

The single central idea in inverse theory is the prescription

$$\text{minimize: } \mathcal{A} + \lambda\mathcal{B} \quad (18.4.12)$$

for various values of  $0 < \lambda < \infty$  along the so-called trade-off curve (see Figure 18.4.1), and then to settle on a “best” value of  $\lambda$  by one or another criterion, ranging from fairly objective (e.g., making  $\chi^2 = N$ ) to entirely subjective. Successful methods, several of which we will now describe, differ as to their choices of  $\mathcal{A}$  and  $\mathcal{B}$ , as to whether the prescription (18.4.12) yields linear or nonlinear equations, as to their recommended method for selecting a final  $\lambda$ , and as to their practicality for computer-intensive two-dimensional problems like image processing.

They also differ as to the philosophical baggage that they (or rather, their proponents) carry. We have thus far avoided the word “Bayesian.” (Courts have consistently held that academic license does not extend to shouting “Bayesian” in a crowded lecture hall.) But it is hard, nor have we any wish, to disguise the fact that  $\mathcal{B}$  has something to do with *a priori* expectation, or knowledge, of a solution, while  $\mathcal{A}$  has something to do with *a posteriori* knowledge. The constant  $\lambda$  adjudicates a delicate compromise between the two. Some inverse methods have acquired a more Bayesian stamp than others, but we think that this is purely an accident of history. An outsider looking only at the equations that are actually solved, and not at the accompanying philosophical justifications, would have a difficult time separating the so-called Bayesian methods from the so-called empirical ones, we think.

The next three sections discuss three different approaches to the problem of inversion, which have had considerable success in different fields. All three fit within the general framework that we have outlined, but they are quite different in detail and in implementation.

#### CITED REFERENCES AND FURTHER READING:

- Craig, I.J.D., and Brown, J.C. 1986, *Inverse Problems in Astronomy* (Bristol, U.K.: Adam Hilger).  
 Twomey, S. 1977, *Introduction to the Mathematics of Inversion in Remote Sensing and Indirect Measurements* (Amsterdam: Elsevier).  
 Tikhonov, A.N., and Arsenin, V.Y. 1977, *Solutions of Ill-Posed Problems* (New York: Wiley).  
 Tikhonov, A.N., and Goncharsky, A.V. (eds.) 1987, *Ill-Posed Problems in the Natural Sciences* (Moscow: MIR).  
 Parker, R.L. 1977, *Annual Review of Earth and Planetary Science*, vol. 5, pp. 35–64.  
 Frieden, B.R. 1975, in *Picture Processing and Digital Filtering*, T.S. Huang, ed. (New York: Springer-Verlag).  
 Tarantola, A. 1987, *Inverse Problem Theory* (Amsterdam: Elsevier).  
 Baumeister, J. 1987, *Stable Solution of Inverse Problems* (Braunschweig, Germany: Friedr. Vieweg & Sohn) [mathematically oriented].  
 Titterton, D.M. 1985, *Astronomy and Astrophysics*, vol. 144, pp. 381–387.  
 Jeffrey, W., and Rosner, R. 1986, *Astrophysical Journal*, vol. 310, pp. 463–472.

## 18.5 Linear Regularization Methods

What we will call *linear regularization* is also called the *Phillips-Twomey method* [1,2], the *constrained linear inversion method* [3], the *method of regularization* [4], and *Tikhonov-Miller regularization* [5-7]. (It probably has other names also,

since it is so obviously a good idea.) In its simplest form, the method is an immediate generalization of zeroth-order regularization (equation 18.4.11, above). As before, the functional  $\mathcal{A}$  is taken to be the  $\chi^2$  deviation, equation (18.4.9), but the functional  $\mathcal{B}$  is replaced by more sophisticated measures of smoothness that derive from first or higher derivatives.

For example, suppose that your *a priori* belief is that a credible  $u(x)$  is not too different from a constant. Then a reasonable functional to minimize is

$$\mathcal{B} \propto \int [\hat{u}'(x)]^2 dx \propto \sum_{\mu=1}^{M-1} [\hat{u}_\mu - \hat{u}_{\mu+1}]^2 \tag{18.5.1}$$

since it is nonnegative and equal to zero only when  $\hat{u}(x)$  is constant. Here  $\hat{u}_\mu \equiv \hat{u}(x_\mu)$ , and the second equality (proportionality) assumes that the  $x_\mu$ 's are uniformly spaced. We can write the second form of  $\mathcal{B}$  as

$$\mathcal{B} = |\mathbf{B} \cdot \hat{\mathbf{u}}|^2 = \hat{\mathbf{u}} \cdot (\mathbf{B}^T \cdot \mathbf{B}) \cdot \hat{\mathbf{u}} \equiv \hat{\mathbf{u}} \cdot \mathbf{H} \cdot \hat{\mathbf{u}} \tag{18.5.2}$$

where  $\hat{\mathbf{u}}$  is the vector of components  $\hat{u}_\mu$ ,  $\mu = 1, \dots, M$ ,  $\mathbf{B}$  is the  $(M - 1) \times M$  first difference matrix

$$\mathbf{B} = \begin{pmatrix} -1 & 1 & 0 & 0 & 0 & 0 & 0 & \dots & 0 \\ 0 & -1 & 1 & 0 & 0 & 0 & 0 & \dots & 0 \\ \vdots & & & \ddots & & & & & \vdots \\ 0 & \dots & 0 & 0 & 0 & 0 & -1 & 1 & 0 \\ 0 & \dots & 0 & 0 & 0 & 0 & 0 & -1 & 1 \end{pmatrix} \tag{18.5.3}$$

and  $\mathbf{H}$  is the  $M \times M$  matrix

$$\mathbf{H} = \mathbf{B}^T \cdot \mathbf{B} = \begin{pmatrix} 1 & -1 & 0 & 0 & 0 & 0 & 0 & \dots & 0 \\ -1 & 2 & -1 & 0 & 0 & 0 & 0 & \dots & 0 \\ 0 & -1 & 2 & -1 & 0 & 0 & 0 & \dots & 0 \\ \vdots & & & \ddots & & & & & \vdots \\ 0 & \dots & 0 & 0 & 0 & -1 & 2 & -1 & 0 \\ 0 & \dots & 0 & 0 & 0 & 0 & -1 & 2 & -1 \\ 0 & \dots & 0 & 0 & 0 & 0 & 0 & -1 & 1 \end{pmatrix} \tag{18.5.4}$$

Note that  $\mathbf{B}$  has one fewer row than column. It follows that the symmetric  $\mathbf{H}$  is degenerate; it has exactly one zero eigenvalue corresponding to the *value* of a constant function, any one of which makes  $\mathcal{B}$  exactly zero.

If, just as in §15.4, we write

$$A_{i\mu} \equiv R_{i\mu}/\sigma_i \quad b_i \equiv c_i/\sigma_i \tag{18.5.5}$$

then, using equation (18.4.9), the minimization principle (18.4.12) is

$$\text{minimize: } \mathcal{A} + \lambda\mathcal{B} = |\mathbf{A} \cdot \hat{\mathbf{u}} - \mathbf{b}|^2 + \lambda \hat{\mathbf{u}} \cdot \mathbf{H} \cdot \hat{\mathbf{u}} \tag{18.5.6}$$

This can readily be reduced to a linear set of *normal equations*, just as in §15.4: The components  $\hat{u}_\mu$  of the solution satisfy the set of  $M$  equations in  $M$  unknowns,

$$\sum_{\rho} \left[ \left( \sum_i A_{i\mu} A_{i\rho} \right) + \lambda H_{\mu\rho} \right] \hat{u}_\rho = \sum_i A_{i\mu} b_i \quad \mu = 1, 2, \dots, M \tag{18.5.7}$$

World Wide Web sample page from NUMERICAL RECIPES IN C: THE ART OF SCIENTIFIC COMPUTING (ISBN 0-521-43108-5)  
Copyright (C) 1988-1992 by Cambridge University Press. Programs Copyright (C) 1988-1992 by Numerical Recipes Software.  
Permission is granted for internet users to make one paper copy for their own personal use. Further reproduction, or any copying of machine-readable files (including this one), to any server computer, is strictly prohibited. To order Numerical Recipes books, diskettes, or CDROMs visit website <http://www.nr.com> or call 1-800-872-7423 (North America only), or send email to [trade@cup.cam.ac.uk](mailto:trade@cup.cam.ac.uk) (outside North America).

or, in vector notation,

$$(\mathbf{A}^T \cdot \mathbf{A} + \lambda \mathbf{H}) \cdot \hat{\mathbf{u}} = \mathbf{A}^T \cdot \mathbf{b} \quad (18.5.8)$$

Equations (18.5.7) or (18.5.8) can be solved by the standard techniques of Chapter 2, e.g., *LU* decomposition. The usual warnings about normal equations being ill-conditioned do not apply, since the whole purpose of the  $\lambda$  term is to cure that same ill-conditioning. Note, however, that the  $\lambda$  term *by itself* is ill-conditioned, since it does not select a preferred constant value. You hope your data can at least do *that!*

Although inversion of the matrix  $(\mathbf{A}^T \cdot \mathbf{A} + \lambda \mathbf{H})$  is not generally the best way to solve for  $\hat{\mathbf{u}}$ , let us digress to write the solution to equation (18.5.8) schematically as

$$\hat{\mathbf{u}} = \left( \frac{1}{\mathbf{A}^T \cdot \mathbf{A} + \lambda \mathbf{H}} \cdot \mathbf{A}^T \cdot \mathbf{A} \right) \mathbf{A}^{-1} \cdot \mathbf{b} \quad (\text{schematic only!}) \quad (18.5.9)$$

where the identity matrix in the form  $\mathbf{A} \cdot \mathbf{A}^{-1}$  has been inserted. This is schematic not only because the matrix inverse is fancifully written as a denominator, but also because, in general, the inverse matrix  $\mathbf{A}^{-1}$  does not exist. However, it is illuminating to compare equation (18.5.9) with equation (13.3.6) for optimal or Wiener filtering, or with equation (13.6.6) for general linear prediction. One sees that  $\mathbf{A}^T \cdot \mathbf{A}$  plays the role of  $S^2$ , the signal power or autocorrelation, while  $\lambda \mathbf{H}$  plays the role of  $N^2$ , the noise power or autocorrelation. The term in parentheses in equation (18.5.9) is something like an optimal filter, whose effect is to pass the ill-posed inverse  $\mathbf{A}^{-1} \cdot \mathbf{b}$  through unmodified when  $\mathbf{A}^T \cdot \mathbf{A}$  is sufficiently large, but to suppress it when  $\mathbf{A}^T \cdot \mathbf{A}$  is small.

The above choices of  $\mathbf{B}$  and  $\mathbf{H}$  are only the simplest in an obvious sequence of derivatives. If your *a priori* belief is that a *linear* function is a good approximation to  $u(x)$ , then minimize

$$\mathcal{B} \propto \int [\hat{u}''(x)]^2 dx \propto \sum_{\mu=1}^{M-2} [-\hat{u}_{\mu} + 2\hat{u}_{\mu+1} - \hat{u}_{\mu+2}]^2 \quad (18.5.10)$$

implying

$$\mathbf{B} = \begin{pmatrix} -1 & 2 & -1 & 0 & 0 & 0 & 0 & \cdots & 0 \\ 0 & -1 & 2 & -1 & 0 & 0 & 0 & \cdots & 0 \\ \vdots & & & & \ddots & & & & \vdots \\ 0 & \cdots & 0 & 0 & 0 & -1 & 2 & -1 & 0 \\ 0 & \cdots & 0 & 0 & 0 & 0 & -1 & 2 & -1 \end{pmatrix} \quad (18.5.11)$$

and

$$\mathbf{H} = \mathbf{B}^T \cdot \mathbf{B} = \begin{pmatrix} 1 & -2 & 1 & 0 & 0 & 0 & 0 & \cdots & 0 \\ -2 & 5 & -4 & 1 & 0 & 0 & 0 & \cdots & 0 \\ 1 & -4 & 6 & -4 & 1 & 0 & 0 & \cdots & 0 \\ 0 & 1 & -4 & 6 & -4 & 1 & 0 & \cdots & 0 \\ \vdots & & & & \ddots & & & & \vdots \\ 0 & \cdots & 0 & 1 & -4 & 6 & -4 & 1 & 0 \\ 0 & \cdots & 0 & 0 & 1 & -4 & 6 & -4 & 1 \\ 0 & \cdots & 0 & 0 & 0 & 1 & -4 & 5 & -2 \\ 0 & \cdots & 0 & 0 & 0 & 0 & 1 & -2 & 1 \end{pmatrix} \quad (18.5.12)$$

This  $\mathbf{H}$  has two zero eigenvalues, corresponding to the two undetermined parameters of a linear function.

If your *a priori* belief is that a *quadratic* function is preferable, then minimize

$$\mathcal{B} \propto \int [\hat{u}'''(x)]^2 dx \propto \sum_{\mu=1}^{M-3} [-\hat{u}_\mu + 3\hat{u}_{\mu+1} - 3\hat{u}_{\mu+2} + \hat{u}_{\mu+3}]^2 \quad (18.5.13)$$

with

$$\mathbf{B} = \begin{pmatrix} -1 & 3 & -3 & 1 & 0 & 0 & 0 & \cdots & 0 \\ 0 & -1 & 3 & -3 & 1 & 0 & 0 & \cdots & 0 \\ \vdots & & & & \ddots & & & & \vdots \\ 0 & \cdots & 0 & 0 & -1 & 3 & -3 & 1 & 0 \\ 0 & \cdots & 0 & 0 & 0 & -1 & 3 & -3 & 1 \end{pmatrix} \quad (18.5.14)$$

and now

$$\mathbf{H} = \begin{pmatrix} 1 & -3 & 3 & -1 & 0 & 0 & 0 & 0 & 0 & \cdots & 0 \\ -3 & 10 & -12 & 6 & -1 & 0 & 0 & 0 & 0 & \cdots & 0 \\ 3 & -12 & 19 & -15 & 6 & -1 & 0 & 0 & 0 & \cdots & 0 \\ -1 & 6 & -15 & 20 & -15 & 6 & -1 & 0 & 0 & \cdots & 0 \\ 0 & -1 & 6 & -15 & 20 & -15 & 6 & -1 & 0 & \cdots & 0 \\ \vdots & & & & \ddots & & & & & & \vdots \\ 0 & \cdots & 0 & -1 & 6 & -15 & 20 & -15 & 6 & -1 & 0 \\ 0 & \cdots & 0 & 0 & -1 & 6 & -15 & 20 & -15 & 6 & -1 \\ 0 & \cdots & 0 & 0 & 0 & -1 & 6 & -15 & 19 & -12 & 3 \\ 0 & \cdots & 0 & 0 & 0 & 0 & -1 & 6 & -12 & 10 & -3 \\ 0 & \cdots & 0 & 0 & 0 & 0 & 0 & -1 & 3 & -3 & 1 \end{pmatrix} \quad (18.5.15)$$

(We'll leave the calculation of cubics and above to the compulsive reader.)

Notice that you can regularize with “closeness to a differential equation,” if you want. Just pick  $\mathbf{B}$  to be the appropriate sum of finite-difference operators (the coefficients can depend on  $x$ ), and calculate  $\mathbf{H} = \mathbf{B}^T \cdot \mathbf{B}$ . You don't need to know the values of your boundary conditions, since  $\mathbf{B}$  can have fewer rows than columns, as above; hopefully, your data will determine them. Of course, if you do know some boundary conditions, you can build these into  $\mathbf{B}$  too.

With all the proportionality signs above, you may have lost track of what actual value of  $\lambda$  to try first. A simple trick for at least getting “on the map” is to first try

$$\lambda = \text{Tr}(\mathbf{A}^T \cdot \mathbf{A}) / \text{Tr}(\mathbf{H}) \quad (18.5.16)$$

where  $\text{Tr}$  is the trace of the matrix (sum of diagonal components). This choice will tend to make the two parts of the minimization have comparable weights, and you can adjust from there.

As for what is the “correct” value of  $\lambda$ , an objective criterion, if you know your errors  $\sigma_i$  with reasonable accuracy, is to make  $\chi^2$  (that is,  $|\mathbf{A} \cdot \hat{\mathbf{u}} - \mathbf{b}|^2$ ) equal to  $N$ , the number of measurements. We remarked above on the twin acceptable choices  $N \pm (2N)^{1/2}$ . A subjective criterion is to pick any value that you like in the

World Wide Web sample page from NUMERICAL RECIPES IN C: THE ART OF SCIENTIFIC COMPUTING (ISBN 0-521-43108-5)  
 Copyright (C) 1988-1992 by Cambridge University Press. Programs Copyright (C) 1988-1992 by Numerical Recipes Software.  
 Permission is granted for internet users to make one paper copy for their own personal use. Further reproduction, or any copying of machine-readable files (including this one), to any server computer, is strictly prohibited. To order Numerical Recipes books, diskettes, or CDROMs visit website <http://www.nr.com> or call 1-800-872-7423 (North America only), or send email to [trade@cup.cam.ac.uk](mailto:trade@cup.cam.ac.uk) (outside North America).

range  $0 < \lambda < \infty$ , depending on your relative degree of belief in the *a priori* and *a posteriori* evidence. (Yes, people actually do that. Don't blame us.)

### Two-Dimensional Problems and Iterative Methods

Up to now our notation has been indicative of a one-dimensional problem, finding  $\hat{u}(x)$  or  $\hat{u}_\mu = \hat{u}(x_\mu)$ . However, all of the discussion easily generalizes to the problem of estimating a two-dimensional set of unknowns  $\hat{u}_{\mu\kappa}$ ,  $\mu = 1, \dots, M$ ,  $\kappa = 1, \dots, K$ , corresponding, say, to the pixel intensities of a measured image. In this case, equation (18.5.8) is still the one we want to solve.

In image processing, it is usual to have the same number of input pixels in a measured "raw" or "dirty" image as desired "clean" pixels in the processed output image, so the matrices  $\mathbf{R}$  and  $\mathbf{A}$  (equation 18.5.5) are square and of size  $MK \times MK$ .  $\mathbf{A}$  is typically much too large to represent as a full matrix, but often it is either (i) sparse, with coefficients blurring an underlying pixel  $(i, j)$  only into measurements  $(i \pm \text{few}, j \pm \text{few})$ , or (ii) translationally invariant, so that  $A_{(i,j)(\mu,\nu)} = A(i-\mu, j-\nu)$ . Both of these situations lead to tractable problems.

In the case of translational invariance, fast Fourier transforms (FFTs) are the obvious method of choice. The general linear relation between underlying function and measured values (18.4.7) now becomes a discrete convolution like equation (13.1.1). If  $\mathbf{k}$  denotes a two-dimensional wave-vector, then the two-dimensional FFT takes us back and forth between the transform pairs

$$A(i-\mu, j-\nu) \iff \tilde{\mathbf{A}}(\mathbf{k}) \quad b_{(i,j)} \iff \tilde{b}(\mathbf{k}) \quad \hat{u}_{(i,j)} \iff \tilde{u}(\mathbf{k}) \quad (18.5.17)$$

We also need a regularization or smoothing operator  $\mathbf{B}$  and the derived  $\mathbf{H} = \mathbf{B}^T \cdot \mathbf{B}$ . One popular choice for  $\mathbf{B}$  is the five-point finite-difference approximation of the Laplacian operator, that is, the difference between the value of each point and the average of its four Cartesian neighbors. In Fourier space, this choice implies,

$$\begin{aligned} \tilde{B}(\mathbf{k}) &\propto \sin^2(\pi k_1/M) \sin^2(\pi k_2/K) \\ \tilde{H}(\mathbf{k}) &\propto \sin^4(\pi k_1/M) \sin^4(\pi k_2/K) \end{aligned} \quad (18.5.18)$$

In Fourier space, equation (18.5.7) is merely algebraic, with solution

$$\tilde{u}(\mathbf{k}) = \frac{\tilde{A}^*(\mathbf{k})\tilde{b}(\mathbf{k})}{|\tilde{A}(\mathbf{k})|^2 + \lambda\tilde{H}(\mathbf{k})} \quad (18.5.19)$$

where asterisk denotes complex conjugation. You can make use of the FFT routines for real data in §12.5.

Turn now to the case where  $\mathbf{A}$  is not translationally invariant. Direct solution of (18.5.8) is now hopeless, since the matrix  $\mathbf{A}$  is just too large. We need some kind of iterative scheme.

One way to proceed is to use the full machinery of the conjugate gradient method in §10.6 to find the minimum of  $\mathcal{A} + \lambda\mathcal{B}$ , equation (18.5.6). Of the various methods in Chapter 10, conjugate gradient is the unique best choice because (i) it does not require storage of a Hessian matrix, which would be infeasible here,

and (ii) it does exploit gradient information, which we can readily compute: The gradient of equation (18.5.6) is

$$\nabla(\mathcal{A} + \lambda\mathcal{B}) = 2[(\mathbf{A}^T \cdot \mathbf{A} + \lambda\mathbf{H}) \cdot \hat{\mathbf{u}} - \mathbf{A}^T \cdot \mathbf{b}] \quad (18.5.20)$$

(cf. 18.5.8). Evaluation of both the function and the gradient should of course take advantage of the sparsity of  $\mathbf{A}$ , for example via the routines `spr sax` and `spr stx` in §2.7. We will discuss the conjugate gradient technique further in §18.7, in the context of the (nonlinear) maximum entropy method. Some of that discussion can apply here as well.

The conjugate gradient method notwithstanding, application of the unsophisticated steepest descent method (see §10.6) can sometimes produce useful results, particularly when combined with projections onto convex sets (see below). If the solution after  $k$  iterations is denoted  $\hat{\mathbf{u}}^{(k)}$ , then after  $k + 1$  iterations we have

$$\hat{\mathbf{u}}^{(k+1)} = [\mathbf{1} - \epsilon(\mathbf{A}^T \cdot \mathbf{A} + \lambda\mathbf{H})] \cdot \hat{\mathbf{u}}^{(k)} + \epsilon\mathbf{A}^T \cdot \mathbf{b} \quad (18.5.21)$$

Here  $\epsilon$  is a parameter that dictates how far to move in the downhill gradient direction. The method converges when  $\epsilon$  is small enough, in particular satisfying

$$0 < \epsilon < \frac{2}{\max \text{eigenvalue}(\mathbf{A}^T \cdot \mathbf{A} + \lambda\mathbf{H})} \quad (18.5.22)$$

There exist complicated schemes for finding optimal values or sequences for  $\epsilon$ , see [7]; or, one can adopt an experimental approach, evaluating (18.5.6) to be sure that downhill steps are in fact being taken.

In those image processing problems where the final measure of success is somewhat subjective (e.g., “how good does the picture look?”), iteration (18.5.21) sometimes produces significantly improved images long before convergence is achieved. This probably accounts for much of its use, since its mathematical convergence is extremely slow. In fact, (18.5.21) can be used with  $\mathbf{H} = 0$ , in which case the solution is not regularized at all, and full convergence would be disastrous! This is called *Van Cittert’s method* and goes back to the 1930s. A number of iterations the order of 1000 is not uncommon [7].

### **Deterministic Constraints: Projections onto Convex Sets**

A set of possible underlying functions (or images)  $\{\hat{\mathbf{u}}\}$  is said to be *convex* if, for any two elements  $\hat{\mathbf{u}}_a$  and  $\hat{\mathbf{u}}_b$  in the set, all the linearly interpolated combinations

$$(1 - \eta)\hat{\mathbf{u}}_a + \eta\hat{\mathbf{u}}_b \quad 0 \leq \eta \leq 1 \quad (18.5.23)$$

are also in the set. Many *deterministic constraints* that one might want to impose on the solution  $\hat{\mathbf{u}}$  to an inverse problem in fact define convex sets, for example:

- positivity
- compact support (i.e., zero value outside of a certain region)

- known bounds (i.e.,  $u_L(x) \leq \hat{u}(x) \leq u_U(x)$  for specified functions  $u_L$  and  $u_U$ ).

(In this last case, the bounds might be related to an initial estimate and its error bars, e.g.,  $\hat{u}_0(x) \pm \gamma\sigma(x)$ , where  $\gamma$  is of order 1 or 2.) Notice that these, and similar, constraints can be either in the image space, or in the Fourier transform space, or (in fact) in the space of any linear transformation of  $\hat{\mathbf{u}}$ .

If  $C_i$  is a convex set, then  $\mathcal{P}_i$  is called a *nonexpansive projection operator* onto that set if (i)  $\mathcal{P}_i$  leaves unchanged any  $\hat{\mathbf{u}}$  already in  $C_i$ , and (ii)  $\mathcal{P}_i$  maps any  $\hat{\mathbf{u}}$  outside  $C_i$  to the *closest* element of  $C_i$ , in the sense that

$$|\mathcal{P}_i\hat{\mathbf{u}} - \hat{\mathbf{u}}| \leq |\hat{\mathbf{u}}_a - \hat{\mathbf{u}}| \quad \text{for all } \hat{\mathbf{u}}_a \text{ in } C_i \quad (18.5.24)$$

While this definition sounds complicated, examples are very simple: A nonexpansive projection onto the set of positive  $\hat{\mathbf{u}}$ 's is "set all negative components of  $\hat{\mathbf{u}}$  equal to zero." A nonexpansive projection onto the set of  $\hat{u}(x)$ 's bounded by  $u_L(x) \leq \hat{u}(x) \leq u_U(x)$  is "set all values less than the lower bound equal to that bound, and set all values greater than the upper bound equal to *that* bound." A nonexpansive projection onto functions with compact support is "zero the values outside of the region of support."

The usefulness of these definitions is the following remarkable theorem: Let  $C$  be the intersection of  $m$  convex sets  $C_1, C_2, \dots, C_m$ . Then the iteration

$$\hat{\mathbf{u}}^{(k+1)} = (\mathcal{P}_1\mathcal{P}_2 \cdots \mathcal{P}_m)\hat{\mathbf{u}}^{(k)} \quad (18.5.25)$$

will converge to  $C$  from all starting points, as  $k \rightarrow \infty$ . Also, if  $C$  is empty (there is no intersection), then the iteration will have no limit point. Application of this theorem is called the *method of projections onto convex sets* or sometimes *POCS* [7].

A generalization of the POCS theorem is that the  $\mathcal{P}_i$ 's can be replaced by a set of  $\mathcal{T}_i$ 's,

$$\mathcal{T}_i \equiv \mathbf{1} + \beta_i(\mathcal{P}_i - \mathbf{1}) \quad 0 < \beta_i < 2 \quad (18.5.26)$$

A well-chosen set of  $\beta_i$ 's can accelerate the convergence to the intersection set  $C$ .

Some inverse problems can be completely solved by iteration (18.5.25) alone! For example, a problem that occurs in both astronomical imaging and X-ray diffraction work is to recover an image given only the *modulus* of its Fourier transform (equivalent to its power spectrum or autocorrelation) and not the *phase*. Here three convex sets can be utilized: the set of all images whose Fourier transform has the specified modulus to within specified error bounds; the set of all positive images; and the set of all images with zero intensity outside of some specified region. In this case the POCS iteration (18.5.25) cycles among these three, imposing each constraint in turn; FFTs are used to get in and out of Fourier space each time the Fourier constraint is imposed.

The specific application of POCS to constraints alternately in the spatial and Fourier domains is also known as the *Gerchberg-Saxton* algorithm [8]. While this algorithm is non-expansive, and is frequently convergent in practice, it has not been proved to converge in all cases [9]. In the phase-retrieval problem mentioned above, the algorithm often "gets stuck" on a plateau for many iterations before making sudden, dramatic improvements. As many as  $10^4$  to  $10^5$  iterations are sometimes

necessary. (For “unsticking” procedures, see [10].) The uniqueness of the solution is also not well understood, although for two-dimensional images of reasonable complexity it is believed to be unique.

Deterministic constraints can be incorporated, via projection operators, into iterative methods of linear regularization. In particular, rearranging terms somewhat, we can write the iteration (18.5.21) as

$$\hat{\mathbf{u}}^{(k+1)} = [\mathbf{1} - \epsilon\lambda\mathbf{H}] \cdot \hat{\mathbf{u}}^{(k)} + \epsilon\mathbf{A}^T \cdot (\mathbf{b} - \mathbf{A} \cdot \hat{\mathbf{u}}^{(k)}) \quad (18.5.27)$$

If the iteration is modified by the insertion of projection operators at each step

$$\hat{\mathbf{u}}^{(k+1)} = (\mathcal{P}_1\mathcal{P}_2 \cdots \mathcal{P}_m)[\mathbf{1} - \epsilon\lambda\mathbf{H}] \cdot \hat{\mathbf{u}}^{(k)} + \epsilon\mathbf{A}^T \cdot (\mathbf{b} - \mathbf{A} \cdot \hat{\mathbf{u}}^{(k)}) \quad (18.5.28)$$

(or, instead of  $\mathcal{P}_i$ 's, the  $\mathcal{T}_i$  operators of equation 18.5.26), then it can be shown that the convergence condition (18.5.22) is unmodified, and the iteration will converge to minimize the quadratic functional (18.5.6) subject to the desired nonlinear deterministic constraints. See [7] for references to more sophisticated, and faster converging, iterations along these lines.

#### CITED REFERENCES AND FURTHER READING:

- Phillips, D.L. 1962, *Journal of the Association for Computing Machinery*, vol. 9, pp. 84–97. [1]  
 Twomey, S. 1963, *Journal of the Association for Computing Machinery*, vol. 10, pp. 97–101. [2]  
 Twomey, S. 1977, *Introduction to the Mathematics of Inversion in Remote Sensing and Indirect Measurements* (Amsterdam: Elsevier). [3]  
 Craig, I.J.D., and Brown, J.C. 1986, *Inverse Problems in Astronomy* (Bristol, U.K.: Adam Hilger). [4]  
 Tikhonov, A.N., and Arsenin, V.Y. 1977, *Solutions of Ill-Posed Problems* (New York: Wiley). [5]  
 Tikhonov, A.N., and Goncharsky, A.V. (eds.) 1987, *Ill-Posed Problems in the Natural Sciences* (Moscow: MIR).  
 Miller, K. 1970, *SIAM Journal on Mathematical Analysis*, vol. 1, pp. 52–74. [6]  
 Schafer, R.W., Mersereau, R.M., and Richards, M.A. 1981, *Proceedings of the IEEE*, vol. 69, pp. 432–450.  
 Biemond, J., Lagendijk, R.L., and Mersereau, R.M. 1990, *Proceedings of the IEEE*, vol. 78, pp. 856–883. [7]  
 Gerchberg, R.W., and Saxton, W.O. 1972, *Optik*, vol. 35, pp. 237–246. [8]  
 Fienup, J.R. 1982, *Applied Optics*, vol. 15, pp. 2758–2769. [9]  
 Fienup, J.R., and Wackerman, C.C. 1986, *Journal of the Optical Society of America A*, vol. 3, pp. 1897–1907. [10]

## 18.6 Backus-Gilbert Method

The *Backus-Gilbert method* [1,2] (see, e.g., [3] or [4] for summaries) differs from other regularization methods in the nature of its functionals  $\mathcal{A}$  and  $\mathcal{B}$ . For  $\mathcal{B}$ , the method seeks to maximize the *stability* of the solution  $\hat{u}(x)$  rather than, in the first instance, its smoothness. That is,

$$\mathcal{B} \equiv \text{Var}[\hat{u}(x)] \quad (18.6.1)$$