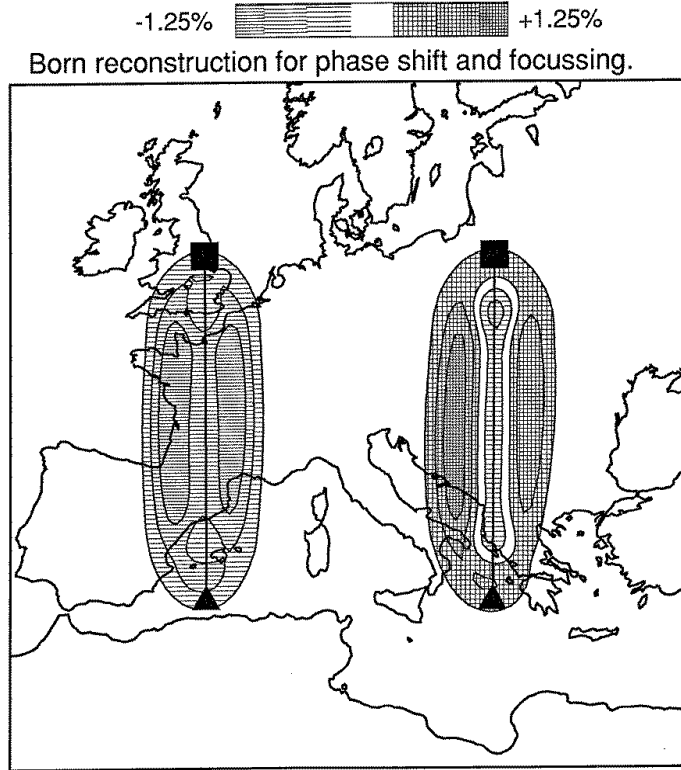


Figure 12. Phase shifted data (left panel) and data with an amplitude anomaly (right panel) before and after inversion. The synthetic data are shown with a solid line, the synthetics before and after inversion with a dashed line.

As an example how this algorithm operates consider the two upper seismograms shown in the two panels of figure 12. The seismogram of the left panel has been shifted in time with respect to the seismogram of a laterally homogeneous Earth model, whereas the amplitude of the seismogram of the right panel has been increased by 20%. Both seismograms are inverted simultaneously, and the corresponding model is shown in figure 13. (The map is shown only to display the scale of the problem but has otherwise no specific meaning.) The fit of the wave-forms after inversion is shown by the lower seismograms in figure 12, it can be seen that both the time shift and the amplitude change are accounted for well.

The triangle and square on the left are the source and receiver respectively of the time-shifted seismogram while the same symbols on the right are for the seismogram with the amplitude change. The velocity



*Figure 13.* The velocity model obtained from the inversion of the seismograms in the previous figure. The source locations are marked with triangles, the receivers with squares. The source-receiver pair on the left is for the phase shifted seismogram, the pair on the right for the seismogram with an amplitude error. The map serves only to fix the scale.

anomaly for the time-shifted seismogram is a negative velocity perturbation straddling the line joining the source and receiver on the left in figure 13. This negative velocity change gives the required time-shift. Note that this velocity perturbation is not confined to the source-receiver line; its finite extent is due to the fact that rays have no physical meaning and that a wave is influenced by the velocity perturbation averaged over the first Fresnel zone [Snieder and Lomax, 1996].

The velocity structure on the right is for the seismogram with the perturbed amplitude. The perturbation is negative on the source-receiver line, but slightly further away from this line the velocity perturbation is positive. Since waves are deflected from areas of high velocity towards regions of low velocity, the wave energy is deflected towards the receiver. This means that the algorithm has realized a fit of the amplitude perturbation by creating a “surface wave lens” that produces just the right

amount of focussing at the receiver to account for the amplitude perturbation.

Note that the model of figure 13 was constructed by numerically solving the linear system of equations (111). The input of the system of equations simple consisted of the real and imaginary components of the perturbation of the surface waves in the frequency domain. At no point in the inversion has the phase or amplitude of the waves been prescribed explicitly. However, the Born approximation contains all the relevant physics needed to translate the perturbation of the real and imaginary components of the waves into the physics of focussing and phase retardation.

#### 4.4. SURFACE WAVE INVERSION OF THE STRUCTURE UNDER NORTH-AMERICA

The surface waveform inversion has been applied to infer the shear-velocity structure under Europe and the Mediterranean by *Snieder* [1988b]. In this section a model for the shear-velocity under North-America is shown that was made by *Alsina et al.* [1996]. The shear-velocity perturbation in three layers with depths between 25 and 300 km is shown in figure 14. Red colours indicate slow velocities that are indicative of high temperature while green colours indicate high velocities correspond to low temperature.<sup>7</sup>

It is instructive to consider the west-coast of North America. Under the Gulf of California the velocity is very low. In this region, the East-Pacific rise meets the continent. (The East-Pacific rise is a spreading center in the ocean bottom comparable to the Mid-Atlantic ridge in the Atlantic Ocean.) Since in a spreading center hot material wells up from the mantle, one expects the temperature to be low, which is consistent with the slow velocity in that region.

Further north at the state of California the velocity perturbation is close to zero. In this area the Pacific plate slides horizontally along the North-American plate. This so called “strike-slip motion” gives rise to earthquakes in California. However, since the plate simply slide along each other, the temperature field is not perturbed very much, which is reflected by the neutral velocity perturbation.

However, northward of Cape Mendocino under the states of Oregon and Washington there is a distinct positive velocity perturbation. In this

<sup>7</sup>The identification of velocity anomalies with temperature perturbations should be treated with care. Apart from temperature perturbations the seismic velocity is also affected by variations in composition, the presence of volatiles such as water and possibly also pressure.

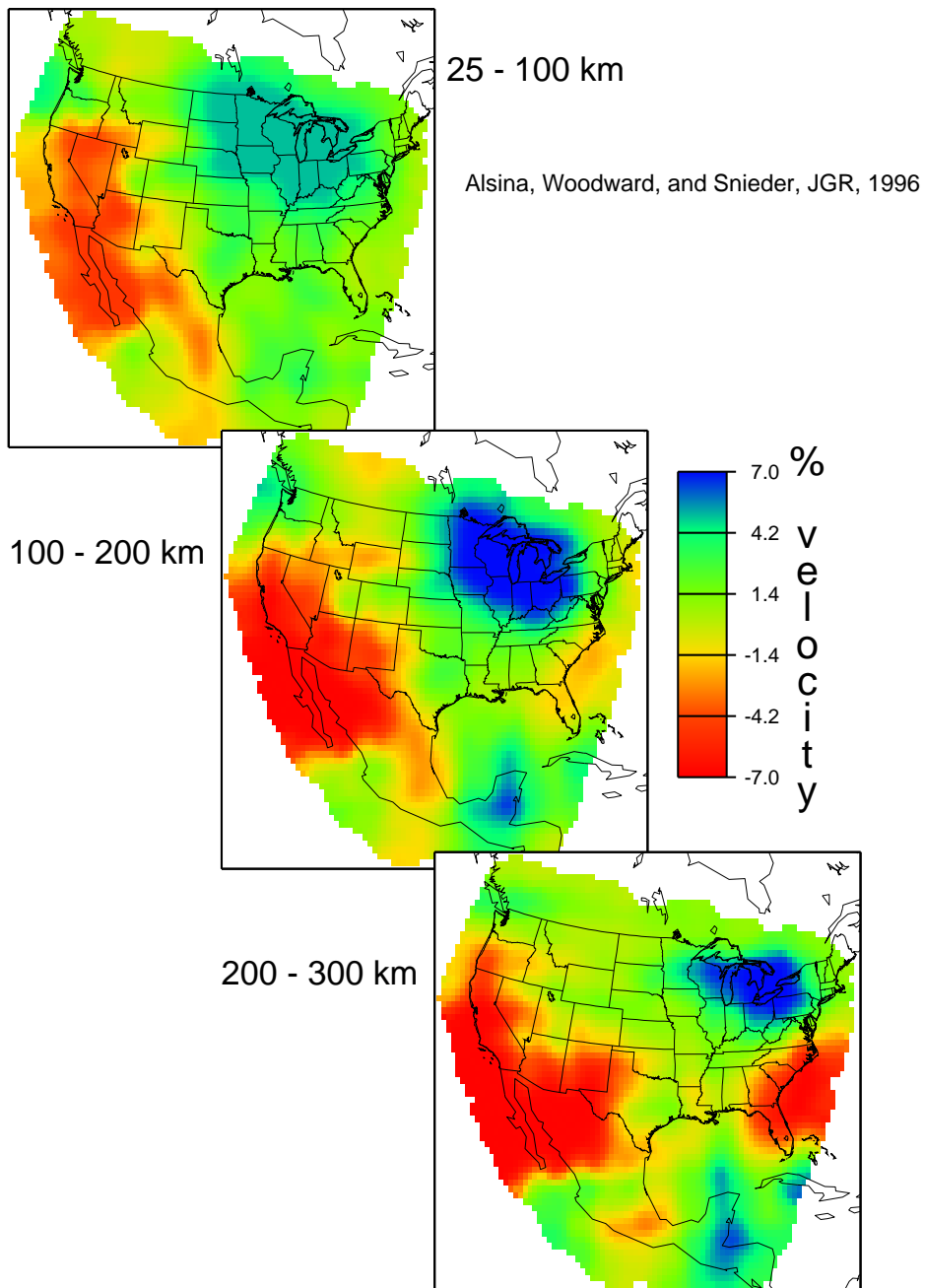


Figure 14. The relative shear velocity perturbation under North America in three layers at depths between 25 and 300 km as determined by [Alsina *et al.*, 1996].

region oceanic plates slide eastward under the continents. This process is called “subduction.” Since the subsiding plate is colder than the surrounding material this is reflected in the image as positive velocity anomalies.

The model shown in figure 14 and it’s relation with the different tectonic regimes at the west coast of North-America show that with tomographic techniques as shown here it is possible to see plate-tectonics in action.

## 5. Rayleigh’s principle and perturbed eigenfrequencies

Important information about the structure of the Earth’s interior follows from observations of perturbations of the frequencies of the free oscillations of the Earth as well as from measurements of the perturbation of the phase or group velocity of surface waves. This information can most easily be retrieved from observations when the relation between the frequency shift (or surface wave phase velocity) and the perturbation of the Earth model can be linearized. This can be achieved by applying Rayleigh-Schrödinger perturbation theory.

### 5.1. RAYLEIGH-SCHRÖDINGER PERTURBATION THEORY

Consider an eigenvalue problem of the following form:

$$Hu_n = \lambda_n r u_n . \quad (112)$$

In this expression  $H$  is an operator, for wave-propagation problems it usually is a differential operator. The eigenfunctions of this operator are denoted by  $u_n$  and the corresponding eigenvalues by  $\lambda_n$ . The function  $r$  denotes a positive weight function. For the Earth the equation of motion is given by

$$\partial_j (c_{ijkl} \partial_k u_l) = -\rho \omega^2 u_i , \quad (113)$$

where  $c_{ijkl}$  is the elasticity tensor. A comparison with (112) shows that  $H(\bullet)_i = \partial_j (c_{ijkl} \partial_k (\bullet)_l)$ , that the weight function  $r$  is given by the density  $\rho$  and that  $\lambda$  corresponds to  $-\omega^2$ . In this case,  $H$  is a Hermitian operator:

$$H^\dagger = H . \quad (114)$$

For the special case of elastic wave propagation, this property stems from the symmetry properties of the elasticity tensor and the stress-free boundary conditions at the Earth’s surface [Nolet, 1981; Dahlen and Tromp, 1998]. For a Hermitian operator the eigenvalues are real and the eigenfunctions are orthonormal with respect to the following inner product

$$\langle u_n | r u_m \rangle = \int u_n^* r u_m dV = \delta_{nm} . \quad (115)$$

Let the operator  $H$  and the weight function  $r$  be decomposed in a reference operator  $H_0$  and weight function  $r_0$  for which we know the eigenfunctions  $u_n^{(0)}$  and eigenvalues  $\lambda_n^{(0)}$  and perturbations  $\varepsilon H_1$  and  $\varepsilon r_1$ :

$$H = H_0 + \varepsilon H_1 \quad , \quad r = r_0 + \varepsilon r_1 . \quad (116)$$

Under this perturbation the eigenfunctions and eigenvalues are perturbed as well:

$$u_n = u_n^{(0)} + \varepsilon u_n^{(1)} + \varepsilon^2 u_n^{(2)} + \dots \quad (117)$$

$$\lambda_n = \lambda_n^{(0)} + \varepsilon \lambda_n^{(1)} + \varepsilon^2 \lambda_n^{(2)} + \dots \quad (118)$$

The goal of this analysis is to find the first order perturbation  $\lambda_n^{(1)}$  of the eigenvalues. This can be achieved by inserting the expressions (116) through (118) in (112) and by collecting the terms of order  $\varepsilon^1$ , this gives:

$$H_1 u_n^{(0)} + H_0 u_n^{(1)} = \lambda_n^{(1)} r_0 u_n^{(0)} + \lambda_n^{(0)} r_1 u_n^{(0)} + \lambda_n^{(0)} r_0 u_n^{(1)} \quad (119)$$

The problem with this expression is that we are only interested in  $\lambda_n^{(1)}$ , but that in order to retrieve this term from (119) we need the first order perturbation  $u_n^{(1)}$  of the eigenfunctions as well. The perturbation of the eigenvalue can be extracted by taking the inner product of (119) with the unperturbed eigenfunction  $u_n^{(0)}$ :

$$\begin{aligned} & \left\langle u_n^{(0)} | H_1 u_n^{(0)} \right\rangle + \underbrace{\left\langle u_n^{(0)} | H_0 u_n^{(1)} \right\rangle}_{(\heartsuit)} \\ &= \lambda_n^{(1)} \left\langle u_n^{(0)} | r_0 u_n^{(0)} \right\rangle + \lambda_n^{(0)} \left\langle u_n^{(0)} | r_1 u_n^{(0)} \right\rangle + \underbrace{\lambda_n^{(0)} \left\langle u_n^{(0)} | r_0 u_n^{(1)} \right\rangle}_{(\spadesuit)} \end{aligned} \quad (120)$$

Note that the perturbed eigenfunctions  $u_n^{(1)}$  only appear in the terms marked  $(\heartsuit)$  and  $(\spadesuit)$ . Using the fact that  $H_0$  is Hermitian and that the eigenvalues  $\lambda_n^{(0)}$  are real one finds that these terms are equal:

$$\begin{aligned} (\heartsuit) &= \left\langle u_n^{(0)} | H_0 u_n^{(1)} \right\rangle = \left\langle H_0^\dagger u_n^{(0)} | u_n^{(1)} \right\rangle = \left\langle H_0 u_n^{(0)} | u_n^{(1)} \right\rangle \\ &= \left\langle \lambda_n^{(0)} u_n^{(0)} | u_n^{(1)} \right\rangle = \lambda_n^{(0)} \left\langle u_n^{(0)} | u_n^{(1)} \right\rangle = (\spadesuit) \end{aligned} \quad (121)$$

This means that the terms containing the perturbed eigenfunctions cancel. Solving the remaining equation for the perturbation of the eigenvalue gives:

$$\lambda_n^{(1)} = \frac{\left\langle u_n^{(0)} | \left( H_1 - \lambda_n^{(0)} r_1 \right) u_n^{(0)} \right\rangle}{\left\langle u_n^{(0)} | r_0 u_n^{(0)} \right\rangle} . \quad (122)$$

The perturbation of the eigenvalue thus follows by evaluating the inner product of the perturbed operators sandwiched between the unperturbed eigenfunctions. This is one form of Rayleigh's principle: the eigenvalues are stationary to first order for perturbations of the eigenfunctions. The crux is that according to 122 the first-order perturbation of the eigenvalues depends linearly on the perturbations of the operator  $H$  and weight  $r$  and hence linear inverse theory as exposed in the previous sections may be applied to infer the perturbations of  $H$  and  $r$  from the measured eigenvalues of the system.

## 5.2. THE PHASE VELOCITY PERTURBATION OF LOVE WAVES

Expression (122) can be used to evaluate the perturbation of the phase velocity of surface waves due to perturbations of the Earth model. This is shown here for the simplest example of Love waves in a plane geometry. The reader is referred to *Aki and Richards* [1980] for the extension to Rayleigh waves, while the analysis for a spherical geometry is treated by *Takeuchi and Saito* [1972]. As shown in *Aki and Richards* [1980] the Love wave eigenfunctions satisfy the following differential equation:

$$\partial_z(\mu \partial_z v) + (\rho \omega^2 - \mu k^2) v = 0, \quad (123)$$

where  $\mu$  is the shear modulus. The surface waves modes vanish at large depth ( $z \rightarrow \infty$ ) and are stress-free at the surface. This gives the following boundary conditions:

$$\partial_z v(z=0) = 0, \quad v(z=\infty) = 0. \quad (124)$$

In this analysis we assume that the frequency  $\omega$  is a given constant. The operator  $H$  is given by  $H = \partial_z \mu \partial_z + \rho \omega^2$ , the weight function  $r$  is given by the shear modulus ( $r(z) = \mu(z)$ ) while  $k^2$  is the eigenvalue ( $\lambda = k^2$ ). An integration by parts using the boundary conditions (124) can be used to show that the operator  $H$  is Hermitian. The theory of section 5.1 can then be used to find the perturbation  $\delta k$  of the wavenumber due to a perturbation  $\mu_1$  in the shear modulus and a perturbation  $\rho_1$  in the density. Using the first order relation  $\delta k^2 = 2k \delta k$  one obtains from (122) that:

$$\delta k = \frac{\int (\rho_1 \omega^2 - \mu_1 k^2) v^2 dz - \int \mu_1 (\partial_z v)^2 dz}{2k \int \mu_0 v^2 dz}. \quad (125)$$

Since the phase velocity is given by  $c = \omega/k$  the phase velocity perturbation follows to leading order from this expression by using the relation  $\delta c/c = -\delta k/k$ . This means that the first order effect of the perturbation of the Earth model on the phase velocity of Love waves can readily be

computed once the unperturbed eigenfunctions  $v(z)$  are known. A similar result can be derived for Rayleigh waves.

In general, the perturbation in the phase velocity due to perturbation in the density, P-wave velocity  $\alpha$  and S-wave velocity  $\beta$  can be written as:

$$\frac{\delta c}{c} = \int K_\beta(z) \frac{\delta \beta(z)}{\beta(z)} dz + \int K_\alpha(z) \frac{\delta \alpha(z)}{\alpha(z)} dz + \int K_\rho(z) \frac{\delta \rho(z)}{\rho(z)} dz. \quad (126)$$

For Love waves the kernel  $K_\alpha(z)$  vanishes because Love waves are independent of the bulk modulus. Examples of the kernels  $K_\beta$ ,  $K_\alpha$  and  $K_\rho$  are shown in figure 15 for the fundamental mode Rayleigh waves of periods of 100s (left panel) and 30s (right panel) respectively. It can be seen that the sensitivity of Rayleigh waves to the P-velocity  $\alpha$  is much less than the sensitivity for changes in the shear velocity  $\beta$  and that the sensitivity kernel  $K_\alpha$  decays more rapidly with depth than the other kernels. Both phenomena are a consequence of the fact that the compressive component of the motion become evanescent at much shallower depth than the shear component.

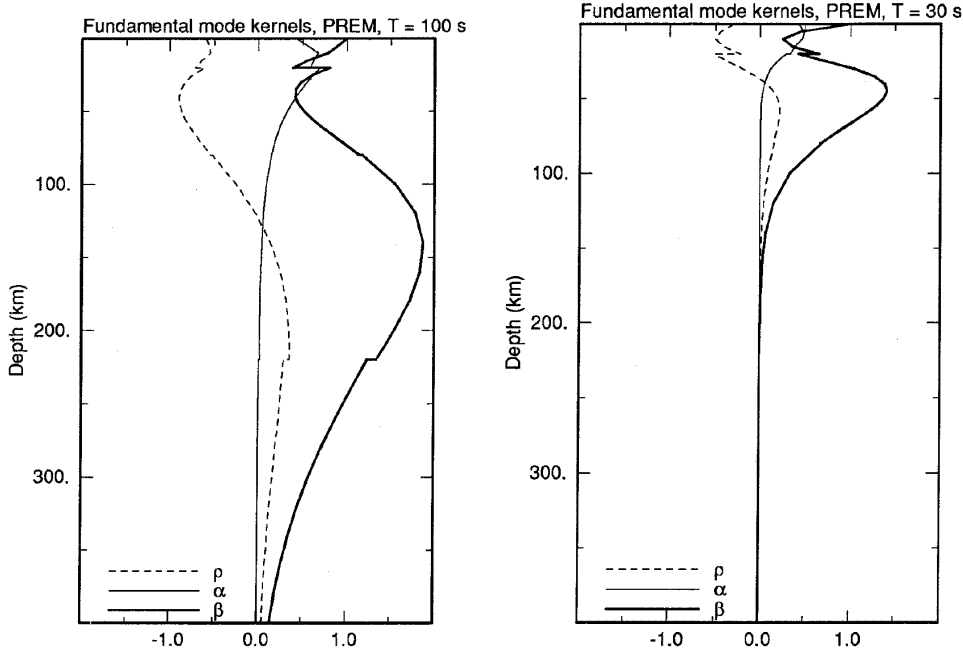


Figure 15. The sensitivity kernels  $K_\alpha$  (thin solid line),  $K_\beta$  (thick solid line) and  $K_\rho$  (dashed line) for the fundamental mode Rayleigh wave for a period of 100s (left panel) and 30s (right panel). The kernels are computed for the PREM model [Dziewonski and Anderson, 1981].



For larger periods (left panel) the sensitivity kernels penetrate over a greater depth range than for shorter periods (right panel). This makes it possible to obtain a depth resolution in surface wave inversions. By measuring the phase velocity of surface waves at different periods, one obtains according to (126) and figure 15 the inner product of the perturbation of the Earth model with different weight functions  $K(z)$ . The theory of linear inversion in section 2 can then be used to determine the perturbation of the Earth model as a function of depth. Of course, only a finite depth resolution can be obtained, but using the expressions (4) or (75) one can determine the depth resolution that can be obtained. The depth resolution can be increased by using higher modes as well as the fundamental mode surface waves and by using both Love waves and Rayleigh waves [Nolet, 1977; Cara, 1978; van Heijst and Woodhouse, 1997].

The theory of this section can be extended to include anisotropic perturbations in the Earth as well [Tanimoto, 1986; Montagner and Nataf, 1986]. The depth-dependence of the sensitivity kernels of the Rayleigh wave phase velocity show a dependence of the anisotropic P-wave velocity that is fundamentally different than that of the isotropic P-wave velocity [Muyzert and Snieder, 1999].

## 6. Fermat's theorem and seismic tomography

Travel-time tomography is a technique where one aims to reconstruct the velocity structure of a body given the measurement of travel times of waves that have propagated through that body. The travel time along a ray is given by

$$T = \int_{\mathbf{r}[u]} u(\mathbf{r}) ds . \quad (127)$$

In this expression,  $u$  is the slowness which is defined as the reciprocal of the velocity:  $u = 1/c$ . The slowness is used rather than the velocity because now the integrand is linear to the quantity we aim to retrieve. It is tempting to conclude from (127) that the relation between the travel time and the slowness is linear. However, this is wrong! The reason for this is that the integration in (127) is along the path on which the waves travel. The rays are curves of stationary travel time, and hence the ray location depends on the slowness as well. Travel time tomography is thus a nonlinear inverse problem: the unknown slowness is present both in the integrand and it determines the ray position  $\mathbf{r}[u]$  in the travel time integral (127). Linear inversion techniques can only be used when the relation between  $T$  and  $u$  is linearized.

Traditionally this is achieved by invoking Fermat's theorem which states that the travel along a ray does not change to first order when this ray is perturbed [e.g. *Nolet*, 1987; *Ben-Menahem and Singh*, 1981]. The proof of Fermat's theorem is based on variational calculus and invokes the equations of kinematic ray tracing. However, the same result can be derived in a much simpler way when one starts from the eikonal equation.

### 6.1. FERMAT'S THEOREM, THE EIKONAL EQUATION AND SEISMIC TOMOGRAPHY

The derivation given in this section was formulated by *Aldridge* [1994] and is based on a perturbation analysis of the eikonal equation. Here, only the first order travel time perturbation will be derived, but *Snieder and Aldridge* [1995] have generalized the analysis to arbitrary order. Starting point is the eikonal equation which governs the propagation of wavefronts:

$$|\nabla T|^2 = u^2(\mathbf{r}) . \quad (128)$$

Consider a reference slowness  $u_0(\mathbf{r})$  that is perturbed by a perturbation  $\varepsilon u_1(\mathbf{r})$ , where the parameter  $\varepsilon$  serves to facilitate a systematic perturbation approach:

$$u(\mathbf{r}) = u_0(\mathbf{r}) + \varepsilon u_1(\mathbf{r}) . \quad (129)$$

Under this perturbation the travel changes, and we assume that the travel time can be written as a regular perturbation series:<sup>8</sup>

$$T = T_0 + \varepsilon T_1 + \varepsilon^2 T_2 + \dots \quad (130)$$

Inserting (129) and (130) in the eikonal equation (128) and collecting terms proportional to  $\varepsilon^0$  and  $\varepsilon^1$  gives:

$$|\nabla T_0|^2 = u_0^2(\mathbf{r}) , \quad (131)$$

$$(\nabla T_0 \cdot \nabla T_1) = u_0 u_1 . \quad (132)$$

The first equation is the eikonal equation for the reference travel time. Let the unit vector  $\hat{\mathbf{t}}_0$  be directed along the gradient of  $T_0$ , then using (131) this gradient can be written as

$$\nabla T_0 = u_0 \hat{\mathbf{t}}_0 . \quad (133)$$

<sup>8</sup>The assumption that the travel time perturbation is regular is not valid when caustics are present and the relation between the slowness perturbation and the travel time surface is not analytic.

Taking the inner product of this expression with  $\hat{\mathbf{t}}_0$  gives  $(\hat{\mathbf{t}}_0 \cdot \nabla T_0) = u_0$ , which can also be written as

$$\frac{dT_0}{ds_0} = u_0 . \quad (134)$$

In this expression  $d/ds_0 = \hat{\mathbf{t}}_0 \cdot \nabla$  is the derivative along the unit vector  $\hat{\mathbf{t}}_0$ . Expression (134) can be integrated to give

$$T_0 = \int_{\mathbf{r}_0[u_0]} u_0(\mathbf{r}_0) ds_0 , \quad (135)$$

where  $\mathbf{r}_0$  is the position of the ray in the reference slowness field.

Using (133), expression (132) for the travel time perturbation can be written as  $(\hat{\mathbf{t}}_0 \cdot \nabla T_1) = u_1$ , but since  $\hat{\mathbf{t}}_0 \cdot \nabla = d/ds_0$  this can also be written as

$$\frac{dT_1}{ds_0} = u_1 . \quad (136)$$

This expression can be integrated to give

$$T_1 = \int_{\mathbf{r}_0[u_0]} u_1(\mathbf{r}_0) ds_0 . \quad (137)$$

The main point of this derivation is that the integration in (137) is along the *reference ray*  $\mathbf{r}_0$  rather than along the true ray in the perturbed medium. Expression (137) constitutes a linearized relation between the travel time perturbation  $T_1$  and the slowness perturbation  $u_1$ . When one divides the model in cells where one assumes the slowness perturbation is constant, then the discretized form of (137) can be written as

$$\delta T_i = \sum_j L_{ij} u_j . \quad (138)$$

In this expression, the subscript  $i$  labels the different travel times that are used in the inversion while  $j$  is the cell index. It follows from figure 16 that  $L_{ij}$  is the length of ray  $\#i$  through cell  $\#j$ . Equation (138) forms a linear system of equations as discussed in section 2.

The matrix  $L_{ij}$  is in general very sparse for tomographic problems because every ray intersects only a small fraction of the cells, see figure 16. This is in particular the case in three dimensions where the relative number of intersected cells is much smaller than in two dimensions. This makes it particularly attractive to use iterative solutions for the linear system of equations as presented in section 2.9.2. It follows from expression

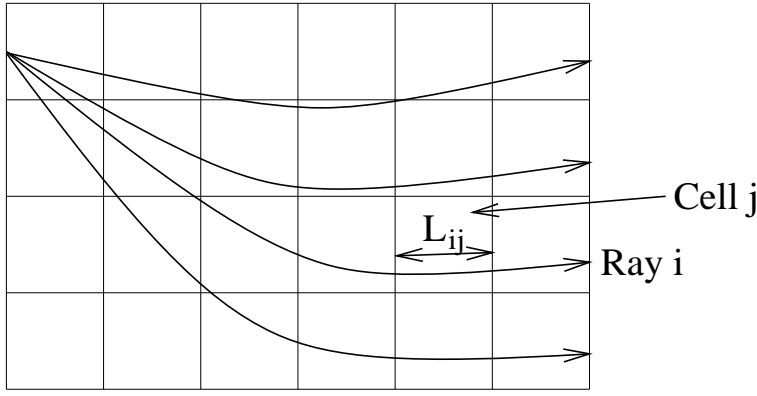


Figure 16. Diagram of a tomographic experiment.

(71) that in this iterative approach the solution is constructed by multiplication with the matrices  $\mathbf{L}$  and  $\mathbf{L}^T$ , but that the inverse  $(\mathbf{L}^T \mathbf{L})^{-1}$  (which is not sparse) is not needed. In such an approach there not even a need to store the matrix  $\mathbf{L}$ , it is much efficient to store only the relatively few nonzero matrix elements and keep track of the indices of these elements. This approach has been developed by *Nolet* [1985] and has been used successfully for extremely large-scale tomographic inversion for the interior of the Earth's mantle [e.g. *Spakman et al.*, 1993; *van der Hilst et al.*, 1997].

It should be remarked that there is no strict need to use cells to discretize the slowness (and the same holds for the treatment of linearized waveform inversion in the sections 4.2 and 4.3). As an alternative the slowness perturbation can be expanded in a finite set of basis functions  $B_j(\mathbf{r})$ :

$$u_1(\mathbf{r}) = \sum_{j=1}^L m_j B_j(\mathbf{r}) . \quad (139)$$

Inserting this in (137) one again arrives at a linear system of the form (138), but the matrix elements are now given by

$$L_{ij} = \int_{ref\ ray\ j} B_j(\mathbf{r}) ds_0 , \quad (140)$$

where the integration now is along reference ray  $\#j$ . However, one should realize that when the basis functions have global character, such as spherical harmonics, the matrix  $\mathbf{L}$  is in general not sparse with this model parameterization. This means one cannot fully exploit the computational efficiency of iterative solutions of linear systems.

## 6.2. SURFACE WAVE TOMOGRAPHY

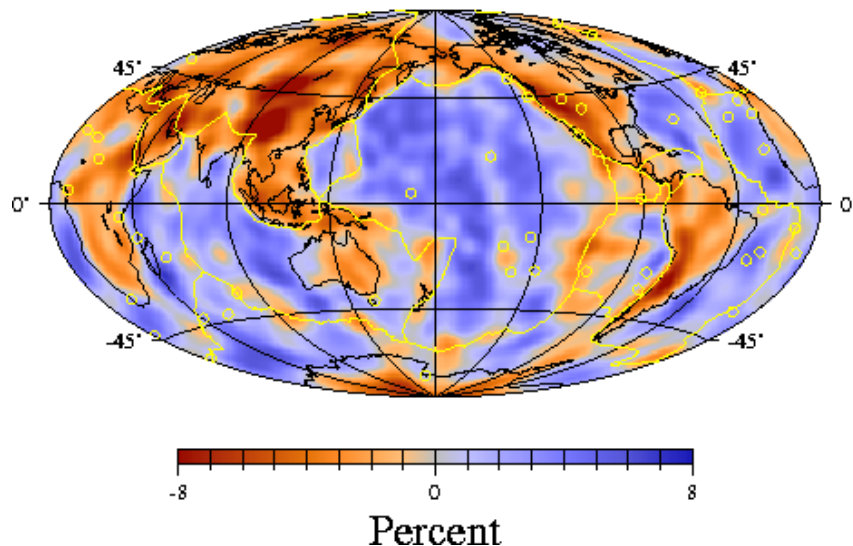
For most people Fermat's principle goes hand in hand with travel time tomography, but surface wave tomography actually relies on it since the early twenties when *Gutenberg* [1924], using data collected by *Tams* [1921], explained the dispersion differences between surface waves propagating along continental and oceanic paths in terms of properties of the Earth's crust.

Surface wave tomography clearly means tomography based on surface waves. There are several ways of doing this. One may directly invert the waveform for structural parameters, provided the sensitivities are known. Generally, the structural parameters are non-linearly related to the waveforms and one way (among others) to linearize the problem has been outlined in section 4.3. A more classical approach, consists of the so-called two-step inversion where one first constructs models of phase or group velocity as a function of frequency. The information contained in these maps is then inverted for depth structure. This is possible because expression (126) gives a linear relation between the phase velocity as a function of frequency and the perturbation of the medium as a function of depth. The mapping from phase velocity to depth then reduces to a linear inverse problem to which the techniques of section 2 can be applied.

It is this first step which commonly assumes Fermat's principle. If we assume that the Earth structure is sufficiently smooth compared to the wavelength of a propagating wave, locally any wave of frequency  $\omega$  may be approximated by a plane wave. We are then in the realm of ray theory and are able to measure the mean phase (or group) velocity along the ray path. If furthermore Fermat's principle holds, the measurements correspond to path integrals along the minor and major arc segments of the great circle connecting station and event. It is undoubtedly the case that there are many examples of off-great-circle propagation and non-theoretical effects. The extent to which such effects corrupt the models constructed under the simple assumptions of Fermat's principle can only be determined by numerical modelling, and is still an open question.

It is commonly accepted that Fermat's principle is valid for surface wave periods longer than 150 seconds roughly, and most work until the mid-nineties concentrated on these longer periods. *Trampert and Woodhouse* [1995, 1996] extend these classical methods in two ways. Firstly, they applied it to much shorter periods in the range of 40 to 150 seconds. Secondly, they exploited the by now tremendous wealth of assembled digital data in a systematic fashion by developing fully automatic methods for extracting path-averaged phase velocities of surface waves from recorded waveforms. The most original feature of these studies is the phase ve-

## Love 40 seconds



*Figure 17.* Love wave phase velocity perturbations at a period of 40 seconds. The variations are given in percent with respect to an average reference model. Yellow lines represent plate boundaries and yellow circles are hotspots.

locity maps at short periods which gives new global information on the uppermost structure of the Earth. As an example, the model for the phase velocity of Love waves at a period of 40 seconds is shown in figure 17. The Love-wave model shows a remarkable correlation with surface topography and bathymetry, and hence with crustal thickness. This is too be expected because these waves sample the surface with an exponentially decaying sensitivity. The reference model used in the inversion has a crustal thickness of 24.4 km. If the true Earth has a thicker crust, mantle material of the reference model is replaced with slower crustal material resulting in a slow velocity perturbation. This is the case for the continents. In oceans, the opposite takes place. The sensitivity of 40 second Love waves decays very rapidly with depth and hence tectonic signatures from the uppermost mantle are less pronounced. Nevertheless, a slow velocity anomaly is observed along most oceanic ridges, where hot, and hence slow, material is put into place.

## 7. Nonlinearity and ill-posedness

Nonlinearity complicates the estimation problem considerably. In many practical problems, nonlinear inversion is treated as a nonlinear optimization problem where a suitably chosen measure of the data misfit is reduced as a function of the model parameters. There is a widespread belief in the inverse problem community that the dominant effect of nonlinearity is the creation of secondary minima in the misfit function that is minimized. This point of view is overly simplistic. Nonlinearity affects both the estimation problem and the appraisal problem. In sections 7.1 and 7.2 it is shown how non-linearity can be a source of ill-posedness of inverse problems. In section 8 the appraisal problem for nonlinear inverse problems is discussed.

### 7.1. EXAMPLE 1, NON-LINEARITY AND THE INVERSE PROBLEM FOR OF THE SCHRÖDINGER EQUATION

The inverse problem of the estimation of a quantum mechanical potential in one dimension from the measurement of the reflection coefficient of waves reflected by the potential is an interesting tool for studying nonlinear inversion because the inverse problem has a stunningly simple solution [Marchenko, 1955; Burridge, 1980]. This inverse problem is of direct relevance in the earth sciences; both the inverse problem of geomagnetic induction [Weidelt, 1972] as well as the seismic reflection problem [Burridge, 1980; Newton, 1981] can be reformulated similar to the problem treated in this section. For the Schrödinger equation the wavefield  $\psi$  satisfies the following differential equation:

$$\psi_{xx} + (k^2 - V(x)) \psi = 0 . \quad (141)$$

Let the reflection coefficient after a Fourier transform to the time domain be denoted by  $R(t)$ . The potential  $V(x)$  at a fixed location  $x$  follows from the reflection coefficient by solving the Marchenko equation:

$$K(x, t) + R(x + t) + \int_{-t}^x K(x, \tau) R(\tau + t) d\tau = 0 , \quad (142)$$

for  $K(x, t)$  and carrying out a differentiation:

$$V(x) = -2 \frac{dK(x, x)}{dx} . \quad (143)$$

In figure 18 the results of a numerical solution of the Marchenko equation is shown. The potential is shown by a thick solid line. For this

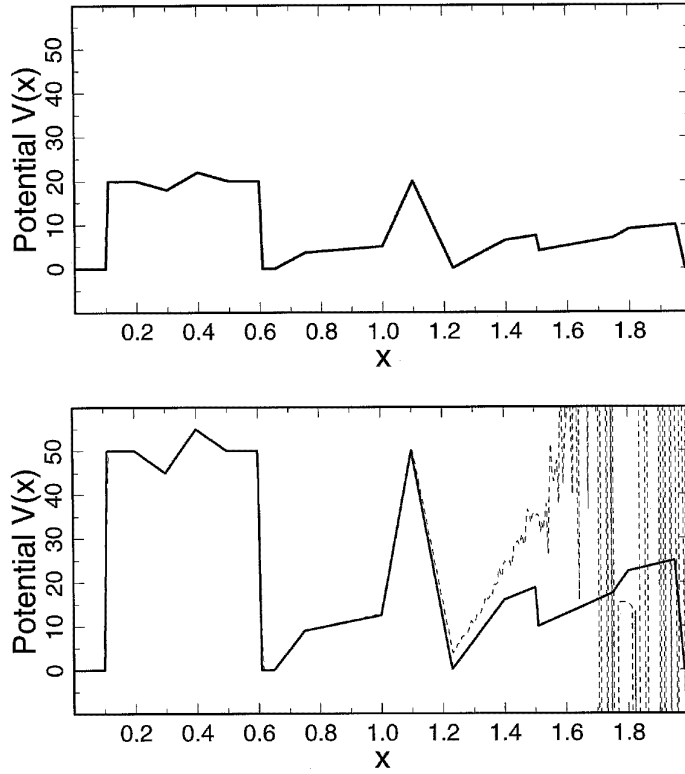


Figure 18. The original potential (solid lines) and the potential reconstructed with the Marchenko equation (dashed lines). The potential in the bottom panel is 2.5 times as strong as the potential in the top panel.

potential the reflection coefficient is computed and the potential is reconstructed by numerically solving the Marchenko equation (142) and carrying out the differentiation (143), details of the calculation can be found in [Dorren *et al.*, 1994]. The reconstructed potential of figure 18 (top panel) is on the scale of the figure indistinguishable from the true potential. In figure 18 (bottom panel) the same synthetic experiment is shown, the only difference is that the potential has been multiplied with a factor 2.5. If the problem would be linear, the reflection coefficients would be 2.5 times as strong as for the potential in the top panel and the reconstructed potential would also be 2.5 times as strong. This would lead to a near-perfect reconstruction of the potential. However, the reconstructed potential is given by the dashed line. It can be seen that the left part of the potential (the side from which the waves are incident) the potential is reconstructed quite well, but that the part of the potential on the right is very poorly reconstructed. According to the reasoning above, the po-



tential would have been reconstructed quite well if the problem had been linear. This implies that the instability in the reconstructed potential is due to the non-linearity of the problem.

The physical reason for this instability can relatively easily be understood. It follows from the Schrödinger equation (141) that the effective wavenumber is given by  $\sqrt{k^2 - V(x)}$ . When  $k^2 < V(x)$  the wavenumber is complex which reflects the fact that the waves are evanescent when the potential energy is larger than the total energy. In that case the wavefield decays exponentially within the potential. For a fixed energy  $k^2$  the wavefield is more evanescent for the potential in the bottom panel of figure 18 than for the potential in the top panel of that figure, simply because the potential energy is 2.5 times as high. This has the result that the wavefield penetrates further in the potential in the top panel than in the potential in the bottom panel of figure 18. Obviously, the potential in a certain region is not constrained by the recorded wavefield if the wavefield does not sample the potential in that region. The strong evanescence of the waves in the potential in the bottom panel of figure 18 implies that that parts of that potential are effectively not sampled by the waves. In that case, the numerical details of the algorithm, including the numerical round-off error determine the reconstructed potential in that region. The essential point is that the values of the model parameters affect the way in which the wavefield interrogates the model.

Physically, the instability in the reconstructed potential in figure 18 can thus be understood. However, what does this imply for the inverse problem? What happens physically if the potential on the left side is high, is that the wavefield is prevented from sampling the potential on the right part. This means that for some values of the model parameters (that describe how high the potential is on the left), other model parameters (that describe the potential on the right) are unconstrained by the data. In terms of a misfit function this implies that the misfit does not depend on the model parameters that describe the right side of the potential (when the left side of the potential is high). In other words, the misfit function does not have a minimum, but it has a broad plateau. Note that as an additional complexity this only occurs for certain values of other model parameters (that describe how high the potential is on the left).

## 7.2. EXAMPLE 2, NON-LINEARITY AND SEISMIC TOMOGRAPHY

A second example of the ill-posedness introduced by the non-linearity in inverse problems is seismic tomography. Consider a cross-borehole tomographic experiment where rays travel from a source in a well to a string of receivers on another well. The case where the velocity is homogeneous is

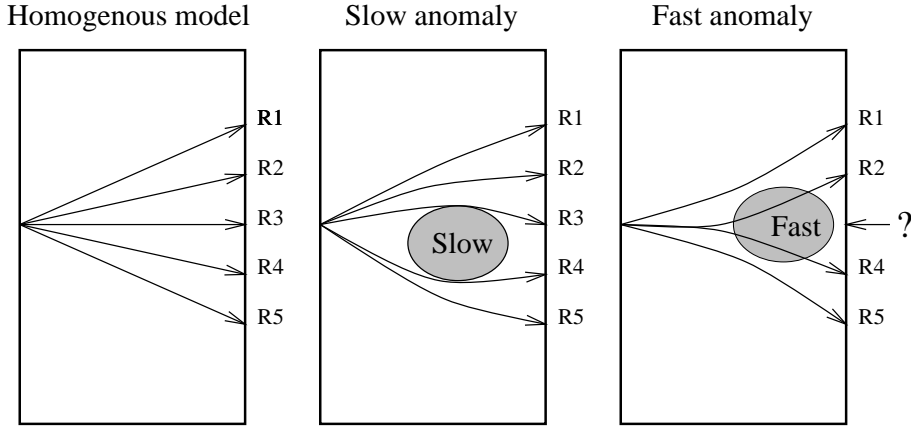


Figure 19. (a) Tomographic experiment where the velocity is homogeneous and the rays are straight. (b) Tomographic experiment where the rays curve around a low-velocity body. (c) Tomographic experiment where a high-velocity anomaly causes a shadow zone at the middle receiver.

shown in the left panel of figure 19. In that case the rays are straight lines that travel from the source on the left to receivers  $R1$  through  $R5$  on the right. If the velocity is not homogeneous, the rays are curved. This implies that the value of the model parameters determine the way in which the rays interrogate the model.

Suppose that a low-velocity anomaly is present, see the middle panel of figure 19. Since the rays of first arrivals are curves of minimal travel time, the rays curve around the slow-velocity anomaly. If the velocity anomaly is sufficiently slow, all the rays may completely curve around the slow-velocity. In that situation the anomaly would not be sampled by any rays and the velocity within the anomaly cannot be determined because it is not sampled by any rays. At best one could derive an upper bound for the velocity within the anomaly. Just as in the previous section, the model parameters affect the way in which the probe (in this case the rays) samples the model. In terms of the misfit function this implies that for a certain range of model parameters, the misfit function does not depend on the model parameters at all, in other words: the misfit function is completely flat over an extended range of parameter values.

Let us now consider the opposite situation where a high-velocity anomaly is present, see the right panel of figure 19. Rays will be defocused by a high velocity body and a shadow-zone is formed behind the anomaly. This may mean that there is no ray that will hit the receiver  $R3$ , see the question mark in figure 19. This means that for some values of the model parameters it is impossible to compute the data because the travel time

cannot be determined when there is no ray that hits the receiver. This means that for some values of the model parameters it is impossible to compute the corresponding data values given the theory that one uses. In this sense one can say that some values of the model parameters are “forbidden.” It should be noted that this is not a complexity created in some exotic thought-experiment; the problem of the “missing rays” is a real problem in seismic travel time tomography [Sambridge, 1990], as well as the tomographic imaging of the temperature field in ovens [Natterer *et al.*, 1997].

The critical reader might remark at this point that the fact that for certain values of the velocity model there are no rays hitting the receiver is due to the fact that ray theory is an approximation to a true theory (wave theory), and that wave theory predicts that some energy will diffract into the shadow-zones. Although this is correct, one should also note that the wavefield in the shadow zones is very weak (that is why they are called shadow zones) and that in practice this diffracted wavefield is usually not detectable.

## 8. Model appraisal for nonlinear inverse problems

In the previous section the effect of non-linearity on model estimation has been discussed. In this section an attempt is made to describe the effect of non-linearity on the appraisal problem where one wants to describe how the estimated model is related to the true model (see figure 2). However, it should be stressed that there is presently no general theory to deal with the appraisal problem for a truly nonlinear inverse problem with infinitely many degrees of freedom. In practice, one often linearizes the problem around the estimated model and then uses linear theory to make inferences about the resolution and reliability of the estimated model. The lack of a general theory for the appraisal problem should be seen as a challenge for theorists! In this section three attempts are described to carry out model assessment for a nonlinear inverse problem. These attempts follow the lines of formal theory (section 8.1), a numerical approach (section 8.2) and a pragmatic approach (section 8.3).

### 8.1. NONLINEAR BACKUS-GILBERT THEORY

Linear Backus-Gilbert theory is based on the equations (72)-(74) for linear inverse problems for continuous models. The model estimate  $\tilde{m}(x)$  at location  $x$  is constructed by making the linear superposition (73) of the data. The resolution kernel in equation (74) specifies the relation between the estimated model and the true model. In the ideal case the resolution kernel is a delta function. *Backus and Gilbert* [1967, 1968] have shown that

the criterion that the resolution kernel should resemble a delta function as much as possible can be used to determine the coefficients  $a_i(x)$  in equation (73) which prescribes how a datum  $d_i$  affects the estimated model at location  $x$ .

Backus-Gilbert theory has been generalized by *Snieder* [1991] for the special case that the forward problem can be written as a perturbation series:

$$d_i = \int G_i^{(1)}(x)m(x)dx + \iint G_i^{(2)}(x_1, x_2)m(x_1)m(x_2)dx_1dx_2 + \dots \quad (144)$$

In a number of applications such a perturbation series arises naturally. Important examples are the Neumann series (101) in scattering theory where the scattering data are written as a sum of integrals that contain successively higher powers of the model perturbation, or ray perturbation theory where the travel time of rays is written as a sum of integrals with increasing powers of the slowness perturbation [e.g. *Snieder and Sambridge*, 1993; *Snieder and Aldridge*, 1995]. When the forward problem is nonlinear, the inverse problem is non-linear as well, this suggest that for a non-linear inverse problem the linear estimator (73) should be generalized to include terms that are nonlinear in the data as well:

$$\tilde{m}(x) = \sum_i a_i^{(1)}(x)d_i + \sum_{i,j} a_{ij}^{(2)}(x)d_id_j + \dots \quad (145)$$

The key of non-linear Backus-Gilbert theory is to insert the expansion of the forward problem (144) in the estimator (145). The result can then be written as:

$$\tilde{m}(x) = \int R^{(1)}(x; x_1)m(x_1)dx_1 + \iint R^{(2)}(x; x_1, x_2)m(x_1)m(x_2)dx_1dx_2 + \dots \quad (146)$$

This expression generalizes the linear resolution kernel of equation (74) to nonlinear inverse problems. The kernel  $R^{(1)}(x; x_1)$  describes to what extent the estimated model is a blurred version of the true model. The higher order kernels such as  $R^{(2)}(x; x_1, x_2)$  can be interpreted as nonlinear resolution kernels that describe to what extent there is a spurious nonlinear mapping from the estimated model onto the true model in the inversion process.

In the ideal case, the estimated model is equal to the true model:  $\tilde{m}(x) = m(x)$ . This is the case when the linear resolution kernel  $R^{(1)}(x; x_1)$  is a delta function  $\delta(x - x_1)$  and when the nonlinear resolution kernels are equal to zero:  $R^{(n)}(x; x_1, \dots, x_n) = 0$  for  $n \geq 2$ . However, as in equation (75) the linear resolution kernel  $R^{(1)}(x; x_1)$  can be written as a sum of a finite amount of data kernels  $G^{(1)}(x_1)$ . Since a delta function cannot

be obtained by summing a finite amount of smooth functions, the linear resolution kernel can never truly be a delta function. This reflects the fact that with a finite amount of data the estimated model will be a blurred version of the true model. *Snieder* [1991] treats the inverse problem of the determination of the mass-density of a vibrating string from the eigenfrequencies of the string. He shows that if only a finite amount of eigenfrequencies are available, the nonlinear resolution kernels cannot be zero. This implies that the finiteness of the data set leads to a spurious nonlinear mapping from the true model to the estimated model. This can be related to the symmetries of the problem and to mode coupling “off the energy shell.” The finite width of the linear resolution kernel and the fact that the nonlinear resolution kernels are nonzero imply not only that the estimated model is a blurred version of the true model, but also that the estimated model is biased. The reader is referred to *Snieder* [1991] for details. In that work it is also described how the coefficients  $a^{(i)}$  in the estimator (145) can be determined.

Although nonlinear Backus-Gilbert theory is a new tool for dealing with the assessment problem for nonlinear inverse problems, one should realize that the theory can only be applied to weakly nonlinear problems where a (very) few orders are sufficient for an accurate description of both the forward and the inverse problem. In addition, the theory of *Snieder* [1991] is so complex that a reformulation is needed to make the theory applicable to the large-scale inverse problems that are being treated in practice.

It is important to realize that any regular (nonlinear) mapping from data  $d_i$  to an estimated model  $\tilde{m}(x)$  can be written in the form of expression (145). The details of the employed algorithm then determine the coefficients  $a_{i_1 \dots i_n}^{(n)}$ . The resolution analysis shown in equation (146) and the subsequent discussion therefore is applicable to the estimated model. The conclusions concerning the linear and nonlinear resolution kernels can thus be used for *any* regular mapping from the data to the estimated model.

## 8.2. GENERATION OF POPULATIONS OF MODELS THAT FIT THE DATA

Another approach to assess the reliability of estimated models is to generate not a single model that fit the data within a certain tolerance but to obtain an ensemble of models that fits the data within a certain tolerance [e.g. *Lomax and Snieder, 1995*]. An alternative approach is to compute the misfit for a very large class of models and to use the data fit, possibly in combination with Bayesian statistics to make inferences about the range

of models that explain the data in a certain sense [e.g. *Mosegaard and Tarantola*, 1995; *Gouveia and Scales*, 1998, *Mosegaard*, 1998]. Obviously, this approach requires a numerical approach to create such ensembles, but with present day computers significance progress has been made.

An important concept in the generation of ensembles of models is the randomness in the search method that one employs. A descent method contains no element of randomness whatsoever, whereas a Monte Carlo search where one randomly samples model space is completely random. In between are algorithms which have both, a random component as well as a mechanism to prefer models that fit the data well. Examples of such algorithms are simulated annealing [*Krikpatrick et al.*, 1983; *Rothman*, 1985] or genetic algorithms [*Sambridge and Drijkoningen*, 1992; *Sen and Stoffa*, 1992; *Lomax and Snieder*, 1995]. A promising technique is the adaptive search [*Mosegaard and Tarantola*, 1995] where in the process of carrying out a random search, information about the misfit function is being built up, and this misfit function is then used to drive the random search in an intelligent way.

The main merit of algorithms that determine an ensemble of models together with information on how well each model explains the data is that this ensemble can be used to make inferences about the model. These inferences may or may not be statistical. The possibility to analyse such ensembles of models in a meaningful way has not yet been fully exploited.

An example is shown from a study of *Douma et al.* [1996]. In their study synthetic group velocity data of the fundamental mode Rayleigh wave that propagates along the earth's surface were computed for periods between 10 and 300 s. The true velocity model is shown as the dashed line in figure 20 as a function of depth. A Monte Carlo search was used to find models of the  $S$ -velocity that were consistent with the data within a realistic measurement error. The resulting population of models is shown in figure 20. In this study, the  $S$ -velocity was deliberately over-parameterized. As a result, the resulting models are highly oscillatory and contain strong trade-offs. The aim of the study of *Douma et al.* [1996] was to extract the robust features of the velocity model from the ensemble of models shown in figure 20. This was achieved by computing "Empirical Orthogonal Functions" (EOF's) from this ensemble. These functions give the patterns with the different degrees of variability within an ensemble.

The EOF's can be used to re-parameterize the model space in an intelligent way that reflects how well model perturbations are constrained by the data. Alternatively, the EOF's could be used to carry out a statistical analysis of the ensemble of models that explains the data. However, as noted by *Douma et al.* [1996] the EOF techniques is only useful for inverse problems that are weakly nonlinear.

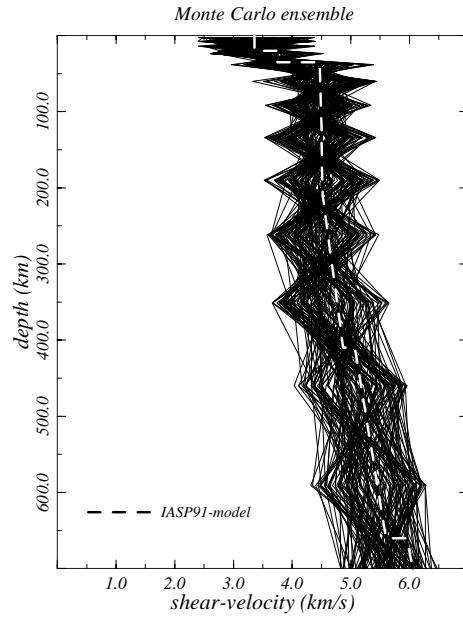


Figure 20. The true velocity model (broken curve) and an ensemble of models generated by a Monte Carlo search that fit surface wave group velocity data with a realistic tolerance (solid lines).

### 8.3. USING DIFFERENT INVERSION METHODS

In the previous sections, a theoretical and a numerical method for model appraisal were presented. Apart from these more formal approaches, “common sense” is a powerful tool for carrying out model assessment. An important way to assess the reliability of a model is to determine the model in different ways. In the ideal case, different data sets are used by different research groups who use different theories to estimate the same properties. The agreement or disagreement between these models can be used as an indicator of the reliability of these models. It is admittedly difficult to quantify the sense of reliability that is thus obtained, but in the absence of an adequate theory to carry out model assessment for nonlinear inverse problems (and remember that we don’t have such a theory) this may be the best approach.

An example of this approach is shown in figure 21. Nonlinear waveform inversion using the Partitioned Waveform Inversion of *Nolet* [1990] has been used by *van der Hilst and Kennett* [1997] to determine a three-dimensional model of the  $S$ -velocity under Australia. A cross section of the model at a depth of 140 km is shown by the colors in figure 21. In such an inversion, averages of the earth structure along long paths are

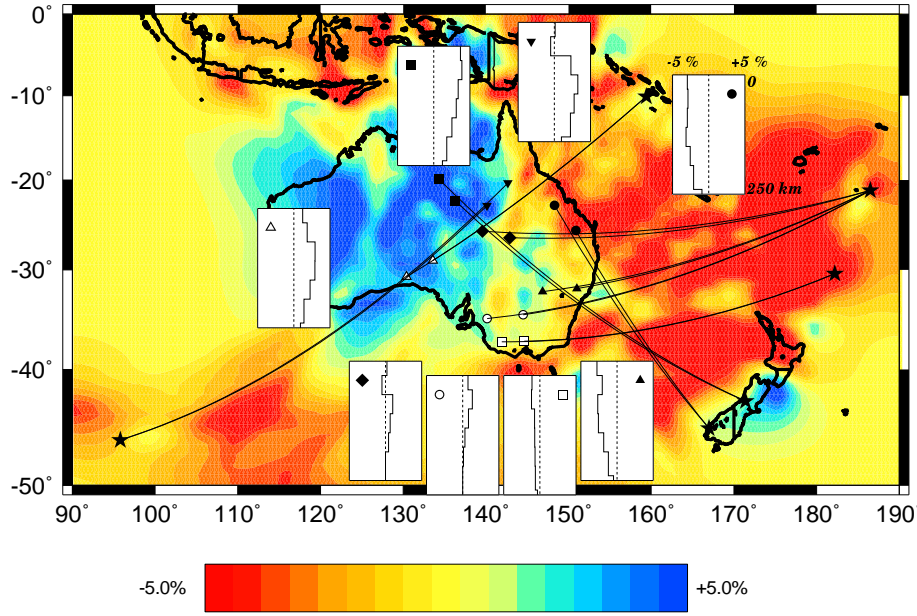


Figure 21. The  $S$ -velocity anomaly under Australia at a depth of 140 km obtained by surface wave tomography by van der Hilst and Kennett [1997] where the relative velocity perturbation is shown by the colour bar at the bottom. The interstation  $S$ -velocity obtained by *Passier et al.* [1997] is shown in the panels. The  $S$ -velocity in panels ranges from -5% to +5% and is shown from the surface down to a depth of 250 km. The employed stations for each panel are indicated that corresponds with the symbol in the upper-right corner of each panel.

used to determine the local variations in the  $S$ -velocity. Effects such as errors in the source mechanisms of the employed earthquakes, seismic anisotropy that is not accounted for in the inversion, variations in the path coverage and the path direction over the domain of inversion can lead to artifacts in the inversion. As a way of verifying the reliability of the solution, *Passier et al.* [1997] determined the local  $S$ -velocity using the waveform inversion method described by *Passier and Snieder* [1995]. In this method, a consistency requirement of the waveforms recorded in two nearby stations is used to estimate the horizontally averaged velocity between the stations as a function of depth. Just as the technique used by *van der Hilst and Kennett* [1997] this technique may lead to models with artifacts, but in general these artifacts will be different, see *Passier et al.* [1997] for a discussion.

The interstation  $S$ -velocity models are shown in the panels in figure 21. In these panels the depth range is from the surface down to a depth of 250 km, whereas the  $S$ -velocity ranges from -5% to +5%. The symbol in each panel (circle, triangle, etc.) corresponds to the employed stations



that are marked on the map with the same symbol. The depth-profile in each panel shows the  $S$ -velocity between the corresponding stations. By comparing the interstation velocity models shown in the panels with the velocity model shown in colors one can see that the locally determined  $S$ -velocity agrees well with the  $S$ -velocity obtained from the 3D tomographic inversions. Consider for example the  $S$ -velocity between the stations in the northeast marked with solid dots that is shown in the panel with the solid dot. At a depth of about 140  $km$  the interstation shear velocity is about -2.5%. This agrees well with the  $S$ -velocity obtained from 3D tomography that is indicated with the colors.

A comparison as presented in this section is very useful for assessing the reliability of models in a qualitative way. In this sense, a healthy competition between different research groups is an important ingredient in the appraisal of models. There is one pitfall in the approach to compare models obtained from different data sets and different inversion methods; it is possible that both models are in error but that they agree with each other. The comparison of models may thus provide a false impression of their reliability. However, in practice this test often is a useful pointer to artefacts in (at least one of) the models.

## 9. Epilogue

Linear inverse problem theory is an extremely powerful tool for solving inverse problems. Much of the information that we currently have on the Earth's interior is based on linear inverse problems. The success of modern-day oil exploration for an affordable price is to a large extent possible because of our ability to use imaging techniques that are based on single scattering theory to map oil reservoirs. One could argue that the modern world, which heavily relies on cheap access to hydrocarbons, might have a drastically different form if single scattering would not explain seismic data so well.

Despite the success of linear inverse theory, one should be aware that for many practical problems our ability to solve inverse problems is largely confined to the estimation problem. This may surprise the reader, because for linear inverse problems the resolution kernel and the statistics of linear error propagation seem to be adequate tools for carrying out the model assessment. However, as noted in section 2, in order to compute the resolution kernel one needs to explicitly define the generalized inverse  $\mathbf{A}^{-g}$ . For many problems, this entails inverting the matrix  $\mathbf{A}^T \mathbf{A}$  or an equivalent matrix that contains regularization terms as well. This task is for many important problem, such as imaging seismic reflection data, not feasible.

The results of section 8 show that for nonlinear inverse problems both the estimation problem and the appraisal problem are much more difficult. In the context of a data-fitting approach, the estimation problem amounts to a problem of nonlinear optimization. Researchers can use the results of this field of mathematics to solve the estimation problem. However, the reader should be aware of the fact that there is presently no theory to describe the appraisal problem for nonlinear inverse problems. Developing a theoretical framework for the appraisal of models obtained from nonlinear inversion of data is a task of considerable complexity which urgently needs the attention of the inverse problem community.

**Acknowledgments:** Discussions with John Scales and Jean-Jacques Lévêque have helped to clarify a number of issues, we appreciate their input very much.

## References

1. Aldridge, D.F., Linearization of the eikonal equation, *Geophysics*, 59, 1631-1632, 1994.
2. Alsina, D., R.L. Woodward, and R.K. Snieder, Shear-Wave Velocity Structure in North America from Large-Scale Waveform Inversions of Surface Waves, *J. Geophys. Res.*, 101, 15969-15986, 1996.
3. Aki, K., and P.G. Richards, *Quantitative Seismology (2 volumes)*, Freeman and Co. New York, 1980.
4. Backus, G., and J.F. Gilbert, Numerical applications of a formalism for geophysical inverse problems, *Geophys. J.R. Astron. Soc.*, 13, 247-276, 1967.
5. Backus, G., and J.F. Gilbert, The resolving power of gross earth data, *Geophys. J.R. Astron. Soc.*, 16, 169-205, 1968.
6. Backus, G. E. and F. Gilbert, Uniqueness in the inversion of inaccurate gross earth data, *Philos. Trans. R. Soc. London, Ser. A*, 266, 123-192, 1970.
7. Ben-Menahem, A. and S.J. Singh, *Seismic waves and sources*, Springer Verlag, New York, 1981.
8. Borg, G., Eine Umkehrung der Sturm-Liouvillischen Eigenwertaufgabe, Bestimmung der Differentialgleichung durch die Eigenwerte, *Acta Math.*, 78, 1-96, 1946.
9. Burridge, R., The Gel'fand-Levitan, the Marchenko and the Gopinath-Sondi integral equations of inverse scattering theory, regarded in the context of the inverse impulse response problems, *Wave Motion*, 2, 305-323, 1980.
10. Cara, M., Regional variations of higher-mode phase velocities: A spatial filtering method, *Geophys. J.R. Astron. Soc.*, 54, 439-460, 1978.
11. Claerbout, J.F., *Fundamentals of Geophysical data processing*, McGraw-Hill, New York, 1976.
12. Claerbout, J.F., *Imaging the Earth's interior*, Blackwell, Oxford, 1985.
13. Clayton, R. W. and R. P. Comer, A tomographic analysis of mantle heterogeneities from body wave travel time data, *EOS, Trans. Am. Geophys. Un.*, 64, 776, 1983.
14. Constable, S.C., R.L. Parker, and C.G. Constable, Occam's inversion: a practical algorithm for generating smooth models from electromagnetic sounding data, *Geophysics*, 52, 289-300, 1987.
15. Dahlen, F.A., and J. Tromp, *Theoretical global seismology*, Princeton University

- Press, Princeton, 1998.
16. Dorren, H.J.S., E.J. Muzert, and R.K. Snieder, The stability of one-dimensional inverse scattering, *Inverse Problems*, 10, 865-880, 1994.
  17. Douma, H., R. Snieder, and A. Lomax, Ensemble inference in terms of Empirical Orthogonal Functions, *Geophys. J. Int.*, 127, 363-378, 1996.
  18. Dziewonski, A.M., and D.L. Anderson, Preliminary Reference Earth Model, *Phys. Earth. Plan. Int.*, 25, 297-356, 1981.
  19. Gerver, M.L. and V. Markushevitch, Determination of a seismic wave velocity from the travel time curve, *Geophys. J. Royal astro. Soc.*, 11 165-173, 1966.
  20. Gilbert, F., Ranking and winnowing gross Earth data for inversion and resolution, *Geophys. J. Royal astro. Soc.*, 23 125-128, 1971.
  21. Gouveia, W.P., and J.A. Scales, Bayesian seismic waveform inversion: parameter estimation and uncertainty analysis, *J. Geophys. Res.*, 103, 2759-2779, 1998.
  22. Gutenberg, B., Dispersion und Extinktion von seismischen Oberflächenwellen und der Aufbau der obersten Erdschichten, *Physikalische Zeitschrift*, 25, 377-382, 1924.
  23. Herglotz, G. Über das Benndorfsche Problem des Fortpflanzungsgeschwindigkeit der Erdbebenstrahlen, *Zeitschrift für Geophys.*, 8 145-147, 1907.
  24. Keller, J.B., I. Kay, and J. Shmoys, Determination of a potential from scattering data, *Phys. Rev.*, 102, 557-559, 1956.
  25. Kirkpatrick, S., C. Gelatt, and M.P. Vecchi, Optimization by simulated annealing, *Science*, 220, 671-680, 1983.
  26. Lanczos, C., *Linear Differential Operators*, Van Nostrand, London, 1961.
  27. Levenberg, K., A method for the solution of certain nonlinear problems in least squares, *Quart. Appl. Math.*, 2, 164-168, 1944.
  28. Lomax, A., and R. Snieder, The contrast in upper-mantle shear-wave velocity between the East European Platform and tectonic Europe obtained with genetic algorithm inversion of Rayleigh-wave group dispersion, *Geophys. J. Int.*, 123, 169-182, 1995.
  29. Marchenko, V.A., The construction of the potential energy from the phases of scattered waves, *Dokl. Akad. Nauk*, 104, 695-698, 1955.
  30. Matsu'ura M. and N. Hirata, Generalized least-squares solutions to quasi-linear inverse problems with a priori information, *J. Phys. Earth*, 30, 451-468, 1982.
  31. Mayer, K., R. Marklein, K.J. Langenberg and T.Kreutter, Three-dimensional imaging system based on Fourier transform synthetic aperture focussing technique, *Ultrasonics*, 28, 241-255, 1990.
  32. Menke, W., *Geophysical data analysis: discrete inverse theory*, Academic Press, San Diego, 1984.
  33. Merzbacher, E., *Quantum mechanics (2nd ed.)*, Wiley, New York, 1970.
  34. Montagner, J.P., and H.C. Nataf, On the inversion of the azimuthal anisotropy of surface waves, *J. Geophys. Res.*, 91, 511-520, 1986.
  35. Mosegaard, K., Resolution analysis of general inverse problems through inverse Monte Carlo sampling, *Inverse Problems*, 14, 405-426, 1998.
  36. Mosegaard, K., and A. Tarantola, Monte Carlo sampling of solutions to inverse problems, *J. Geophys. Res.*, 100, 12431-12447, 1995.
  37. Muzert, E., and R. Snieder, An alternative parameterization for surface waves in a transverse isotropic medium, *Phys. Earth Planet Int. (submitted)*, 1999.
  38. Natterer, F., H. Sielschott, and W. Derichs, Schallpyrometrie, in *Mathematik - Schlüsseltechnologie für die Zukunft*, edited by K.H. Hoffmann, W. Jäger, T. Lochmann and H. Schunk, 435-446, Springer Verlag, Berlin, 1997.
  39. Newton, R.G., Inversion of reflection data for layered media: A review of exact methods, *Geophys. J.R. Astron. Soc.*, 65, 191-215, 1981.
  40. Newton, R.G., *Inverse Schrödinger scattering in three dimensions*, Springer Verlag, Berlin, 1989.

41. Nolet, G., The upper mantle under Western-Europe inferred from the dispersion of Rayleigh wave modes, *J. Geophys.*, *43*, 265-285, 1977.
42. Nolet, G., Linearized inversion of (teleseismic) data, in *The Solution of the Inverse Problem in Geophysical Interpretation*, edited by R.Cassinis, Plenum Press, New York, 1981.
43. Nolet, G., Solving or resolving inadequate and noisy tomographic systems, *J. Comp. Phys.*, *61*, 463-482, 1985.
44. Nolet, G., Seismic wave propagation and seismic tomography, in *Seismic Tomography*, edited by G.Nolet, pp. 1-23, Reidel, Dordrecht, 1987.
45. Nolet, G., Partitioned waveform inversion and two-dimensional structure under the Network of Autonomously Recording Seismographs, *J. Geophys. Res.*, *95*, 8499-8512, 1990.
46. Nolet, G., S.P. Grand, and B.L.N. Kennett, Seismic heterogeneity in the upper mantle, *J. Geophys. Res.*, *99*, 23753-23766, 1994.
47. Nolet, G., and R. Snieder, Solving large linear inverse problems by projection, *Geophys. J. Int.*, *103*, 565-568, 1990.
48. Paige, C.G., and M.A. Saunders, LSQR: An algorithm for sparse linear equations and sparse least-squares, *ACM Trans. Math. Software*, *8*, 43-71, 1982.
49. Paige, C.G., and M.A. Saunders, LSQR: Sparse linear equations and least-squares problems, *ACM Trans. Math. Software*, *8*, 195-209, 1982.
50. Parker, R.L., *Geophysical Inverse Theory*, Princeton University Press, Princeton, New Jersey, 1994.
51. Passier, M.L., and R.K. Snieder, Using differential waveform data to retrieve local S-velocity structure or path-averaged S-velocity gradients, *J. Geophys. Res.*, *100*, 24061 - 24078, 1995.
52. Passier, M.L., and R.K. Snieder, Correlation between shear wave upper mantle structure and tectonic surface expressions: Application to central and southern Germany, *J. Geophys. Res.*, *101*, 25293-25304, 1996.
53. Passier, T.M., R.D. van der Hilst, and R.K. Snieder, Surface wave waveform inversions for local shear-wave velocities under eastern Australia, *Geophys. Res. Lett.*, *24*, 1291-1294, 1997.
54. Press, W.H., Flannery, B.P., Teukolsky, S.A. and W.T. Vetterling, *Numerical Recipes*, Cambridge University Press, Cambridge, 1989.
55. Rothman, D.H., Nonlinear inversion, statistical mechanics and residual statics estimation, *Geophysics*, *50*, 2784-2796, 1985.
56. Sabatier, P.C., Discrete ambiguities and equivalent potentials, *Phys. Rev. A*, *8*, 589-601, 1973.
57. Sambridge, M., Non-linear arrival time inversion: constraining velocity anomalies by seeking smooth models in 3-D, *Geophys. J.R. Astron. Soc.*, *102*, 653-677, 1990.
58. Sambridge, M., and G. Drijkoningen, Genetic algorithms in seismic waveform inversion, *Geophys. J. Int.*, *109*, 323-342, 1992.
59. Scales, J., and R. Snieder, To Bayes or not to Bayes?, *Geophysics*, *62*, 1045-1046, 1997.
60. Scales, J., and R. Snieder, What is noise?, *Geophysics*, *63*, 1122-1124, 1998.
61. Sen, M.K., and P.L. Stoffa, Rapid sampling of model space using genetic algorithms: examples of seismic wave from inversion, *Geophys. J. Int.*, *198*, 281-292, 1992.
62. Sluis, A. van der, and H.A. van der Vorst, Numerical solution of large, sparse linear algebraic systems arising from tomographic problems, in *Seismic tomography, with applications in global seismology and exploration geophysics*, edited by G. Nolet, Reidel, Dordrecht, 1987.
63. Snieder, R., 3D Linearized scattering of surface waves and a formalism for surface wave holography, *Geophys. J. R. astron. Soc.*, *84*, 581-605, 1986a.
64. Snieder, R., The influence of topography on the propagation and scattering of

- surface waves, *Phys. Earth Planet. Inter.*, 44, 226-241, 1986b.
65. van Heijst, H.J. and J.H. Woodhouse, Measuring surface wave overtone phase velocities using a mode-branch stripping technique, *Geophys. J. Int.*, 131, 209-230, 1997.
  66. Snieder, R., Surface wave holography, in *Seismic tomography, with applications in global seismology and exploration geophysics*, edited by G. Nolet, pp. 323-337, Reidel, Dordrecht, 1987.
  67. Snieder, R., Large-Scale Waveform Inversions of Surface Waves for Lateral Heterogeneity, 1, Theory and Numerical Examples, *J. Geophys. Res.*, 93, 12055-12065, 1988.
  68. Snieder, R., Large-Scale Waveform Inversions of Surface Waves for Lateral Heterogeneity, 2, Application to Surface Waves in Europe and the Mediterranean, *J. Geophys. Res.*, 93, 12067-12080, 1988.
  69. Snieder, R., A perturbative analysis of nonlinear inversion, *Geophys. J. Int.*, 101, 545-556, 1990.
  70. Snieder, R., The role of the Born-approximation in nonlinear inversion, *Inverse Problems*, 6, 247-266, 1990.
  71. Snieder, R., An extension of Backus-Gilbert theory to nonlinear inverse problems, *Inverse Problems*, 7, 409-433, 1991.
  72. Snieder, R., Global inversions using normal modes and long-period surface waves, in *Seismic tomography*, edited by H.M. Iyer and K. Hirahara, pp. 23-63, Prentice-Hall, London, 1993.
  73. Snieder, R., and D.F. Aldridge, Perturbation theory for travel times, *J. Acoust. Soc. Am.*, 98, 1565-1569, 1995.
  74. Snieder, R.K., J. Beckers, and F. Neele, The effect of small-scale structure on normal mode frequencies and global inversions, *J. Geophys. Res.*, 96, 501-515, 1991.
  75. Snieder, R., and A. Lomax, Wavefield smoothing and the effect of rough velocity perturbations on arrival times and amplitudes, *Geophys. J. Int.*, 125, 796-812, 1996.
  76. Snieder, R., and G. Nolet, Linearized scattering of surface waves on a spherical Earth, *J. Geophys.*, 61, 55-63, 1987.
  77. Snieder, R., and M. Sambridge, The ambiguity in ray perturbation theory, *J. Geophys. Res.*, 98, 22021-22034, 1993.
  78. Spakman, W., S. Van der Lee, and R.D. van der Hilst, Travel-time tomography of the European-Mediterranean mantle down to 1400 km, *Phys. Earth Planet. Int.*, 79, 3-74, 1993.
  79. Strang, *Linear algebra and its applications*, Harbourt Brace Jovanovich Publishers, Fort Worth, 1988.
  80. Takeuchi, H. and M. Saito, Seismic surface waves, in *Seismology: Surface waves and earth oscillations*, (Methods in computational physics, 11), Ed. B.A. Bolt, Academic Press, New York, 1972.
  81. Tams, E., 1921. Über Fortplanzungsgeschwindigkeit der seismischen Oberflächenwellen längs kontinentaler und ozeanischer Wege, *Centralblatt für Mineralogie, Geologie und Paläontologie*, 2-3, 44-52, 1921.
  82. Tanimoto, T., Free oscillations in a slightly anisotropic earth, *Geophys. J.R. Astron. Soc.*, 87, 493-517, 1986.
  83. Tarantola, A., Linearized inversion of seismic reflection data, *Geophys. Prosp.*, 32, 998-1015, 1984.
  84. Tarantola, A., *Inverse problem theory*, Elsevier, Amsterdam, 1987.
  85. Tarantola, A. and B. Valette, Inverse problems = quest for information, *J. Geophys.*, 50, 159-170, 1982a.
  86. Tarantola, A., and B. Valette, Generalized nonlinear inverse problems solved using the least squares criterion, *Rev. Geophys. Space Phys.*, 20, 219-232, 1982b.

87. Trampert, J., Global seismic tomography: the inverse problem and beyond, *Inverse Problems*, 14, 371-385, 1998.
88. Trampert, J., and J.J. L  v  que, Simultaneous Iterative Reconstruction Technique: Physical interpretation based on the generalized least squares solution, *J. Geophys. Res.*, 95, 12553-12559, 1990.
89. Trampert, J., J.J. L  v  que, and M. Cara, Inverse problems in seismology, in *Inverse problems in scattering and imaging*, edited by M. Bertero and E.R. Pike, pp. 131-145, Adam Hilger, Bristol, 1992.
90. Trampert, J., and R. Snieder, Model estimations based on truncated expansions: Possible artifacts in seismic tomography, *Science*, 271, 1257-1260, 1996.
91. Trampert, J., and J.H. Woodhouse, Global phase velocity maps of Love and Rayleigh waves between 40 and 150 seconds, *Geophys. J. Int.*, 122, 675-690, 1995.
92. Trampert, J., and J.H. Woodhouse, High resolution global phase velocity distributions, *Geophys. Res. Lett.*, 23, 21-24, 1996.
93. VanDecar, J.C., and R. Snieder, Obtaining smooth solutions to large linear inverse problems, *Geophysics*, 59, 818-829, 1994.
94. van der Hilst, R.D., S. Widiyantoro, and E.R. Engdahl, Evidence for deep mantle circulation from global tomography, *Nature*, 386, 578-584, 1997.
95. van der Hilst, R.D., and B.L.N. Kennett, Upper mantle structure beneath Australia from portable array deployments, *American Geophysical Union 'Geodynamics Series'*, 38, 39-57, 1998.
96. van Heijst, H.J. and J.H. Woodhouse, Measuring surface wave overtone phase velocities using a mode-branch stripping technique, *Geophys. J. Int.*, 131, 209-230, 1997.
97. Weidelt, P., The inverse problem of geomagnetic induction, *J. Geophys.*, 38, 257-289, 1972.
98. Wiechert, E.,   ber Erdbebenwellen. I. Theoretisches   ber die Ausbreitung der Erdbebenwellen, *Nachr. Ges. Wiss. G  ttingen*, Math.-Phys. Klasse, 415-529 1907.
99. Woodhouse, J. H. and A. M. Dziewonski, Mapping the upper mantle: Three dimensional modelling of Earth structure by inversion of seismic waveforms, *J. Geophys. Res.*, 89, 5953-5986, 1984.
100. Yilmaz, O., Seismic data processing, *Investigations in geophysics*, 2, Society of Exploration Geophysicists, Tulsa, 1987.