

Chapter 2

The Inverse Problem

2.1 Introduction

Given a list of effects, the problem of determining cause has intrigued philosophers, mathematicians and engineers throughout recorded history. Problems of this type are formally referred to as *inverse problems*. Inverse problems pose a particularly difficult challenge: no solution is guaranteed to be unique or stable. The solution is unique only if for some reason *known to the observer* the given list of effects can be due to one and only one cause.

We are concerned here with the inverse problem as it relates to signal and image restoration. In this context of linear time-invariant (LTI) systems, it is common to use the terms *inverse problem* and *deconvolution* interchangeably. The problem here may be stated as that of estimating the true signal given a distorted and noisy version of the true signal.

2.2 Signal Restoration

In general, the goal of signal recovery is to find the best estimate of a signal that has been distorted. Although the mathematics is the same, we would like to distinguish between signal restoration and signal reconstruction. In the first problem, the research is concerned with obtaining a signal that has been distorted by a measuring device whose transfer function is available. Such a problem arises in image processing, wherein the distorting apparatus could be a lens or an image grabber. In the second problem, the scientist is faced with the challenge of reconstructing a signal from a set of its projections, generally corrupted by noise. This problem arises in spectral estimation, tomography, and image compression. In the image-compression problem, a finite subset of projections of the original signal are given, perhaps on the orthonormal cosine basis, and the original signal is desired.

Generally, to go about the problem of signal recovery, a mathematical model of the signal-formation system is needed. Different models are available; simple linear models are easy to work with but do not reflect the real world. More realistic models are complex and may be used at some additional computational cost.¹

Once a model is specified, a recovery criterion must be selected. Many such criteria exist—ME, minimum-mean-squared error (MMSE), maximum likelihood

(ML), and maximum *a posteriori* (MAP) probability are but a few criteria that have proved useful. Mathematically, the problem of signal recovery is referred to as solving an inverse problem. Generally, an inverse problem will be characterized as being either *well posed* or *ill posed*.^{2,3} We clarify these notions below.

Typically, signal restoration/reconstruction belongs to the class of ill-posed problems. That is, we are often concerned with inverting a singular or nearly singular operator. Our goal is to convert a real problem into one which is well posed in the sense that the statement of the problem gives just enough information to determine one unique solution. However, the new reformulation of the problem is often unrealistic due to the assumptions made to choose a particular model or prior information used to tackle the problem. Many attempts have been made to deal with such problems by inventing *ad-hoc* algorithms that imitate this direct mathematical inversion that one would like to carry out. We shall examine these problems in detail.

As stated earlier, signal-restoration techniques seek the best estimate of the true signal given the observed signal and other *a priori* data such as the noise variance, noise PDF, or positivity constraints on the signal itself. The performance of the Gerchberg–Saxton algorithm⁴ illustrates the degree of success one can achieve in reconstructing the original signal given only the magnitude of its Fourier transform along with the knowledge that the signal is nonnegative.

In recent years deconvolution has become a science in itself. Extremely sophisticated mathematical techniques have been applied to the solution of this problem. Each of the resulting algorithms has its advantages and pitfalls. Experimentation has confirmed that the success of a given algorithm is intimately related to the characteristics of the data. Thus the need for newer and more general algorithms remains.

The objective of this work is to present and analyze a new generalized formulation for iterative signal restoration. The generalized mapping function (GMF) is presented, and its convergence is studied both in the general formulation and for specific cases. The van Cittert algorithm is a special case of this mapping function. The convergence of the van Cittert algorithm has been discussed by Hill and Ioup⁵ and by Jansson.⁶ We present a novel and elegant method of obtaining the criteria for convergence of this algorithm. This also serves as a check to establish the validity of the general formulation. Further, we demonstrate that some popular algorithms are special cases of the GMF for most practical purposes. The convergence of these algorithms is analyzed using the structure developed in this work. A few examples are presented, and the direction of our future work is described.

2.3 Well-Posed and Ill-Posed Problems

Many modern experimental devices for investigating physical phenomena and objects of different kinds are complicated. The results of observations are to be processed and interpreted to extract the necessary information about the characteristics

of the phenomenon or object to be studied.

Most often, what is measured in a physical experiment is not the desired parameter, represented here by the vector \mathbf{x} , but instead a certain effect, $\mathbf{y} = \mathbf{H}\mathbf{x}$. Therefore, the interpretation problem usually reduces to solving an algebraic equation of the form

$$\mathbf{H} \cdot \mathbf{x} = \mathbf{y} \quad (2.1)$$

for the unknown vector \mathbf{x} of length N . Usually, \mathbf{H} represents the apparatus function matrix ($N \times N$), often termed the impulse transfer function. If the measuring device is linear, then the functional relationship between \mathbf{x} and \mathbf{y} is given by*

$$\int_a^b H(t, s)x(s)ds = y(t) \quad \forall t \in [c, d], \quad (2.2)$$

where the kernel H represents the measuring device and is assumed known. The integral equation is a Fredholm integral of the first kind.

For the following discussion, let us assume that the unknown function $x(s)$ belongs to a metric space[†] F and the known function $y(t)$ to a metric space U . Also, assume the kernel $H(s, t)$ is continuous with respect to t , and that it has a continuous partial derivative $\partial H/\partial t$. Usually, we measure changes in both spaces with the L_2 metric defined by Avriel:⁷

$$\rho_c(\omega_1(t), \omega_2(t)) = \left(\int_c^d [\omega_1(t) - \omega_2(t)]^2 \right)^{0.5} \quad (2.3)$$

in the continuous domain. For the discrete field, we have

$$\rho_D(\omega_1(t), \omega_2(t)) = \left(\frac{1}{N} \sum_{i=1}^N [\omega_1(t) - \omega_2(t)]^2 \right)^{0.5}. \quad (2.4)$$

In the classical sense, solving for x is equivalent to finding the inverse operator \mathbf{H}^{-1} , which leads to:

$$\mathbf{x} = \mathbf{H}^{-1}\mathbf{y}. \quad (2.5)$$

Obviously, Eq. (2.1) has solutions for functions \mathbf{y} that lie in the image space $\mathbf{H}\mathbf{F}$. Since the right-hand member $y(t)$ is usually obtained experimentally, only an approximation is available and the apparatus function is only known to some given accuracy. Thus, we are faced with the challenging problem of solving for \mathbf{x} when only partial or approximate information is available. Hence, we are dealing with a system:

$$\tilde{\mathbf{H}}\mathbf{x} = \tilde{\mathbf{y}}, \quad (2.6)$$

*The notation used in this work is described in Appendix D.

†metric space: see Appendix C for definition.

which deviates from the initial equation given in Eq. (2.1). Specifically,

$$\|\tilde{H} - H\| \leq \delta, \quad \|\tilde{y} - y\| \leq \tau, \quad (2.7)$$

where the norm is arbitrary and δ, τ are some positive numbers. The question then arises: is the approximate system solvable? Frequently, the operator H is not invertible or its inverse is not continuous (when H is everywhere continuous). Then, the problem on hand is termed “ill posed.”

Definition 2.3.1. The problem of determining the solution \mathbf{x} in the space F from the “initial data” \mathbf{y} in the space U is said to be ill posed on the pair of metric spaces (F, U) if at least one of the following three conditions is violated:

1. For every element y in U there exists a solution x in the space F .
2. The solution is unique.
3. The problem is stable in the spaces (F, U) .

The property of stability is defined as follows.

Definition 2.3.2. The problem of determining the solution $\mathbf{x} = \mathbf{R}(\mathbf{y})$ in the space F from the initial data \mathbf{y} in U is said to be stable on the spaces (F, U) if, for every positive number ϵ , there exists a positive number $\delta(\epsilon)$ such that the inequality $\rho_U(y_1, y_2) \leq \delta(\epsilon)$ implies $\rho_F(x_1, x_2) \leq \epsilon$ where $x_i = R(y_i)$ with y_i in U and x_i in F for $i = 1, 2$.

Finally, throughout this work we will consider the case when the transformer or operator H is homogeneous. That is, the general integral equation becomes

$$\int_{-\infty}^{\infty} H(t - \tau)x(\tau)d\tau = y(t), \quad (2.8)$$

known as the Fredholm integral equation of the first kind.⁸ Any scanning measurement device leads to this form of the integral equation in a noise-free situation.

In the presence of noise, the models presented above are modified to account for the effect of the noise. Also, we will take as a given that the distorting processes are time (or space) invariant. With these changes the model using signal notation is as in Eqs. (2.9) (2.10) below

$$y(t) = \int_{-\infty}^{\infty} h(t - \tau)x(\tau) d\tau + n(t), \quad (2.9)$$

which may be written more compactly as

$$y = h * x + n. \quad (2.10)$$

In Eq. (2.10), y, x, h , and n could be continuous functions of time or their sampled versions. We develop algorithms in this book under the assumption that the

noise n may follow any distribution. On the other hand, when verifying these algorithms and the applications they admit, we assume a normal distribution on the noise. We use a minimax argument to justify this assumption. As argued in Ref. 9, this assumption provides the largest lower bound for the variance of any unbiased estimator of x for a general class of linear models. Consequently, assuming noise normality provides the worst-case scenario under which to investigate our estimation and restoration algorithms.

The same model may be written in matrix notation as

$$\mathbf{y} = \mathbf{H} \cdot \mathbf{x} + \mathbf{n}, \quad (2.11)$$

where \mathbf{y} is the sample observation vector, \mathbf{x} is a sample vector representing the true signal, \mathbf{H} is the distorting function matrix, and \mathbf{n} is the sample noise vector. The distortion process is the standard convolution integral. In the matrix model, \mathbf{H} is a circulant Toeplitz matrix and the matrix product in Eq. (2.11) specifies discrete convolution. The noise variance is assumed to be σ^2 .

In the case of two-dimensional signals such as images, the models are the same except that the functions are now over two variables. There is a slight modification to the structure of \mathbf{H} and the vectors \mathbf{x} , \mathbf{y} in the matrix model; \mathbf{H} is a block Toeplitz matrix, and the vectors \mathbf{x} , \mathbf{y} are lexicographically ordered.¹⁰ Lexicographic ordering of image matrices is described in Appendix A.

2.4 Naïve Approaches to Inverse Problems

Our interest is in estimating the process x given the measurement y . The naïve approach is to find the inverse of the distorting function and ignore the effect of noise. In the matrix case this involves computing the inverse of the distorting matrix, i.e.,

$$\mathbf{x} = \mathbf{H}^{-1} \cdot \mathbf{y}. \quad (2.12)$$

This approach runs into difficulties for any but the most trivial cases. The problem lies in determining the inverse matrix. There is, as we have stated before, no guarantee that the inverse exists. In cases where the matrix is noninvertible there is scope for preconditioning it to obtain an inverse, but such approaches may be unstable.

Define h to be the sampled distortion function vector. Let x , y , and n be the true sampled signal vector, measurement vector, and sampled noise vector respectively. Then, instead of using the Toeplitz form Eq. (2.11), the degradation process can be written as

$$y = h * x + n. \quad (2.13)$$

In the frequency domain this is represented as

$$Y = H \cdot X + N. \quad (2.14)$$