

# Multidimensional autoregression

*Jon Claerbout*

Occam's razor says we should try understand the world by the simplest explanation. So, how do we decompose a complicated thing into its essential parts? That's far too difficult a question, but the word "covariance" points the way. If things are found statistically connected (covary), the many might be explainable by a few. For example a one-dimensional waveform can excite a wave equation filling a 3-D space. The values in that space will have a lot of covariance. In this chapter we take multidimensional spaces full of numbers and answer the question, "what causal differential (difference) equation might have created these numbers?" Our answer here, an autoregressive filter, does the job imperfectly, but it is a big step away from complete ignorance. As the book progresses we find three kinds of uses: (1) filling in missing data and uncontrolled parts of models, (2) preparing residuals for data fitting, (3) providing "prior" models for preconditioning an estimation.

Recall that residuals (and preconditioning variables) should be Independent, and Identically Distributed (IID). In practice the "ID" means all residuals should have the same variance, and the preceding "I" means likewise in Fourier space (whiteness). This is the "I" chapter. Conceptually we might jump in and out of Fourier space, but here we learn processes in physical space that whiten in Fourier space. In earlier chapters we transformed from a physical space to something more like an IID space when we said, "Topography is smooth, so let us estimate and view instead its derivative."

The branch of mathematics introduced here is young. Physicists seem to know nothing of it, perhaps because it begins with time not being a continuous variable. About 100 years ago people looked at market prices and wondered why they varied from day to day. To try to make money from the market fluctuations they schemed to try to predict prices. That is a good place to begin. The subject is known as "**time-series analysis.**" In this chapter we define the *autoregression* filter, also known as the **prediction-error filter (PEF)**. It gathers statistics for us. It gathers not the autocorrelation or the spectrum directly but it gathers them indirectly as the inverse of the amplitude spectrum of its input. Although time-series analysis is a one dimensional study, we naturally use the helix to broaden it to multidimensional space. The PEF leads us to the "inverse-covariance matrix" of statistical estimation theory. Theoreticians tell us we need this before we can properly find a solution. Here we see how to go after it.

## Time domain versus frequency domain

In the simplest applications, solutions can be most easily found in the frequency domain. When complications arise, it is better to use the time domain, to directly apply the convolution operator and the method of least squares.

A first complicating factor in the frequency domain is a required boundary in the time domain, such as that between past and future, or requirements that a filter be nonzero in a stated time interval. Another factor that attracts us to the time domain rather than the frequency domain is **weighting functions.**

Weighting functions are appropriate whenever a signal or image amplitude varies from place to place. Much of the literature on **time-series analysis** applies to the limited case of uniform weighting functions. Such time series are said to be “stationary.” This means that their statistical properties do not change in time. In real life, particularly in the analysis of echoes, signals are never stationary in time and space. A **stationarity** assumption is a reasonable starting assumption, but we should know how to go beyond it so we can take advantage of the many opportunities that do arise. In order of increasing difficulty in the frequency domain are the following complications:

1. A time boundary such as between past and future.
2. More time boundaries such as delimiting a filter.
3. More time boundaries such as erratic locations of missing data.
4. Nonstationary signal, i.e., time-variable weighting.
5. Time-axis stretching such as normal moveout.

We will not have difficulty with any of these complications here, because we will stay in the time domain and set up and solve optimization applications by use of the conjugate-direction method. Thus we will be able to cope with great complexity in goal formulation and get the right answer without approximations. By contrast, analytic or partly analytic methods can be more economical, but they generally solve somewhat different applications than those given to us by nature.

## SOURCE WAVEFORM, MULTIPLE REFLECTIONS

Here we devise a simple mathematical model for deep **water bottom** multiple reflections.<sup>1</sup> There are two unknown waveforms, the source waveform  $S(\omega)$  and the ocean-floor reflection  $F(\omega)$ . The water-bottom primary reflection  $P(\omega)$  is the convolution of the source waveform with the water-bottom response; so  $P(\omega) = S(\omega)F(\omega)$ . The first multiple reflection  $M(\omega)$  sees the same source waveform, the ocean floor, a minus one for the free surface, and the ocean floor again. Thus the observations  $P(\omega)$  and  $M(\omega)$  as functions of the physical parameters are

$$P(\omega) = S(\omega)F(\omega) \tag{1}$$

$$M(\omega) = -S(\omega)F(\omega)^2 \tag{2}$$

Algebraically the solutions of equations (1) and (2) are

$$F(\omega) = -M(\omega)/P(\omega) \tag{3}$$

$$S(\omega) = -P(\omega)^2/M(\omega) \tag{4}$$

These solutions can be computed in the Fourier domain by simple division. The difficulty is that the divisors in equations (3) and (4) can be zero, or small. This difficulty can be

---

<sup>1</sup> For this short course I am omitting here many interesting examples of multiple reflections shown in my 1992 book, PVI.

attacked by use of a positive number  $\epsilon$  to **stabilize** it. For example, multiply equation (3) on top and bottom by  $P(\omega)^T$  and add  $\epsilon > 0$  to the denominator. This gives

$$F(\omega) = - \frac{M(\omega)P^T(\omega)}{P(\omega)P(\omega)^T + \epsilon} \quad (5)$$

where  $P^T(\omega)$  is the complex conjugate of  $P(\omega)$ . Although the  $\epsilon$  stabilization seems nice, it apparently produces a nonphysical model. For  $\epsilon$  large or small, the time-domain response could turn out to be of much greater duration than is physically reasonable. This should not happen with perfect data, but in real life, data always has a limited spectral band of good quality.

Functions that are rough in the frequency domain will be long in the time domain. This suggests making a short function in the time domain by local smoothing in the frequency domain. Let the notation  $\langle \cdots \rangle$  denote smoothing by local averaging. Thus, to specify filters whose time duration is not unreasonably long, we can revise equation (5) to

$$F(\omega) = - \frac{\langle M(\omega)P^T(\omega) \rangle}{\langle P(\omega)P^T(\omega) \rangle} \quad (6)$$

where instead of deciding a size for  $\epsilon$  we need to decide how much smoothing. I find that smoothing has a simpler physical interpretation than choosing  $\epsilon$ . The goal of finding the filters  $F(\omega)$  and  $S(\omega)$  is to best model the multiple reflections so that they can be subtracted from the data, and thus enable us to see what primary reflections have been hidden by the multiples.

These frequency-duration difficulties do not arise in a time-domain formulation. Unlike in the frequency domain, in the time domain it is easy and natural to limit the duration and location of the nonzero time range of  $F(\omega)$  and  $S(\omega)$ . First express (3) as

$$0 = P(\omega)F(\omega) + M(\omega) \quad (7)$$

To imagine equation (7) as a fitting goal in the time domain, instead of scalar functions of  $\omega$ , think of vectors with components as a function of time. Thus  $\mathbf{f}$  is a column vector containing the unknown sea-floor filter,  $\mathbf{m}$  contains the “multiple” portion of a seismogram, and  $\mathbf{P}$  is a matrix of down-shifted columns, each column being the “primary”.

$$\mathbf{0} \approx \mathbf{r} = \begin{bmatrix} r_1 \\ r_2 \\ r_3 \\ r_4 \\ r_5 \\ r_6 \\ r_7 \\ r_8 \end{bmatrix} = \begin{bmatrix} p_1 & 0 & 0 \\ p_2 & p_1 & 0 \\ p_3 & p_2 & p_1 \\ p_4 & p_3 & p_2 \\ p_5 & p_4 & p_3 \\ p_6 & p_5 & p_4 \\ 0 & p_6 & p_5 \\ 0 & 0 & p_6 \end{bmatrix} \begin{bmatrix} f_1 \\ f_2 \\ f_3 \end{bmatrix} + \begin{bmatrix} m_1 \\ m_2 \\ m_3 \\ m_4 \\ m_5 \\ m_6 \\ m_7 \\ m_8 \end{bmatrix} \quad (8)$$

## TIME-SERIES AUTOREGRESSION

Given  $y_t$  and  $y_{t-1}$ , you might like to predict  $y_{t+1}$ . Earliest application of the ideas in this chapter came in the predictions of markets. Prediction of a signal from its past is called

“**autoregression**”, because a signal is regressed on itself “auto”. To find the scale factors you would optimize the fitting goal below, for the **prediction filter**  $(f_1, f_2)$ :

$$\mathbf{0} \approx \mathbf{r} = \begin{bmatrix} y_1 & y_0 \\ y_2 & y_1 \\ y_3 & y_2 \\ y_4 & y_3 \\ y_5 & y_4 \end{bmatrix} \begin{bmatrix} f_1 \\ f_2 \end{bmatrix} - \begin{bmatrix} y_2 \\ y_3 \\ y_4 \\ y_5 \\ y_6 \end{bmatrix} \quad (9)$$

(In practice, of course the system of equations would be much taller, and perhaps somewhat wider.) A typical row in the matrix (9) says that  $y_{t+1} \approx y_t f_1 + y_{t-1} f_2$  hence the description of  $f$  as a “prediction” filter. The error in the prediction is simply the residual. Define the residual to have opposite polarity and merge the column vector into the matrix, so you get

$$\begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} \approx \mathbf{r} = \begin{bmatrix} y_2 & y_1 & y_0 \\ y_3 & y_2 & y_1 \\ y_4 & y_3 & y_2 \\ y_5 & y_4 & y_3 \\ y_6 & y_5 & y_4 \end{bmatrix} \begin{bmatrix} 1 \\ -f_1 \\ -f_2 \end{bmatrix} \quad (10)$$

which is a standard form for autoregressions and prediction error.

**Multiple reflections** are predictable. It is the unpredictable part of a signal, the prediction residual, that contains the primary information. The output of the filter  $(1, -f_1, -f_2) = (a_0, a_1, a_2)$  is the unpredictable part of the input. This filter is a simple example of a “prediction-error” (PE) filter. It is one member of a family of filters called “error filters.”

The error-filter family are filters with one coefficient constrained to be unity and various other coefficients constrained to be zero. Otherwise, the filter coefficients are chosen to have minimum power output. Names for various error filters follow:

$(1, a_1, a_2, a_3, \dots, a_n)$	<b>prediction-error (PE) filter</b>
$(1, 0, 0, a_3, a_4, \dots, a_n)$	gapped PE filter with a gap
$(a_{-m}, \dots, a_{-2}, a_{-1}, 1, a_1, a_2, a_3, \dots, a_n)$	<b>interpolation-error (IE) filter</b>

We introduce a **free-mask matrix**  $\mathbf{K}$  which “passes” the freely variable coefficients in the filter and “rejects” the constrained coefficients (which in this first example is merely the first coefficient  $a_0 = 1$ ).

$$\mathbf{K} = \begin{bmatrix} 0 & \cdot & \cdot \\ \cdot & 1 & \cdot \\ \cdot & \cdot & 1 \end{bmatrix} \quad (11)$$

To compute a simple prediction error filter  $\mathbf{a} = (1, a_1, a_2)$  with the CD method, we write (9) or (10) as

$$\mathbf{0} \approx \mathbf{r} = \begin{bmatrix} y_2 & y_1 & y_0 \\ y_3 & y_2 & y_1 \\ y_4 & y_3 & y_2 \\ y_5 & y_4 & y_3 \\ y_6 & y_5 & y_4 \end{bmatrix} \begin{bmatrix} 0 & \cdot & \cdot \\ \cdot & 1 & \cdot \\ \cdot & \cdot & 1 \end{bmatrix} \begin{bmatrix} 1 \\ a_1 \\ a_2 \end{bmatrix} + \begin{bmatrix} y_2 \\ y_3 \\ y_4 \\ y_5 \\ y_6 \end{bmatrix} \quad (12)$$

Let us move from this specific fitting goal to the general case. (Notice the similarity of the free-mask matrix  $\mathbf{K}$  in this filter estimation application with the free-mask matrix  $\mathbf{J}$  in missing data goal (??).) The fitting goal is,

$$\mathbf{0} \approx \mathbf{Y}\mathbf{a} \quad (13)$$

$$\mathbf{0} \approx \mathbf{Y}(\mathbf{I} - \mathbf{K} + \mathbf{K})\mathbf{a} \quad (14)$$

$$\mathbf{0} \approx \mathbf{YK}\mathbf{a} + \mathbf{Y}(\mathbf{I} - \mathbf{K})\mathbf{a} \quad (15)$$

$$\mathbf{0} \approx \mathbf{YK}\mathbf{a} + \mathbf{Y}\mathbf{a}_0 \quad (16)$$

$$\mathbf{0} \approx \mathbf{YK}\mathbf{a} + \mathbf{y} \quad (17)$$

$$\mathbf{0} \approx \mathbf{r} = \mathbf{YK}\mathbf{a} + \mathbf{r}_0 \quad (18)$$

which means we initialize the residual with  $\mathbf{r}_0 = \mathbf{y}$ . and then iterate with

$$\Delta\mathbf{a} \leftarrow \mathbf{K}^T\mathbf{Y}^T \mathbf{r} \quad (19)$$

$$\Delta\mathbf{r} \leftarrow \mathbf{YK} \Delta\mathbf{a} \quad (20)$$

## PREDICTION-ERROR FILTER OUTPUT IS WHITE

In Chapter ?? we learned that least squares residuals should be IID (Independent, Identically Distributed) which in practical terms means in both Fourier space and physical space they should have a uniform variance. Further, not only should residuals have the IID property, but we should choose a preconditioning transformation so that our unknowns have the same IID nature. It is easy enough to achieve flattening in physical space by means of weighting functions. Here we see that Prediction-error filters (PEFs) enable us to flatten in fourier space.

PEFs transform signals and images to whiteness. Residuals and preconditioned models should be white.

### *The relationship between spectrum and PEF*

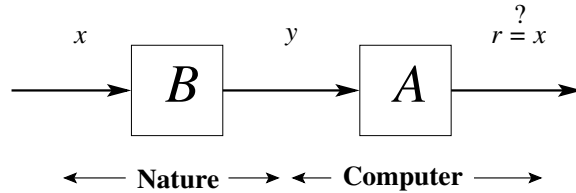
Knowledge of an autocorrelation function is equivalent to knowledge of a spectrum. The two are simply related by Fourier transform. A spectrum or an autocorrelation function encapsulates an important characteristic of a signal or an image. Generally the spectrum changes slowly from place to place although it could change rapidly. Of all the assumptions we could make to fill empty bins, one that people usually find easiest to agree with is that the spectrum should be the same in the empty-bin regions as where bins are filled. In practice we deal with neither the spectrum nor its autocorrelation but with a third object. This third object is the Prediction Error Filter (PEF), the filter in equation (10).

Take equation (10) for  $\mathbf{r}$  and multiply it by the adjoint  $\mathbf{r}^T$  getting a quadratic form in the PEF coefficients. Minimizing this quadratic form determines the PEF. This quadratic form depends only on the autocorrelation of the original data  $y_t$ , not on the data  $y_t$  itself. Clearly the PEF is unchanged if the data has its polarity reversed or its time axis reversed. Indeed, we'll see here that knowledge of the PEF is equivalent to knowledge of the autocorrelation or the spectrum.

### *Undoing convolution in nature*

Prediction-error filtering is called “**blind deconvolution**”. In the exploration industry it is simply called “**deconvolution**”. This word goes back to very basic models and concepts. In this model one envisions a random white-spectrum excitation function  $\mathbf{x}$  existing in nature, and this excitation function is somehow filtered by unknown natural processes, with a filter operator  $\mathbf{B}$  producing an *output*  $\mathbf{y}$  in nature that becomes the *input*  $\mathbf{y}$  to our computer programs. This is sketched in Figure 1. Then we design a prediction-error filter  $\mathbf{A}$  on  $\mathbf{y}$ ,

Figure 1: Flow of information from nature, to observation, into computer.



which yields a white-spectrum residual  $\mathbf{r}$ . Because  $\mathbf{r}$  and  $\mathbf{x}$  theoretically have the same spectrum, the tantalizing prospect is that maybe  $\mathbf{r}$  equals  $\mathbf{x}$ , meaning that the PEF  $\mathbf{A}$  has *deconvolved* the unknown convolution  $\mathbf{B}$ .

### *Causal with causal inverse*

Theoretically, a PEF is a causal filter with a causal inverse. This adds confidence to the likelihood that deconvolution of natural processes with a PEF might get the correct phase spectrum as well as the correct amplitude spectrum. Naturally, the PEF does not give the correct phase to an “all-pass” filter. That is a filter with a phase shift but a constant amplitude spectrum. (I think most migration operators are in this category.)

Theoretically we should be able to use a PEF in either convolution or polynomial division. There are some dangers though, mainly connected with dealing with data in small windows. Truncation phenomena might give us PEF estimates that are causal, but whose inverse is not, so they cannot be used in polynomial division. This is a lengthy topic in the classic literature. This old, fascinating subject is examined in my books, FGDP and PVI. A classic solution is one by John Parker Burg. We should revisit the Burg method in light of the helix.

### *PEF output tends to whiteness*

The most important property of a **prediction-error filter** or **PEF** is that its output tends to a **white spectrum** (to be proven here). No matter what the input to this filter, its output tends to whiteness as the number of the coefficients  $n \rightarrow \infty$  tends to infinity. Thus, the **PE filter** adapts itself to the input by absorbing all its **color**. This has important statistical implications and important geophysical implications.

### *Spectral estimation*

The PEF’s output being white leads to an important consequence: To specify a spectrum, we can give the spectrum (of an input) itself, give its autocorrelation, or give its PEF

coefficients. Each is transformable to the other two. Indeed, an effective mechanism of spectral estimation, developed by John P. **Burg** and described in **FGDP**, is to compute a PE filter and look at the inverse of its spectrum.

### *Short windows*

The power of a PE filter is that a short filter can often extinguish, and thereby represent, the information in a long resonant filter. If the input to the PE filter is a sinusoid, it is exactly predictable by a three-term recurrence relation, and all the color is absorbed by a three-term PE filter (see exercises). Burg's spectral estimation is especially effective in short windows.

### *Weathered layer resonance*

That the output spectrum of a PE filter is **white** is also useful geophysically. Imagine the reverberation of the **soil** layer, highly variable from place to place, as the resonance between the surface and shallow more-consolidated soil layers varies rapidly with surface location because of geologically recent fluvial activity. The spectral **color** of this erratic variation on surface-recorded seismograms is compensated for by a PE filter. Usually we do not want PE-filtered seismograms to be white, but once they all have the same spectrum, it is easy to postfilter them to any desired spectrum.

## PEF whiteness proof in 1-D

The basic idea of least-squares fitting is that the residual is orthogonal to the fitting functions. Applied to the PE filter, this idea means that the output of a PE filter is orthogonal to lagged inputs. The **orthogonality** applies only for lags in the past, because prediction knows only the past while it aims to the future. What we want to show here is different, namely, that the output is uncorrelated with *itself* (as opposed to the input) for lags in *both* directions; hence the output spectrum is **white**.

In (21) are two separate and independent autoregressions,  $\mathbf{0} \approx \mathbf{Y}_a \mathbf{a}$  for finding the filter  $\mathbf{a}$ , and  $\mathbf{0} \approx \mathbf{Y}_b \mathbf{b}$  for finding the filter  $\mathbf{b}$ . By noticing that the two matrices are really the same (except a row of zeros on the bottom of  $\mathbf{Y}_a$  is a row in the top of  $\mathbf{Y}_b$ ) we realize that the two regressions must result in the same filters  $\mathbf{a} = \mathbf{b}$ , and the residual  $\mathbf{r}_b$  is a shifted version of  $\mathbf{r}_a$ . In practice, I visualize the matrix being a thousand components tall (or a million) and a hundred components wide.

$$\mathbf{0} \approx \mathbf{r}_a = \begin{bmatrix} y_1 & 0 & 0 \\ y_2 & y_1 & 0 \\ y_3 & y_2 & y_1 \\ y_4 & y_3 & y_2 \\ y_5 & y_4 & y_3 \\ y_6 & y_5 & y_4 \\ 0 & y_6 & y_5 \\ 0 & 0 & y_6 \\ 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} 1 \\ a_1 \\ a_2 \end{bmatrix}; \quad \mathbf{0} \approx \mathbf{r}_b = \begin{bmatrix} 0 & 0 & 0 \\ y_1 & 0 & 0 \\ y_2 & y_1 & 0 \\ y_3 & y_2 & y_1 \\ y_4 & y_3 & y_2 \\ y_5 & y_4 & y_3 \\ y_6 & y_5 & y_4 \\ 0 & y_6 & y_5 \\ 0 & 0 & y_6 \end{bmatrix} \begin{bmatrix} 1 \\ b_1 \\ b_2 \end{bmatrix} \quad (21)$$

When the energy  $\mathbf{r}^T \mathbf{r}$  of a residual has been minimized, the residual  $\mathbf{r}$  is orthogonal to the fitting functions. For example, choosing  $a_2$  to minimize  $\mathbf{r}^T \mathbf{r}$  gives  $0 = \partial \mathbf{r}^T \mathbf{r} / \partial a_2 = 2\mathbf{r}^T \partial \mathbf{r} / \partial a_2$ . This shows that  $\mathbf{r}^T$  is perpendicular to  $\partial \mathbf{r} / \partial a_2$  which is the rightmost column of the  $\mathbf{Y}_a$  matrix. Thus the vector  $\mathbf{r}_a$  is orthogonal to all the columns in the  $\mathbf{Y}_a$  matrix except the first (because we do not minimize with respect to  $a_0$ ).

Our goal is a different theorem that is imprecise when applied to the three coefficient filters displayed in (21), but becomes valid as the filter length tends to infinity  $\mathbf{a} = (1, a_1, a_2, a_3, \dots)$  and the matrices become infinitely wide. Actually, all we require is the last component in  $\mathbf{b}$ , namely  $b_n$  tend to zero. This generally happens because as  $n$  increases,  $y_{t-n}$  becomes a weaker and weaker predictor of  $y_t$ .

Here's a mathematical fact we soon need: For any vectors  $\mathbf{u}$  and  $\mathbf{v}$ , if  $\mathbf{r} \cdot \mathbf{u} = \mathbf{0}$  and  $\mathbf{r} \cdot \mathbf{v} = \mathbf{0}$ , then  $\mathbf{r} \cdot (\mathbf{u} + \mathbf{v}) = \mathbf{0}$  and  $\mathbf{r} \cdot (6\mathbf{u} - 3\mathbf{v}) = \mathbf{0}$  and  $\mathbf{r} \cdot (a_1\mathbf{u} + a_2\mathbf{v}) = \mathbf{0}$  for any  $a_1$  and  $a_2$ .

The matrix  $\mathbf{Y}_a$  contains all of the columns that are found in  $\mathbf{Y}_b$  except the last (and the last one is not important). This means that  $\mathbf{r}_a$  is not only orthogonal to all of  $\mathbf{Y}_a$ 's columns (except the first) but  $\mathbf{r}_a$  is also orthogonal to all of  $\mathbf{Y}_b$ 's columns except the last. Although  $\mathbf{r}_a$  isn't really perpendicular to the last column of  $\mathbf{Y}_b$ , it doesn't matter because that column has hardly any contribution to  $\mathbf{r}_b$  since  $|b_n| \ll 1$ . Because  $\mathbf{r}_a$  is (effectively) orthogonal to all the components of  $\mathbf{r}_b$ ,  $\mathbf{r}_a$  is also orthogonal to  $\mathbf{r}_b$  itself.

Here is a detail: In choosing the example of equation (21), I have shifted the two fitting problems by only one lag. We would like to shift by more lags and get the same result. For this we need more filter coefficients. By adding many more filter coefficients we are adding many more columns to the right side of  $\mathbf{Y}_b$ . That's good because we'll be needing to neglect more columns as we shift  $\mathbf{r}_b$  further from  $\mathbf{r}_a$ . Neglecting these columns is commonly justified by the experience that "after short range regressors have had their effect, long range regressors generally find little remaining to predict." (Recall that the damped harmonic oscillator from physics, the finite difference equation that predicts the future from the past, uses only two lags.)

Here is the main point: Since  $\mathbf{r}_b$  and  $\mathbf{r}_a$  both contain the same signal  $\mathbf{r}$  but time-shifted, the orthogonality at all shifts means that the autocorrelation of  $\mathbf{r}$  vanishes at all lags. An exception, of course, is at zero lag. The autocorrelation does not vanish there because  $\mathbf{r}_a$  is not orthogonal to its first column (because we did not minimize with respect to  $a_0$ ).

As we redraw  $\mathbf{0} \approx \mathbf{r}_b = \mathbf{Y}_b \mathbf{b}$  for various lags, we may shift the columns only downward because shifting them upward would bring in the first column of  $\mathbf{Y}_a$  and the residual  $\mathbf{r}_a$  is not orthogonal to that. Thus we have only proven that one side of the autocorrelation of  $\mathbf{r}$  vanishes. That is enough however, because autocorrelation functions are symmetric, so if one side vanishes, the other must also.

If  $\mathbf{a}$  and  $\mathbf{b}$  were two-sided filters like  $(\dots, b_{-2}, b_{-1}, 1, b_1, b_2, \dots)$  the proof would break. If  $\mathbf{b}$  were two-sided,  $\mathbf{Y}_b$  would catch the nonorthogonal column of  $\mathbf{Y}_a$ . Not only is  $\mathbf{r}_a$  not proven to be perpendicular to the first column of  $\mathbf{Y}_a$ , but it cannot be orthogonal to it because a signal cannot be orthogonal to itself.

The implications of this theorem are far reaching. The residual  $\mathbf{r}$ , a convolution of  $\mathbf{y}$  with  $\mathbf{a}$  has an autocorrelation that is an impulse function. The Fourier transform of an impulse is a constant. Thus the spectrum of the residual is "white". Thus  $\mathbf{y}$  and  $\mathbf{a}$  have



mutually inverse spectra.

Since the output of a PEF is white, the PEF itself has a spectrum inverse to its input.

An important application of the PEF is in missing data interpolation. We'll see examples later in this chapter. My third book, PVI<sup>2</sup> has many examples<sup>3</sup> in one dimension with both synthetic data and field data including the **gap** parameter. Here we next extend these ideas to two (or more) dimensions.

## Simple dip filters

Convolution in two dimensions is just like convolution in one dimension except that convolution is done on two axes. The input and output data are planes of numbers and the filter is also a plane. A two-dimensional filter is a small plane of numbers that is convolved over a big data plane of numbers.

Suppose the data set is a collection of seismograms uniformly sampled in space. In other words, the data is numbers in a  $(t, x)$ -plane. For example, the following filter destroys any wavefront aligned along the direction of a line containing both the “+1” and the “−1”.

$$\begin{array}{ccc} -1 & \cdot & \\ \cdot & \cdot & \\ \cdot & & 1 \end{array} \quad (22)$$

The next filter destroys a wave with a slope in the opposite direction:

$$\begin{array}{ccc} \cdot & & 1 \\ -1 & \cdot & \end{array} \quad (23)$$

To convolve the above two filters, we can reverse either on (on both axes) and correlate them, so that you can get

$$\begin{array}{ccc} \cdot & -1 & \cdot \\ 1 & \cdot & \cdot \\ \cdot & \cdot & 1 \\ \cdot & -1 & \cdot \end{array} \quad (24)$$

which destroys waves of both slopes.

A **two-dimensional filter** that can be a **dip-rejection filter** like (22) or (23) is

$$\begin{array}{ccc} a & \cdot & \\ b & \cdot & \\ c & & 1 \\ d & \cdot & \\ e & \cdot & \end{array} \quad (25)$$

where the coefficients  $(a, b, c, d, e)$  are to be estimated by least squares in order to minimize the power out of the filter. (In the filter table, the time axis runs vertically.)

<sup>2</sup> [http://sepwww.stanford.edu/sep/prof/pvi/toc\\_html/index.html](http://sepwww.stanford.edu/sep/prof/pvi/toc_html/index.html)

<sup>3</sup> [http://sepwww.stanford.edu/sep/prof/pvi/tsa/paper\\_html/node1.html](http://sepwww.stanford.edu/sep/prof/pvi/tsa/paper_html/node1.html)

Fitting the filter to two neighboring traces that are identical but for a time shift, we see that the filter coefficients  $(a, b, c, d, e)$  should turn out to be something like  $(-1, 0, 0, 0, 0)$  or  $(0, 0, -0.5, -0.5, 0)$ , depending on the dip (stepout) of the data. But if the two channels are not fully coherent, we expect to see something like  $(-0.9, 0, 0, 0, 0)$  or  $(0, 0, -0.4, -0.4, 0)$ . To find filters such as (24), we adjust coefficients to minimize the power out of filter shapes, as in

$$\begin{array}{rcc} v & a & \cdot \\ w & b & \cdot \\ x & c & 1 \\ y & d & \cdot \\ z & e & \cdot \end{array} \quad (26)$$

With 1-dimensional filters, we think mainly of power spectra, and with 2-dimensional filters we can think of temporal spectra and spatial spectra. What is new, however, is that in two dimensions we can think of dip spectra (which is when a 2-dimensional spectrum has a particularly common form, namely when energy organizes on radial lines in the  $(\omega, k_x)$ -plane). As a short (three-term) 1-dimensional filter can devour a sinusoid, we have seen that simple 2-dimensional filters can devour a small number of dipoles.

### PEF whiteness proof in 2-D

A well-known property (see FGDP or PVI) of a 1-D PEF is that its energy clusters immediately after the impulse at zero delay time. Applying this idea to the helix in Figure shows us that we can consider a 2-D PEF to be a small halfplane with an impulse along a side. These shapes are what we see here in Figure 2.

Figure 2: A 2-D whitening filter template, and itself lagged. At output locations “A” and “B,” the filter coefficient is constrained to be “1”. When the semicircles are viewed as having infinite radius, the B filter is contained in the A filter. Because the output at A is orthogonal to all its inputs, which include all inputs of B, the output at A is orthogonal to the output of B.

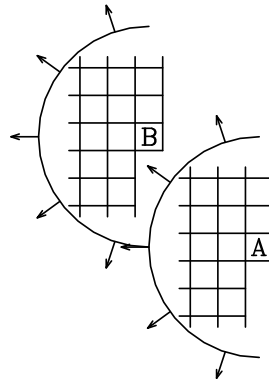


Figure 2 shows the input plane with a 2-D filter on top of it at two possible locations. The filter shape is a semidisk, which you should imagine being of infinitely large radius. Notice that semidisk A includes all the points in B. The output of disk A will be shown to be orthogonal to the output of disk B. Conventional least squares theory says that the coefficients of the filter are designed so that the output of the filter is orthogonal to each of the inputs to that filter (except for the input under the “1,” because any nonzero signal cannot be orthogonal to itself). Recall that if a given signal is orthogonal to each in a given group of signals, then the given signal is orthogonal to all linear combinations within that group. The output at B is a linear combination of members of its input group, which is

included in the input group of A, which are already orthogonal to A. Therefore the output at B is orthogonal to the output at A. In summary,

residual	$\perp$	fitting function
output at A	$\perp$	each input to A
output at A	$\perp$	each input to B
output at A	$\perp$	linear combination of each input to B
output at A	$\perp$	output at B

The essential meaning is that a particular lag of the output **autocorrelation** function vanishes.

Study Figure 2 to see for what lags all the elements of the B filter are wholly contained in the A filter. These are the lags where we have shown the output autocorrelation to be vanishing. Notice another set of lags where we have proven nothing (where B is moved to the right of A). Autocorrelations are centrosymmetric, which means that the value at any lag is the same as the value at the negative of that lag, even in 2-D and 3-D where the lag is a vector quantity. Above we have shown that a halfplane of autocorrelation values vanishes. By the centrosymmetry, the other half must vanish too. Thus the autocorrelation of the PEF output is an impulse function, so its 2-D spectrum is white.

The helix tells us why the proper filter form is not a square with the “1” on the corner. Before I discovered the helix, I understood it another way (that I learned from John P. Burg): For a spectrum to be white, *all* nonzero autocorrelation lags must be zero-valued. If the filter were a quarter-plane, then the symmetry of autocorrelations would only give us vanishing in another quarter, so there would be two remaining quarter-planes where the autocorrelation was not zero.

Fundamentally, the white-output theorem requires a one-dimensional ordering to the values in a plane or volume. The filter must contain a halfplane of values so that symmetry gives the other half.

You will notice some nonuniqueness. We could embed the helix with a 90° rotation in the original physical application. Besides the difference in side boundaries, the 2-D PEF would have a different orientation. Both PEFs should have an output that tends to whiteness as the filter is enlarged. It seems that we could design whitening autoregression filters for 45° rotations also, and we could also design them for hexagonal coordinate systems. In some physical applications, you might find the nonuniqueness unsettling. Does it mean the “final solution” is nonunique? Usually not, or not seriously so. Recall even in one dimension, the time reverse of a PEF has the same spectrum as the original PEF. When a PEF is used for regularizing a fitting application, it is worth noticing that the quadratic form minimized is the PEF times its adjoint so the phase drops out. Likewise, a missing data restoration also amounts to minimizing a quadratic form so the phase again drops out.

## BASIC BLIND DECONVOLUTION

Here are the basic definitions of blind deconvolution: If a model  $m_t$  (with FT  $M$ ) is made of random numbers and convolved with a “source waveform” (having FT)  $F^{-1}$  it creates data  $D$ . From data  $D$  you find the model  $M$  by  $M = FD$ . Trouble is, you typically do not know  $F$  and need to estimate (guess) it hence the word “blind.”

Suppose we have many observations or many channels of  $D$  so we label them  $D_j$ . We can define a model  $M_j$  as

$$M_j = \frac{D_j}{\sqrt{\sum_j D^* D}} \quad (27)$$

so blind deconvolution removes the average spectrum.

Sometimes we have only a single signal  $D$  but it is quite long. Because the signal is long, the magnitude of its Fourier transform is rough, so we smooth it over frequency, and denote it thus:

$$M = \frac{D}{\sqrt{\ll D^* D \gg}} \quad (28)$$

Smoothing the spectrum makes the time function shorter. Indeed, the amount of smoothing may be chosen by the amount of shortness wanted.

These preliminary models are the most primitive forms of deconvolved data. They deal only with the amplitude spectrum. Most deconvolutions involve also the phase. The generally chosen phase is one with a causal filter. A causal filter  $f_t$  (vanishes before  $t = 0$ ) with FT  $F$  is chosen so that  $M = FD$  is white. Finding this filter is a serious undertaking, normally done in a one-dimensional space. Here, taking advantage of the helix, we do it in space of any number of dimensions.

For reasons explained later, this is equivalent to minimizing the energy output of a filter beginning with a one,  $(1, f_1, f_2, f_3, \dots)$ . The inverse of this filter  $1/F$  is often called “the impulse response”, or “the source waveform”. Whether it actually is a source waveform depends on the physical setup as well as some mathematical assumptions we will learn.

## Examples of modeling and deconvolving with a 2-D PEF

Here we examine elementary signal-processing applications of 2-D prediction-error filters (PEFs) on both everyday 2-D textures and on seismic data. Some of these textures are easily modeled with prediction-error filters (PEFs) while others are not. All figures used the same  $10 \times 10$  filter shape. No attempt was made to optimize filter size or shape or any other parameters.

Results in Figures 3-9 are shown with various familiar textures<sup>4</sup> on the left as training data sets. From these training data sets, a prediction-error filter (PEF) is estimated using module `pef` on page 25. The center frame is simulated data made by deconvolving (polynomial division) random numbers by the estimated PEF. The right frame is the more familiar process, convolving the estimated PEF on the training data set.

Theoretically, the right frame tends towards a white spectrum. Earlier you could notice the filter size by knowing that the output was taken to be zero where the filter is only partially on the data. This was annoying on real data where we didn’t want to throw away any data around the sides. Now the filtering is done without a call to the boundary module so we have typical helix wraparound.

---

<sup>4</sup> I thank Morgan Brown for finding these textures.

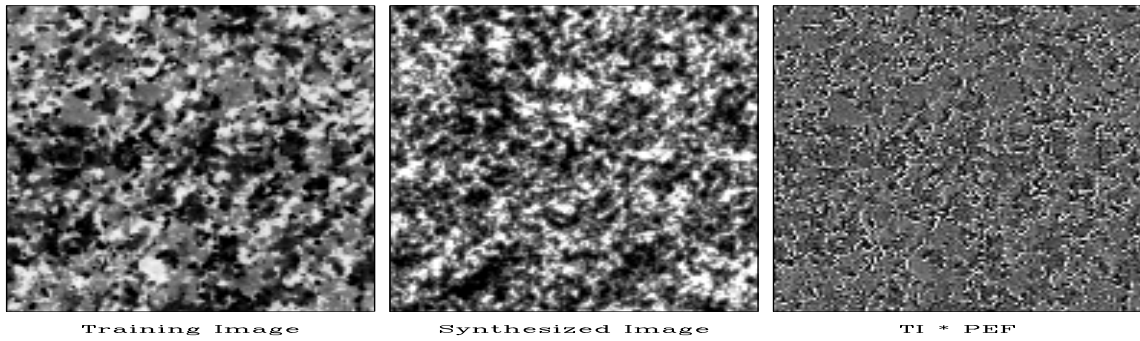


Figure 3: Synthetic granite matches the training image quite well. The prediction error (PE) is large at grain boundaries so it almost seems to outline the grains.

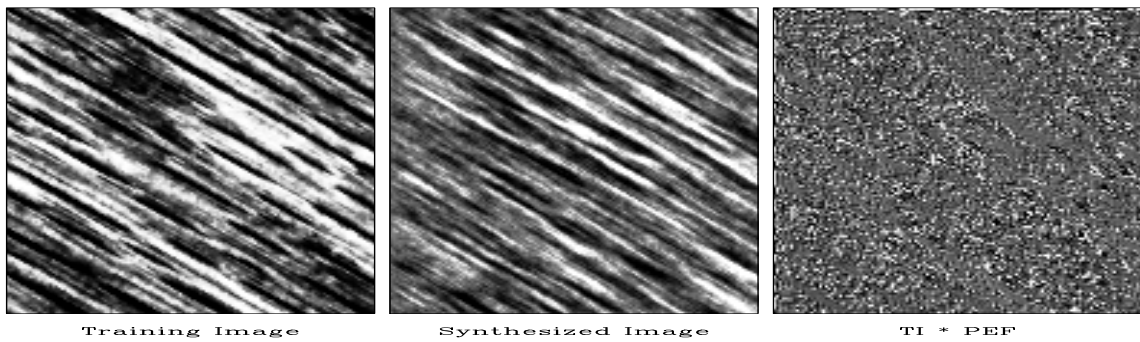


Figure 4: Synthetic wood grain has too little white. This is because of the nonsymmetric brightness histogram of natural wood. Again, the PEF output looks random as expected.

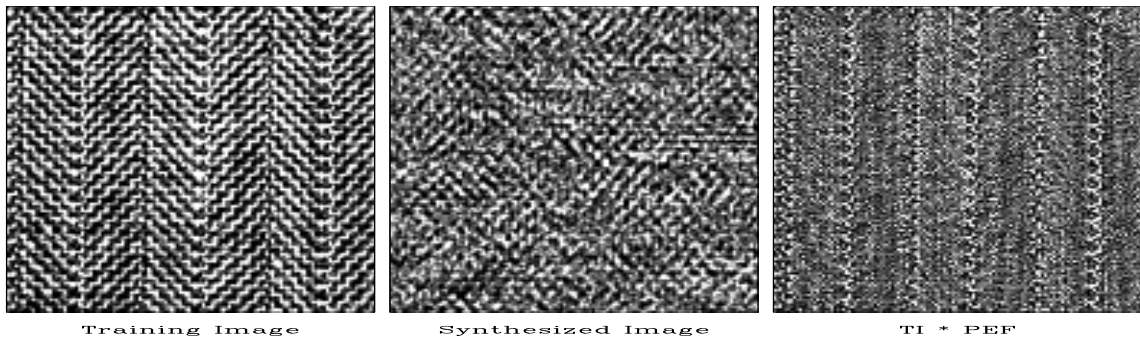


Figure 5: A banker's suit (left). A student's suit (center). My suit (right). The prediction error is large where the weave changes direction.



Figure 6: Basket weave. The simulated data fails to segregate the two dips into a checkerboard pattern. The PEF output looks structured perhaps because the filter is too small.

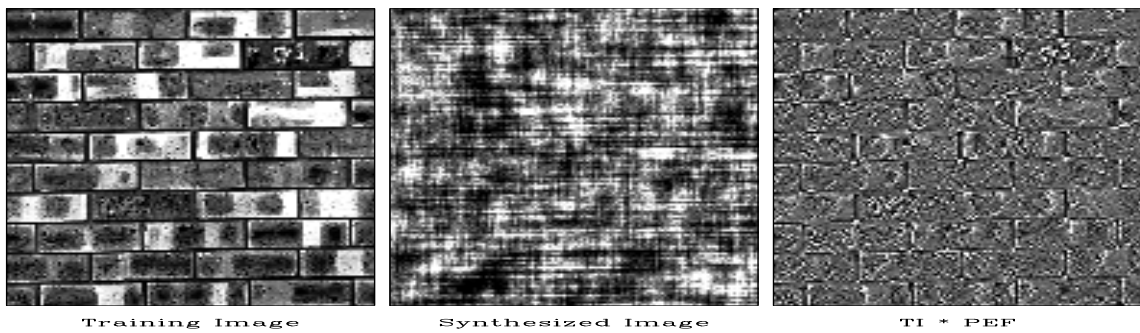


Figure 7: Brick. Synthetic brick edges are everywhere and do not enclose blocks containing a fixed color. PEF output highlights the mortar.

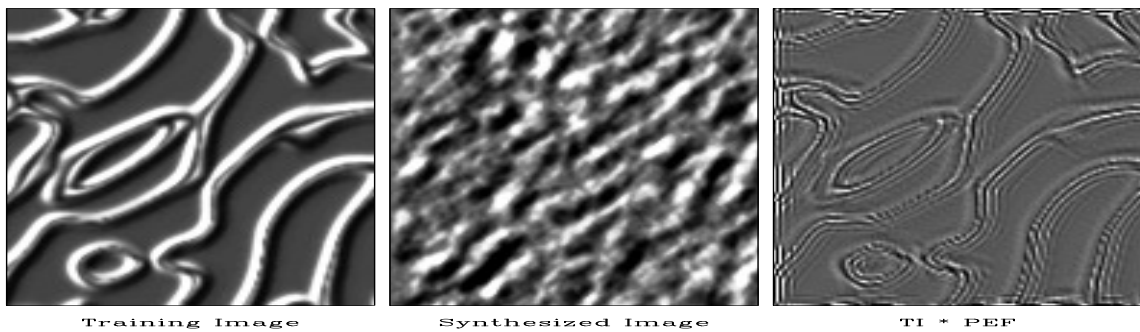


Figure 8: Ridges. A spectacular failure of the stationarity assumption. All dips are present but in different locations. Never-the-less, the ridges have been sharpened by the deconvolution.

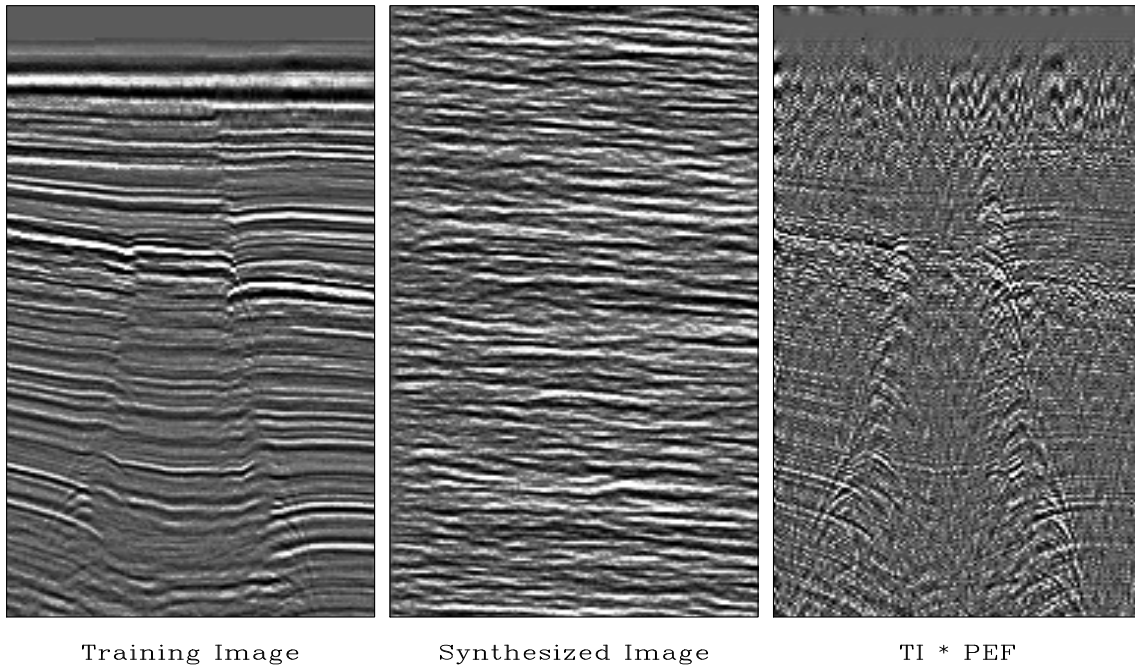


Figure 9: Gulf of Mexico seismic section, modeled, and deconvolved. Do you see any drilling prospects in the simulated data? In the deconvolution, the strong horizontal layering is suppressed giving a better view of the hyperbolas. The decon filter is the same  $10 \times 10$  used on the everyday textures.

Since a PEF tends to the inverse of the spectrum of its input, results similar to these could probably be found using Fourier transforms, smoothing spectra, etc. We used PEFs because of their flexibility. The filters can be any shape. They can dodge around missing data, or we can use them to estimate missing data. We avoid periodic boundary assumptions inherent to FT. The PEF's are designed only internal to known data, not off edges so they are readily adaptable to nonstationarity. Thinking of these textures as seismic time slices, the textures could easily be required to pass thru specific values at well locations.

### Seismic field data examples

Figures 10-13 are based on exploration seismic data from the Gulf of Mexico deep water. A ship carries an air gun and tows a streamer with some hundreds of geophones. First we look at a single pop of the gun. We use all the hydrophone signals to create a single 1-D PEF for the time axis. This changes the average temporal frequency spectrum as shown in Figure 10. Signals from 60 Hz to 120 Hz are boosted substantially. The raw data has evidently been prepared with strong filtering against signals below about 8 Hz. The PEF attempts to recover these signals, mostly unsuccessfully, but it does boost some energy near the 8 Hz cutoff. Choosing a longer filter would flatten the spectrum further. The big question is, "Has the PEF improved the appearance of the data?"

The data itself from the single pop, both before and after PE-filtering is shown in Figure 11. For reasons of aesthetics of human perception I have chosen to display a mirror

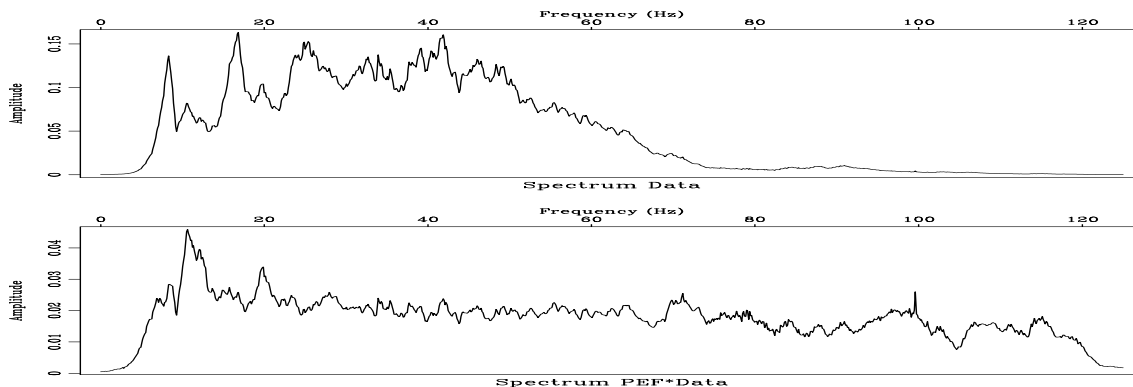


Figure 10:  $\omega$  spectrum of a shot gather of Figure 11 before and after 1-D decon with a 30 point filter.

image of the PEF’ed data. To see a blink movie of superposition of before-and-after images you need the electronic book. We notice that signals of high temporal frequencies indeed have the expected hyperbolic behavior in space. Thus, these high-frequency signals are wavefields, not mere random noise.

Given that all visual (or audio) displays have a bounded range of amplitudes, increasing the frequency content (bandwidth) means that we will need to turn down the amplification so we do not wish to increase the bandwidth unless we are adding signal.

Increasing the spectral bandwidth always requires us to diminish the gain.

The same ideas but with a two-dimensional PEF are in Figure 12 (the same data but with more of it squeezed onto the page.) As usual, the raw data is dominated by events arriving later at greater distances. After the PEF, we tend to see equal energy in dips in all directions. We have strongly enhanced the “backscattered” energy, those events that arrive later at *shorter* distances.

Figure 13 shows echos from the all shots, the nearest receiver on each shot. This picture of the earth is called a “near-trace section.” This earth picture shows us why there is so much backscattered energy in Figure 12 (which is located at the left side of Figure 13). The backscatter comes from any of the many of near-vertical faults.

We have been thinking of the PEF as a tool for shaping the spectrum of a display. But does it have a physical meaning? What might it be? Referring back to the beginning of the chapter we are inclined to regard the PEF as the convolution of the source waveform with some kind of water-bottom response. In Figure 12 we used many different shot-receiver separations. Since each different separation has a different response (due to differing moveouts) the water bottom reverberation might average out to be roughly an impulse. Figure 12 is a different story. Here for each shot location, the distance to the receiver is constant. Designing a single channel PEF we can expect the PEF to contain both the shot waveform and the water bottom layers because both are nearly identical in all the shots. We would rather have a PEF that represents only the shot waveform (and perhaps a radiation pattern).



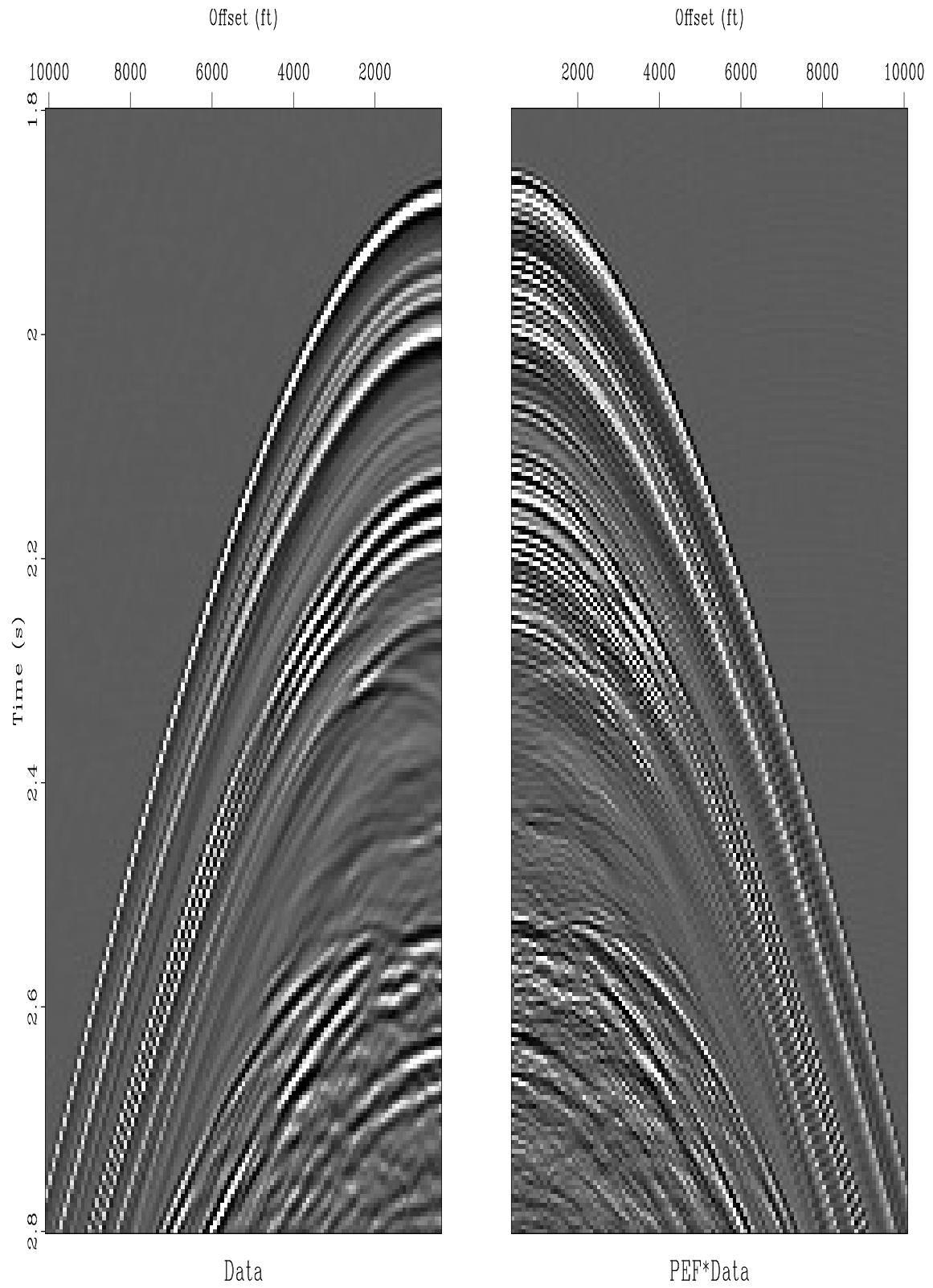


Figure 11: Raw data with its mirror. Mirror had 1-D PEF applied, 30 point filter.

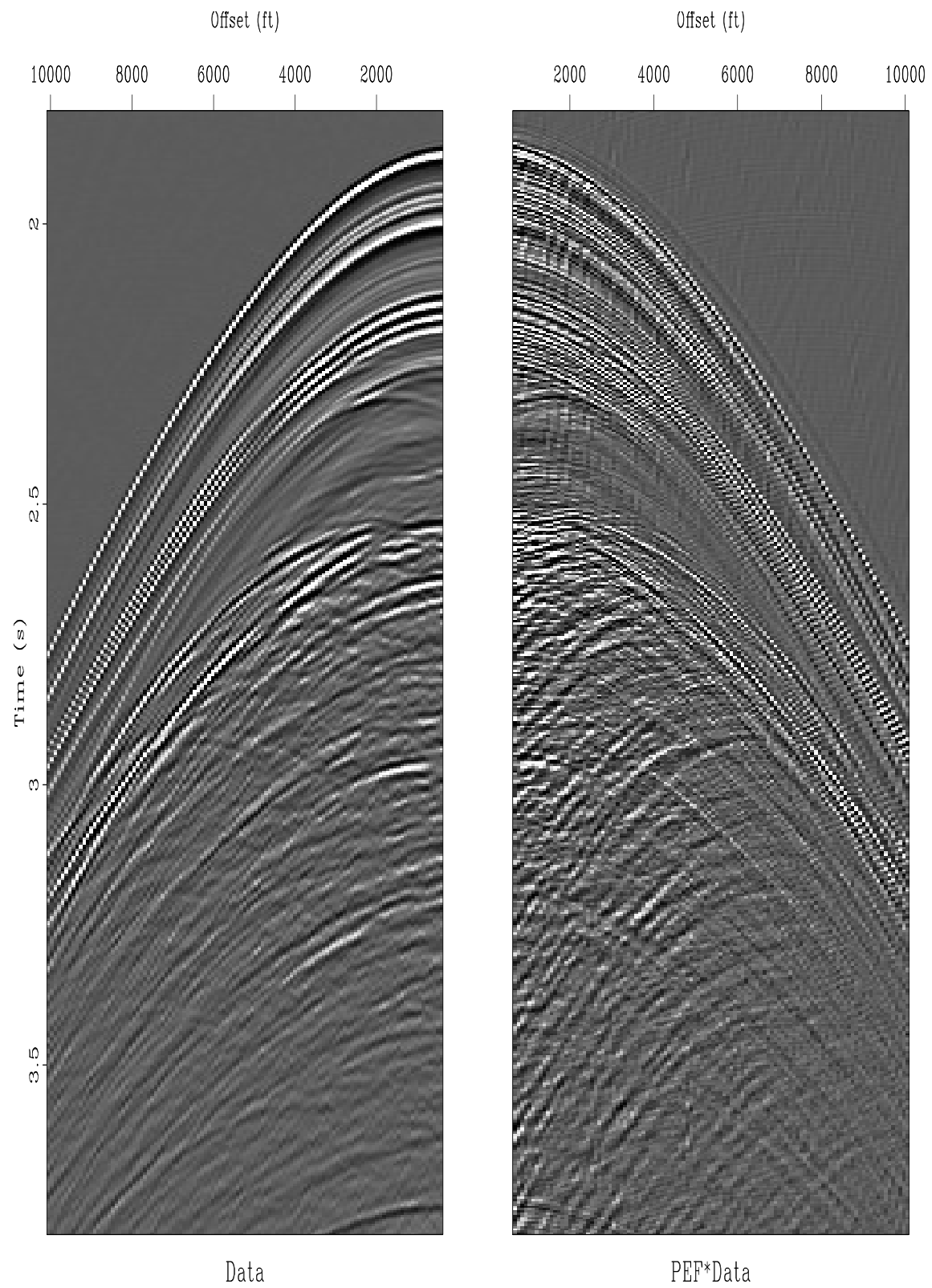


Figure 12: A 2-D filter (here  $20 \times 5$ ) brings out the backscattered energy.

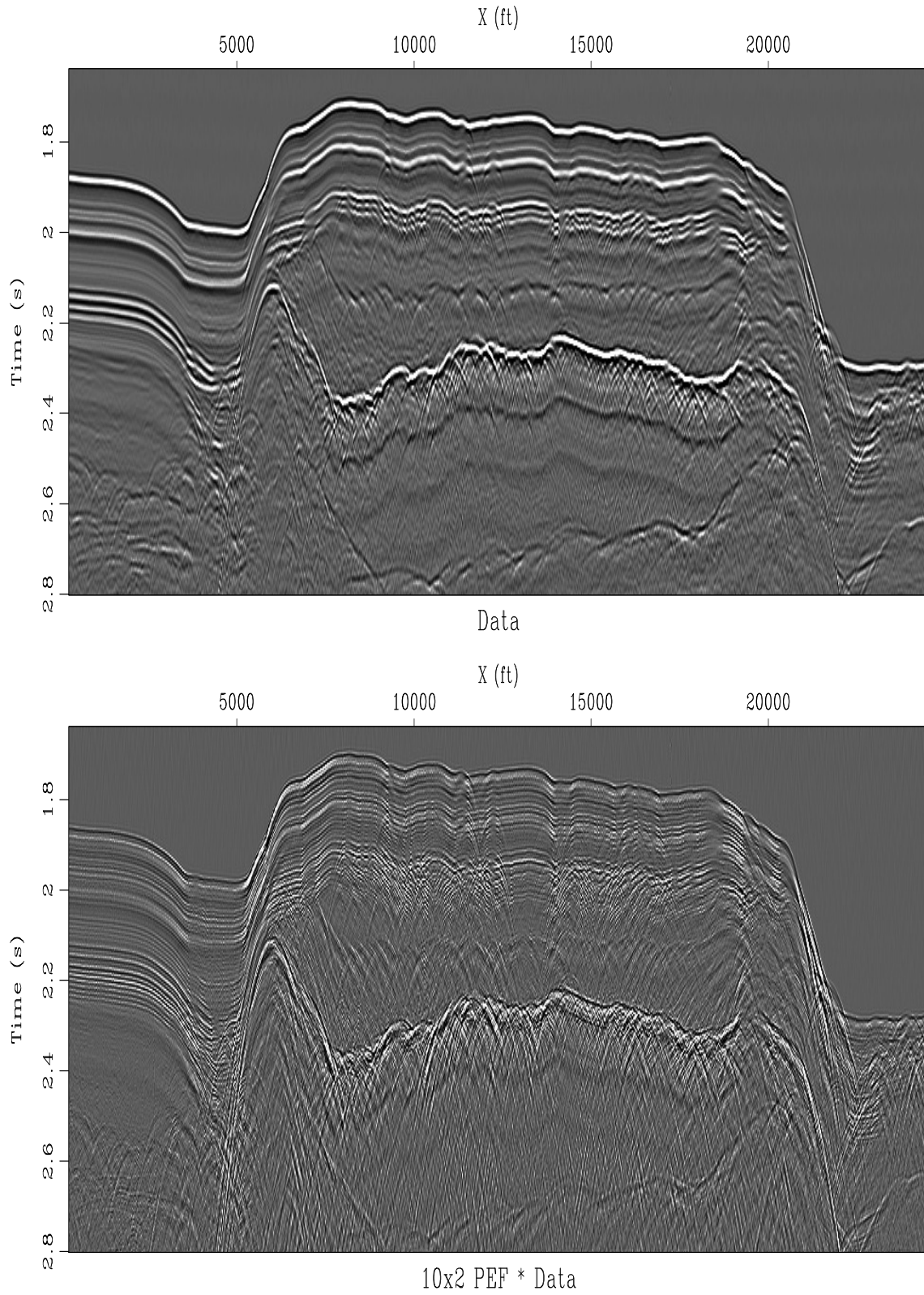


Figure 13: Raw data, near-trace section (top). Filtered with a two-channel PEF (bottom). The movie has other shaped filters.

Let us consider how we might work to push the water-bottom reverberation out of the PEF. This data is recorded in water 600 meters deep. A consequence is that the sea bottom is made of fine-grained sediments that settled very slowly and rather similarly from place to place. In shallow water the situation is different. The sands near estuaries are always shifting. Sedimentary layers thicken and thin. They are said to “on-lap and off-lap.” Here I do notice where the water bottom is sloped the layers do thin a little. To push the water bottom layers out of the PEF our idea is to base its calculation not on the raw data, but on the spatial prediction error of the raw data. On a perfectly layered earth a perfect spatial prediction error filter would zero all traces but the first one. Since a 2-D PEF includes spatial prediction as well as temporal prediction, we can expect it to contain much less of the sea-floor layers than the 1-D PEF. If you have access to the electronic book, you can blink the figure back and forth with various filter shapes.

### PEF ESTIMATION WITH MISSING DATA

If we are not careful, our calculation of the PEF could have the pitfall that it would try to use the missing data to find the PEF, and hence it would get the wrong PEF. To avoid this pitfall, imagine a PEF finder that uses weighted least squares where the weighting function vanishes on those fitting equations that involve missing data. The weighting would be unity elsewhere. Instead of weighting bad results by zero, we simply will not compute them. The residual there will be initialized to zero and never changed. Likewise for the adjoint, these components of the residual will never contribute to a gradient. So now we need a convolution program that produces no outputs where missing inputs would spoil it.

Recall there are two ways of writing convolution, equation (??) when we are interested in finding the filter *inputs*, and equation (??) when we are interested in finding the *filter itself*. We have already coded equation (??), operator `helicon` on page ???. That operator was useful in missing data applications. Now we want to find a prediction-error filter so we need the other case, equation (??), and we need to ignore the outputs that will be broken because of missing inputs. The operator module `hconest` does the job.

user/gee/hconest.c

```

1   for (ia = 0; ia < na; ia++) {
2       for (iy = aa->lag[ia]; iy < ny; iy++) {
3           if(aa->mis[iy]) continue;
4
5           ix = iy - aa->lag[ia];
6
7           if( adj) a[ia] -= y[iy] * x[ix];
8           else   y[iy] -= a[ia] * x[ix];
9       }
10  }
```

We are seeking a prediction error filter  $(1, a_1, a_2)$  but some of the data is missing. The data is denoted  $\mathbf{y}$  or  $y_i$  above and  $x_i$  below. Because some of the  $x_i$  are missing, some of the regression equations in (29) are worthless. When we figure out which are broken, we will put zero weights on those equations.

$$\mathbf{0} \approx \mathbf{r} = \mathbf{W}\mathbf{X}\mathbf{a} = \begin{bmatrix} w_1 & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & w_2 & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & w_3 & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & w_4 & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & w_5 & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & w_6 & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & w_7 & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & w_8 & \cdot \end{bmatrix} \begin{bmatrix} x_1 & 0 & 0 \\ x_2 & x_1 & 0 \\ x_3 & x_2 & x_1 \\ x_4 & x_3 & x_2 \\ x_5 & x_4 & x_3 \\ x_6 & x_5 & x_4 \\ 0 & x_6 & x_5 \\ 0 & 0 & x_6 \end{bmatrix} \begin{bmatrix} 1 \\ a_1 \\ a_2 \end{bmatrix} \quad (29)$$

Suppose that  $x_2$  and  $x_3$  were missing or known bad. That would spoil the 2nd, 3rd, 4th, and 5th fitting equations in (29). In principle, we want  $w_2, w_3, w_4$  and  $w_5$  to be zero. In practice, we simply want those components of  $\mathbf{r}$  to be zero.

What algorithm will enable us to identify the regression equations that have become defective, now that  $x_2$  and  $x_3$  are missing? Take filter coefficients  $(a_0, a_1, a_2, \dots)$  to be all ones. Let  $\mathbf{d}_{\text{free}}$  be a vector like  $\mathbf{x}$  but containing 1's for the missing (or "freely adjustable") data values and 0's for the known data values. Recall our very first definition of filtering showed we can put the filter in a vector and the data in a matrix or vice versa. Thus  $\mathbf{X}\mathbf{a}$  above gives the same result as  $\mathbf{A}\mathbf{x}$  below.

$$\begin{bmatrix} m_1 \\ m_2 \\ m_3 \\ m_4 \\ m_5 \\ m_6 \\ m_7 \\ m_8 \end{bmatrix} = \begin{bmatrix} 0 \\ 1 \\ 2 \\ 2 \\ 1 \\ 0 \\ 0 \\ 0 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 1 & 0 & 0 & 0 & 0 \\ 1 & 1 & 1 & 0 & 0 & 0 \\ 0 & 1 & 1 & 1 & 0 & 0 \\ 0 & 0 & 1 & 1 & 1 & 0 \\ 0 & 0 & 0 & 1 & 1 & 1 \\ 0 & 0 & 0 & 0 & 1 & 1 \\ 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 0 \\ 1 \\ 1 \\ 0 \\ 0 \\ 0 \end{bmatrix} = \mathbf{A}\mathbf{d}_{\text{free}} \quad (30)$$

The numeric value of each  $m_i$  tells us how many of its inputs are missing. Where none are missing, we want unit weights  $w_i = 1$ . Where any are missing, we want zero weights  $w_i = 0$ . The desired residual under partially missing inputs is computed by module `misinput` on the following page.

## Internal boundaries to multidimensional convolution

Sometimes we deal with small patches of data. In order that boundary phenomena not dominate the calculation intended in the central region, we need to take care that input data is not assumed to be zero beyond the interval that the data is given.

The two little triangular patches of zeros in the convolution matrix in equation (29) describe end conditions where it is assumed that the data  $y_t$  vanishes before  $t = 1$  and after  $t = 6$ . Alternately we might not wish to make that assumption. Thus the triangles filled with zeros could be regarded as missing data. In this one-dimensional example, it is easy to see that the filter, say `yy->mis` should be set to `true` at the ends so no output would ever be computed there. We would like to find a general multidimensional algorithm to correctly

user/gee/misinput.c

```

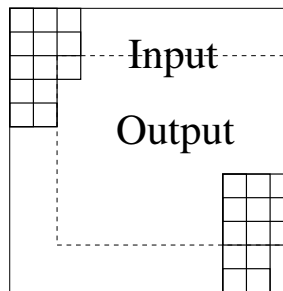
1 void find_mask(int n          /* data size */,
2               const int *known /* mask for known data [n] */,
3               sf_filter aa    /* helical filter */)
4 /*< create a filter mask >*/
5 {
6     int i, ih;
7     float *rr, *dfre;
8
9     rr = sf_floatalloc(n);
10    dfre = sf_floatalloc(n);
11
12    for (i=0; i < n; i++) {
13        dfre[i] = known[i]? 0.:1.;
14    }
15
16    sf_helicon_init(aa);
17
18    for (ih=0; ih < aa->nh; ih++) {
19        aa->flt[ih] = 1.;
20    }
21
22    sf_helicon_lop(false, false, n, n, dfre, rr);
23
24    for (ih=0; ih < aa->nh; ih++) {
25        aa->flt[ih] = 0.;
26    }
27
28    for (i=0; i < n; i++) {
29        if ( rr[i] > 0.) aa->mis[i] = true;
30    }
31
32    free (rr);
33    free (dfre);
34 }

```

specify `yy->mis` around the multidimensional boundaries. This proceeds like the missing data algorithm, i.e. we apply a filter of all ones to a data space template that is taken all zeros except ones at the locations of missing data, in this case  $y_0, y_{-1}$  and  $y_7, y_8$ . This amounts to surrounding the original data set with some missing data. We need padding the size of the filter on all sides. The padded region would be filled with ones (designating missing inputs). Where the convolution output is nonzero, there `yy->mis` is set to `true` denoting an output with missing inputs.

The two-dimensional case is a little more cluttered than the 1-D case but the principle is about the same. Figure 14 shows a larger input domain, a  $5 \times 3$  filter, and a smaller output domain. There are two things to notice. First, sliding the filter everywhere inside the outer

Figure 14: Domain of inputs and outputs of a two-dimensional filter like a PEF.



box, we get outputs (under the 1 location) only in the inner box. Second, (the adjoint idea) crosscorrelating the inner and outer boxes gives us the  $3 \times 5$  patch of information we use to build the filter coefficients. We need to be careful not to assume that signals vanish outside the region where they are defined. In a later chapter we will break data spaces into overlapping patches, separately analyze the patches, and put everything back together. We do this because crosscorrelations change with time and they are handled as constant in short time windows. There we must be particularly careful that zero signal values not be presumed outside of the small volumes; otherwise the many edges and faces of the many small volumes can overwhelm the interior that we want to study.

In practice, the input and output are allocated equal memory, but the output residual is initialized to zero everywhere and then not computed except where shown in figure 14. Below is module `bound` to build a selector for filter outputs that should never be examined or even computed (because they need input data from outside the given data space). Inputs are a filter `aa` and the size of its cube `na = (na(1), na(2), ...)`. Also input are two cube dimensions, that of the data last used by the filter `no1d` and that of the filter's next intended use `nd`. (`no1d` and `nd` are often the same). Module `bound` begins by defining a bigger data space with room for a filter surrounding the original data space `nd` on all sides. It does this by the line `nb=nd+2*na`. Then we allocate two data spaces `xx` and `yy` of the bigger size `nb` and pack many ones in a frame of width `na` around the outside of `xx`. The filter `aa` is also filled with ones. The filter `aa` must be regridded for the bigger `nb` data space (regridding merely changes the lag values of the ones). Now we filter the input `xx` with `aa` getting `yy`. Wherever the output is nonzero, we have an output that has been affected by the boundary. Such an output should not be computed. Thus we allocate the logical mask `aa->mis` (a part of the helix filter definition in module `helix` on page ??) and wherever we see a nonzero value of `yy` in the output, we designate the output as depending on missing inputs by setting `aa->mis` to `true`.

In reality one would set up the boundary conditions with module `bound` before identify-

user/gee/bound.c

```

1 void bound (int dim          /* number of dimensions */,
2             const int *nold /* old data coordinates [dim] */,
3             const int *nd   /* new data coordinates [dim] */,
4             const int *na   /* filter box size [dim] */,
5             const sf_filter aa /* helix filter */)
6 /*< Mark helix filter outputs where input is off data. >*/
7 {
8     int iy, my, ib, mb, i, nb[SF_MAX_DIM], ii[SF_MAX_DIM];
9     float *xx, *yy;
10
11     my = mb = 1;
12     for (i=0; i < dim; i++) {
13         nb[i] = nd[i] + 2*na[i]; /* nb is a bigger space. */
14         mb *= nb[i];
15         my *= nd[i];
16     }
17
18     xx = sf_floatalloc(mb); yy = sf_floatalloc(mb);
19
20     for (ib=0; ib < mb; ib++) {
21         sf_line2cart(dim, nb, ib, ii);
22         xx[ib] = 0.;
23         for (i=0; i < dim; i++)
24             if(ii[i]+1 <= na[i] || ii[i]+1 > nb[i]-na[i]) {
25                 xx[ib] = 1.;
26                 break;
27             }
28     }
29     sf_helicon_init(aa);
30     regrid(dim, nold, nb, aa);
31     for (i=0; i < aa->nh; i++) aa->flt[i] = 1.;
32     /* apply filter */
33     sf_helicon_lop(false, false, mb, mb, xx, yy);
34     regrid(dim, nb, nd, aa);
35     for (i=0; i < aa->nh; i++) aa->flt[i] = 0.;
36
37     aa->mis = sf_boolalloc(my); /* attach missing designation */
38     for (iy = 0; iy < my; iy++) { /* map to padded space */
39         sf_line2cart(dim, nd, iy, ii);
40         for (i=0; i < dim; i++) ii[i] += na[i];
41         ib = sf_cart2line(dim, nb, ii);
42         aa->mis[iy] = (bool) (yy[ib] > 0.);
43     }
44
45     free (xx); free (yy);
46 }

```



ing locations of missing data with module `misinput`. Both modules are based on the same concept, but the boundaries are more cluttered and confusing which is why we examined them later.

## Finding the prediction-error filter

The first stage of the least-squares estimation is computing the **prediction-error filter**. The second stage will be using it to find the missing data. The input data space contains a mixture of known data values and missing unknown ones. For the first stage of finding the filter, we generally have many more fitting equations than we need so we can proceed by ignoring the fitting equations that involve missing data values. We ignore them everywhere that the missing inputs hit the filter.

The codes here do not address the difficulty that maybe too much data is missing so that all weights are zero. To add stabilization we could supplement the data volume with a “training dataset” or by a “prior filter”. With things as they are, if there is not enough data to specify a prediction-error filter, you should encounter the error exit from `cgstep()` on page ??.

```

                                     user/gee/pef.c
1  void find_pef(int nd           /* data size */,
2                float* dd       /* data [nd] */,
3                sf_filter aa    /* estimated filter */,
4                int niter       /* number of iterations */)
5  /*< find PEF >*/
6  {
7      honest_init( dd, aa);
8      sf_solver(honest_lop, sf_cgstep, aa->nh, nd, aa->flt, dd,
9                niter, "x0", aa->flt, "end");
10     sf_cgstep_close();
11 }

```

## TWO-STAGE LINEAR LEAST SQUARES

In Chapter ?? and Chapter ?? we filled empty bins by minimizing the energy output from the filtered mesh. In each case there was arbitrariness in the choice of the filter. Here we find and use the optimum filter, the PEF.

The first stage is that of the previous section, finding the optimal PEF while carefully avoiding using any regression equations that involve boundaries or missing data. For the second stage, we take the PEF as known and find values for the empty bins so that the power out of the prediction-error filter is minimized. To do this we find missing data with module `mis2()` on page ??.

This two-stage method avoids the nonlinear problem we would otherwise face if we included the fitting equations containing both free data values and free filter values. Pre-

sumably, after two stages of linear least squares we are close enough to the final solution that we could switch over to the full nonlinear setup described near the end of this chapter.

The synthetic data in Figure 15 is a superposition of two plane waves of different directions, each with a random (but low-passed) waveform. After punching a hole in the data, we find that the lost data is pleasingly restored, though a bit weak near the side boundary. This imperfection could result from the side-boundary behavior of the operator or from an insufficient number of missing-data iterations.

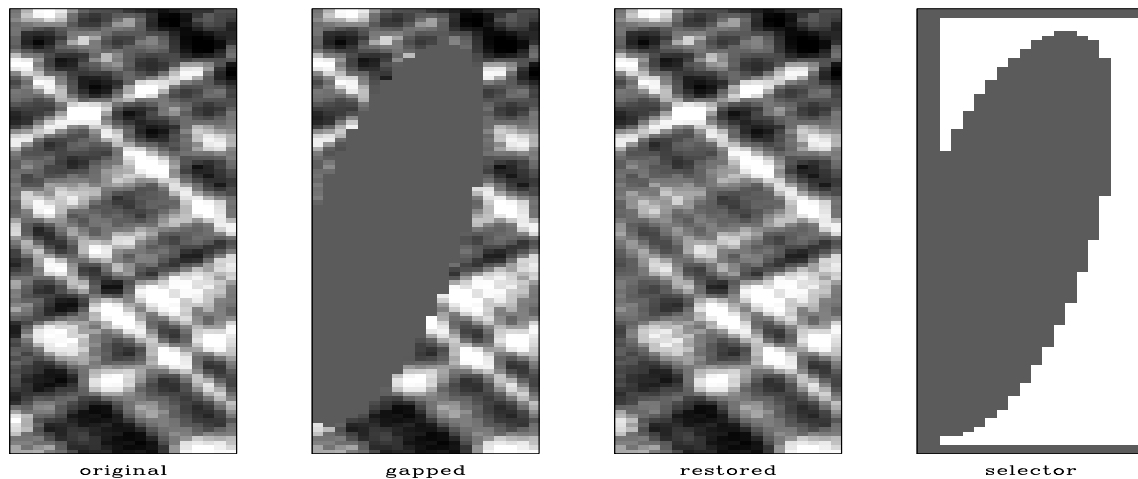


Figure 15: Original data (left), with a zeroed hole, restored, residual selector (right).

The residual selector in Figure 15 shows where the filter output has valid inputs. From it you can deduce the size and shape of the filter, namely that it matches up with Figure 14. The ellipsoidal hole in the residual selector is larger than that in the data because we lose regression equations not only at the hole, but where any part of the filter overlaps the hole.

The results in Figure 15 are essentially perfect representing the fact that that synthetic example fits the conceptual model perfectly. Before we look at the many examples in Figures 16-19 we will examine another gap-filling strategy.

### Adding noise (Geostat)

In chapter ?? we restored missing data by adopting the philosophy of minimizing the energy in filtered output. In this chapter we learned about an optimum filter for this task, the prediction-error filter (PEF). Let us name this method the “minimum noise” method of finding missing data.

A practical application with the minimum-noise method is evident in a large empty hole such as in Figures 16- 17. In such a void the interpolated data diminishes greatly. Thus we have not totally succeeded in the goal of “hiding our data acquisition footprint” which we would like to do if we are trying to make pictures of the earth and not pictures of our data acquisition footprint.

What we will do next is useful in some applications but not in others. Misunderstood

or misused it is rightly controversial. We are going to fill the empty holes with something that looks like the original data but really isn't. I will distinguish the words “**synthetic data**” (that derived from a physical model) from “**simulated data**” (that manufactured from a statistical model). We will fill the empty holes with simulated data like what you see in the center panels of Figures 3-9. We will add just enough of that “wall paper noise” to keep the variance constant as we move into the void.

Given some data  $\mathbf{d}$ , we use it in a filter operator  $\mathbf{D}$ , and as described with equation (29) we build a weighting function  $\mathbf{W}$  that throws out the broken regression equations (ones that involve missing inputs). Then we find a PEF  $\mathbf{a}$  by using this regression.

$$\mathbf{0} \approx \mathbf{r} = \mathbf{W}\mathbf{D}\mathbf{a} \quad (31)$$

Because of the way we defined  $\mathbf{W}$ , the “broken” components of  $\mathbf{r}$  vanish. We need to know the variance  $\sigma$  of the nonzero terms. It can be expressed mathematically in a couple different ways. Let  $\mathbf{1}$  be a vector filled with ones and let  $\mathbf{r}^2$  be a vector containing the squares of the components of  $\mathbf{r}$ .

$$\sigma = \sqrt{\frac{1}{N} \sum_i r_i^2} = \sqrt{\frac{\mathbf{1}'\mathbf{W}\mathbf{r}^2}{\mathbf{1}'\mathbf{W}\mathbf{1}}} \quad (32)$$

Let us go to a random number generator and get a noise vector  $\mathbf{n}$  filled with random numbers of variance  $\sigma$ . We'll call this the “added random noise”. Now we solve this new regression for the data space  $\mathbf{d}$  (both known and missing)

$$\mathbf{0} \approx \mathbf{r} = \mathbf{A}\mathbf{d} - \mathbf{n} \quad (33)$$

keeping in mind that known data is constrained (as detailed in chapter ??).

To understand why this works, consider first the training image, a region of known data. Although we might think that the data defines the white noise residual by  $\mathbf{r} = \mathbf{A}\mathbf{d}$ , we can also imagine that the white noise determines the data by  $\mathbf{d} = \mathbf{A}^{-1}\mathbf{r}$ . Then consider a region of wholly missing data. This data is determined by  $\mathbf{d} = \mathbf{A}^{-1}\mathbf{n}$ . Since we want the data variance to be the same in known and unknown locations, naturally we require the variance of  $\mathbf{n}$  to match that of  $\mathbf{r}$ .

A very minor issue remains. Regression equations may have all of their required input data, some of it, or none of it. Should the  $\mathbf{n}$  vector add noise to every regression equation? First, if a regression equation has all its input data that means there are no free variables so it doesn't matter if we add noise to that regression equation because the constraints will overcome that noise. I don't know if I should worry about how *many* inputs are missing for each regression equation.

It is fun making all this interesting “wall paper” noticing where it is successful and where it isn't. We cannot help but notice that it seems to work better with the genuine geophysical data than it does with many of the highly structured patterns. Geophysical data is expensive to acquire. Regrettably, we have uncovered a technology that makes counterfeiting much easier.

Examples are in Figures 16-19. In the electronic book, the right-side panel of each figure is a movie, each panel being derived from different random numbers.

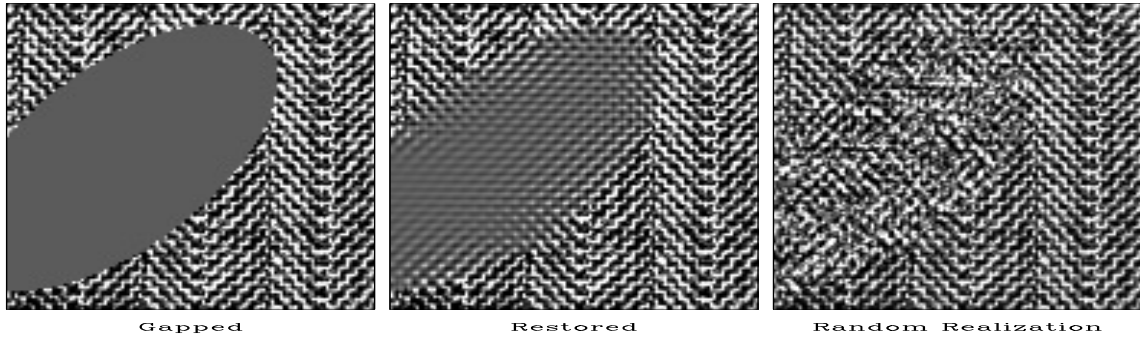


Figure 16: The herringbone texture is a patchwork of two textures. We notice that data missing from the hole tends to fill with the texture at the edge of the hole. The spine of the herring fish, however, is not modeled at all.

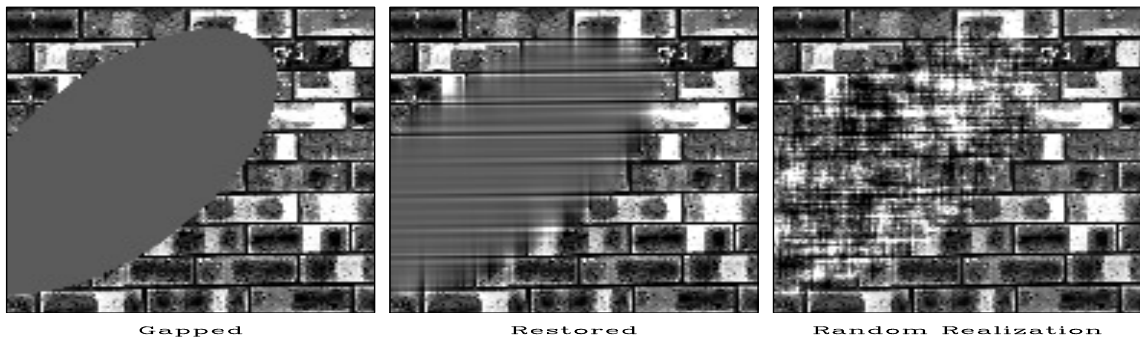


Figure 17: The brick texture has a mortar part (both vertical and horizontal joins) and a brick surface part. These three parts enter the empty area but do not end where they should.

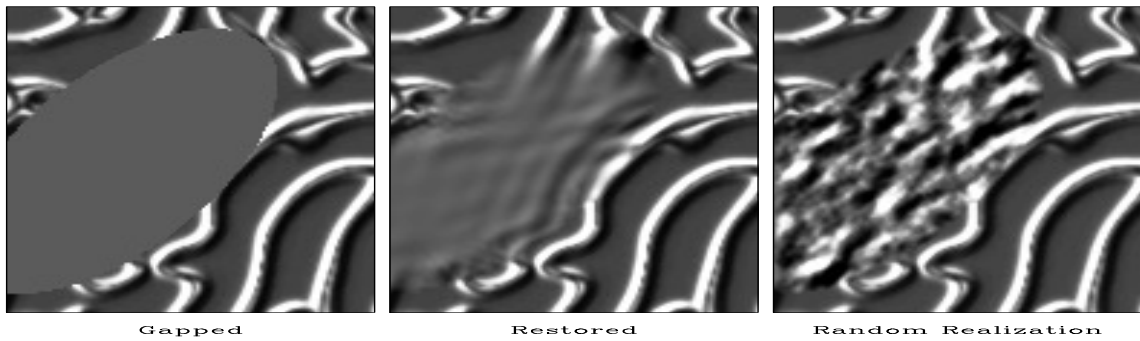


Figure 18: The theoretical model is a poor fit to the ridge data since the prediction must try to match ridges of all possible orientations. This data requires a broader theory which incorporates the possibility of nonstationarity (space variable slope).

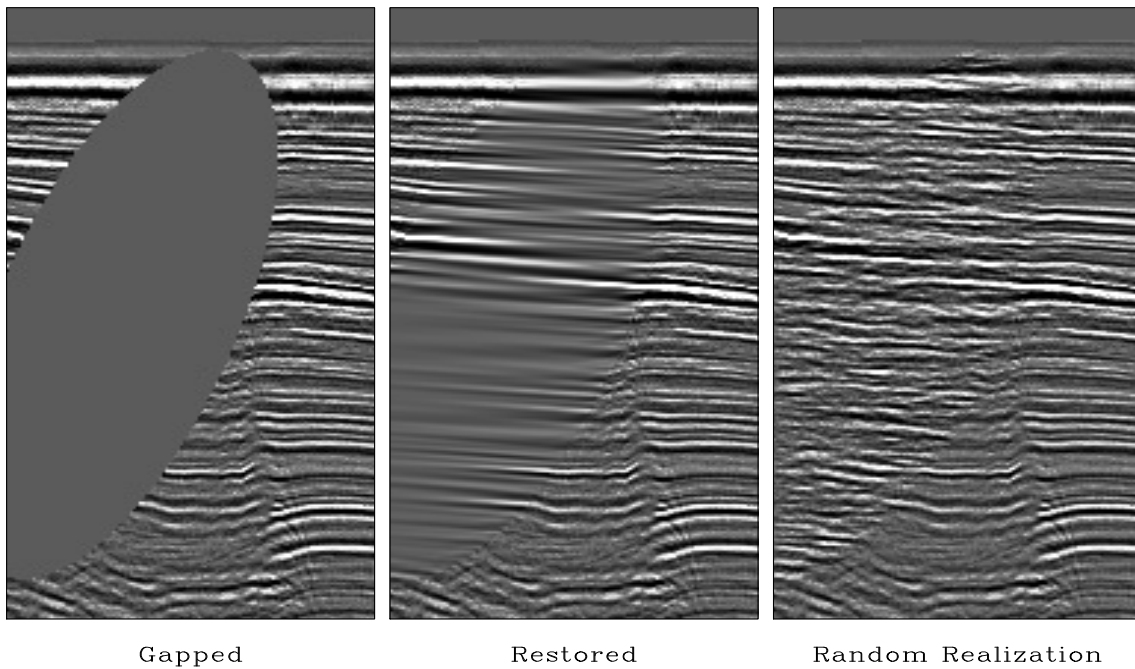


Figure 19: Filling the missing seismic data. The imaging process known as “migration” would suffer diffraction artifacts in the gapped data that it would not suffer on the restored data.

The seismic data in Figure 19 illustrates a fundamental principle: In the restored hole we do not see the same spectrum as we do on the other panels. This is because the hole is filled, not with all frequencies (or all slopes) but with those that are most predictable. The filled hole is devoid of the unpredictable noise that is a part of all real data.

### Inversions with geostat

In geophysical estimation (inversion) we use model styling (regularization) to handle the portion of the model that is not determined by the data. This results in the addition of minimal noise. Alternately, like in Geostatistics, we could make an assumption of statistical stationarity and add much more noise so the signal variance in poorly determined regions matches that in well determined regions. Here is how to do this. Given the usual data fitting and model styling goals

$$\mathbf{0} \approx \mathbf{Lm} - \mathbf{d} \quad (34)$$

$$\mathbf{0} \approx \mathbf{Am} \quad (35)$$

We introduce a sample of random noise  $\mathbf{n}$  and fit instead these regressions

$$\mathbf{0} \approx \mathbf{Lm} - \mathbf{d} \quad (36)$$

$$\mathbf{0} \approx \mathbf{Am} - \mathbf{n} \quad (37)$$

Of course you get a different solution for each different realization of the random noise. You also need to be a little careful to use noise  $\mathbf{n}$  of the appropriate variance. Figure 20 shows

a result on the SeaBeam data. **Bob Clapp** developed this idea at SEP and also applied it

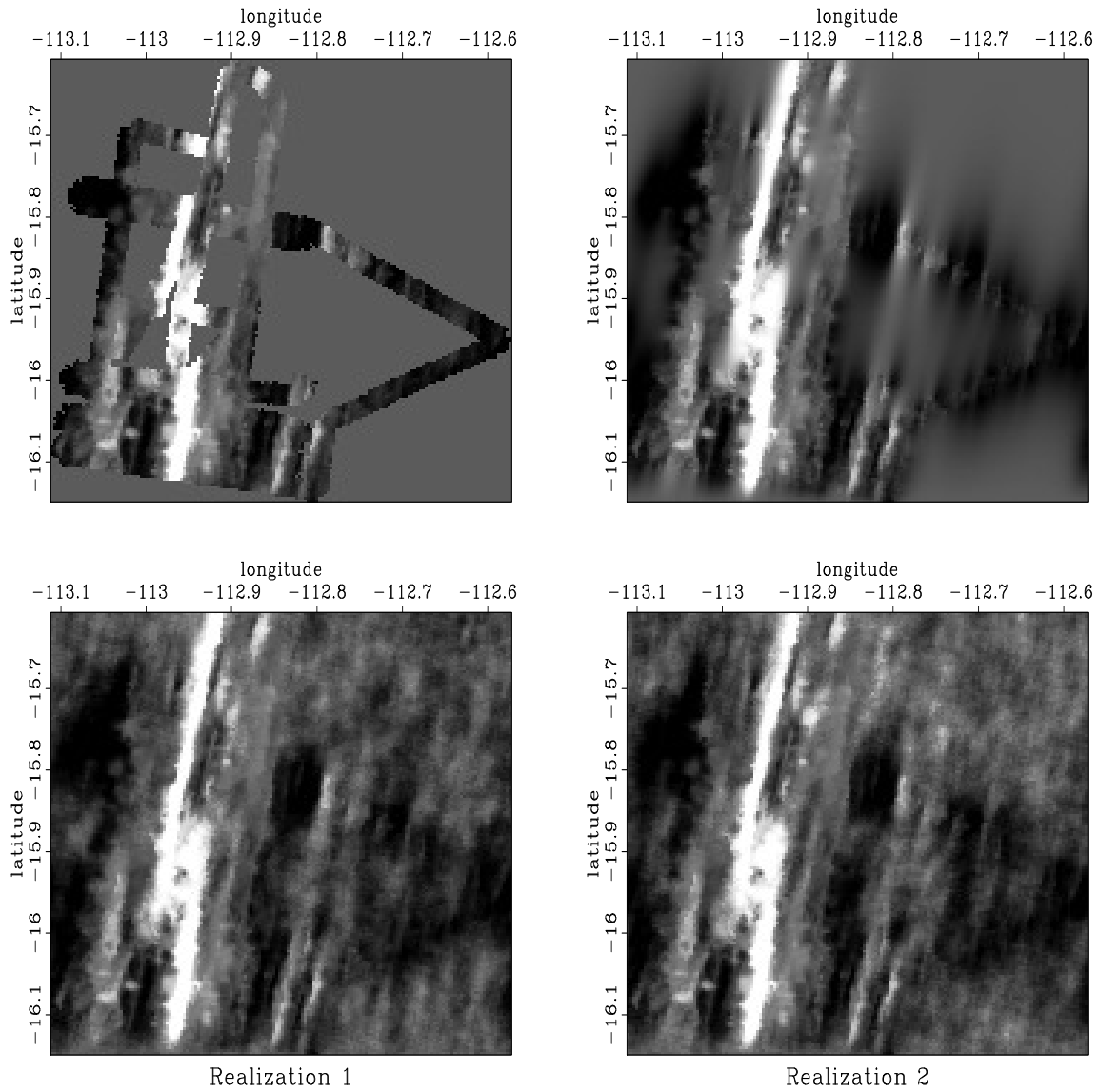


Figure 20: Top left is binned data. Top right extends the data with a PEF. The bottom two panels add appropriately colored random noise in the regions of missing data.

to interval velocity estimation, the example of Figures ??-??.

### Infill of 3-D seismic data from a quarry blast

Finding **missing data** (filling empty bins) requires use of a filter. Because of the helix, the codes work in spaces of all dimensions.

An open question is how many conjugate-direction iterations are needed in missing-data programs. When estimating filters, I set the **iteration count** `niter` at the number of free filter parameters. Theoretically, this gives me the exact solution but sometimes I run double

the number of iterations to be sure. The missing-data estimation, however is a completely different story. The number of free parameters in the missing-data estimation, could be very large. This often implies impractically long compute times for the exact solution. In practice I experiment carefully with `niter` and hope for the best. I find that where gaps are small, they fill in quickly. Where the gaps are large, they don't, and more iterations are required. Where the gaps are large is where we must experiment with preconditioning.

Figure 21 shows an example of replacing missing data by values predicted from a 3-D PEF. The data was recorded at **Stanford University** with a  $13 \times 13$  array of independent recorders. The figure shows 12 of the 13 lines each of length 13. Our main goal was to measure the ambient night-time noise. By morning about half the recorders had dead batteries but the other half recorded a wave from a quarry blast. The raw data was distracting to look at because of the many missing traces so I interpolated it with a small 3-D filter. That filter was a PEF.

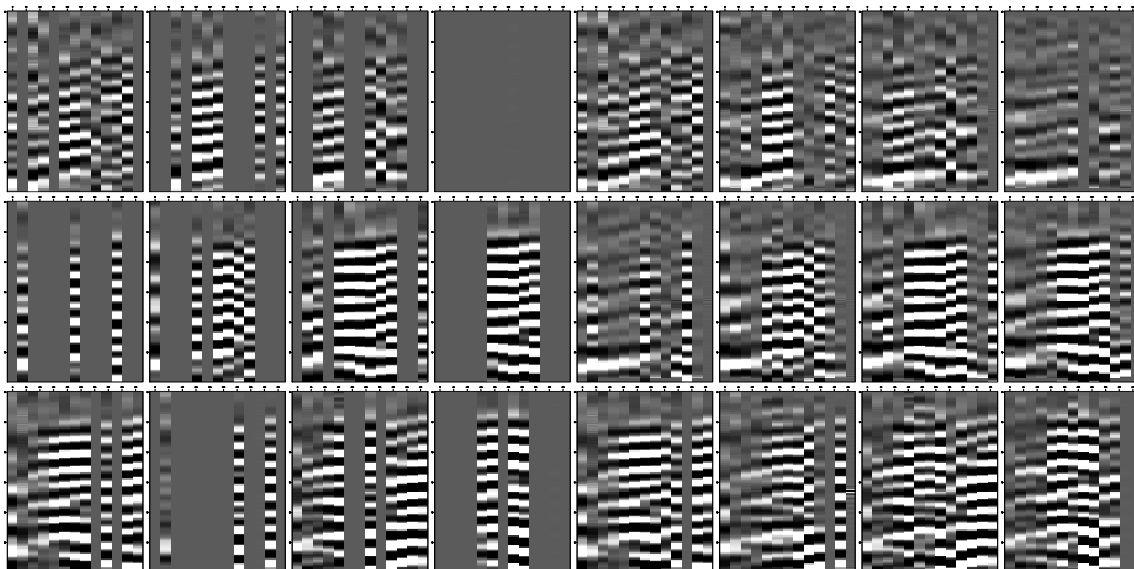


Figure 21: The left 12 panels are the inputs. The right 12 panels are outputs.

### Imposing prior knowledge of symmetry

Reversing a signal in time does not change its autocorrelation. In the analysis of stationary time series, it is well known (**FGDP**) that the filter for predicting forward in time should be the same as that for “predicting” backward in time (except for time reversal). When the data samples are short, however, a different filter may be found for predicting forward than for backward. Rather than average the two filters directly, the better procedure is to find the filter that minimizes the sum of power in two residuals. One is a filtering of the original signal, and the other is a filtering of a time-reversed signal, as in equation (38), where the top half of the equations represent prediction-error predicting forward in time

and the second half is prediction backward.

$$\begin{bmatrix} r_1 \\ r_2 \\ r_3 \\ r_4 \\ r_5 \\ r_6 \\ r_7 \\ r_8 \end{bmatrix} = \begin{bmatrix} y_3 & y_2 & y_1 \\ y_4 & y_3 & y_2 \\ y_5 & y_4 & y_3 \\ y_6 & y_5 & y_4 \\ y_1 & y_2 & y_3 \\ y_2 & y_3 & y_4 \\ y_3 & y_4 & y_5 \\ y_4 & y_5 & y_6 \end{bmatrix} \begin{bmatrix} 1 \\ a_1 \\ a_2 \end{bmatrix} \quad (38)$$

To get the bottom rows from the top rows, we simply reverse the order of all the components within each row. That reverses the input time function. (Reversing the order within a column would reverse the output time function.) Instead of the matrix being diagonals tipping  $45^\circ$  to the right, they tip to the left. We could make this matrix from our old familiar convolution matrix and a time-reversal matrix

$$\begin{bmatrix} 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 \end{bmatrix}$$

It is interesting to notice how time-reversal symmetry applies to Figure 15. First of all, with time going both forward and backward the residual space gets twice as big. The time-reversal part gives a selector for Figure 15 with a gap along the right edge instead of the left edge. Thus, we have acquired a few new regression equations.

Some of my research codes include these symmetries, but I excluded them here. Nowhere did I see that the reversal symmetry made noticeable difference in results, but in coding, it makes a noticeable clutter by expanding the residual to a two-component *residual array*.

Where a data sample grows exponentially towards the boundary, I expect that extrapolated data would diverge too. You can force it to go to zero (or any specified value) at some distance from the body of the known data. To do so, surround the body of data by missing data and surround that by specification of “enough” zeros. “Enough” is defined by the filter length.

## Hexagonal coordinates

In a two-dimensional plane it seems that the one-sidedness of the PEF could point in any direction. Since we usually have a rectangular mesh, however, we can only do the calculations along the axes so we have only two possibilities, the helix can wrap around the 1-axis, or it can wrap around the 2-axis.

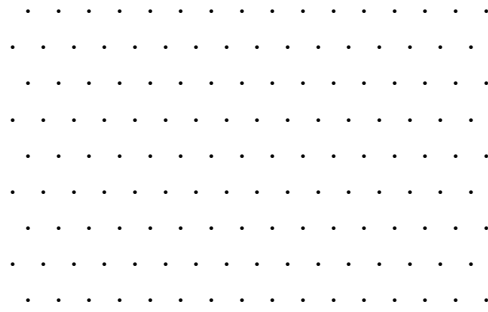
Suppose you acquire data on a hexagonal mesh as below

```

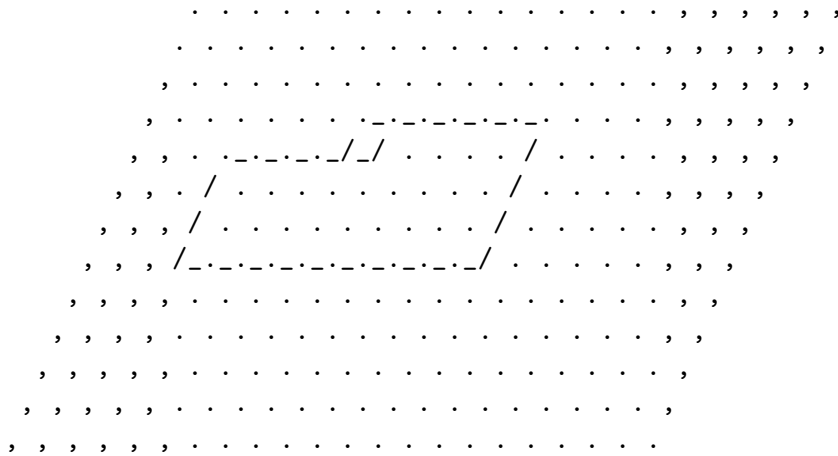
. . . . .
. . . . .
. . . . .

```





and some of the data values are missing. How can we apply the methods of this chapter? The solution is to append the given data by more missing data shown by the commas below.



Now we have a familiar two-dimensional coordinate system in which we can find missing values, as well as perform signal and noise separations as described in a later chapter.

### BOTH MISSING DATA AND UNKNOWN FILTER

Recall the missing-data figures beginning with Figure ???. There the filters were taken as known, and the only unknowns were the missing data. Now, instead of having a predetermined filter, we will solve for the filter along with the missing data. The principle we will use is that the output power is minimized while the filter is constrained to have one nonzero coefficient (else all the coefficients would go to zero). We will look first at some results and then see how they were found.

In Figure 22 the filter is constrained to be of the form  $(1, a_1, a_2)$ . The result is pleasing in that the interpolated traces have the same general character as the given values. The filter came out slightly different from the  $(1, 0, -1)$  that I guessed and tried in Figure ??. Curiously, constraining the filter to be of the form  $(a_{-2}, a_{-1}, 1)$  in Figure 23 yields the same interpolated missing data as in Figure 22. I understand that the sum squared of the coefficients of  $A(Z)P(Z)$  is the same as that of  $A(1/Z)P(Z)$ , but I do not see why that would imply the same interpolated data; never the less, it seems to.

Figure 22: Top is known data. Middle includes the interpolated values. Bottom is the filter with the left-most point constrained to be unity and other points chosen to minimize output power.

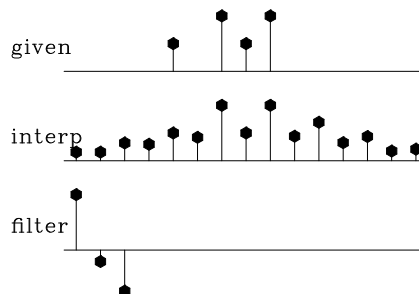
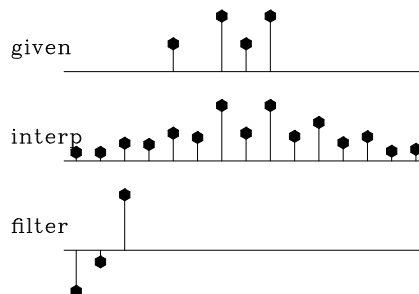


Figure 23: The filter here had its rightmost point constrained to be unity—i.e., this filtering amounts to backward prediction. The interpolated data seems to be identical to that of forward prediction.



## Objections to interpolation error

In any data interpolation or extrapolation, we want the extended data to behave like the original data. And, in regions where there is no observed data, the extrapolated data should drop away in a fashion consistent with its **spectrum** determined from the known region.

My basic idea is that the spectrum of the missing data should match that of the known data. This is the idea that the spectrum should be unchanging from a known region to an unknown region. A technical word to express the idea of spectra not changing is “**stationary**.” This happens with the PEF (one-sided filter) because its spectrum tends to the inverse of that of the known data while that of the unknown data tends to the inverse of that of the PEF. Thus the spectrum of the missing data is the “inverse of the inverse” of the spectrum of the known. The PEF enables us to fill in the missing area with the spectral shape of the known area. (In regions far away or unpredictable, the spectral shape may be the same, but the energy drops to zero.)

On the other hand, the **interpolation-error filter**, a filter like  $(a_{-2}, a_{-1}, 1, a_1, a_2)$ , should fail to do the job because it has the wrong spectrum. (I am stating this fact without proof).

To confirm and show these concepts, I prepared synthetic data consisting of a fragment of a damped exponential, and off to one side of it an impulse function. Most of the energy is in the damped exponential. Figure 24 shows that the spectrum and the extended data are about what we would expect. From the extrapolated data, it is impossible to see where the given data ends.

For comparison, I prepared Figure 25. It is the same as Figure 24, except that the filter is constrained in the middle. Notice that the extended data does *not* have the spectrum of the given data—the wavelength is much shorter. The boundary between real data and extended data is not nearly as well hidden as in Figure 24.

Figure 24: Top is synthetic data with missing portions. Middle includes the interpolated values. Bottom is the filter, a *prediction-error* filter which may look symmetric but is not quite.

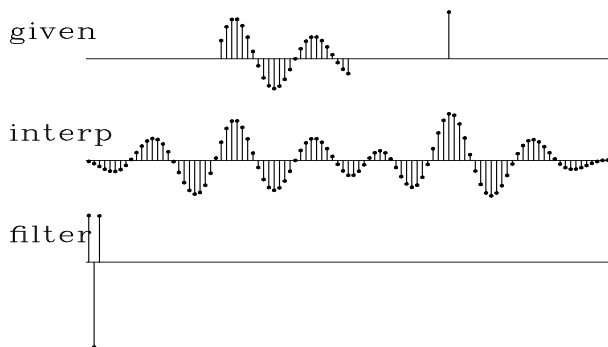
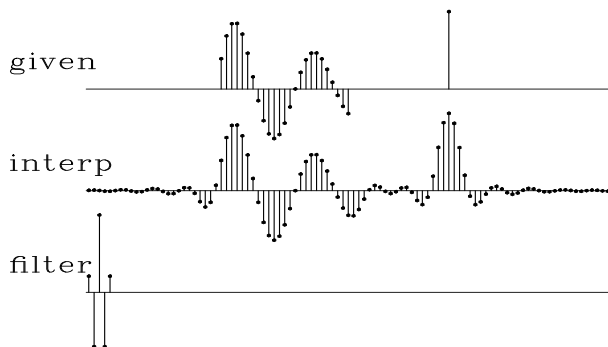


Figure 25: Top is the same synthetic data. Middle includes the interpolated values. Bottom is the filter, an *interpolation-error* filter.



## Packing both missing data and filter into a vector

Now let us examine the theory and coding behind the above examples. Define a roughening filter  $A(\omega)$  and a data signal  $Y(\omega)$  at some stage of interpolation. The fitting goal is  $0 \approx A(\omega)Y(\omega)$  where the filter  $A(\omega)$  has at least one time-domain coefficient constrained to be nonzero and the data contains both known and missing values. Think of perturbations  $\Delta A$  and  $\Delta Y$ . We neglect the nonlinear term  $\Delta A \Delta Y$  as follows:

$$0 \approx (A + \Delta A)(Y + \Delta Y) \quad (39)$$

$$0 \approx A \Delta Y + Y \Delta A + AY + \Delta A \Delta Y \quad (40)$$

$$0 \approx A \Delta Y + Y \Delta A + AY \quad (41)$$

Let us use matrix algebraic notation to rewrite the fitting goals (41). For this we need mask matrices (diagonal matrices with ones on the diagonal where variables are free and zeros where they are constrained i.e., where  $\Delta a_i = 0$  and  $\Delta y_i = 0$ ). The free-mask matrix for missing data is denoted  $\mathbf{J}$  and that for the PE filter is  $\mathbf{K}$ . The fitting goal (41) becomes

$$\mathbf{0} \approx \mathbf{A}\mathbf{J}\Delta\mathbf{y} + \mathbf{Y}\mathbf{K}\Delta\mathbf{a} + (\mathbf{A}\mathbf{y} \text{ or } \mathbf{Y}\mathbf{a}) \quad (42)$$

Defining the original residual as  $\bar{\mathbf{r}} = \mathbf{A}\mathbf{y}$  this becomes

$$\mathbf{0} \approx \begin{bmatrix} \mathbf{A}\mathbf{J} & \mathbf{Y}\mathbf{K} \end{bmatrix} \begin{bmatrix} \Delta\mathbf{y} \\ \Delta\mathbf{a} \end{bmatrix} + \bar{\mathbf{r}} \quad (43)$$

For a 3-term filter and a 7-point data signal, the fitting goal (42) becomes

$$\begin{bmatrix} a_0 & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & y_0 & \cdot & \cdot \\ a_1 & a_0 & \cdot & \cdot & \cdot & \cdot & \cdot & y_1 & y_0 & \cdot \\ a_2 & a_1 & a_0 & \cdot & \cdot & \cdot & \cdot & y_2 & y_1 & y_0 \\ \cdot & a_2 & a_1 & a_0 & \cdot & \cdot & \cdot & y_3 & y_2 & y_1 \\ \cdot & \cdot & a_2 & a_1 & a_0 & \cdot & \cdot & y_4 & y_3 & y_2 \\ \cdot & \cdot & \cdot & a_2 & a_1 & a_0 & \cdot & y_5 & y_4 & y_3 \\ \cdot & \cdot & \cdot & \cdot & a_2 & a_1 & a_0 & y_6 & y_5 & y_4 \\ \cdot & \cdot & \cdot & \cdot & \cdot & a_2 & a_1 & \cdot & y_6 & y_5 \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & a_2 & \cdot & \cdot & y_6 \end{bmatrix} \begin{bmatrix} \mathbf{J} & \mathbf{0} \\ \mathbf{0} & \mathbf{K} \end{bmatrix} \begin{bmatrix} \Delta y_0 \\ \Delta y_1 \\ \Delta y_2 \\ \Delta y_3 \\ \Delta y_4 \\ \Delta y_5 \\ \Delta y_6 \\ \Delta a_0 \\ \Delta a_1 \\ \Delta a_2 \end{bmatrix} + \begin{bmatrix} \bar{r}_0 \\ \bar{r}_1 \\ \bar{r}_2 \\ \bar{r}_3 \\ \bar{r}_4 \\ \bar{r}_5 \\ \bar{r}_6 \\ \bar{r}_7 \\ \bar{r}_8 \end{bmatrix} \approx \mathbf{0} \quad (44)$$

Recall that  $\bar{r}_t$  is the convolution of  $a_t$  with  $y_t$ , namely,  $\bar{r}_0 = y_0 a_0$  and  $\bar{r}_1 = y_0 a_1 + y_1 a_0$ , etc. To optimize this fitting goal we first initialize  $\mathbf{a} = (1, 0, 0, \dots)$  and then put zeros in for missing data in  $\mathbf{y}$ . Then we iterate over equations (45) to (49).

$$\mathbf{r} \leftarrow \mathbf{A}\mathbf{y} \quad (45)$$

$$\begin{bmatrix} \Delta \mathbf{y} \\ \Delta \mathbf{a} \end{bmatrix} \leftarrow \begin{bmatrix} \mathbf{J}^T \mathbf{A}^T \\ \mathbf{K}^T \mathbf{Y}^T \end{bmatrix} \mathbf{r} \quad (46)$$

$$\Delta \mathbf{r} \leftarrow \begin{bmatrix} \mathbf{A}\mathbf{J} & \mathbf{Y}\mathbf{K} \end{bmatrix} \begin{bmatrix} \Delta \mathbf{y} \\ \Delta \mathbf{a} \end{bmatrix} \quad (47)$$

$$\mathbf{y} \leftarrow \text{cgstep}(\mathbf{y}, \Delta \mathbf{y}) \quad (48)$$

$$\mathbf{a} \leftarrow \text{cgstep}(\mathbf{a}, \Delta \mathbf{a}) \quad (49)$$

This is the same idea as all the linear fitting goals we have been solving, except that now we recompute the residual  $\mathbf{r}$  inside the iteration loop so that as convergence is achieved (*if* it is achieved), the neglected nonlinear term  $\Delta \mathbf{A} \Delta \mathbf{Y}$  tends to zero.

My initial research proceeded by linearization like (41). Although I ultimately succeeded, I had enough difficulties that I came to realize that linearization is dangerous. When you start “far enough” from the correct solution the term  $\Delta \mathbf{A} \Delta \mathbf{Y}$  might not actually be small enough. You don’t know how small is small, because these are not scalars but operators. Then the solution may not converge to the minimum you want. Your solution will depend on where you start from. I no longer exhibit the nonlinear solver `missif` until I find a real data example where it produces noticeably better results than multistage linear-least squares.

The alternative to linearization is two-stage linear least squares. In the first stage you estimate the PEF; in the second you estimate the missing data. If need be, you can re-estimate the PEF using all the data both known and missing (downweighted if you prefer).

If you don’t have enough regression equations because your data is irregularly distributed, then you can use binning. Still not enough? Try coarser bins. The point is that nonlinear solvers will not work unless you begin close enough to the solution, and the way to

get close is by arranging first to solve a sensible (though approximate) linearized problem. Only as a last resort, after you have gotten as near as you can, should you use the nonlinear least-squares techniques.

## LEVELED INVERSE INTERPOLATION

Eighteenth- and nineteenth- century mathematics literature gives us many methods of interpolating functions. These classical methods are generally based on polynomials. The user specifies some order of polynomial and the theory works out the coefficients. Today our interest is in both interpolating and extrapolating wavefields (which are solutions to low order differential equations) and we use methods that are much better behaved than polynomials when extrapolating data, methods which behave acceptably when faced with contradictory data values, and methods which also apply in two and three dimensions.

In Chapter ??, subroutine `invint1()` on page ?? solved the problem of inverse linear interpolation, which is, given scattered data points, to find a function on a uniform mesh from which linear interpolation gives the scattered data points. To cope with regions having no data points, the subroutine requires an input roughening filter. This is a bit like specifying a differential equation to be satisfied between the data points. The question is, how should we choose a roughening filter? The importance of the roughening filter grows as the data gets sparser or as the mesh is refined.

Figures 22-25 suggest that the choice of the roughening filter need not be subjective, nor a priori, but that the prediction-error filter (PEF) is the ideal roughening filter. Spectrally, the PEF tends to the inverse of its input hence its output tends to be “level”. Missing data that is interpolated with this “leveler” tends to have the spectrum of given data.

### Test results for leveled inverse interpolation

Figures 26 and 27 show the same example as in Figures ?? and ??. What is new here is that the proper PEF is not given but is determined from the data. Figure 26 was made with a three-coefficient filter  $(1, a_1, a_2)$  and Figure 27 was made with a five-coefficient filter  $(1, a_1, a_2, a_3, a_4)$ . The main difference in the figures is where the data is sparse. The data points in Figures ??, 26 and 27 are samples from a sinusoid.

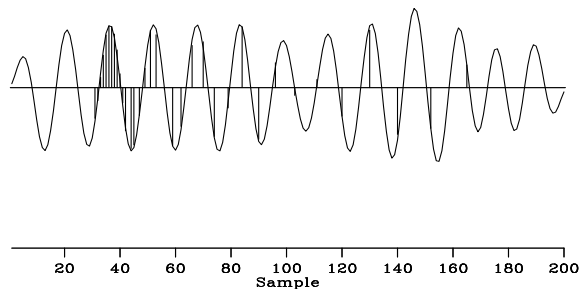


Figure 26: Interpolating with a three-term filter. The interpolated signal is fairly monofrequency.

Comparing Figures ?? and ?? to Figures 26 and 27 we conclude that by finding and imposing the prediction-error filter while finding the model space, we have interpolated beyond **aliasing** in data space.



for  $\mathbf{m}$  and  $\mathbf{a}$ . The guess  $\mathbf{m} = 0$  is satisfactory (as explained later). For the first guess of the filter, I suggest you load it up with  $\mathbf{a} = (1, -2, 1)$  as I did for the examples here.

### Seabeam: theory to practice

I provide here a more fundamental theory for dealing with the Seabeam data. I originally approached the data in this more fundamental way, but with time, I realized that I paid a high price in code complexity, computational speed, and reliability. The basic problem is that the elegant theory requires a good starting model which can only come from the linearized theory. I briefly recount the experience here, because the fundamental theory is interesting and because in other applications, you will face the challenge of sorting out the fundamental features from the essential features.

The linear-interpolation operator carries us from a uniform mesh to irregularly distributed data. Fundamentally we seek to solve the inverse problem to go the other direction. A nonlinear approach to filling in the missing data is suggested by the one-dimensional examples in Figures 26–27, where the PEF and the missing data are estimated simultaneously. The nonlinear approach has the advantage that it allows for completely arbitrary data positioning, whereas the two-stage linear approach forces the data to be on a uniform mesh and requires there not be too many empty mesh locations.

For the 2-D nonlinear application, we follow the same approach we used in one dimension, equations (50) and (51), except that the filtering and the linear interpolations are two dimensional.

I have had considerable experience with this problem on this data set and I can report that bin filling is easier and works much more quickly and reliably. Eventually I realized that the best way to start the nonlinear iteration (50-51) is with the final result of bin filling. Then I learned that the extra complexity of the nonlinear iteration (50-51) offers little apparent improvement to the quality of the SeaBeam result. (This is not to say that we should not try more variations on the idea).

Not only did I find the binning method faster, but I found it to be *much* faster (compare a minute to an hour). The reasons for being faster (most important first) are,

1. Binning reduces the amount of data handled in each iteration by a factor of the average number of points per bin.
2. The 2-D linear interpolation operator adds many operations per data point.
3. Using two fitting goals seems to require more iterations.

(Parenthetically, I later found that helix preconditioning speeds the Seabeam interpolation from minutes to seconds.)

The most serious criticism of the nonlinear approach is that it does not free us from the linearized approaches. We need them to get a “close enough” starting solution to the nonlinear problem. I learned that the iteration (50-51), like most nonlinear sequences, behaves unexpectedly and badly when you start too far from the desired solution. For example, I often began from the assumed PEF being a Laplacian and the original map

being fit from that. Oddly, from this starting location I sometimes found myself stuck. The iteration (50-51) would not move towards the map we humans consider a better one.

Having said all those bad things about iteration (50-51), I must hasten to add that with a different type of data set, you might find the results of (50-51) to be significantly better.

### Risky ways to do nonlinear optimization

I have noticed that some geophysicists have adopted a risky method of nonlinear optimization, which is not advocated in the professional optimization literature. This risky method is to linearize a goal (with a multiparameter model space), then optimize the linearized goal, then relinearize, etc. The safer method is to relinearize after each step of CD.

An instructive example I learned about many years ago was earthquake epicenter location. Model space is latitude, longitude, and origin time. When people added a new variable, the depth, the solutions went wild until they learned to restrict the depth to zero until the other three parameters were stabilized. Apparently the instability stems from the fact that depth and origin time affect distant receivers in a similar way.

### The bane of PEF estimation

This is the place where I would like to pat myself on the back for having “solved” the problem of missing data. Actually, an important practical problem remains. I’ve been trying to coax younger, more energetic people to think about it. The problem arises when there is too much missing data.

The bane of PEF estimation is too much missing data.

Then *all* the regression equations disappear. The nonlinear methods are particularly bad because if they don’t have a good enough starting location, they can and do go crazy. My only suggestion is to begin with a linear PEF estimator. Shrink the PEF and coarsen the mesh in model space until you do have enough equations. Starting from there, hopefully you can refine this crude solution without dropping into a local minimum.

Another important practical problem remains, that of nonstationarity. We’ll see the beginnings of the solution to that problem in chapter ??.

## MULTIVARIATE SPECTRUM

A common **spectrum** is the Fourier spectrum. More fundamentally, a spectrum is a decomposition of a model space or data space into components. The components are in some sense independent; more specifically, the components are orthogonal to one another. Another well-known spectrum is provided by eigenvectors and eigenvalues. In statistical signal processing we handle a third type of spectrum, the multivariate spectrum.

Working in an optimization application, we begin from residuals between theory and practice. These residuals can be scaled to make new optimization residuals before we start



minimizing their energy. What scaling should we use? The scaling can be a simple weighting function or a filter. A filter is simply a weighting function in Fourier space.

The basic idea of common sense, which also comes to us as results proven by Gauss or from the theory of statistical signal processing, is this: The optimization residuals should be roughly of equal scale. This makes sense because squaring magnifies scale, and anything small will be ignored while anything large will dominate. Scaling optimization residuals to be in a common range makes them all equally influential on the final solution. Not only should optimization residuals be of like scale in physical space, they should be of like scale in Fourier space or eigenvector space, or any other space that we might use to represent the optimization residuals. This implies that the optimization residuals should be uncorrelated. If the optimization residuals were correlated, they would have a spectrum that was not white. Not white means of differing sizes in Fourier space. Residuals should be the same size as one another in physical space, likewise in Fourier space. Thus the optimization residuals should be orthogonal and of unit scale, much like Fourier components or as eigenvectors are orthonormal.

Let us approach the problem backwards. Suppose we have two random variables that we take to be the ideal optimization residuals  $x_1$  and  $x_2$ . In reality the two may be few or trillions. In the language of statistics, the optimization residuals are expected to have zero mean, an idea that is formalized by writing  $E(x_1) = 0$  and  $E(x_2) = 0$ . Likewise these ideal optimization residuals have equal energy,  $E(x_1^2) = 1$  and  $E(x_2^2) = 1$ . Finally, these two optimization residuals are uncorrelated, a condition which is written as  $E(x_1x_2) = 0$ . The expectation symbol  $E()$  is like a summation over many instances of the random variable.

Now suppose there exists a transformation  $\mathbf{B}$  from these ideal optimization residuals to two experimental residuals  $y_1$  and  $y_2$ , say  $\mathbf{y} = \mathbf{B}\mathbf{x}$  where

$$\begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = \begin{bmatrix} b_{11} & b_{12} \\ b_{21} & b_{22} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \quad (53)$$

The experimental residuals  $y_1$  and  $y_2$  are likely to be neither orthogonal nor equal in energy. From the column vector  $\mathbf{y}$ , the experimenter can form a square matrix. Let us also allow the experimenter to write the symbol  $E()$  to denote summation over many trials or over many sections of data, ranges over time or space, over soundings or over receiver locations. The experimenter writes

$$\mathbf{R} = E(\mathbf{y}\mathbf{y}^T) \quad (54)$$

$$\mathbf{R} = E(\mathbf{B}\mathbf{x}\mathbf{x}^T\mathbf{B}^T) \quad (55)$$

Given a random variable  $r$ , the expectation of  $2r$  is simply  $E(2r) = 2E(r)$ . The  $E()$  symbol is a summation on random variables, but constants like the coefficients of  $\mathbf{B}$  pass right through it. Thus,

$$\mathbf{R} = \mathbf{B} E(\mathbf{x}\mathbf{x}^T) \mathbf{B}^T \quad (56)$$

$$\mathbf{R} = \mathbf{B} E\left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \begin{bmatrix} x_1 & x_2 \end{bmatrix}\right) \mathbf{B}^T \quad (57)$$

$$\mathbf{R} = \mathbf{B} \begin{bmatrix} E(x_1x_1) & E(x_1x_2) \\ E(x_2x_1) & E(x_2x_2) \end{bmatrix} \mathbf{B}^T \quad (58)$$

$$\mathbf{R} = \mathbf{B}\mathbf{B}^T \quad (59)$$

Given a matrix  $\mathbf{R}$ , there is a simple well-known method called the **Cholesky factorization** method that will factor  $\mathbf{R}$  into two parts like  $\mathbf{B}$  and  $\mathbf{B}^T$ . The method creates for us either an upper or a lower triangular matrix (our choice) for  $\mathbf{B}$ . You can easily reinvent the Cholesky method if you multiply the symbols for two triangular matrices like  $\mathbf{B}$  and  $\mathbf{B}^T$  and notice the procedure that works backwards from  $\mathbf{R}$  to  $\mathbf{B}$ . The experimenter seeks not  $\mathbf{B}$ , however, but its inverse, the matrix that takes us from the experimental residuals to the ideal optimization residuals that are uncorrelated and of equal energies. The Cholesky factorization costs  $N^3$  computations, which is about the same as the cost of the matrix inversion of  $\mathbf{R}$  or  $\mathbf{B}$ . For geophysical maps and other functions on Cartesian spaces, the Prediction Error Filter (PEF) accomplishes the same general goal and has the advantage that we have already learned how to perform the operation using operators instead of matrices.

The multivariate spectrum of experimental residuals  $\mathbf{y}$  is the matrix  $\mathbf{R} = E(\mathbf{y}\mathbf{y}^T)$ . For optimum model finding, the experimental residuals (squared) should be weighted inversely (matrix inverse) by their multivariate spectrum.

If I were a little stronger at analysis (or rhetoric) I would tell you that the optimizers preconditioned variable  $\mathbf{p}$  is the statisticians IID (Independent Identically Distributed) random variable. For stationary (statistically constant) signals and images,  $\mathbf{A}_m$  is the model-space PEF. Echo soundings and **interval velocity** have statistical properties that change with depth. There  $\mathbf{A}_m$  is a diagonal weighting matrix (perhaps before or after a PEF).

## What should we optimize?

Least-squares applications often present themselves as fitting goals such as

$$\mathbf{0} \approx \mathbf{F}\mathbf{m} - \mathbf{d} \quad (60)$$

$$\mathbf{0} \approx \mathbf{m} \quad (61)$$

To balance our possibly contradictory goals we need weighting functions. The quadratic form that we should minimize is

$$\min_m (\mathbf{F}\mathbf{m} - \mathbf{d})^T \mathbf{A}_n^T \mathbf{A}_n (\mathbf{F}\mathbf{m} - \mathbf{d}) + \mathbf{m}^T \mathbf{A}_m^T \mathbf{A}_m \mathbf{m} \quad (62)$$

where  $\mathbf{A}_n^T \mathbf{A}_n$  is the inverse multivariate spectrum of the noise (data-space residuals) and  $\mathbf{A}_m^T \mathbf{A}_m$  is the inverse multivariate spectrum of the model. In other words,  $\mathbf{A}_n$  is a leveler on the data fitting error and  $\mathbf{A}_m$  is a leveler on the model. There is a curious unresolved issue: What is the most suitable constant scaling ratio of  $\mathbf{A}_n$  to  $\mathbf{A}_m$ ?

## Confusing terminology for data covariance

Confusion often stems from the mean of the data  $E(\mathbf{d})$ .

An experimentalist would naturally believe that the expectation of the data is solely a function of the data, that it can be estimated by averaging data.

On the other hand, a theoretician’s idea of the expectation of the observational data  $E(\mathbf{d})$  is that it is the theoretical data  $\mathbf{Fm}$ , that the expectation of the data  $E(\mathbf{d}) = \mathbf{Fm}$  is a function of the model. The theoretician thinks this way because of the idea of noise  $\mathbf{n} = \mathbf{Fm} - \mathbf{d}$  as having zero mean.

Seismological data is highly complex but also highly reproducible. In studies like seismology, the world is deterministic but more complicated than our ability to model. Thus, as a practical matter, the discrepancy between observational data and theoretical data is more realistically attributed to the theoretical data. It is not adequately modeled and computed.

This superficial difference in viewpoint becomes submerged to a more subtle level by statistical textbooks that usually define weighting functions in terms of variances instead of spectra. This is particularly confusing with the noise spectrum  $(\mathbf{A}_n^T \mathbf{A}_n)^{-1}$ . It is often referred to as the “**data covariance**” defined as  $E[(\mathbf{d} - E(\mathbf{d}))(\mathbf{d} - E(\mathbf{d}))^T]$ . Clearly, the noise spectrum is the same as the data covariance only if we accept the theoretician’s definition that  $E(\mathbf{d}) = \mathbf{Fm}$ .

There is no ambiguity and no argument if we drop the word “variance” and use the word “spectrum”. Thus, (1) the “inverse **noise spectrum**” is the appropriate weighting for data-space residuals; and (2) the “inverse **model spectrum**” is the appropriate model-space weighting. Theoretical expositions generally require these spectra to be given as “**prior information**.” In this book we see how, when the model space is a map, we can solve for the “prior information” along with everything else.

The statistical words “covariance matrix” are suggestive and appealing, but I propose not to use them because of the ambiguity of  $E(\mathbf{d})$ . For example, we understand that people who say “data covariance” intend the “multivariate noise spectrum” but we cannot understand their meaning of “model covariance”. They should intend the “multivariate model spectrum” but that implies that  $E(\mathbf{m}) = \mathbf{0}$ , which seems wrong. Avoiding the word “covariance” avoids the problem.

## Hermeneutics

In seismology the data is usually better than the theory.

Hermeneutics is the study of the methodological principles of interpretation. Historically, it refers to bible study. Never-the-less, it seems entirely appropriate for Geophysical Estimation. If Albert’s book is “Inverse Problem Theory” and mine is “Inverse Problem Practice”, and if the difference between theory and practice is smaller in theory than it is in practice, then there are two fundamental questions:

1. In theory, what is the difference between theory and practice? In theory, the difference is data error.
2. In practice, what is the difference between theory and practice? One suggestion is that the discrepancy is entirely due to inadequate modeling. It is well known that geophysical data is highly repeatable. The problem is that the modeling neglects far too much.

Here is a perspective drawn from analysis of the human genome: “The problem is that it is possible to use empirical data to calibrate a model that generates simulated data that is similar to the empirical data. The point of using such a calibrated model is to be able to show how strange certain regions are if they don’t fit the simulated distribution, which is based on the empirical distribution.” In other words, “inversion” is just the process of calibrating a model. To learn something new we study the *failures* of such models.