

Discrete Signal Processing Notes

Charles L. Byrne

January 25, 2006

Contents

1	Farfield Propagation	3
1.1	The Solar-Emission Problem	4
1.2	The One-Dimensional Case	4
1.2.1	The Plane-Wave Model	5
1.3	Fourier-Transform Pairs	6
1.3.1	The Fourier Transform	6
1.3.2	Sampling	6
1.3.3	Reconstructing from Fourier-Transform Data	7
1.3.4	An Example	7
1.4	The Dirac Delta	8
1.5	Practical Limitations	8
1.5.1	Convolution Filtering	9
1.5.2	Low-Pass Filtering	10
1.6	Point Sources as Dirac Deltas	11
1.7	The Limited-Aperture Problem	11
1.7.1	Resolution	12
1.7.2	The Solar-Emission Problem Revisited	13
1.8	Discrete Data	14
1.8.1	Reconstruction from Samples	15
1.9	The Finite-Data Problem	15
1.10	Functions of Several Variables	16
1.10.1	Two-Dimensional Farfield Object	16
1.10.2	Two-Dimensional Fourier Transforms	17
1.10.3	Two-Dimensional Fourier Inversion	18
1.10.4	Limited Apertures in Two Dimensions	19
1.11	Broadband Signals	19
1.12	The Laplace Transform and the Ozone Layer	20
1.12.1	The Laplace Transform	20
1.12.2	Scattering of Ultraviolet Radiation	20
1.12.3	Measuring the Scattered Intensity	21
1.12.4	The Laplace Transform Data	21
1.13	Summary	22

2	Reconstruction from Line-Integral Data	23
2.1	Ocean Acoustic Tomography	23
2.1.1	Obtaining Line-Integral Data	24
2.1.2	The Difficulties	24
2.1.3	Why “Tomography”?	25
2.1.4	An Algebraic Approach	25
2.2	X-ray Transmission Tomography	26
2.2.1	The Exponential-Decay Model	26
2.2.2	Difficulties to be Overcome	27
2.3	Positron Emission Tomography	27
2.3.1	The Coincidence-Detection Model	28
2.3.2	Line-Integral Data	28
2.4	Single-Photon Emission Tomography	29
2.4.1	The Line-Integral Model	29
2.4.2	Problems with the Line-Integral Model	30
2.4.3	The Stochastic Model: Discrete Poisson Emitters	31
2.4.4	Reconstruction as Parameter Estimation	32
2.5	Reconstruction from Line Integrals	32
2.5.1	The Radon Transform	33
2.5.2	The Central Slice Theorem	33
2.5.3	Ramp Filter, then Backproject	34
2.5.4	Backproject, then Ramp Filter	35
2.5.5	Radon’s Inversion Formula	36
2.5.6	Practical Issues	36
2.6	Summary	37
3	Discrete Signal Processing	39
3.1	Discrete Signals	39
3.2	Notation	40
3.3	Operations on Discrete Signals	40
3.3.1	Linear Operators	40
3.3.2	Shift-invariant Operators	41
3.3.3	Convolution Operators	41
3.3.4	LSI Filters are Convolutions	42
3.4	Special Types of Discrete Signals	42
3.5	The Frequency-Response Function	43
3.5.1	The Response of a LSI System to $x = e_\omega$	44
3.5.2	Relating $H(\omega)$ to $h = T(\delta)$	45
3.6	The Discrete Fourier Transform	46
3.7	The Convolution Theorem	47
3.8	Sampling and Aliasing	48
3.9	Important Problems in Discrete Signal Processing	49
3.9.1	Low-pass Filtering	49
3.9.2	The Finite-Data Problem	50

<i>CONTENTS</i>	1
3.9.3 The Extrapolation Problem	50
3.10 Discrete Signals from Finite Data	52
3.10.1 Zero-extending the Data	52
3.10.2 Periodically Extending the Data	53
3.10.3 A Third Way to Extend the Data	54
3.10.4 A Fourth Way: Bandlimited Extrapolation	54
3.11 Is this Analysis or Representation?	56
3.12 Oversampling	58
3.13 Finite Data and the Fast Fourier Transform	59
4 Randomness in Signal Processing	63
4.1 Randomness in Farfield Propagation	63
4.2 Random Variables as Models	65
4.3 Discrete Random Signal Processing	67
4.3.1 The Simplest Random Sequence	67
4.4 Random Discrete Functions or Discrete Random Processes .	68
4.5 Correlation Functions and Power Spectra	71
4.6 Random Sinusoidal Sequences	72
4.7 Spread-Spectrum Communication	74
4.8 Stochastic Difference Equations	74
Bibliography	75
Index	85

Chapter 1

Farfield Propagation

A basic problem in remote sensing is to determine the nature of a distant object by measuring signals transmitted by or reflected from that object. If the object of interest is sufficiently remote, that is, is in the *farfield*, it can be assumed that the data we obtain by sampling the propagating spatio-temporal field is related to what we want by *Fourier transformation*. The problem is then to estimate a function from finitely many (usually noisy) values of its *Fourier transform*. Although there are many important mathematical tools employed to solve signal-processing problems, the Fourier transform is the most important. Our discussion of farfield propagation will serve to motivate the Fourier transform, not only as a useful mathematical device, but also as an object having actual physical significance.

We shall begin our discussion of farfield propagation by considering an extended object transmitting or reflecting a single-frequency, or *narrowband*, signal. Later, we shall move to the problem of a distant point source whose location we wish to ascertain, as well as to signals involving multiple frequencies, the so-called *broadband-signal* case. The narrowband, extended-object case is a good place to begin, since a point object is simply a limiting case of an extended object, and broadband received signals can always be filtered to reduce their frequency band.

The application we consider here is a common one of remote-sensing of transmitted or reflected waves propagating from distant sources. Examples include optical imaging of planets and asteroids using reflected sunlight, radio-astronomy imaging of distant sources of radio waves, active and passive sonar, and radar imaging.

1.1 The Solar-Emission Problem

In [1] Bracewell discusses the *solar-emission* problem. In 1942, it was observed that radio-wave emissions in the one-meter wavelength range were arriving from the sun. Were they coming from the entire disk of the sun or were the sources more localized, in sunspots, for example? The problem then was to view each location on the sun's surface as a potential source of these radio waves and to determine the intensity of emission corresponding to each location. The sun has an angular diameter of 30 min. of arc, or one-half of a degree, when viewed from earth, but the needed resolution was more like 3 min. of arc. As we shall see shortly, such resolution requires a radio telescope 1000 wavelengths across, which means a diameter of 1km at a wavelength of 1 meter; in 1942 the largest military radar antennas were less than 5 meters across. A solution was found, using the method of reconstructing an object from line-integral data, a technique that surfaced again in tomography. The problem here is inherently two-dimensional, but, for simplicity, we shall begin with the one-dimensional case.

1.2 The One-Dimensional Case

Because our purpose is to motivate the Fourier transform by showing how it arises naturally in a discussion of farfield propagation, we begin with the more tractable *narrowband-signal* case. We assume that each of the signals being transmitted or reflected is a single-frequency complex sinusoid, having the form $Ae^{i\omega t}$, with complex amplitude A that varies as a function of position within the distant object. The Fourier transform enters the picture when we make the *farfield assumption* that the distance from the object to the sensors is much larger than the distance between sensors. Equivalently, we assume that we are far enough away from the sources that the spherically spreading waves they have generated appear to the sensors as planewave fronts.

Suppose that $D > 0$ represents a large distance from our sensors. Imagine each point $(x, D, 0)$ along an axis parallel to the x -axis in three-dimensional space transmitting or reflecting the sinusoidal signal $g(x)e^{i\omega t}$, where ω is the common frequency of these signals and the $g(x)$ is the complex amplitude associated with each particular x . Our objective is to determine the values $g(x)$, for each x . In the sun-spot problem such information will help us decide where the transmitted radio waves are coming from. In a radar problem, determining the $g(x)$, the amplitudes of the reflected radio wave, will tell us something about the nature of the extended object, since different materials reflect the waves differently. We calculate the signal received at the point $(s, 0, 0)$, under the assumption that $D > 0$ is much, much larger than $|s|$.

1.2.1 The Plane-Wave Model

Let θ denote the angle between the x -axis and the line from $(0, 0, 0)$ to $(x, D, 0)$. Because D is so much larger than $|s|$, the angle θ remains the same, when observed from any other point $(s, 0, 0)$, that is, there is no parallax, and the use of the point $(0, 0, 0)$ is merely a convenience. Again, because D is so large, the spherically spreading field originating at $(x, D, 0)$ is essentially a plane surface as it reaches the sensors. The planes of constant value are normal to the direction vector $\theta = (\cos \theta, \sin \theta)$. Let $b(s, t)$ be the signal from $(x, D, 0)$ that is received at location $(s, 0, 0)$ at time t . For reference, let us suppose that

$$b(0, t) = e^{i\omega(t - \frac{D}{c})} g(x).$$

Because the planewaves travel at a speed c , we have

$$b(s, t) = u(0, t + \frac{s \cos \theta}{c}) = e^{i\omega(t - \frac{D}{c})} e^{i\frac{\omega s \cos \theta}{c}} g(x).$$

Of course, the signal received at $(s, 0, 0)$ does not come only from a single point $(x, D, 0)$, but from all the points $(x, D, 0)$, so the combined signal received at $(s, 0, 0)$ is

$$B(s, t) = e^{i\omega(t - \frac{D}{c})} \int e^{i\frac{\omega s \cos \theta}{c}} g(x) dx. \quad (1.1)$$

Since θ is a one-to-one function of x , we can view $g(x)$ as a function of θ , and write $g(\theta)$ in place of $g(x)$. We then introduce the new variable $k = \frac{\omega}{c} \cos \theta$ and write the integral

$$\int e^{i\frac{\omega s \cos \theta}{c}} g(x) dx$$

as

$$\frac{c}{\omega} \int_{-\frac{\omega}{c}}^{\frac{\omega}{c}} f(k) e^{isk} dk, \quad (1.2)$$

where $f(k)$ is the function $g(\theta)/\sin \theta$, written as a function of the variable k . Since, in most applications, the distant object has a small angular diameter when viewed from a great distance, the sun's is 30 minutes of arc, the angle θ will be restricted to a small interval centered at $\theta = \frac{\pi}{2}$. Therefore, $\sin \theta$ is bounded away from zero and $f(k)$ is well defined.

The integral

$$\int_{-\frac{\omega}{c}}^{\frac{\omega}{c}} f(k) e^{isk} dk$$

is the familiar one that defines the Fourier transform of the function $f(k)$. Using the approximations permitted under the farfield assumption, the received signal $B(s, t)$ is easily shown to provide the Fourier transform of the object function $f(k)$.

1.3 Fourier-Transform Pairs

We consider now the Fourier transform of a function of a single real variable. In the previous section it was reasonable to denote the Fourier transform of $f(k)$ by $F(s)$, with s denoting location in sensor space and k denoting wave vectors associated with given angles. However, in discussing the more general case, it is better to use more conventional notation. Therefore, we shall consider a function $f(x)$ having Fourier transform $F(\gamma)$. The variable x has no relation to the variable of the same name used to describe the spatial extent of the distant object being imaged in our previous example.

1.3.1 The Fourier Transform

Let $f(x)$ be defined for real variable x in $(-\infty, \infty)$. The *Fourier transform* of $f(x)$ is the function of the real variable γ given by

$$F(\gamma) = \int_{-\infty}^{\infty} f(x)e^{i\gamma x} dx. \quad (1.3)$$

In our example of farfield propagation, the signal received at $(s, 0, 0)$, as given by Equation (1.1), can be rewritten as

$$B(s, t) = \frac{c}{\omega} e^{i\omega(t - \frac{D}{c})} F(s), \quad (1.4)$$

where $F(s)$ is the Fourier transform of $f(k)$. Consequently, we can say that the data measured at the sensor locations $(s, 0, 0)$ give us (noisy) values of the Fourier transform of $f(k)$.

1.3.2 Sampling

Because the function $f(k)$ is zero outside the interval $[-\frac{\omega}{c}, \frac{\omega}{c}]$, the function $F(s)$ is *band-limited*. The *Nyquist spacing* in the variable s is therefore

$$\Delta_s = \frac{\pi c}{\omega}.$$

The wavelength λ associated with the frequency ω is defined to be

$$\lambda = \frac{2\pi c}{\omega},$$

so that

$$\Delta_s = \frac{\lambda}{2}.$$

The significance of the Nyquist spacing comes from Shannon's Sampling Theorem, which says that if we have the values $F(m\Delta_s)$, for all integers

m , then we have enough information to recover $f(k)$ exactly. In practice, of course, this is never the case.

Notice that $B(s, t)$ is not just $F(s)$, but

$$B(s, t) = \frac{c}{\omega} e^{i\omega(t - \frac{D}{c})} F(s).$$

To extract $F(s)$ from $B(s, t)$, we need to remove the factor $e^{i\omega(t - \frac{D}{c})}$. When the frequency ω is large, as in optical remote sensing, for example, determining this value accurately may be impossible. What we then have is the *phase problem*; that is, we can measure only $|F(s)|$, and not the phase of the complex numbers $F(s)$.

1.3.3 Reconstructing from Fourier-Transform Data

As illustrated by the farfield propagation example, our goal is often to reconstruct the function $f(x)$ from measurements of its Fourier transform $F(\gamma)$. But, how?

If we have $F(\gamma)$ for all real γ , then we can recover the function $f(x)$ using the *Fourier Inversion Formula*:

$$f(x) = \frac{1}{2\pi} \int F(\gamma) e^{-i\gamma x} d\gamma. \quad (1.5)$$

The functions $f(x)$ and $F(\gamma)$ are called a *Fourier-transform pair*.

1.3.4 An Example

For example, consider an extended object of finite length, with uniform amplitude function $f(x) = \frac{1}{2X}$, for $|x| \leq X$, and $f(x) = 0$, otherwise. The Fourier transform of this $f(x)$ is

$$F(\gamma) = \frac{\sin(X\gamma)}{X\gamma},$$

for all real $\gamma \neq 0$, and $F(0) = 1$. Note that $F(\gamma)$ is nonzero throughout the real line, except for isolated zeros, but that it goes to zero as we go to the infinities. This is typical behavior. Notice also that the smaller the X , the slower $F(\gamma)$ dies out; the first zeros of $F(\gamma)$ are at $|\gamma| = \frac{\pi}{X}$, so the main lobe widens as X goes to zero.

It may seem paradoxical that when X is larger, its Fourier transform dies off more quickly. The Fourier transform $F(\gamma)$ goes to zero faster for larger X because of destructive interference. Because of differences in their complex phases, the magnitude of the sum of the signals received from various parts of the object is much smaller than we might expect, especially when X is large. For smaller X the signals received at a sensor are much

more *in phase* with one another, and so the magnitude of the sum remains large. A more quantitative statement of this phenomenon is provided by the *uncertainty principle* (see [19]).

1.4 The Dirac Delta

Consider what happens in the limit, as $X \rightarrow 0$. Then we have an infinitely high point source at $x = 0$; we denote this by $\delta(x)$, the *Dirac delta*. The Fourier transform approaches the constant function with value 1, for all γ ; the Fourier transform of $f(x) = \delta(x)$ is the constant function $F(\gamma) = 1$, for all γ . The Dirac delta $\delta(x)$ has the *sifting property*:

$$\int h(x)\delta(x)dx = h(0),$$

for each function $h(x)$ that is continuous at $x = 0$.

Because the Fourier transform of $\delta(x)$ is the function $F(\gamma) = 1$, the Fourier inversion formula tells us that

$$\delta(x) = \frac{1}{2\pi} \int_{-\infty}^{\infty} e^{-i\omega x} d\omega. \quad (1.6)$$

Obviously, this integral cannot be understood in the usual way. The integral in Equation (1.6) is a symbolic way of saying that

$$\int h(x) \left(\frac{1}{2\pi} \int_{-\infty}^{\infty} e^{-i\omega x} d\omega \right) dx = \int h(x)\delta(x)dx = h(0), \quad (1.7)$$

for all $h(x)$ that are continuous at $x = 0$; that is, the integral in Equation (1.6) has the sifting property, so it acts like $\delta(x)$. Interchanging the order of integration in Equation (1.7), we obtain

$$\begin{aligned} \int h(x) \left(\frac{1}{2\pi} \int_{-\infty}^{\infty} e^{-i\omega x} d\omega \right) dx &= \frac{1}{2\pi} \int_{-\infty}^{\infty} \left(\int h(x)e^{-i\omega x} dx \right) d\omega \\ &= \frac{1}{2\pi} \int_{-\infty}^{\infty} H(-\omega) d\omega = \frac{1}{2\pi} \int_{-\infty}^{\infty} H(\omega) d\omega = h(0). \end{aligned}$$

We shall return to the Dirac delta when we consider farfield point sources.

1.5 Practical Limitations

In actual remote-sensing problems, antennas cannot be of infinite extent. In digital signal processing, moreover, there are only finitely many sensors. We never measure the entire Fourier transform of $f(x)$, but, at best, just

part of it. In fact, the data we are able to measure is almost never exact values of the Fourier transform of $f(x)$, but rather, values of some distorted or blurred version. To describe such situations, we usually resort to *convolution-filter* models.

1.5.1 Convolution Filtering

Imagine that what we measure are not values of $F(\gamma)$, but of $F(\gamma)H(\gamma)$, where $H(\gamma)$ is a function that describes the limitations and distorting effects of the measuring process, including any blurring due to the medium through which the signals have passed, such as refraction of light as it passes through the atmosphere. If we apply the Fourier Inversion Formula to $F(\gamma)H(\gamma)$, instead of to $F(\gamma)$, we get

$$g(x) = \frac{1}{2\pi} \int F(\gamma)H(\gamma)e^{-i\gamma x} d\gamma. \quad (1.8)$$

The function $g(x)$ that results is $g(x) = (f * h)(x)$, the *convolution* of the functions $f(x)$ and $h(x)$, with the latter given by

$$h(x) = \frac{1}{2\pi} \int H(\gamma)e^{-i\gamma x} d\gamma.$$

Note that, if $f(x) = \delta(x)$, then $g(x) = h(x)$; that is, our reconstruction of the object from distorted data is the function $h(x)$ itself. For that reason, the function $h(x)$ is called the *point-spread function* of the imaging system.

Convolution filtering refers to the process of converting any given function, say $f(x)$, into a different function, say $g(x)$, by convolving $f(x)$ with a fixed function $h(x)$. Since this process can be achieved by multiplying $F(\gamma)$ by $H(\gamma)$ and then inverse Fourier transforming, such convolution filters are studied in terms of the properties of the function $H(\gamma)$, known in this context as the *system transfer function*, or the *optical transfer function* (OTF); when γ is a frequency, rather than a spatial frequency, $H(\gamma)$ is called the *frequency-response function* of the filter. The magnitude of $H(\gamma)$, $|H(\gamma)|$, is called the *modulation transfer function* (MTF). The study of convolution filters is a major part of signal processing. Such filters provide both reasonable models for the degradation signals undergo, and useful tools for reconstruction.

Let us rewrite Equation (1.8), replacing $F(\gamma)$ and $H(\gamma)$ with their definitions, as given by Equation (1.3). Then we have

$$g(x) = \int \left(\int f(t)e^{i\gamma t} dt \right) \left(\int h(s)e^{i\gamma s} ds \right) e^{-i\gamma x} d\gamma.$$

Interchanging the order of integration, we get

$$g(x) = \int \int f(t)h(s) \left(\int e^{i\gamma(t+s-x)} d\gamma \right) ds dt.$$

Now using Equation (1.6) to replace the inner integral with $\delta(t + s - x)$, the next integral becomes

$$\int h(s)\delta(t + s - x)ds = h(x - t).$$

Finally, we have

$$g(x) = \int f(t)h(x - t)dt; \quad (1.9)$$

this is the definition of the convolution of the functions f and h .

1.5.2 Low-Pass Filtering

A major problem in image reconstruction is the removal of blurring, which is often modelled using the notion of convolution filtering. In the one-dimensional case, we describe blurring by saying that we have available measurements not of $F(\gamma)$, but of $F(\gamma)H(\gamma)$, where $H(\gamma)$ is the frequency-response function describing the blurring. If we know the nature of the blurring, then we know $H(\gamma)$, at least to some degree of precision. We can try to remove the blurring by taking measurements of $F(\gamma)H(\gamma)$, dividing these numbers by the value of $H(\gamma)$, and then inverse Fourier transforming. The problem is that our measurements are always noisy, and typical functions $H(\gamma)$ have many zeros and small values, making division by $H(\gamma)$ dangerous, except where the values of $H(\gamma)$ are not too small. These values of γ tend to be the smaller ones, centered around zero, so that we end up with estimates of $F(\gamma)$ itself only for the smaller values of γ . The result is a *low-pass filtering* of the object $f(x)$.

To investigate such low-pass filtering, we suppose that $H(\gamma) = 1$, for $|\gamma| \leq \Gamma$, and is zero, otherwise. Then the filter is called the ideal Γ -lowpass filter. In the farfield propagation model, the variable x is spatial, and the variable γ is spatial frequency, related to how the function $f(x)$ changes spatially, as we move x . Rapid changes in $f(x)$ are associated with values of $F(\gamma)$ for large γ . For the case in which the variable x is time, the variable γ becomes frequency, and the effect of the low-pass filter on $f(x)$ is to remove its higher-frequency components.

One effect of low-pass filtering in image processing is to smooth out the more rapidly changing features of an image. This can be useful if these features are simply unwanted oscillations, but if they are important detail, the smoothing presents a problem. Restoring such wanted detail is often viewed as removing the unwanted effects of the low-pass filtering; in other words, we try to recapture the missing high-spatial-frequency values that have been zeroed out. Such an approach to image restoration is called *frequency-domain extrapolation*. How can we hope to recover these missing spatial frequencies, when they could have been anything? To have

some chance of estimating these missing values we need to have some prior information about the image being reconstructed.

1.6 Point Sources as Dirac Deltas

Television signals reflected from satellites are picked up using antennas in the shape of parabolic dishes. The idea here is to point the dish at the satellite, so that signals from other sources are discriminated against and the one from the satellite is reinforced. In applications such as sonar surveillance, it is often the case that the array of sensors cannot be moved. In such cases electronic steering using phase shifts replaces the physical turning of the antenna. A common practice in sonar is to place sensors at equal intervals along a straight line; such an arrangement is called a *linear array*. If our sensor array is linear, along the line making the angle $\phi = 0$ with the horizontal axis, then each sensor in the linear array receives the same signal. If the line of the array corresponds to an angle ϕ that is not zero, then two sensors a distance Δ apart along the line receive the signal with time delays that differ by $\frac{\Delta \sin \phi}{c}$, that is, with a phase difference of $\frac{\omega \Delta \sin \phi}{c}$. Therefore, the data we measure along this linear array contains, in the phase differences, information about the direction of the farfield point source, relative to the line of the array. This forms the basis for sonar direction-of-arrival estimation and detection; for further details see [19].

As we shall see, if we had available an infinite number of sensors, properly spaced along the line of the array, we could determine the direction of the distant point source with perfect accuracy. In the real world, we must make due with finitely many imperfect sensors. In addition, it is rarely the case that the received signal comes from a single point source; there will always be background noise, other point sources, and so on. Limitations on the number of sensors, and on where they can be placed, make it harder to separate closely-spaced distant point sources. If we know *a priori* that we are looking at point sources, and not extended objects, the *resolution problem* can be partly overcome, using nonlinear *high-resolution* techniques. We shall consider high-resolution methods, such as *entropy maximization* and *likelihood maximization*, in subsequent chapters.

1.7 The Limited-Aperture Problem

In the farfield propagation model, our measurements in the farfield give us the values $F(s)$. Suppose now that we are able to take measurements only for limited values of s , say for $|s| \leq A$; then $2A$ is the *aperture* of our antenna or array of sensors. We describe this, in the general case, by saying that we have available measurements of $F(\gamma)H(\gamma)$, where $H(\gamma) = \chi_{\Gamma}(\gamma) = 1$,

for $|\gamma| \leq \Gamma$, and zero otherwise. So, in addition to describing blurring and low-pass filtering, the convolution-filter model can also be used to model the limited-aperture problem. As in the low-pass case, the limited-aperture problem can be attacked using extrapolation, but with the same sort of risks described for the low-pass case. A much different approach is to increase the aperture by physically moving the array of sensors, as in *synthetic aperture radar* (SAR).

Returning to the farfield propagation model, if we have Fourier transform data only for $|s| \leq A$, then we have $F(s)$ for $|s| \leq A$. Using $H(s) = \chi_A(s)$ to describe the limited aperture of the system, the point-spread function is $h(k) = \frac{\sin(Ak)}{\pi k}$. The first zeros of the numerator occur at $|k| = \frac{\pi}{A}$, so the main lobe of the point-spread function has width $\frac{2\pi}{A}$. For this reason, the resolution of such a limited-aperture imaging system is said to be on the order of $\frac{1}{A}$. Because the distant object, expressed as a function of k in the interval $[-\frac{\omega}{c}, \frac{\omega}{c}]$, the resolution achieved in imaging the distant object will depend on the frequency ω , as well. For that reason, it is common practice to measure the aperture A in units of wavelength λ , rather than, say, in units of meters; an aperture of $A = 5$ meters may be acceptable if the frequency is high, but not if the radiation is in the one-meter-wavelength range.

1.7.1 Resolution

If $f(x) = \delta(x)$ and $H(\gamma) = \chi_\Gamma(\gamma)$ describes the aperture-limitation of the imaging system, then the point-spread function is $h(x) = \frac{\sin \Gamma x}{\pi x}$. The maximum of $h(x)$ still occurs at $x = 0$, but the main lobe of $h(x)$ extends from $-\frac{\pi}{\Gamma}$ to $\frac{\pi}{\Gamma}$; the point source has been spread out. If the point-source object shifts, so that $f(x) = \delta(x - a)$, then the reconstructed image of the object is $h(x - a)$, so the peak is still in the proper place. If we know *a priori* that the object is a single point source, but we do not know its location, the spreading of the point poses no problem; we simply look for the maximum in the reconstructed image. Problems arise when the object contains several point sources, or when we do not know *a priori* what we are looking at, or when the object contains no point sources, but is just a continuous distribution.

Suppose that $f(x) = \delta(x - a) + \delta(x - b)$; that is, the object consists of two point sources. Then Fourier inversion of the aperture-limited data leads to the reconstructed image

$$g(x) = \frac{\sin \Gamma(x - a)}{\pi(x - a)} + \frac{\sin \Gamma(x - b)}{\pi(x - b)}.$$

If $|b - a|$ is large enough, $g(x)$ will have two distinct maxima, at approximately $x = a$ and $x = b$, respectively. However, if $|b - a|$ is too small, the distinct maxima merge into one, at $x = \frac{a+b}{2}$ and resolution will be lost.

How small is too small will depend on Γ , which, of course, depends on both A and ω .

Suppose now that $f(x) = \delta(x - a)$, but we do not know *a priori* that the object is a single point source. We calculate

$$g(x) = h(x - a) = \frac{\sin \Gamma(x - a)}{\pi(x - a)}$$

and use this function as our reconstructed image of the object, for all x . What we see when we look at $g(x)$ for some $x = b \neq a$ is $g(b)$, which is the same thing we see when the point source is at $x = b$ and we look at $x = a$. Point-spreading is, therefore, more than a cosmetic problem. When the object is a point source at $x = a$, but we do not know *a priori* that it is a point source, the spreading of the point causes us to believe that the object function $f(x)$ is nonzero at values of x other than $x = a$. When we look at, say, $x = b$, we see a nonzero value that is caused by the presence of the point source at $x = a$.

Suppose now that the object function $f(x)$ contains no point sources, but is simply an ordinary function of x . If the aperture A is very small, then the function $h(x)$ is nearly constant over the entire extent of the object. The convolution of $f(x)$ and $h(x)$ is essentially the integral of $f(x)$, so the reconstructed object is $g(x) = \int f(x)dx$, for all x .

Let's see what this means for the solar-emission problem discussed earlier.

1.7.2 The Solar-Emission Problem Revisited

The wavelength of the radiation is $\lambda = 1$ meter. Therefore, $\frac{\omega}{c} = 2\pi$, and k in the interval $[-2\pi, 2\pi]$ corresponds to the angle θ in $[0, \pi]$. The sun has an angular diameter of 30 minutes of arc, which is about 10^{-2} radians. Therefore, the sun subtends the angles θ in $[\frac{\pi}{2} - (0.5) \cdot 10^{-2}, \frac{\pi}{2} + (0.5) \cdot 10^{-2}]$, which corresponds roughly to the variable k in the interval $[-3 \cdot 10^{-2}, 3 \cdot 10^{-2}]$. Resolution of 3 minutes of arc means resolution in the variable k of $3 \cdot 10^{-3}$. If the aperture is $2A$, then to achieve this resolution, we need

$$\frac{\pi}{A} = 3 \cdot 10^{-3},$$

or

$$A = \frac{\pi}{3} \cdot 10^3$$

meters, or about 1000 meters.

The radio-wave signals emitted by the sun are focused, using a parabolic radio-telescope. The telescope is pointed at the center of the sun. Because the sun is a great distance from the earth and the subtended arc is small (30 min.), the signals from each point on the sun's surface arrive at the parabola

head-on, that is, parallel to the line from the vertex to the focal point, and are reflected to the receiver located at the focal point of the parabola. The effect of the parabolic antenna is not to discriminate against signals coming from other directions, since there are none, but to effect a summation of the signals received at points $(s, 0, 0)$, for $|s| \leq A$, where $2A$ is the diameter of the parabola. When the aperture is large, the function $H(s)$ is nearly one for all s and the signal received at the focal point is essentially

$$\int F(s)ds = f(0);$$

we are now able to distinguish between $f(0)$ and other values $f(k)$. When the aperture is small, $H(s)$ is essentially $\delta(s)$ and the signal received at the focal point is essentially

$$\int F(s)\delta(s)d\gamma = F(0) = \int f(k)dk;$$

now all we get is the contribution from all the k , superimposed, and all resolution is lost.

Since the solar emission problem is clearly two-dimensional, and we need 3 min. resolution in both dimensions, it would seem that we would need a circular antenna with a diameter of about one kilometer, or a rectangular antenna roughly one kilometer on a side. We shall return to this problem later, once when we discuss multi-dimensional Fourier transforms, and then again when we consider tomographic reconstruction of images from line integrals.

1.8 Discrete Data

A familiar topic in signal processing is the passage from functions of continuous variables to discrete sequences. This transition is achieved by *sampling*, that is, extracting values of the continuous-variable function at discrete points in its domain. Our example of farfield propagation can be used to explore some of the issues involved in sampling.

Imagine an infinite *uniform line array* of sensors formed by placing receivers at the points $(n\Delta, 0, 0)$, for some $\Delta > 0$ and all integers n . Then our data are the values $F(n\Delta)$. Because we defined $k = \frac{\omega}{c} \cos \theta$, it is clear that the function $f(k)$ is zero for k outside the interval $[-\frac{\omega}{c}, \frac{\omega}{c}]$.

Exercise 1.1 Show that our discrete array of sensors cannot distinguish between the signal arriving from θ and a signal with the same amplitude, coming from an angle α with

$$\frac{\omega}{c} \cos \alpha = \frac{\omega}{c} \cos \theta + \frac{2\pi}{\Delta} m,$$

where m is an integer.

To avoid the ambiguity described in Exercise 1.1, we must select $\Delta > 0$ so that

$$-\frac{\omega}{c} + \frac{2\pi}{\Delta} \geq \frac{\omega}{c},$$

or

$$\Delta \leq \frac{\pi c}{\omega} = \frac{\lambda}{2}.$$

The sensor spacing $\Delta_s = \frac{\lambda}{2}$ is the *Nyquist spacing*.

In the sunspot example, the object function $f(k)$ is zero for k outside of an interval much smaller than $[-\frac{\omega}{c}, \frac{\omega}{c}]$. Knowing that $f(k) = 0$ for $|k| > K$, for some $0 < K < \frac{\omega}{c}$, we can accept ambiguities that confuse θ with another angle that lies outside the angular diameter of the object. Consequently, we can redefine the Nyquist spacing to be

$$\Delta_s = \frac{\pi}{K}.$$

This tells us that when we are imaging a distant object with a small angular diameter, the Nyquist spacing is greater than $\frac{\lambda}{2}$. If our sensor spacing has been chosen to be $\frac{\lambda}{2}$, then we have *oversampled*. In the oversampled case, band-limited extrapolation methods can be used to improve resolution (see [20]).

1.8.1 Reconstruction from Samples

From the data gathered at our infinite array we have extracted the Fourier transform values $F(n\Delta)$, for all integers n . The obvious question is whether or not the data is sufficient to reconstruct $f(k)$. We know that, to avoid ambiguity, we must have $\Delta \leq \frac{\pi c}{\omega}$. The good news is that, provided this condition holds, $f(k)$ is uniquely determined by this data and formulas exist for reconstructing $f(k)$ from the data; this is the content of the *Shannon Sampling Theorem*. Of course, this is only of theoretical interest, since we never have infinite data. Nevertheless, a considerable amount of traditional signal-processing exposition makes use of this infinite-sequence model. The real problem, of course, is that our data is always finite.

1.9 The Finite-Data Problem

Suppose that we build a *uniform line array* of sensors by placing receivers at the points $(n\Delta, 0, 0)$, for some $\Delta > 0$ and $n = -N, \dots, N$. Then our data are the values $F(n\Delta)$, for $n = -N, \dots, N$. Suppose, as previously, that the object of interest, the function $f(k)$, is nonzero only for values of k in the interval $[-K, K]$, for some $0 < K < \frac{\omega}{c}$. Once again, we must have $\Delta \leq \frac{\pi c}{\omega}$ to avoid ambiguity; but this is not enough, now. The finite Fourier data is no longer sufficient to determine a unique $f(k)$. The best we can hope

to do is to estimate the true $f(k)$, using both our measured Fourier data and whatever prior knowledge we may have about the function $f(k)$, such as where it is nonzero, if it consists of Dirac delta point sources, or if it is nonnegative. The data is also noisy, and that must be accounted for in the reconstruction process. We shall return later to this important problem of reconstructing a general function $f(x)$ from finitely many noisy values of its Fourier transform.

In certain applications, such as sonar array processing, the sensors are not necessarily arrayed at equal intervals along a line, or even at the grid points of a rectangle, but in an essentially arbitrary pattern in two, or even three, dimensions. In such cases, we have values of the Fourier transform of the object function, but at essentially arbitrary values of the variable. How best to reconstruct the object function in such cases is not obvious.

1.10 Functions of Several Variables

Fourier transformation applies, as well, to functions of several variables. As in the one-dimensional case, we can motivate the multi-dimensional Fourier transform using the farfield propagation model. As we noted earlier, the solar emission problem is inherently a two-dimensional problem.

1.10.1 Two-Dimensional Farfield Object

Consider the case of a distant two-dimensional transmitting or reflecting object. Let each point (x, D, z) in the x, z -plane send out the signal $g(x, z)e^{i\omega t}$. As in the one-dimensional case, D is so large that the spherically spreading wave from (x, D, z) is essentially a plane surface when it reaches the plane $y = 0$ of the sensors. Let θ be the unit vector along the line from $(0, 0, 0)$ to (x, D, z) . Then θ is normal to the planes of constant value of the field originating at (x, D, z) . As before, we assume that D is so large that the direction of (x, D, z) as measured from $(0, 0, 0)$ is the same as would have been measured at any other location $(u, 0, v)$ at which we may locate a sensor.

Let $b(u, v, t)$ be the signal from (x, D, z) that is received at location $(u, 0, v)$ at time t . For reference, let us suppose that

$$u(0, 0, t) = e^{i\omega(t - \frac{D}{c})} g(x, z).$$

Because the planewaves travel at a speed c , we have

$$b(u, v, t) = u(0, t + \frac{\mathbf{s} \cdot \theta}{c}) = e^{i\omega(t - \frac{D}{c})} e^{i\frac{\omega \mathbf{s} \cdot \theta}{c}} g(x),$$

where $\mathbf{s} = (u, 0, v)$.

Of course, the signal received at $(u, 0, v)$ does not come only from a single point (x, D, z) , but from all the points (x, D, z) , so the combined signal received at $(u, 0, v)$ is

$$B(u, v, t) = e^{i\omega(t - \frac{D}{c})} \int \int e^{i\frac{\omega s \cdot \theta}{c}} g(x, z) dx dz. \quad (1.10)$$

Since there is a one-to-one relationship between the direction vectors θ and the points (x, D, z) , we can view $g(x, z)$ as a function of θ , and write $g(\theta)$ in place of $g(x, z)$. We then introduce the new variable $\mathbf{k} = \frac{\omega}{c}\theta$ and write the integral

$$\int \int e^{i\frac{\omega s \cdot \theta}{c}} g(x, z) dx dz$$

as

$$\frac{c}{\omega} \int \int f(\mathbf{k}) e^{i\mathbf{s} \cdot \mathbf{k}} d\mathbf{k}, \quad (1.11)$$

where the integral is over all three-dimensional vectors having length $\frac{\omega}{c}$, $f(\mathbf{k})$ is the function obtained from $g(\theta)$ and the Jacobian of the transformation of the variables of integration. Since, in most applications, the distant object has a small angular diameter when viewed from a great distance - the sun's is 30 minutes of arc - the direction vector θ will be restricted to a small subset of vectors centered at $\theta = (0, D, 0)$.

The integral

$$\int \int f(\mathbf{k}) e^{i\mathbf{s} \cdot \mathbf{k}} d\mathbf{k}$$

is the familiar one that defines the Fourier transform of the function $f(\mathbf{k})$. Using the approximations permitted under the farfield assumption, the received signal $B(u, v, t)$ provides the Fourier transform of the object function $f(\mathbf{k})$.

1.10.2 Two-Dimensional Fourier Transforms

Generally, we consider a function $f(x, z)$ of two real variables. Its Fourier transformation is

$$F(\alpha, \beta) = \int \int f(x, z) e^{i(x\alpha + z\beta)} dx dz. \quad (1.12)$$

For example, suppose that $f(x, z) = 1$ for $\sqrt{x^2 + z^2} \leq R$, and zero, otherwise. Then we have

$$F(\alpha, \beta) = \int_{-\pi}^{\pi} \int_0^R e^{-i(\alpha r \cos \theta + \beta r \sin \theta)} r dr d\theta.$$

In polar coordinates, with $\alpha = \rho \cos \phi$ and $\beta = \rho \sin \phi$, we have

$$F(\rho, \phi) = \int_0^R \int_{-\pi}^{\pi} e^{ir\rho \cos(\theta-\phi)} d\theta r dr.$$

The inner integral is well known;

$$\int_{-\pi}^{\pi} e^{ir\rho \cos(\theta-\phi)} d\theta = 2\pi J_0(r\rho),$$

where J_0 denotes the 0th order Bessel function. Using the identity

$$\int_0^z t^n J_{n-1}(t) dt = z^n J_n(z),$$

we have

$$F(\rho, \phi) = \frac{2\pi R}{\rho} J_1(\rho R).$$

Notice that, since $f(x, z)$ is a radial function, that is, dependent only on the distance from $(0, 0, 0)$ to $(x, 0, z)$, its Fourier transform is also radial.

The first positive zero of $J_1(t)$ is around $t = 4$, so when we measure F at various locations and find $F(\rho, \phi) = 0$ for a particular (ρ, ϕ) , we can estimate $R \approx 4/\rho$. So, even when a distant spherical object, like a star, is too far away to be imaged well, we can sometimes estimate its size by finding where the intensity of the received signal is zero.

1.10.3 Two-Dimensional Fourier Inversion

Just as in the one-dimensional case, the Fourier transformation that produced $F(\alpha, \beta)$ can be inverted to recover the original $f(x, y)$. The Fourier Inversion Formula in this case is

$$f(x, y) = \frac{1}{4\pi^2} \int \int F(\alpha, \beta) e^{-i(\alpha x + \beta y)} d\alpha d\beta. \quad (1.13)$$

It is important to note that this procedure can be viewed as two one-dimensional Fourier inversions: first, we invert $F(\alpha, \beta)$, as a function of, say, β only, to get the function of α and y

$$g(\alpha, y) = \frac{1}{2\pi} \int F(\alpha, \beta) e^{-i\beta y} d\beta;$$

second, we invert $g(\alpha, y)$, as a function of α , to get

$$f(x, y) = \frac{1}{2\pi} \int g(\alpha, y) e^{-i\alpha x} d\alpha.$$

If we write the functions $f(x, y)$ and $F(\alpha, \beta)$ in polar coordinates, we obtain alternative ways to implement the two-dimensional Fourier inversion. We shall consider these other ways when we discuss the tomography problem of reconstructing a function $f(x, y)$ from line-integral data.

1.10.4 Limited Apertures in Two Dimensions

Suppose we have the values of the Fourier transform, $F(\alpha, \beta)$, for $|\alpha| \leq A$, $|\beta| \leq B$. We describe this limited-data problem using the function $H(\alpha, \beta)$ that is one for $|\alpha| \leq A$, $|\beta| \leq B$, and zero, otherwise. Then the point-spread function is the inverse Fourier transform of this $H(\alpha, \beta)$, given by

$$h(x, z) = \frac{\sin Ax}{\pi x} \frac{\sin Bz}{\pi z}.$$

The resolution in the horizontal (x) direction is on the order of $\frac{1}{A}$, and $\frac{1}{B}$ in the vertical.

Suppose our aperture is circular, with radius A . Then we have Fourier transform values $F(\alpha, \beta)$ for $\sqrt{\alpha^2 + \beta^2} \leq A$. Let $H(\alpha, \beta)$ equal one, for $\sqrt{\alpha^2 + \beta^2} \leq A$, and zero, otherwise. Then the point-spread function of this limited-aperture system is the inverse Fourier transform of $H(\alpha, \beta)$, given by $h(x, z) = \frac{A}{2\pi r} J_1(rA)$, with $r = \sqrt{x^2 + z^2}$. The resolution of this system is roughly the distance from the origin to the first null of the function $J_1(rA)$, which means that $rA = 4$, roughly.

For the solar emission problem, this says that we would need a circular aperture with radius approximately one kilometer to achieve 3 minutes of arc resolution. But this holds only if the antenna is stationary; a moving antenna is different! The solar emission problem was solved by using a rectangular antenna with a large A , but a small B , and exploiting the rotation of the earth. The resolution is then good in the horizontal, but bad in the vertical, so that the imaging system discriminates well between two distinct vertical lines, but cannot resolve sources within the same vertical line. Because B is small, what we end up with is essentially the integral of the function $f(x, z)$ along each vertical line. By tilting the antenna, and waiting for the earth to rotate enough, we can get these integrals along any set of parallel lines. The problem then is to reconstruct $f(x, z)$ from such line integrals. This is also the main problem in tomography, as we shall see.

1.11 Broadband Signals

We have spent considerable time discussing the case of a distant point source or an extended object transmitting or reflecting a single-frequency signal. If the signal consists of many frequencies, the so-called broadband case, we can still analyze the received signals at the sensors in terms of time delays, but we cannot easily convert the delays to phase differences, and thereby make good use of the Fourier transform. One approach is to filter each received signal, to remove components at all but a single frequency, and then to proceed as previously discussed. In this way we can process

one frequency at a time. The object now is described in terms of a function of both x and ω , with $f(x, \omega)$ the complex amplitude associated with the spatial variable x and the frequency ω . In the case of radar, the function $f(x, \omega)$ tells us how the material at $(x, 0, 0)$ reflects the radio waves at the various frequencies ω , and thereby gives information about the nature of the material making up the object near the point $(x, 0, 0)$.

There are times, of course, when we do not want to decompose a broadband signal into single-frequency components. A satellite reflecting a TV signal is a broadband point source. All we are interested in is receiving the broadband signal clearly, free of any other interfering sources. The direction of the satellite is known and the antenna is turned to face the satellite. Each location on the parabolic dish reflects the same signal. Because of its parabolic shape, the signals reflected off the dish and picked up at the focal point have exactly the same travel time from the satellite, so they combine coherently, to give us the desired TV signal.

1.12 The Laplace Transform and the Ozone Layer

In the farfield propagation example just considered, we found the measured data to be related to the desired object function by a Fourier transformation. The image reconstruction problem then became one of estimating a function from finitely many noisy values of its Fourier transform. In this section we consider an inverse problem involving the Laplace transform. The example is taken from Twomey's book [86].

1.12.1 The Laplace Transform

The Laplace transform of the function $f(x)$ defined for $0 \leq x < +\infty$ is the function

$$F(s) = \int_0^{+\infty} f(x)e^{-sx} dx.$$

1.12.2 Scattering of Ultraviolet Radiation

The sun emits ultraviolet (UV) radiation that enters the Earth's atmosphere at an angle θ_0 that depends on the sun's position, and with intensity $I(0)$. Let the x -axis be vertical, with $x = 0$ at the top of the atmosphere and x increasing as we move down to the Earth's surface, at $x = X$. The intensity at x is given by

$$I(x) = I(0)e^{-kx/\cos\theta_0}.$$

Within the ozone layer, the amount of UV radiation scattered in the direction θ is given by

$$S(\theta, \theta_0)I(0)e^{kx/\cos\theta_0}\Delta p,$$

where $S(\theta, \theta_0)$ is a known parameter, and Δp is the change in the pressure of the ozone within the infinitesimal layer $[x, x + \Delta x]$, and so is proportional to the concentration of ozone within that layer.

1.12.3 Measuring the Scattered Intensity

The radiation scattered at the angle θ then travels to the ground, a distance of $X - x$, weakened along the way, and reaches the ground with intensity

$$S(\theta, \theta_0)I(0)e^{-kx/\cos\theta_0}e^{-k(X-x)/\cos\theta}\Delta p.$$

The total scattered intensity at angle θ is then a superposition of the intensities due to scattering at each of the thin layers, and is then

$$S(\theta, \theta_0)I(0)e^{-kX/\cos\theta_0}\int_0^X e^{-x\beta} dp,$$

where

$$\beta = k\left[\frac{1}{\cos\theta_0} - \frac{1}{\cos\theta}\right].$$

This superposition of intensity can then be written as

$$S(\theta, \theta_0)I(0)e^{-kX/\cos\theta_0}\int_0^X e^{-x\beta}p'(x)dx.$$

1.12.4 The Laplace Transform Data

Using integration by parts, we get

$$\int_0^X e^{-x\beta}p'(x)dx = p(X)e^{-\beta X} - p(0) + \beta\int_0^X e^{-\beta x}p(x)dx.$$

Since $p(0) = 0$ and $p(X)$ can be measured, our data is then the Laplace transform value

$$\int_0^{+\infty} e^{-\beta x}p(x)dx;$$

note that we can replace the upper limit X with $+\infty$ if we extend $p(x)$ as zero beyond $x = X$.

The variable β depends on the two angles θ and θ_0 . We can alter θ as we measure and θ_0 changes as the sun moves relative to the earth. In this way we get values of the Laplace transform of $p(x)$ for various values of β . The problem then is to recover $p(x)$ from these values. Because the Laplace transform involves a smoothing of the function $p(x)$, recovering $p(x)$ from its Laplace transform is more ill-conditioned than is the Fourier transform inversion problem.

1.13 Summary

Our goal in this chapter has been to introduce the Fourier transform through the use of the example of farfield propagation. For a more detailed discussion of Fourier transforms and Fourier series, see [20]. As our example of farfield propagation shows, the Fourier transform arises naturally in remote sensing and measured data is often related by Fourier transformation to what we really want. The theory also connects the Fourier transform to the important class of convolution filters, which are used to model various types of signal degradation, such as blurring and point-spreading, as well as the limitations on the aperture of the array of sensors.

Chapter 2

Reconstruction from Line-Integral Data

In many tomographic reconstruction problems, the data we have are not Fourier transform values, but line integrals associated with the function of interest. However, such data can, in principle, be used to obtain Fourier transform values, so that reconstruction can be achieved by Fourier inversion. For reasons that we shall explore, this approach is not usually practical. However, it does suggest approximate solution methods, involving convolution filtering and backprojection, that lead to useful algorithms.

We saw earlier that the solar emission problem was solved by formulating it as a problem of reconstruction from line-integral data. We begin here with several other signal-processing problems that require reconstruction of a function from its line integrals, including ocean acoustic tomography, x-ray transmission tomography, and positron- and single-photon emission tomography. Then we establish the connection between the tomography problem and Fourier-transform inversion. Finally, we consider several approaches to Fourier inversion that lead to practical algorithms.

2.1 Ocean Acoustic Tomography

Sound travels in the ocean at approximately $c = 1500$ mps, with deviations from this figure due to water temperature, depth at which the sound is travelling, salinity of the water, and so on. If c is constant, sound emitted at point A at time t will reach point B at time $t + d/c$, where d is the distance from A to B . If we know d and measure the delay in receiving the signal, we can find c . The sound speed is not truly constant, however, but is a function $c(x, y, z)$ of position. In fact, it may depend on time, as well, due, for example, to changing seasons of the year; because temporal changes

are much slower to occur, we usually ignore time-dependence. Determining the spatial sound-speed profile, the function $c(x, y, z)$, is the objective of ocean acoustic tomography.

2.1.1 Obtaining Line-Integral Data

Since the sound speed is not constant, the sound travelling from point A to point B can now take a curved path; the shortest-time route may not be the shortest-distance route. To keep things from getting too complicated in this example, we consider the situation in which the sound still moves from A to B along the straight line segment joining them, but does not travel at a constant speed. We parameterize this line segment with the variable s , with $s = 0$ corresponding to the point A and $s = d$ the point B . We denote by $c(s)$ the sound speed at the point along the line having parameter value s . The time required for the sound to travel from s to $s + \Delta s$ is approximately $\Delta t = \frac{\Delta s}{c(s)}$, so that the signal reaches point B after a delay of $\int_0^d \frac{1}{c(s)} ds$ seconds. Ocean acoustic tomography has as its goal the estimation of the sound speed profile $c(x, y, z)$ from finitely many such line integrals. Because the sound speed is closely related to ocean temperature, ocean acoustic tomography has important applications in weather prediction, as well as in sonar imaging and active and passive sonar detection and surveillance.

2.1.2 The Difficulties

Now let's consider the various obstacles that we face as we try to solve this problem. First of all, we need to design a signal to be transmitted. It must be one from which we can easily and unambiguously determine the delays. When the delayed signal is received, it will not be the only sound in the ocean and must be clearly distinguished from the acoustic background. The processing of the received signals will be performed digitally, which means that we will have to convert the analog functions of the continuous time variable into discrete samples. These vectors of discrete samples will then be processed mathematically to obtain estimates of the line integrals. Once we have determined the line integrals, we must estimate the function $c(x, y, z)$ from these line integrals. We will know the line integrals only approximately and will have only finitely many of them, so the best we can hope to do is to approximate the function $c(x, y, z)$. How well we do will depend on which pairs of sources and receivers we have chosen to use. On the bright side, we have good prior information about the behavior of the sound speed in the ocean, and can specify *a priori* upper and lower bounds on the possible deviations from the nominal speed of 1500 mps. Even so, we need good algorithms that incorporate our prior information. As we shall see later, the Fourier transform will provide an important tool for solving these problems.

2.1.3 Why “Tomography”?

Although the sound-speed profile $c(x, y, z)$ is a function of the three spatial variables, accurate reconstruction of such a three-dimensional function from line integrals would require a large number of lines. In ocean acoustic tomography, as well as in other applications, such as x-ray transmission tomography, the three-dimensional object of interest is studied one slice at a time, so that the function is reduced to a two-dimensional distribution. In fact, the term *tomography*, coming as it does from the Greek word for *part* or *slice*, and thereby related to the word *atom* (“no parts”), is used to describe such problems, because of the early emphasis placed on computationally tractable slice-by-slice reconstruction.

2.1.4 An Algebraic Approach

There is a more algebraic way to reconstruct a function from line integrals. Suppose that we transmit our signal from points A_i , $i = 1, \dots, I$ and receive them at points B_j , $j = 1, \dots, J$. Then we have $N = IJ$ transmitter-receiver pairs, so we have N line integrals, corresponding to N line segments, which we denote L_n , $n = 1, \dots, N$. Imagine the part of the ocean involved to be discretized into M cubes or *voxels*, or, in the slice-by slice approach, two-dimensional squares, or *pixels*, and suppose that within the m th voxel the sound speed is equal to c_m ; also let $x_m = 1/c_m$. For each line segment L_n let P_{nm} be the length of the intersection of line segment L_n with the m th voxel. The time it takes for the acoustic signal to traverse line segment L_n is then approximately

$$(P\mathbf{x})_n = \sum_{m=1}^M P_{nm}x_m,$$

where P denotes the matrix with entries P_{nm} and \mathbf{x} denotes the vector with entries x_m . Our problem now is to solve the system of linear equations $P\mathbf{x} = \mathbf{t}$, where the entries of the vector \mathbf{t} are the travel times we have measured for each line segment. This system can be solved by any number of well known algorithms. Notice that the entries of P , \mathbf{x} and \mathbf{t} are all nonnegative. This suggests that algorithms designed specifically to deal with nonnegative problems may work better. In many cases, both M and N are large, making some algorithms, such as Gauss elimination, impractical, and iterative algorithms competitive.

Although we have presented tomography within the context of ocean acoustics, most of what we have discussed in this section carries over, nearly unchanged, to a number of medical imaging problems.

2.2 X-ray Transmission Tomography

Computer-assisted tomography (CAT) scans have revolutionized medical practice. One example of CAT is x-ray transmission tomography. The goal here is to image the spatial distribution of various matter within the body, by estimating the distribution of x-ray attenuation. Once again, the data are line integrals of the function of interest.

2.2.1 The Exponential-Decay Model

As an x-ray beam passes through the body, it encounters various types of matter, such as soft tissue, bone, ligaments, air, each weakening the beam to a greater or lesser extent. If the intensity of the beam upon entry is I_{in} and I_{out} is its lower intensity after passing through the body, then

$$I_{out} = I_{in} e^{-\int_L f},$$

where $f = f(x, y) \geq 0$ is the *attenuation function* describing the two-dimensional distribution of matter within the slice of the body being scanned and $\int_L f$ is the integral of the function f over the line L along which the x-ray beam has passed. To see why this is the case, imagine the line L parameterized by the variable s and consider the intensity function $I(s)$ as a function of s . For small $\Delta s > 0$, the drop in intensity from the start to the end of the interval $[s, s + \Delta s]$ is approximately proportional to the intensity $I(s)$, to the attenuation $f(s)$ and to Δs , the length of the interval; that is,

$$I(s) - I(s + \Delta s) \approx f(s)I(s)\Delta s.$$

Dividing by Δs and letting Δs approach zero, we get

$$\frac{dI}{ds} = -f(s)I(s).$$

Exercise 2.1 Show that the solution to this differential equation is

$$I(s) = I(0) \exp\left(-\int_{u=0}^{u=s} f(u)du\right).$$

Hint: Use an integrating factor.

From knowledge of I_{in} and I_{out} , we can determine $\int_L f$. If we know $\int_L f$ for every line in the x, y -plane we can reconstruct the attenuation function f . In the real world we know line integrals only approximately and only for finitely many lines. The goal in x-ray transmission tomography is to estimate the attenuation function $f(x, y)$ in the slice, from finitely many noisy measurements of the line integrals. As in the case of ocean acoustic tomography, we usually have prior information about the values that

$f(x, y)$ can take on. We also expect to find sharp boundaries separating regions where the function $f(x, y)$ varies only slightly. Therefore, we need algorithms capable of providing such images.

2.2.2 Difficulties to be Overcome

Once again, there are hurdles to be overcome. X-ray beams are not exactly straight lines; the beams tend to spread out. The x-rays are not monochromatic, and their various frequency components are attenuated at different rates. The beams consist of photons obeying statistical laws, so our algorithms probably should be based on these laws. How we choose the line segments is determined by the nature of the problem; in certain cases we are somewhat limited in our choice of these segments. Patients move; they breathe, their heart beats, and, occasionally, they shift position during the scan. Compensating for these motions is an important, and difficult, aspect of the image reconstruction process. Finally, to be practical in a clinical setting, the processing that leads to the reconstructed image must be completed in a short time, usually around fifteen minutes. This time constraint is what motivates viewing the three-dimensional attenuation function in terms of its two-dimensional slices.

The mathematical similarities between x-ray transmission tomography and ocean acoustic tomography suggest that the reconstruction algorithms used will be similar, and this is the case. As we shall see later, the Fourier transform and the associated theory of convolution filters play important roles.

The data we actually obtain at the detectors are counts of detected photons. These counts are not the line integrals; they are random quantities whose means, or expected values, are related to the line integrals. The Fourier inversion methods for solving the problem ignore its statistical aspects; in contrast, other methods, such as likelihood maximization, are based on a statistical model that involves Poisson-distributed emissions.

2.3 Positron Emission Tomography

In positron emission tomography (PET) and single photon emission tomography (SPECT) the patient inhales, or is injected with, chemicals to which radioactive material has been chemically attached [88]. The chemicals are designed to accumulate in that specific region of the body we wish to image. For example, we may be looking for tumors in the abdomen, weakness in the heart wall, or evidence of brain activity in a selected region. The patient is placed on a table surrounded by detectors that count the number of emitted photons. On the basis of where the various counts were obtained, we wish to determine the concentration of radioactivity at

various locations throughout the region of interest within the patient.

2.3.1 The Coincidence-Detection Model

In PET the radionuclide emits individual positrons, which travel, on average, between 4 mm and 2.5 cm (depending on their kinetic energy) before encountering an electron. The resulting annihilation releases two gamma-ray photons that then proceed in essentially opposite directions. Detection in the PET case means the recording of two photons at nearly the same time at two different detectors. The locations of these two detectors then provide the end points of the line segment passing, more or less, through the site of the original positron emission. Therefore, each possible pair of detectors determines a *line of response* (LOR). When a LOR is recorded, it is assumed that a positron was emitted somewhere along that line. The PET data consists of a chronological list of LOR that are recorded. Because the two photons detected at either end of the LOR are not detected at exactly the same time, the time difference can be used in *time-of-flight* PET to further localize the site of the emission to a smaller segment of perhaps 8 cm in length.

2.3.2 Line-Integral Data

Let the LOR be parameterized by the variable s , with $s = 0$ and $s = L$ denoting the two ends, and L the distance from one end to the other. For a fixed value $s = s_0$, let $P(s)$ be the probability of reaching s for a photon resulting from an emission at s_0 . For small $\Delta s > 0$ the probability that a photon that reached s is absorbed in the interval $[s, s + \Delta s]$ is approximately $\mu(s)\Delta s$, where $\mu(s) \geq 0$ is the photon attenuation density at s . Then $P(s + \Delta s) \approx P(s)[1 - \mu(s)\Delta s]$, so that

$$P(s + \Delta s) - P(s) \approx -P(s)\mu(s)\Delta s.$$

Dividing by Δs and letting Δs go to zero, we get

$$P'(s) = -P(s)\mu(s).$$

It follows that

$$P(s) = e^{-\int_{s_0}^s \mu(t)dt}.$$

The probability that the photon will reach $s = L$ and be detected is then

$$P(L) = e^{-\int_{s_0}^L \mu(t)dt}.$$

Similarly, we find that the probability that a photon will succeed in reaching $s = 0$ from s_0 is

$$P(0) = e^{-\int_0^{s_0} \mu(t)dt}.$$

Since having one photon reach $s = 0$ and the other reach $s = L$ are independent events, their probabilities multiply, so that the probability of a coincident detection along the LOR, due to an emission at s_0 , is

$$e^{-\int_0^L \mu(t)dt}.$$

The expected number of coincident detections along the LOR is then proportional to

$$\int_0^L f(s)e^{-\int_0^L \mu(t)dt} ds = e^{-\int_0^L \mu(t)dt} \int_0^L f(s)ds,$$

where $f(s)$ is the intensity of radionuclide at s . Assuming we know the attenuation function $\mu(s)$, we can estimate the line integral $\int_0^L f(s)ds$ from the number of coincident detections recorded along the LOR. So, once again, we have line-integral data pertaining to the function of interest.

2.4 Single-Photon Emission Tomography

Single-photon emission tomography (SPECT) is similar to PET and has the same objective: to image the distribution of a radionuclide within the body of the patient. In SPECT the radionuclide emits single photons, which then travel through the body of the patient and, in some fraction of the cases, are detected. Detections in SPECT correspond to individual sensor locations outside the body. The data in SPECT are the photon counts at each of the finitely many detector locations. Lead collimators are used in front of the gamma-camera detectors to eliminate photons arriving at oblique angles. While this helps us narrow down the possible sources of detected photons, it also reduces the number of detected photons and thereby decreases the signal-to-noise ratio.

2.4.1 The Line-Integral Model

To solve the reconstruction problem we need a model that relates the count data to the radionuclide density function. A somewhat unsophisticated, but computationally attractive, model is to view the count at a particular detector as the line integral of the radionuclide density function along the line from the detector that is perpendicular to the camera face. The count data then provide many such line integrals and the reconstruction problem becomes the familiar one of estimating a function from noisy measurements of line integrals. Viewing the data as line integrals allows us to use the Fourier transform in reconstruction. The resulting *filtered backprojection* (FBP) algorithm is a commonly used method for medical imaging in clinical settings.

2.4.2 Problems with the Line-Integral Model

It is not really accurate, however, to view the photon counts at each detector as line integrals. Consequently, applying filtered backprojection to the counts at each detector can lead to distorted reconstructions. There are at least three degradations that need to be corrected before FBP can be successfully applied [64]: attenuation, scatter, and spatially dependent resolution.

Some photons never reach the detectors because they are absorbed in the body. As in the PET case, correcting for attenuation requires knowledge of the patient's body; this knowledge can be obtained by performing a transmission scan at the same time. In contrast to the PET case, the attenuation due to absorption is difficult to correct, since it does not involve merely the line integral of the attenuation function, but a half-line integral that depends on the distribution of matter between each photon source and each detector.

As in the PET case previously discussed, the probability that a photon emitted at the point on the line corresponding to the variable $s = s_0$ will reach $s = L$ and be detected is then

$$P(s_0) = e^{-\int_{s_0}^L \mu(t) dt}.$$

If $f(s)$ is the expected number of photons emitted from point s during the scanning, then the expected number of photons detected at L is proportional to

$$\int_0^L f(s) e^{-\int_s^L \mu(t) dt} ds.$$

This quantity varies with the line being considered; the resulting function of lines is called the *attenuated Radon transform*. If the attenuation function μ is constant, then the attenuated Radon transform is called the *exponential Radon transform*.

While some photons are absorbed within the body, others are first deflected and then detected; this is called *scatter*. Consequently, some of the detected photons do not come from where we think they come from. The scattered photons often have reduced energy, compared to *primary*, or unscattered, photons, and scatter-correction can be based on this energy difference; see [64].

Finally, even if there were no attenuation and no scatter, it would be incorrect to view the detected photons as having originated along a straight line from the detector. The detectors have a cone of acceptance that widens as it recedes from the detector. This results in spatially varying resolution. There are mathematical ways to correct for both spatially varying resolution and uniform attenuation [84]. Correcting for the more realistic non-uniform and patient-specific attenuation is more difficult and is the subject of on-going research.

Spatially varying resolution complicates the quantitation problem, which is the effort to determine the exact amount of radionuclide present within a given region of the body, by introducing the *partial volume effect* and *spill-over* (see [88]). To a large extent, these problems are shortcomings of reconstruction based on the line-integral model. If we assume that all photons detected at a particular detector came from points within a narrow strip perpendicular to the camera face, and we reconstruct the image using this assumption, then photons coming from locations outside this strip will be incorrectly attributed to locations within the strip (spill-over), and therefore not correctly attributed to their true source location. If the true source location also has its counts raised by spill-over, the net effect may not be significant; if, however, the true source is a hot spot surrounded by cold background, it gets no spill-over from its neighbors and its true intensity value is underestimated, resulting in the partial-volume effect. The term “partial volume” indicates that the hot spot is smaller than the region that the line-integral model offers as the sources of the emitted photons. One way to counter these effects is to introduce a description of the spatially dependent blur into the reconstruction, which is then performed by iterative methods [80].

In the SPECT case, as in most such inverse problems, there is a trade-off to be made between careful modelling of the physical situation and computational tractability. The FBP method slights the physics in favor of computational simplicity and speed. In recent years, iterative methods that incorporate more of the physics have become competitive.

2.4.3 The Stochastic Model: Discrete Poisson Emitters

In iterative reconstruction we begin by *discretizing* the problem; that is, we imagine the region of interest within the patient to consist of finitely many tiny squares, called *pixels* for two-dimensional processing or cubes, called *voxels* for three-dimensional processing. In what follows we shall not distinguish the two cases, but as a linguistic shorthand, we shall refer to ‘pixels’ indexed by $j = 1, \dots, J$. The detectors are indexed by $i = 1, \dots, I$, the count obtained at detector i is denoted y_i , and the vector $\mathbf{y} = (y_1, \dots, y_I)^T$ is our data. In practice, for the fully three-dimensional case, I and J can be several hundred thousand.

We imagine that each pixel j has its own level of concentration of radioactivity and these concentration levels are what we want to determine. Proportional to these concentration levels are the average rates of emission of photons; the average rate for j we denote by x_j . The goal is to determine the vector $\mathbf{x} = (x_1, \dots, x_J)^T$ from \mathbf{y} .

To achieve our goal we must construct a model that relates \mathbf{y} to \mathbf{x} . The standard way to do this is to adopt the model of *independent Poisson*

emitters. For $i = 1, \dots, I$ and $j = 1, \dots, J$, denote by Z_{ij} the random variable whose value is to be the number of photons emitted from pixel j , and detected at detector i , during the scanning time. We assume that the members of the collection $\{Z_{ij} | i = 1, \dots, I, j = 1, \dots, J\}$ are independent. In keeping with standard practice in modelling radioactivity, we also assume that the Z_{ij} are Poisson-distributed.

We assume that Z_{ij} is a Poisson random variable whose mean value (and variance) is $\lambda_{ij} = P_{ij}x_j$. Here the $x_j \geq 0$ is the average rate of emission from pixel j , as discussed previously, and $P_{ij} \geq 0$ is the probability that a photon emitted from pixel j will be detected at detector i . We then define the random variables $Y_i = \sum_{j=1}^J Z_{ij}$, the total counts to be recorded at detector i ; our actual count y_i is then the observed value of the random variable Y_i . Note that the actual values of the individual Z_{ij} are not observable.

2.4.4 Reconstruction as Parameter Estimation

The goal is to estimate the distribution of radionuclide intensity by calculating the vector \mathbf{x} . The entries of \mathbf{x} are parameters and the data are instances of random variables, so the problem looks like a fairly standard parameter estimation problem of the sort studied in beginning statistics. One of the basic tools for statistical parameter estimation is likelihood maximization, which is playing an increasingly important role in medical imaging. There is several problems, however. One is that the number of parameters is quite large, as large as the number of data values, in most cases. Standard statistical parameter estimation usually deals with the estimation of a handful of parameters. Another problem is that we do not know what the P_{ij} are. These values will vary from one patient to the next, since whether or not a photon makes it from a given pixel to a given detector depends on the geometric relationship between detector i and pixel j , as well as what is in the patient's body between these two locations. If there are ribs or skull getting in the way, the probability of making it goes down. If there are just lungs, the probability goes up. These values can change during the scanning process, when the patient moves. Some motion is unavoidable, such as breathing and the beating heart. Determining good values of the P_{ij} in the absence of motion, and correcting for the effects of motion, are important parts of SPECT image reconstruction.

2.5 Reconstruction from Line Integrals

As we have just seen, a wide variety of applications involve the determination of a function of several variables from knowledge of line integrals of that function. We turn now to the underlying problem of reconstructing

such functions from line-integral data.

2.5.1 The Radon Transform

Our goal is to reconstruct the function $f(x, y)$ from line-integral data. Let θ be a fixed angle in the interval $[0, \pi)$, and consider the rotation of the x, y -coordinate axes to produce the t, s -axis system, where

$$t = x \cos \theta + y \sin \theta,$$

and

$$s = -x \sin \theta + y \cos \theta.$$

We can then write the function f as a function of the variables t and s . For each fixed value of t , we compute the integral $\int f(x, y) ds$, obtaining the integral of $f(x, y) = f(t \cos \theta - s \sin \theta, t \sin \theta + s \cos \theta)$ along the single line L corresponding to the fixed values of θ and t . We repeat this process for every value of t and then change the angle θ and repeat again. In this way we obtain the integrals of f over every line L in the plane. We denote by $r_f(\theta, t)$ the integral

$$r_f(\theta, t) = \int_L f(x, y) ds.$$

The function $r_f(\theta, t)$ is called the *Radon transform* of f .

2.5.2 The Central Slice Theorem

For fixed θ the function $r_f(\theta, t)$ is a function of the single real variable t ; let $R_f(\theta, \omega)$ be its Fourier transform. Then,

$$R_f(\theta, \omega) = \int \left(\int f(x, y) ds \right) e^{i\omega t} dt,$$

which we can write as

$$R_f(\theta, \omega) = \iint f(x, y) e^{i\omega(x \cos \theta + y \sin \theta)} dx dy = F(\omega \cos \theta, \omega \sin \theta),$$

where $F(\omega \cos \theta, \omega \sin \theta)$ is the two-dimensional Fourier transform of the function $f(x, y)$, evaluated at the point $(\omega \cos \theta, \omega \sin \theta)$; this relationship is called the *central slice theorem*. For fixed θ , as we change the value of ω , we obtain the values of the function F along the points of the line making the angle θ with the horizontal axis. As θ varies in $[0, \pi)$, we get all the values of the function F . Once we have F , we can obtain f using the formula for the two-dimensional inverse Fourier transform. We conclude that we are able to determine f from its line integrals.

The Fourier-transform inversion formula for two-dimensional functions tells us that the function $f(x, y)$ can be obtained as

$$f(x, y) = \frac{1}{4\pi^2} \int \int F(u, v) e^{-i(xu+yv)} du dv. \quad (2.1)$$

We now derive alternative inversion formulas.

2.5.3 Ramp Filter, then Backproject

Expressing the double integral in Equation (2.1) in polar coordinates (ω, θ) , with $\omega \geq 0$, $u = \omega \cos \theta$, and $v = \omega \sin \theta$, we get

$$f(x, y) = \frac{1}{4\pi^2} \int_0^{2\pi} \int_0^\infty F(u, v) e^{-i(xu+yv)} \omega d\omega d\theta,$$

or

$$f(x, y) = \frac{1}{4\pi^2} \int_0^\pi \int_{-\infty}^\infty F(u, v) e^{-i(xu+yv)} |\omega| d\omega d\theta.$$

Now write

$$F(u, v) = F(\omega \cos \theta, \omega \sin \theta) = R_f(\theta, \omega),$$

where $R_f(\theta, \omega)$ is the FT with respect to t of $r_f(\theta, t)$, so that

$$\int_{-\infty}^\infty F(u, v) e^{-i(xu+yv)} |\omega| d\omega = \int_{-\infty}^\infty R_f(\theta, \omega) |\omega| e^{-i\omega t} d\omega.$$

The function $g_f(\theta, t)$ defined for $t = x \cos \theta + y \sin \theta$ by

$$g_f(\theta, x \cos \theta + y \sin \theta) = \frac{1}{2\pi} \int_{-\infty}^\infty R_f(\theta, \omega) |\omega| e^{-i\omega t} d\omega \quad (2.2)$$

is the result of a linear filtering of $r_f(\theta, t)$ using a *ramp filter* with transfer function $H(\omega) = |\omega|$. Then,

$$f(x, y) = \frac{1}{2\pi} \int_0^\pi g_f(\theta, x \cos \theta + y \sin \theta) d\theta \quad (2.3)$$

gives $f(x, y)$ as the result of a *backprojection operator*; for every fixed value of (θ, t) add $g_f(\theta, t)$ to the current value at the point (x, y) for all (x, y) lying on the straight line determined by θ and t by $t = x \cos \theta + y \sin \theta$. The final value at a fixed point (x, y) is then the average of all the values $g_f(\theta, t)$ for those (θ, t) for which (x, y) is on the line $t = x \cos \theta + y \sin \theta$. It is therefore said that $f(x, y)$ can be obtained by *filtered backprojection* (FBP) of the line-integral data.

Knowing that $f(x, y)$ is related to the complete set of line integrals by filtered backprojection suggests that, when only finitely many line integrals are available, a similar ramp filtering and backprojection can be used to estimate $f(x, y)$; in the clinic this is the most widely used method for the reconstruction of tomographic images.

2.5.4 Backproject, then Ramp Filter

There is a second way to recover $f(x, y)$ using backprojection and filtering, this time in the reverse order; that is, we backproject the Radon transform and then ramp filter the resulting function of two variables. We begin again with the relation

$$f(x, y) = \frac{1}{4\pi^2} \int_0^{2\pi} \int_0^\infty F(u, v) e^{-i(xu+yv)} \omega d\omega d\theta,$$

which we write as

$$\begin{aligned} f(x, y) &= \frac{1}{4\pi^2} \int_0^{2\pi} \int_0^\infty \frac{F(u, v)}{\sqrt{u^2 + v^2}} \sqrt{u^2 + v^2} e^{-i(xu+yv)} \omega d\omega d\theta \\ &= \frac{1}{4\pi^2} \int_0^{2\pi} \int_0^\infty G(u, v) \sqrt{u^2 + v^2} e^{-i(xu+yv)} \omega d\omega d\theta, \end{aligned} \quad (2.4)$$

using

$$G(u, v) = \frac{F(u, v)}{\sqrt{u^2 + v^2}}$$

for $(u, v) \neq (0, 0)$. Equation (2.4) expresses $f(x, y)$ as the result of performing a two-dimensional ramp filtering of $g(x, y)$, the inverse Fourier transform of $G(u, v)$. We show now that $g(x, y)$ is the backprojection of the function $r_f(\omega, t)$; that is, we show that

$$g(x, y) = \frac{1}{2\pi} \int_0^\pi r_f(\theta, x \cos \theta + y \sin \theta) d\theta.$$

From the central slice theorem we know that $g(x, y)$ can be written as

$$g(x, y) = \frac{1}{2\pi} \int_0^\pi h_g(\theta, x \cos \theta + y \sin \theta) d\theta,$$

where

$$h_g(\theta, x \cos \theta + y \sin \theta) = \frac{1}{2\pi} \int_{-\infty}^\infty R_g(\theta, \omega) |\omega| e^{-i\omega(x \cos \theta + y \sin \theta)} d\omega.$$

Since

$$R_g(\theta, \omega) = G(\omega \cos \theta, \omega \sin \theta),$$

we have

$$g(x, y) = \frac{1}{4\pi^2} \int_0^\pi \int_{-\infty}^\infty G(\omega \cos \theta, \omega \sin \theta) |\omega| e^{-i\omega(x \cos \theta + y \sin \theta)} d\omega d\theta$$

$$\begin{aligned}
&= \frac{1}{4\pi^2} \int_0^\pi \int_{-\infty}^\infty F(\omega \cos \theta, \omega \sin \theta) e^{-i\omega(x \cos \theta + y \sin \theta)} d\omega d\theta \\
&= \frac{1}{4\pi^2} \int_0^\pi \int_{-\infty}^\infty R_f(\theta, \omega) e^{-i\omega(x \cos \theta + y \sin \theta)} d\omega d\theta \\
&= \frac{1}{2\pi} \int_0^\pi r_f(\theta, x \cos \theta + y \sin \theta) d\theta,
\end{aligned}$$

as required.

2.5.5 Radon's Inversion Formula

To get Radon's inversion formula, we need two basic properties of the Fourier transform. First, if $f(x)$ has Fourier transform $F(\gamma)$ then the derivative $f'(x)$ has Fourier transform $-i\gamma F(\gamma)$. Second, if $F(\gamma) = \text{sgn}(\gamma)$, the function that is $\frac{\gamma}{|\gamma|}$ for $\gamma \neq 0$, and equal to zero for $\gamma = 0$, then its inverse Fourier transform is $f(x) = \frac{1}{i\pi x}$.

Writing equation (2.3) as

$$h_f(\theta, t) = \frac{1}{2\pi} \int_{-\infty}^\infty \omega R_f(\theta, \omega) \text{sgn}(\omega) e^{-i\omega t} d\omega,$$

we see that h_f is the inverse Fourier transform of the product of the two functions $\omega R_f(\theta, \omega)$ and $\text{sgn}(\omega)$. Consequently, h_f is the convolution of their individual inverse Fourier transforms, $i \frac{\partial}{\partial t} r_f(\theta, t)$ and $\frac{1}{i\pi t}$; that is,

$$h_f(\theta, t) = \frac{1}{\pi} \int_{-\infty}^\infty \frac{\partial}{\partial t} r_f(\theta, s) \frac{1}{t-s} ds,$$

which is the Hilbert transform of the function $\frac{\partial}{\partial t} r_f(\theta, t)$, with respect to the variable t . Radon's inversion formula is then

$$f(x, y) = \frac{1}{2\pi} \int_0^\pi HT\left(\frac{\partial}{\partial t} r_f(\theta, t)\right) d\theta.$$

2.5.6 Practical Issues

Of course, we never have the Radon transform $r_f(\theta, t)$ for all values of its variables. Only finitely many angles θ are used, and, for each θ , we will have (approximate) values of line integrals for only finitely many t . Therefore, taking the Fourier transform of $r_f(\theta, t)$, as a function of the single variable t , is not something we can actually do. At best, we can approximate $R_f(\theta, \omega)$ for finitely many θ . From the Central Slice Theorem, we can then say that we have approximate values of $F(\omega \cos \theta, \omega \sin \theta)$, for finitely many θ . This means that we have (approximate) Fourier transform values for $f(x, y)$ along finitely many lines through the origin, like the spokes of a wheel. The

farther from the origin we get, the fewer values we have, so the *coverage* in Fourier space is quite uneven. The low-spatial-frequencies are much better estimated than higher ones, meaning that we have a low-pass version of the desired $f(x, y)$. The filtered backprojection approaches we have just discussed both involve ramp filtering, in which the higher frequencies are increased, relative to the lower ones. This too can only be implemented approximately, since the data is noisy and careless ramp filtering will cause the reconstructed image to be unacceptably noisy.

2.6 Summary

We have seen how the problem of reconstructing a function from line integrals arises in a number of applications. The Central Slice Theorem connects the line integrals and the Radon transform to the Fourier transform of the desired distribution. Various approaches to implementing the Fourier Inversion Formula lead to filtered backprojection algorithms for the reconstruction. In x-ray tomography and PET, viewing the data as line integrals ignores the statistical aspects of the problem, and in SPECT, it ignores, as well, the important physical effects of attenuation. To incorporate more of the physics of the problem, iterative algorithms based on statistical models have been developed. We shall consider some of these algorithms later.

Chapter 3

Discrete Signal Processing

Although we usually model real-world distributions as functions of continuous variables, while the data we actually obtain are finite, it is standard practice to develop signal processing fundamentals within the context of infinite sequences, or functions of discrete variables. Infinite sequences arise when we sample functions of continuous variables, or when we extend finite data. Within the context of discrete signal processing, Fourier series replace Fourier transforms as the key mathematical tool. The Shannon sampling theorem provides the link between these two branches of Fourier analysis.

3.1 Discrete Signals

A discrete signal is a function $x = \{x(n)\}$ defined for all integers n . In signal processing, such discrete signals are often the result of *sampling* a function of a continuous variable. In our discussion of farfield propagation, we saw that the data gathered at each sensor effected a sampling of the Fourier transform, $F(\gamma)$, of the distant distribution $f(x)$. In the theoretical situation in which we had available an infinite discrete set of sensors, we would have an infinite sequence, obtained by sampling the function $F(\gamma)$. In many applications, the function that is being sampled is a function of time, say $f(t)$; we shall use this example in our discussion here.

In the most common case, that of equispaced sampling, we have $x(n) = f(n\Delta)$, where $\Delta > 0$ is the sampling interval. Generally, such discrete signals are neither a realistic model of the physical situation nor an accurate description of what we have actually obtained through measurement. Nevertheless, discrete signals provide the most convenient framework within which to study the the basic tools of signal processing coming from Fourier analysis.

3.2 Notation

It is common practice to denote functions of a discrete variable by the letters x, y or z , as well as f, g or h . So we speak of the discrete signals $x = \{x(n) = 2n - 1, -\infty < n < \infty\}$ or $y = \{y(n) = -n^3 + n, -\infty < n < \infty\}$. For convenience, we often just say $x(n) = 2n - 1$ or $y(n) = n^3 + n$ when we mean the whole function x or y . However, if k is regarded as a fixed, but unspecified, integer, $x(k)$ means the value of the function x at k . This is really the same thing that we do in calculus, when we define a function $f(x) = x^2 - 6$ and then speak about the value of this function at the point $x = a$, denoted $f(a)$. Speaking more precisely, in the first instance, n is a variable, while k is a parameter, and in the second instance, x is a variable, while a is a parameter; variables change their values during the course of the problem, while parameters have values that are chosen at the outset and retain their chosen values throughout the problem.

There are two special discrete signals with *reserved names*, δ and u : $\delta(0) = 1$ and $\delta(n) = 0$, for $n \neq 0$; $u(n) = 1$, for $n \geq 0$ and $u(n) = 0$ for $n < 0$. When we say that their names are reserved we mean that whenever you see these names you can (usually) assume that they refer to the same functions as just defined; in calculus e^x and $\sin x$ are reserved names, while in signal processing δ and u are reserved names.

3.3 Operations on Discrete Signals

Because discrete signals are functions, we can perform on them many of the operations we perform on functions of a continuous variable. For instance, we can add discrete signals x and y , to get the discrete signal $x + y$, we can multiply x by a real number c to get the discrete signal cx , we can multiply x and y to get xy , and so on. We can *shift* x to the right k units to get y with $y(n) = x(n - k)$. Notice that, if we shift $x = \delta$ to the right k units, we have y with $y(k) = 1$ and $y(n) = 0$ for $n \neq k$; we call this function δ_k , so we sometimes say that $\delta = \delta_0$.

In general, an operation, or, to use the official word, an *operator*, T works on a discrete signal x to produce another discrete signal y ; we describe this situation by writing $y = T(x)$. For example, the operator $T = S_k$ shifts any x to the right by k units; for example, $S_3(\delta) = \delta_3$. We are particularly interested in operators that possess certain nice properties.

3.3.1 Linear Operators

An operator T is called *linear* if, for any x and z and numbers a and b we have $T(ax + bz) = aT(x) + bT(z)$; for example, the operator $T = S_k$ is linear.

Exercise 3.1 Which of the following operators are linear?

- $T(x)(n) = x(n-1) + x(n)$;
- $T(x)(n) = nx(n)$;
- $T(x)(n) = x(n)^2$.

3.3.2 Shift-invariant Operators

Notice that operators are also functions, although not the sort that we usually study; their domains and ranges consist of functions. We have seen such operator-type functions in calculus class- the operator that transforms a function into its derivative is an operator-type function. Therefore we can combine operators using composition, in the same way we compose functions. The composition of operators T and S is the operator that first performs S and then performs T on the result; that is, the composition of T and S begins with x and ends with $y = T(S(x))$. Notice that, just as with ordinary functions, the order of the operators in the composition matters; $T(S(x))$ and $S(T(x))$ need not be the same discrete signal. We say that operators T and S *commute* if $T(S(x)) = S(T(x))$, for all x ; in that case we write $TS = ST$.

An operator T is said to be *shift-invariant* if $TS_k = S_kT$ for all integers k . This means that if y is the output of the system described by T when the input is x , then when we shift the input by k , from x to S_kx , all that happens to the output is that the y is also shifted by k , from y to S_ky . For example, suppose that T is the squaring operator, defined by $T(x) = y$ with $y(n) = x(n)^2$. Then T is shift-invariant. On the other hand, the operator T with $y = T(x)$ such that $y(n) = x(-n)$ is not shift-invariant.

Exercise 3.2 Which of the following operators are shift-invariant?

- $T(x)(n) = x(0) + x(n)$;
- $T(x)(n) = x(n) + x(-n)$;
- $T(x)(n) = \sum_{k=-2}^2 x(n+k)$.

We are most interested in operators T that are both linear and shift-invariant; these are called LSI operators. An LSI operator T is often viewed as a linear system having inputs called x and outputs called y , where $y = T(x)$, and we speak of a LSI system.

3.3.3 Convolution Operators

Let h be a fixed discrete signal. For any discrete signal x define $y = T(x)$ by

$$y(n) = \sum_{k=-\infty}^{\infty} h(k)x(n-k),$$

for any integer n . We then say that y is the *convolution* of x with h and write $y = x * h$. Notice that $x * h = h * x$; that is,

$$\sum_{k=-\infty}^{\infty} h(k)x(n-k) = \sum_{k=-\infty}^{\infty} x(k)h(n-k).$$

The operator T is then the *convolution with h* operator. Any such T is linear.

3.3.4 LSI Filters are Convolutions

The operator T that is convolution with h is linear and shift-invariant. The most important fact in signal processing is that every T that is *linear and shift-invariant* (LSI) must be convolution with h , for some fixed discrete signal h .

Because of the importance of this result we give a proof now. First, we must find the h . To do this we let $x = \delta$; the h we seek is then the output $h = T(\delta)$. Now we must show that, for any other input x , we have $T(x) = x * h$. Note that for any k we have $\delta_k = S_k(\delta)$, so that

$$T(\delta_k) = T(S_k(\delta)) = S_k(T(\delta)) = S_k(h),$$

and so

$$T(\delta_k)(n) = S_k(h)(n) = h(n-k).$$

We can write an arbitrary x in terms of the δ_k as

$$x = \sum_{k=-\infty}^{\infty} x(k)\delta_k.$$

Then

$$T(x)(n) = T\left(\sum_{k=-\infty}^{\infty} x(k)\delta_k\right)(n) = \sum_{k=-\infty}^{\infty} x(k)T(\delta_k)(n) = \sum_{k=-\infty}^{\infty} x(k)h(n-k).$$

Therefore, $T(x) = x * h$, as we claimed. Because the h associated with the operator T is $h = T(\delta)$, the discrete signal h is called the *impulse-response function* of the system.

3.4 Special Types of Discrete Signals

Some of our calculations, such as convolution, involve infinite sums. In order for these sums to make sense we would need to impose certain restrictions on the discrete signals involved. Some books consider only discrete

signals x that are *absolutely summable*, that is, for which

$$\sum_{n=-\infty}^{\infty} |x(n)| < \infty,$$

or, at least, x that are *bounded*, which means that there is a positive constant $b > 0$ with $|x(n)| \leq b$ for all n . Sometimes the condition of absolute summability is imposed only on discrete functions h that are to be associated with LSI operators. Operators T whose h is absolutely summable have the desirable property of *stability*; that is, if the input function x is bounded, so is the output function $y = T(x)$. This property is also called the *bounded in, bounded out* (BIBO) property.

Exercise 3.3 Show that the operator T is a stable operator if and only if its associated h is absolutely summable. Hint: If h is not absolutely summable, consider the input sequence with $x(n) = \overline{h(-n)}/|h(n)|$.

In order to make use of the full power of Fourier methods some texts require that discrete signals x be *absolutely square-summable*, that is,

$$\sum_{n=-\infty}^{\infty} |x(n)|^2 < \infty.$$

Exercise 3.4 Show that the discrete signal $x(n) = \frac{1}{|n|+1}$ is absolutely square-summable, but not absolutely summable.

Our approach will be to avoid discussing specific requirements, with the understanding that some requirements will usually be needed to make the mathematics rigorous.

3.5 The Frequency-Response Function

Just as sine and cosine functions play important roles in calculus, so do their discrete counterparts in signal processing. The discrete sine function with frequency ω is the discrete signal \sin_{ω} with

$$\sin_{\omega}(n) = \sin(\omega n),$$

for each integer n . Similarly, the discrete cosine function with frequency ω is \cos_{ω} with

$$\cos_{\omega}(n) = \cos(\omega n).$$

It is convenient to include in the discussion the complex exponential e_{ω} defined by

$$e_{\omega}(n) = \cos_{\omega}(n) + i \sin_{\omega}(n) = e^{i\omega n}.$$

Since these discrete signals are the same for ω and $\omega + 2\pi$ we assume that ω lies in the interval $[-\pi, \pi)$.

3.5.1 The Response of a LSI System to $x = e_\omega$

Let T denote a LSI system and let ω be fixed. We show now that

$$T(e_\omega) = He_\omega,$$

for some constant H . Since the H can vary as we change ω it is really a function of ω , so we denote it $H = H(\omega)$.

Let $v = \{v(n)\}$ be the signal $v = e_\omega - S_1(e_\omega)$. Then we have

$$v(n) = e^{in\omega} - e^{i(n-1)\omega} = (1 - e^{-i\omega})e^{in\omega}.$$

Therefore, we can write

$$v = (1 - e^{-i\omega})e_\omega,$$

from which it follows that

$$T(v) = (1 - e^{-i\omega})T(e_\omega). \quad (3.1)$$

But we also have

$$T(v) = T(e_\omega - S_1(e_\omega)) = T(e_\omega) - TS_1(e_\omega),$$

and, since T is shift-invariant, $TS_1 = S_1T$, we know that

$$T(v) = T(e_\omega) - S_1T(e_\omega). \quad (3.2)$$

Combining Equations (3.1) and (3.2), we get

$$(1 - e^{-i\omega})T(e_\omega) = T(e_\omega) - S_1T(e_\omega).$$

Therefore,

$$S_1T(e_\omega) = e^{-i\omega}T(e_\omega),$$

or

$$T(e_\omega)(n-1) = S_1T(e_\omega)(n) = e^{-i\omega}T(e_\omega)(n).$$

We conclude from this that

$$e^{in\omega}T(e_\omega)(0) = T(e_\omega)(n),$$

for all n . Finally, we let $H(\omega) = T(e_\omega)(0)$.

It is useful to note that we did not use here the fact that T is a convolution operator. However, since we do know that $T(x) = x * h$, for $h = T(\delta)$, we can relate the function $H(\omega)$ to h .

3.5.2 Relating $H(\omega)$ to $h = T(\delta)$

Since T is a LSI operator, T operates by convolving with $h = T(\delta)$. Consider what happens when we select for the input the discrete signal $x = e_\omega$. Then the output is $y = T(e_\omega)$ with

$$y(n) = \sum_{k=-\infty}^{\infty} h(k)e^{i\omega(n-k)} = H(e^{i\omega})e^{i\omega n},$$

where

$$H(e^{i\omega}) = \sum_{k=-\infty}^{\infty} h(k)e^{-i\omega k} \quad (3.3)$$

is the value, at ω , of the *frequency-response function* of T . The point here is that when the input is $x = e_\omega$ the output is a multiple of e_ω , the multiplier being the (possibly complex) number $H(e^{i\omega})$. Linear, shift-invariant systems T do not alter the frequency of the input, but just change its amplitude and/or phase. The constant $H(e^{i\omega})$ is the same as $H(\omega)$ obtained earlier; having two different notations for the same function is an unfortunate, but common, occurrence in the signal-processing literature.

It is important to note that the infinite sum in Equation (3.3) need not converge for arbitrary $h = \{h(k)\}$. It does converge, obviously, whenever h is finitely nonzero; it will also converge for infinitely nonzero sequences that are suitably restricted.

A common problem in signal processing is to design a LSI filter with a desired frequency-response function $H(e^{i\omega})$. To determine $h(m)$, given $H(e^{i\omega})$, we multiply both sides of Equation (3.3) by $e^{i\omega m}$, multiply by $\frac{1}{2\pi}$, integrate over the interval $[-\pi, \pi]$, and use the helpful fact that

$$\int_{-\pi}^{\pi} e^{i(m-k)\omega} d\omega = 0,$$

for $m \neq k$. The result is

$$h(m) = \frac{1}{2\pi} \int_{-\pi}^{\pi} H(e^{i\omega})e^{i\omega m} d\omega. \quad (3.4)$$

It is useful to extend the definition of $H(e^{i\omega})$ to permit $e^{i\omega}$ to be replaced by any complex number z . Then we get the z -transform of h , given by

$$H(z) = \sum_{k=-\infty}^{\infty} h(k)z^{-k}.$$

We can study the working of the system T on more general inputs x by representing x as a sum of complex-exponential discrete signals e_ω .

The representation, in Equation (3.4), of the infinite sequence $h = \{h(k)\}$ as a superposition of complex-exponential discrete signals suggests the possibility that such a representation is available for general infinite discrete signals, a notion we take up in the next section.

3.6 The Discrete Fourier Transform

A common theme running through mathematics is the representation of complicated objects in terms of simpler ones. Taylor-series expansion enables us to view quite general functions as infinite versions of polynomials by representing them as infinite sums of the power functions. Fourier-series expansions give representations of quite general functions as infinite sums of sines and cosines. Here we obtain similar representation for discrete signals, as infinite sums of the complex exponentials, e_{ω} , for ω in $[-\pi, \pi)$.

Our goal is to represent a general discrete signal x as a sum of the e_{ω} , for ω in the interval $[-\pi, \pi)$. Such a sum is, in general, an integral over ω . So we seek to represent x as

$$x(n) = \frac{1}{2\pi} \int_{-\pi}^{\pi} X(\omega) e^{i\omega n} d\omega, \quad (3.5)$$

where $X(\omega)$ is a function to be determined. As we shall see, the function we seek is the *discrete Fourier transform* (DFT) of x , defined by

$$X(\omega) = \sum_{m=-\infty}^{\infty} x(m) e^{-i\omega m}. \quad (3.6)$$

This follows from the discussion leading up to Equation (3.4). Notice that in the case $x = h$ the function $H(\omega)$ is the same as the frequency-response function $H(e^{i\omega})$ defined earlier. For this reason the notation $X(\omega)$ and $X(e^{i\omega})$ are used interchangeably. The DFT of the discrete signal x is sometimes called the *discrete-time Fourier transform* (DTFT).

The sum in Equation (3.6) is the *Fourier-series expansion* for the function $X(\omega)$, over the interval $[-\pi, \pi)$; the $x(n)$ are its *Fourier coefficients*.

The infinite series in Equation (3.4) that is used to define $X(\omega)$ may not converge. For example, suppose that x is an exponential signal, with $x(n) = e^{i\omega_0 n}$. Then the infinite sum would be

$$\sum_{m=-\infty}^{\infty} e^{i(\omega_0 - \omega)m},$$

which obviously does not converge, at least in any ordinary sense. Consider, though, what happens when we put this sum inside an integral and reverse

the order of integration and summation. Specifically, consider

$$\begin{aligned} & \frac{1}{2\pi} \int_{-\pi}^{\pi} F(\omega) \sum_{m=-\infty}^{\infty} e^{i(\omega_0-\omega)m} d\omega, \\ &= \sum_{m=-\infty}^{\infty} \left(\frac{1}{2\pi} \int_{-\pi}^{\pi} F(\omega) e^{i(\omega_0-\omega)m} d\omega \right), \\ &= \sum_{m=-\infty}^{\infty} e^{i\omega_0 m} f(m) = F(\omega_0). \end{aligned}$$

So, the infinite sum acts like the Dirac delta $\delta(\omega - \omega_0)$. This motivates the following definition of the infinite sum:

$$\sum_{m=-\infty}^{\infty} e^{i(\omega_0-\omega)m} = \delta(\omega - \omega_0). \quad (3.7)$$

A different approach to the infinite sum is to consider

$$\lim_{N \rightarrow +\infty} \frac{1}{2N+1} \sum_{m=-N}^N e^{i(\omega_0-\omega)m}.$$

According to Equation (??), we have

$$\sum_{n=-N}^N e^{i\omega n} = \frac{\sin(\omega(N + \frac{1}{2}))}{\sin(\frac{\omega}{2})}.$$

Therefore,

$$\lim_{N \rightarrow +\infty} \frac{1}{2N+1} \sum_{m=-N}^N e^{i(\omega_0-\omega)m} = 1, \quad (3.8)$$

for $\omega = \omega_0$, and zero, otherwise.

3.7 The Convolution Theorem

Once again, let $y = T(x)$, where T is a LSI operator with associated filter $h = \{h(k)\}$. Because we can write

$$x(n) = \frac{1}{2\pi} \int_{-\pi}^{\pi} X(\omega) e_{\omega}(n) d\omega,$$

or, in shorthand, leaving out the n , as

$$x = \frac{1}{2\pi} \int_{-\pi}^{\pi} X(\omega) e_{\omega} d\omega,$$

we have

$$\begin{aligned} y &= T(x) = \frac{1}{2\pi} \int_{-\pi}^{\pi} X(\omega)T(e_{\omega})d\omega, \\ &= \frac{1}{2\pi} \int_{-\pi}^{\pi} X(\omega)H(\omega)e_{\omega}d\omega, \end{aligned}$$

or

$$y(n) = \frac{1}{2\pi} \int_{-\pi}^{\pi} X(\omega)H(\omega)e_{\omega}(n)d\omega.$$

But we also have

$$y(n) = \frac{1}{2\pi} \int_{-\pi}^{\pi} Y(\omega)e_{\omega}(n)d\omega,$$

from which we conclude that

$$Y(\omega) = X(\omega)H(\omega), \quad (3.9)$$

for each ω in $[-\pi, \pi)$.

Equation (3.9) is the most important equation in signal processing. It describes the activity of an LSI system by telling us that the system simply multiplies the DFT of the input x by the DFT of the h , the frequency-response function of the system, to produce the DFT of the output y . Since $y = x * h$ it also tells us that whenever y is formed by convolving two discrete signals x and h , its DFT is the product of the DFT of x and the DFT of h .

3.8 Sampling and Aliasing

The term *sampling* refers to the transition from a function $f(t)$ of a continuous variable to a discrete signal x , defined by $x(n) = f(n\Delta)$, where $\Delta > 0$ is the *sample spacing*. For example, suppose that $f(t) = \sin(\gamma t)$ for some frequency $\gamma > 0$. Then $x(n) = \sin(\gamma n\Delta) = \sin(\omega n)$, where $\omega = \gamma\Delta$. We define $X(\omega)$, the DFT of the discrete signal x , for $|\omega| \leq \pi$, so we need $|\gamma|\Delta \leq \pi$. This means we must select Δ so that $\Delta \leq \pi/|\gamma|$. In general, if the function $f(t)$ has sinusoidal components with frequencies γ such that $|\gamma| \leq \Gamma$ then we should select $\Delta \leq \pi/\Gamma$.

If we select Δ too large, then a frequency component of $f(t)$ corresponding to $|\gamma| > \pi/\Delta$ will be mistaken for a frequency with smaller magnitude. This is *aliasing*. For example, if $f(t) = \sin(3t)$, but $\Delta = \pi/2$, then the frequency $\gamma = 3$ is mistaken for the frequency $\gamma = -1$, which lies in $[-2, 2]$. When we sample we get

$$x(n) = \sin(3\Delta n) = \sin(-\Delta n + 4\Delta n) = \sin(-\Delta n + 2\pi n) = \sin(-\Delta n),$$

for each n .

3.9 Important Problems in Discrete Signal Processing

A number of important problems in signal processing involve the relation between a discrete signal and its DFT. One problem is the design of a system to achieve a certain desired result, such as low-pass filtering. A second problem is to estimate the $X(\omega)$ when we do not have all the values $x(n)$, but only finitely many of them.

3.9.1 Low-pass Filtering

When we represent a discrete signal x using

$$x(n) = \frac{1}{2\pi} \int_{-\pi}^{\pi} X(\omega) e_{\omega}(n) d\omega,$$

we take the point of view that the function x is made up of the various discrete sinusoids, the functions e_{ω} , each contributing in the amount $\frac{1}{2\pi} X(\omega)$. Since $X(\omega)$ is usually complex we must interpret this in terms of both an amplitude modulation and a phase change. Suppose that, for some fixed Ω in the interval $(0, \pi)$, we wish to design a system that will leave $X(\omega)$ unchanged for those ω in the interval $[-\Omega, \Omega]$ and change $X(\omega)$ to zero otherwise; such a system is called the (ideal) Ω -low-pass filter. To achieve this result we need to take $H(\omega)$ to be $\chi_{\Omega}(\omega)$, the characteristic function of the interval $[-\Omega, \Omega]$, with $\chi_{\Omega}(\omega) = 1$, for $|\omega| \leq \Omega$, and $\chi_{\Omega}(\omega) = 0$, otherwise. We find the $h(k)$ using

$$h(k) = \frac{1}{2\pi} \int_{-\pi}^{\pi} \chi_{\Omega}(\omega) e^{i\omega k} d\omega.$$

Performing the integration, we find that $h(0) = \Omega/\pi$ and, for $k \neq 0$,

$$h(k) = \frac{\sin \Omega k}{\pi k}.$$

To calculate the low-pass output

$$y(n) = \sum_{k=-\infty}^{\infty} \frac{\sin \Omega k}{\pi k} x(n-k)$$

we need infinitely many values $x(m)$ for $m > n$, as well as infinitely many values for $m < n$. If we think of n as time, then to calculate the value of y at time n we need to know the values of x for the entire infinite past before time n , as well as the values for the entire infinite future after time n . Clearly, this is inconvenient if we wish to perform the filtering in real-time. One goal of signal processing is to approximate such filters with ones that are more convenient, using, say, only finitely many past and future values of the input.

3.9.2 The Finite-Data Problem

In practice we have finite data obtained from measurements. We view these data as values $x(n)$ for finitely many values of n , say $n = 0, 1, \dots, N - 1$. The function $X(\omega)$ often is an important object in the problem and must be estimated from the data. One possible estimate is

$$\hat{X}(\omega) = \sum_{n=0}^{N-1} x(n)e^{-i\omega n}.$$

To distinguish this from the DFT, which involves the infinite sum, we shall call $\hat{X}(\omega)$ the DFT of the vector $\mathbf{x} = (x(0), \dots, x(N-1))^T$. If N is large, the DFT of \mathbf{x} will usually be a satisfactory approximation of $X(\omega)$. However, in many applications N is not large and the DFT of \mathbf{x} is not adequate. The *finite-data problem* is how to find better estimates of $X(\omega)$ from the limited data we have.

Because the finite-data problem involves approximating one function of a continuous variable by another, we need some way to measure how far apart two such functions are. The way most commonly used in signal processing is the so-called *Hilbert-space distance*, given by

$$\|X(\omega) - Y(\omega)\| = \sqrt{\int_{-\pi}^{\pi} |X(\omega) - Y(\omega)|^2 d\omega}.$$

We shall return later to the problem of describing best approximations in Hilbert space.

3.9.3 The Extrapolation Problem

If $x(n)$ is obtained from $f(t)$ by sampling, that is, $x(n) = f(n\Delta)$, we have

$$f(n\Delta) = \frac{1}{2\pi} \int_{-\pi}^{\pi} X(\omega)e^{in\omega} d\omega. \quad (3.10)$$

Changing to the variable $\gamma = \omega/\Delta$, and defining $\Gamma = \pi/\Delta$, we can write

$$f(n\Delta) = \frac{\Delta}{2\pi} \int_{-\Gamma}^{\Gamma} X(\gamma\Delta)e^{i(n\Delta)\gamma} d\gamma, \quad (3.11)$$

which makes clearer the use of the sampling time $t = n\Delta$.

The representation in Equation (3.11) is suggestive! Let us define $g(t)$ for all t by the formula

$$g(t) = \frac{\Delta}{2\pi} \int_{-\Gamma}^{\Gamma} X(\gamma\Delta)e^{it\gamma} d\gamma. \quad (3.12)$$

Do we have $g(t) = f(t)$ for all t ? On the face of it, it would seem that the answer is clearly no. How could a function of a continuous variable be completely determined by such a sequence of its values? How can we capture all of a function $f(t)$ from discrete samples? In fact, under certain conditions, the answer is yes. Let us investigate what those conditions might be.

Let $\epsilon > 0$ and let $h_\epsilon(t) = \sin((\Gamma + \epsilon)t) - \sin((-\Gamma + \epsilon)t)$. Then $h_\epsilon(n\Delta) = 0$ for each integer n . From the data we have, we cannot decide if $f(t) = g(t)$ or $f(t) = g(t) + h_\epsilon(t)$, or, perhaps, $f(t) = g(t) + h_\epsilon(t)$ for some other ϵ . Notice that, in order to construct $h_\epsilon(t)$ we need a sine function with a frequency outside the interval $[-\Gamma, \Gamma]$.

On the other hand, if $F(\gamma)$, the Fourier transform of $f(t)$, is zero outside $[-\Gamma, \Gamma]$, then $f(t) = g(t)$. This is because the function $F(\gamma)$ has a Fourier-series representation

$$F(\gamma) = \sum_{n=-\infty}^{\infty} a_n e^{i\gamma n\Delta},$$

where, as in our discussion of the DFT, we have

$$a_n = \frac{1}{2\Gamma} \int_{-\Gamma}^{\Gamma} F(\gamma) e^{-i\gamma n\Delta} d\gamma.$$

But the expression on the right side of this equation equals $\Delta f(n\Delta)$, according to the Fourier Inversion Formula. Therefore

$$\begin{aligned} F(\gamma) &= \Delta \sum_{n=-\infty}^{\infty} f(n\Delta) e^{i\gamma n\Delta} \\ &= \Delta \sum_{n=-\infty}^{\infty} x(n) e^{i\gamma n\Delta} \\ &= \Delta \sum_{n=-\infty}^{\infty} x(n) e^{i\omega n} = \Delta X(-\gamma\Delta). \end{aligned}$$

So, we can write

$$\begin{aligned} f(t) &= \frac{1}{2\pi} \int_{-\Gamma}^{\Gamma} F(\gamma) e^{-it\gamma} d\gamma, \\ &= \frac{\Delta}{2\pi} \int_{-\Gamma}^{\Gamma} X(\gamma\Delta) e^{it\gamma} d\gamma = g(t). \end{aligned}$$

For an arbitrary function $f(t)$ we seek a representation of $f(t)$ as a superposition of complex exponential functions, that is,

$$f(t) = \frac{1}{2\pi} \int_{-\infty}^{\infty} A(\gamma) e^{it\gamma} d\gamma, \quad (3.13)$$

for some function $A(\gamma)$. The function $A(\gamma)$ that does the job is $A(\gamma) = F(-\gamma)$, where $F(\gamma)$ is the *Fourier transform* of $f(t)$. If $F(\gamma) = 0$ for $|\gamma| > \Gamma$, then $f(t)$ is said to be Γ -*bandlimited*; in this case $F(\gamma) = \Delta X(-\gamma\Delta)$, as discussed previously.

It is important to note that we cannot tell from the samples $x(n) = f(n\Delta)$ whether or not $f(t)$ is Γ -bandlimited. If $f(t)$ is not Γ -bandlimited, but we assume that it is, there will be components of $f(t)$ with frequencies outside the band $[-\Gamma, \Gamma]$ that will be mistaken for sinusoids having frequencies inside the band; this is aliasing.

3.10 Discrete Signals from Finite Data

In problems involving actual data obtained from measurements we may have a vector $\mathbf{x} = (x_1, \dots, x_N)^T$ that we wish to associate with a discrete function x . There are, of course, any number of ways to do this. Two of the most commonly used ways employ *zero extension* and *periodic extension*.

3.10.1 Zero-extending the Data

We define $x(n)$ to be x_{n+1} , for $n = 0, \dots, N - 1$ and $x(n) = 0$ otherwise. Then x is a discrete function that extends the data. The DFT of x is now

$$X(\omega) = \sum_{n=0}^{N-1} x(n)e^{-in\omega}, \quad (3.14)$$

for $|\omega| \leq \pi$ and, from the fact that

$$0 = \int_{-\pi}^{\pi} e^{i(m-n)\omega} d\omega$$

for $m \neq n$, we have

$$x(m) = \frac{1}{2\pi} \int_{-\pi}^{\pi} X(\omega)e^{im\omega} d\omega,$$

for all integers m .

The DFT of x obtained by zero-extending the data provides a way to represent the data as a (continuous) sum, or integral, of the discrete exponential functions e_ω :

$$x_n = \frac{1}{2\pi} \int_{-\pi}^{\pi} X(\omega)e^{i(n-1)\omega} d\omega,$$

for $n = 1, \dots, N$.

3.10.2 Periodically Extending the Data

Another way to associate a discrete function \tilde{x} with the data vector \mathbf{x} is by extending the data periodically. For $n = 0, \dots, N - 1$ let $\tilde{x}(n) = x_{n+1}$ and for any integer n define $\tilde{x}(n) = \tilde{x}(n \bmod N)$. Then \tilde{x} extends the data and is N -periodic; that is, $\tilde{x}(n + N) = \tilde{x}(n)$ for all integers n .

Now we want to represent the N -periodic \tilde{x} as a sum of the discrete exponential functions e_ω . Notice, however, that most of the e_ω are not N -periodic; in order for $e^{i(n+N)\omega} = e^{in\omega}$ for all integers n we need $e^{iN\omega} = 1$. This means that $\omega = 2\pi k/N$, for some integer k . Therefore, we shall seek to represent \tilde{x} as a sum of the discrete exponential functions e_ω only for $\omega = 2\pi k/N$. Let us denote such functions as e_k . Notice also that e_{k+N} and e_k are the same function, for any integer k . Therefore, we seek to represent \tilde{x} as a sum of the discrete exponential functions e_k , for $k = 0, 1, \dots, N - 1$; that is, we want

$$\tilde{x}(n) = \sum_{k=0}^{N-1} X_k e^{2\pi i k n / N}, \quad (3.15)$$

for some choice of numbers X_k .

To determine the X_k we multiply both sides of Equation (3.15) by $e^{-2\pi i j n / N}$ and sum over n . Using the fact that

$$\sum_{n=0}^{N-1} e^{2\pi i (k-j)n / N} = 0,$$

if $k \neq j$, it follows that

$$X_j = \frac{1}{N} \sum_{n=0}^{N-1} x(n) e^{-2\pi i j n / N}, \quad (3.16)$$

for $j = 0, \dots, N - 1$.

We began with a finite vector $\mathbf{x} = (x_1, \dots, x_N)^T$, which we chose to write as $\mathbf{x} = (x(0), \dots, x(N-1))^T$, and ended with a finite set of numbers X_j , $j = 0, \dots, N - 1$, which we used to form the vector $\mathbf{X} = (X_0, \dots, X_{N-1})^T$. It is common practice to call the vector \mathbf{X} the DFT of the vector \mathbf{x} , but to avoid confusion, we shall refer to the vector \mathbf{X} as the *vector* DFT (vDFT) of the vector \mathbf{x} , leaving the terminology DFT of \mathbf{x} to refer to the DFT of the zero-extended discrete function x in equation (3.14). Notice, though, that the vDFT and the DFT are related; for $0 \leq k \leq N/2$ we have $X_k = X(2\pi k/N)$ and for $N/2 < k \leq N - 1$ we have $X_k = X(-\pi + 2\pi k/N)$. The vector DFT plays an important role in signal processing because, as we shall see later, there is a fast algorithm for calculating it from the data, called the *fast Fourier transform* (FFT).

3.10.3 A Third Way to Extend the Data

Another way to extend the data vector to a discrete function is to *zero-pad* and then to perform periodic extension. Given the data $x(n)$, $n = 0, \dots, N-1$, let $x(n) = 0$, $n = N, N+1, \dots, M-1$. Then extend these M numbers M -periodically, so that $\tilde{x}(n) = x(n \bmod M)$, for each integer n . Then $\tilde{x}(n+M) = \tilde{x}(n)$, for all n .

Now, when we represent \tilde{x} as a sum of sinusoids we have

$$\tilde{x}(n) = \sum_{k=0}^{M-1} X_k e^{2\pi i k n / M}, \quad (3.17)$$

for some choice of numbers X_k . Arguing as before, we find that now we have

$$X_k = \frac{1}{M} \sum_{n=0}^{M-1} x(n) e^{-2\pi i k n / M}, \quad (3.18)$$

for $k = 0, \dots, M-1$.

3.10.4 A Fourth Way: Bandlimited Extrapolation

Suppose that $f(t)$ is Γ -bandlimited, so that

$$f(t) = \frac{\Delta}{2\pi} \int_{-\Gamma}^{\Gamma} X(\gamma\Delta) e^{it\gamma} d\gamma. \quad (3.19)$$

Inserting $X(\gamma\Delta)$ as in Equation (3.6) into Equation (3.19) and performing the indicated integration, we obtain

$$f(t) = \Delta \sum_{n=-\infty}^{\infty} f(n\Delta) \frac{\sin \Gamma(t - n\Delta)}{\pi(t - n\Delta)}. \quad (3.20)$$

This formula illustrates Shannon's sampling theorem, by showing how to reconstruct the Γ -bandlimited function $f(t)$ from the infinite sequence of samples $\{f(n\Delta)\}$, for any $\Delta < \frac{\pi}{\Gamma}$. We shall use this formula to extend our finite data to obtain a Γ -bandlimited function that is consistent with the finite data. It is not required that the data be equispaced.

Arbitrarily Spaced Data

Now suppose that our data are the values $f(t_m)$, $m = 1, \dots, N$, where the t_m are arbitrary. From Equation (3.20) we have

$$f(t_m) = \Delta \sum_{n=-\infty}^{\infty} f(n\Delta) \frac{\sin \Gamma(t_m - n\Delta)}{\pi(t_m - n\Delta)}, \quad (3.21)$$

for each t_m . In this case, however, we do not know the $f(n\Delta)$. Can we find a sequence $\{f(n\Delta)\}$ for which Equation (3.21) is satisfied for each m ? The answer is yes; in fact, there are infinitely many ways to do this, as we shall see shortly. But, first, we need a useful identity concerning Γ -bandlimited functions.

A Useful Identity

The function $G(\gamma) = \chi_\Gamma(\gamma)$ that is one for $|\gamma| \leq \Gamma$ and is zero otherwise is the Fourier transform of the function $g(x) = \frac{\sin \Gamma x}{\pi x}$. Therefore, its sequence of Fourier coefficients is $\{\Delta g(n\Delta) = \Delta \frac{\sin \Gamma n\Delta}{\pi n\Delta}\}$. For any fixed t , the function $H_t(\gamma) = G(\gamma)e^{i\gamma t}$ has, for its sequence of Fourier coefficients $h_t = \{\Delta \frac{\sin \Gamma(n\Delta - t)}{\pi(n\Delta - t)}\}$. Since $H_t(\gamma)H_{-s}(\gamma) = H_{t-s}(\gamma)$, we have $h_t * h_{-s} = h_{t-s}$. Writing this out, we get

$$\frac{\sin \Gamma(n\Delta - t + s)}{\pi(n\Delta - t + s)} = \Delta \sum_{k=-\infty}^{\infty} \frac{\sin \Gamma(k\Delta - t)}{\pi(k\Delta - t)} \frac{\sin \Gamma((n-k)\Delta + s)}{\pi((n-k)\Delta + s)}. \quad (3.22)$$

Minimum-Norm Extrapolation

One possibility is to provide a finite-parameter model for the desired sequence $\{f(n\Delta)\}$, as

$$f(n\Delta) = \sum_{j=1}^N z_j \frac{\sin \Gamma(t_j - n\Delta)}{\pi(t_j - n\Delta)}. \quad (3.23)$$

Inserting this $f(n\Delta)$ into Equation (3.21), reversing the order of summation, and using the identity in Equation (3.22), we obtain

$$f(t_m) = \Delta \sum_{j=1}^N z_j \frac{\sin \Gamma(t_j - t_m)}{\pi(t_j - t_m)}. \quad (3.24)$$

This system of N equations in N unknowns can be solved uniquely for the z_j . Placing these z_j into Equation (3.23) to get the $f(n\Delta)$ and then using these $f(n\Delta)$ in Equation (3.20), we obtain a Γ -bandlimited function $\hat{f}(t)$ that extrapolates the finite data. The function $\hat{f}(t)$ can be written explicitly as

$$\hat{f}(t) = \Delta \sum_{j=1}^N z_j \frac{\sin \Gamma(t_j - t)}{\pi(t_j - t)}. \quad (3.25)$$

It can be shown that this choice of $\hat{f}(t)$ is the Γ -bandlimited function extrapolating the data for which the energy $\sum_{n=-\infty}^{\infty} |\hat{f}(n\Delta)|^2$ is the smallest.

Estimating the Fourier Transform

We take the Fourier transform of $\hat{f}(t)$ in Equation (3.25), to obtain an explicit formula for $\hat{F}(\gamma)$, our estimate of the Fourier transform of $f(t)$:

$$\hat{F}(\gamma) = \Delta \chi_{\Gamma}(\gamma) \sum_{j=1}^J z_j e^{it_j \gamma}.$$

When $t_j = j\Delta$, with $\Delta = \frac{\pi}{\Gamma}$, we find that $\Delta z_j = f(j\Delta)$, so that our estimate of $F(\gamma)$ becomes

$$\hat{F}(\gamma) = \sum_{j=1}^J f(j\Delta) e^{ij\Delta\gamma}.$$

So our estimate of $X(\omega)$ is

$$\hat{X}(\omega) = \hat{F}\left(-\frac{\omega}{\Delta}\right) = \sum_{j=1}^J f(j\Delta) e^{-ij\omega},$$

which is the DFT we get when we zero-extend the finite data.

Note that if $f(t)$ is known to be Γ -bandlimited, then $f(t)$ is $(\Gamma + \epsilon)$ -bandlimited, for any $\epsilon > 0$. Therefore, we can use $\Gamma + \epsilon$ in place of Γ , in the calculations above, to achieve a bandlimited extrapolation of the finite data. So there are infinitely many different ways to extend the finite data as samples of a bandlimited function. Each of these ways leads to a different estimate for the Fourier transform.

3.11 Is this Analysis or Representation?

As we just saw, we can represent the finite data $x(n)$, $n = 0, \dots, N - 1$, in any number of different ways as sums of discrete exponential functions. In the first way we have

$$x(n) = \frac{1}{2\pi} \int_{-\pi}^{\pi} X(\omega) e^{in\omega} d\omega, \quad (3.26)$$

in the second way

$$x(n) = \sum_{k=0}^{N-1} X_k e^{2\pi i kn/N}, \quad (3.27)$$

and in yet a third way

$$x(n) = \sum_{k=0}^{M-1} X_k e^{2\pi i kn/M}. \quad (3.28)$$

Using the bandlimited extrapolation approach, we can also write

$$x(n) = \frac{1}{2\pi} \int_{-\Gamma}^{\Gamma} \hat{F}(\gamma) e^{-in\Delta\gamma} d\gamma. \quad (3.29)$$

In each of these cases it would appear that the data contains certain sinusoidal components, and yet in each of these ways the sinusoidal frequencies involved are different. How can this be?

By *analysis* we mean the identification of the components of the data, in this case, the complex-exponential components or complex sinusoids, that are really there in the data. When we have at least two different ways to represent the data as a sum of such complex exponentials, can either of these be said to provide true analysis of the data? Equation (3.26) seems to say that the data is made up of complex exponentials whose frequencies encompass the entire interval $[-\pi, \pi)$, while Equation (3.27) exhibits the same data as consisting only of N complex exponentials, with frequencies equispaced through the interval $[-\pi, \pi)$, and Equation (3.28) employs a whole new set of M frequencies, equispaced through the interval $[-\pi, \pi)$. Equation (3.29) says the frequencies are spread over the interval $[-\Gamma, \Gamma]$. Which one is correct? This is not really the right question to ask. The proper response depends on the context, that is, on what the problem is that we are trying to solve.

If the goal is to perform some operation on the data, it may not matter greatly how it is represented. However, as we saw in our discussion of farfield propagation, the data can be finitely many samples of an underlying continuous-variable function $f(t)$ or a discrete function x , for which the frequency-space representation has real physical significance. In the discrete case, the DFT of x can have physical significance beyond simply providing a way to represent the x as a sum of exponential functions. For example, in sonar and radar array processing, the arguments ω may correspond to a direction of a distant object of interest, and ω may take on any value in $[-\pi, \pi)$. In such cases we would like to have all of x , but must settle for the finite data vector \mathbf{x} . The goal then is to use the finite data to approximate or estimate $X(\omega)$, the DFT of x . The DFT of the data is then a finite Fourier-series approximation of the infinite Fourier series that is $X(\omega)$. The vector DFT \mathbf{X} of the data gives us N equispaced values of this approximation, which can be calculated efficiently using the FFT.

There is an added twist to the story, however. Given only the data, we have no way of knowing the complete x ; there are infinitely many x

that extend the data. Which one is the correct one? In most applications we have some prior information about the nature of the function $X(\omega)$ that we seek to estimate from the data. Effective estimation procedures make use of this additional information to obtain better estimates when the data, by itself, is insufficient. Our fourth way to extend the finite data includes, in the extrapolation process, the prior knowledge that $f(t)$ is Γ -bandlimited. Later, we shall consider other ways to employ prior knowledge to extrapolate the data.

3.12 Oversampling

In many applications, we are essentially free to take as many samples as we wish, but are required to take those samples from within some finite region. In the model of farfield propagation, for example, there may be physical limitations on length of our array of sensors, but within that length, we may place as many sensors as seems reasonable. In synthetic-aperture radar, the array of sensors is moving, simulating a longer array, the length of which is limited, in practice, by the need to correct for time differences in the receipt of the signals. In sampling a function of time, the signal being sampled may only last for a short while, but while it lasts, we may take as many samples as we wish; this is the case in seismic exploration, magnetic resonance imaging, and speech processing. In our discussion previously, we saw that if the function $f(t)$ is Γ -bandlimited, then we must sample at a spacing $\Delta \leq \frac{\pi}{\Gamma}$. If we are required to take all our samples from within the time interval $[0, T]$, and if we use $\Delta = \frac{\pi}{\Gamma}$, we may not be able to take a large number of samples. Would it be better, under these circumstances, to *oversample*, that is, to use, say $\frac{\Delta}{2}$, in order to generate more data? Is there any limit on how small the spacing should be?

Suppose we begin with the samples $f(n\Delta)$, for $n = 0, 1, \dots, N - 1$, $\Delta = \frac{\pi}{\Gamma}$, and $T = N\Delta$. The DFT of the zero-extended data,

$$\hat{F}(\gamma) = \Delta \sum_{n=0}^{N-1} f(n\Delta) e^{in\Delta\gamma},$$

for $|\gamma| \leq \Gamma$, is then a reasonable estimate of the Fourier transform, $F(\gamma)$. Now let us take samples at spacing $\frac{\Delta}{2}$; that is, we take $f(\frac{m\Delta}{2})$, for $m = 0, \dots, 2N - 1$. The DFT of the zero-extension of this data is

$$\tilde{F}(\gamma) = \frac{\Delta}{2} \sum_{m=0}^{2N-1} f(m\frac{\Delta}{2}) e^{im\frac{\Delta}{2}\gamma}.$$

But now the interval outside of which the sum repeats itself is no longer $[-\Gamma, \Gamma]$, but $[-2\Gamma, 2\Gamma]$; $\tilde{F}(\gamma)$ is an estimate of $F(\gamma)$ for γ in this larger

interval. If we consider $\tilde{F}(\gamma)$ only for γ within the smaller interval $[-\Gamma, \Gamma]$, we find that $\tilde{F}(\gamma)$ is not much different from $\hat{F}(\gamma)$ for those values of γ . What has happened is that, when we chose to sample faster, the DFT estimation “believes” that our function $f(t)$ is 2Γ -bandlimited, which is true, but not precise. We do get twice as many data points, but we then are forced to use them to estimate the Fourier transform over an interval that is twice as wide as before.

There is a way out of this predicament, however. The bandlimited extrapolation method discussed earlier permits us to use any finite set of samples, $t_j, j = 1, \dots, J$. Therefore, we can take $t_j = (j - 1)\frac{\Delta}{2}, j = 1, \dots, J = 2N$. Then our estimate of $F(\gamma)$ has the form

$$\hat{F}(\gamma) = \Delta \chi_{\Gamma}(\gamma) \sum_{m=0}^{2N-1} z_{m+1} e^{im\frac{\Delta}{2}\gamma},$$

but, unlike for $\tilde{F}(\gamma)$, the z_{m+1} are not $\frac{1}{\Delta}f(m\frac{\Delta}{2})$.

Simulation experiments show that this method of estimating the Fourier transform from oversampled data does lead to improved estimates, but becomes increasingly sensitive to noise in the data, as the sample spacing gets smaller. The signal-to-noise ratio in the data provides the ultimate limitation on how small we can make the sample spacing.

3.13 Finite Data and the Fast Fourier Transform

Given the finite measurements x_1, \dots, x_N , we chose to write these as samples of a function $x(t)$, so that $x_n = x(n-1)$, for $n = 1, \dots, N$. We then analyzed the vector $\mathbf{x} = (x(0), \dots, x(N-1))^T$ in an attempt to uncover interesting components of the function $x(t)$. One approach involved estimating the Fourier transform $X(\omega)$ of $x(t)$ by means of the DFT,

$$\hat{X}(\omega) = \sum_{n=0}^{N-1} x(n-1) e^{-in\omega},$$

for $|\omega| \leq \pi$. As we noted previously, the Fast Fourier Transform algorithm can be used to calculate finitely many equi-spaced values of $\hat{X}(\omega)$.

There is another way to view the problem. Our data consists of the vector \mathbf{x} and we choose to write \mathbf{x} as a linear combination of other vectors, in the hope of discovering information that lies within the data. There are infinitely many ways to do this, however.

One way is to select N arbitrary distinct frequencies $\omega_m, m = 0, 1, \dots, N-1$ in $[-\pi, \pi)$ and define the vectors \mathbf{e}_{ω_m} by

$$\mathbf{e}_{\omega_m}(n) = e^{in\omega_m},$$

for $n = 0, \dots, N - 1$. We then write

$$\mathbf{x} = \sum_{m=0}^{N-1} a_m \mathbf{e}_{\omega_m},$$

where the coefficients a_m are found by solving the system of linear equations

$$x(n) = \sum_{m=0}^{N-1} a_m \mathbf{e}_{\omega_m}(n),$$

$n = 0, \dots, N - 1$.

We write the system of linear equations in matrix form as

$$\mathbf{x} = E\mathbf{a}, \tag{3.30}$$

for $\mathbf{a} = (a_0, \dots, a_{N-1})^T$ and E the N by N matrix with the entries

$$E_{nm} = e^{in\omega_m}.$$

Such a system will have a unique solution, and we will always be able to write the data vector as a finite sum of N arbitrarily chosen sinusoidal vectors \mathbf{e}_{ω_m} .

In general, the matrix E is invertible, but solving the system in Equation (3.30) when N is large can be computationally expensive. Since we are choosing the frequencies ω_m arbitrarily, why not select them so that the matrix E is easily inverted. This is motivation for the vector DFT of the data.

We now select the frequencies ω_m more carefully. We take

$$\omega_m = -\pi + \frac{2\pi}{N}m,$$

for $m = 0, \dots, N - 1$.

Exercise 3.5 Show that, for this choice of the ω_m , the inverse of the matrix E is

$$E^{-1} = \frac{1}{N}E^\dagger.$$

Using the result of this exercise, we find that the coefficient vector \mathbf{a} has entries

$$a_m = \frac{1}{N} \sum_{n=0}^{N-1} x(n) e^{-2\pi imn/N},$$

for $m = 0, \dots, N - 1$. These are the entries of the vector DFT, \mathbf{X} , as given in Equation (3.16). These a_m are what the FFT calculates.

When we consider the problem from this viewpoint, we see that the representation of the data vector \mathbf{x} as a superposition of sinusoidal vectors involves a completely arbitrary selection of the frequencies ω_m to be used, and yet, once the a_m are found, the data vector is completely described as that superposition. The equi-spaced frequencies used in the previous paragraph were chosen merely to facilitate the inversion of E . What does it mean to say that the data actually contains the components with frequencies ω_m , when we are free to select whichever ones we wish? What does it mean to say that the function $x(t)$ that was sampled to get the data actually contains sinusoids at these frequencies?

Chapter 4

Randomness in Signal Processing

We begin our discussion of randomness in signal processing with the example of farfield propagation.

4.1 Randomness in Farfield Propagation

In our earlier discussion of farfield propagation in the one-dimensional case, we imagined that each point $(x, 0, 0)$ on the x-axis transmitting or reflecting a sinusoidal signal $f(x)e^{i\omega t}$, where the complex number $f(x)$ denotes the magnitude and phase of the signal associated with this point. Our goal was to determine $f(x)$, in order to learn something about the extent of the object and what it is made of. The data collected in the farfield turned out to be, essentially, values of the Fourier transform of $f(x)$. The goal then became reconstructing a function from finitely many noisy values of its Fourier transform. Except for the additive noise, this model is completely deterministic and, as such, a somewhat unrealistic model for many situations.

We suppose now that, for each x , the complex number $f(x) = |f(x)|e^{i\theta(x)}$ is a random variable, with both its magnitude, $|f(x)|$, and its phase, $\theta(x)$ real-valued random variables. The randomness is often introduced to account for perturbations caused by the medium through which the signals must travel, such as light passing through the atmosphere, or sound through the ocean with changing sound speed.

One simple way to model such a complex random variable is to write $f(x) = a(x) + ib(x)$, with both $a(x)$ and $b(x)$ real-valued random variables. If $a(x)$ and $b(x)$ are independent Gaussian random variables with the same variance, then $f(x)$ is called a *complex Gaussian* random variable. In that

case, its magnitude has the Rayleigh distribution and its phase is uniformly distributed on the interval $[-\pi, \pi]$. For some applications we assume that the random variables $f(x)$ corresponding to different x are independent; in other cases, the $f(x)$ for neighboring values of x may be assumed correlated. Our goal is no longer to determine a particular value of $f(x)$, for each x , but to estimate the mean value of $|f(x)|$, which is then the average intensity of the signal from the point x .

Suppose, for example, that the object function $f(x)$ consists of finitely many point sources, that is,

$$f(x) = \sum_{j=1}^J f(x_j) \delta(x - x_j),$$

where the amplitudes $f(x_j)$ are independent complex random variables. We assume that the random phase of each $f(x_j)$ is distributed uniformly over the interval $[-\pi, \pi]$, so that the expected value of $f(x_j)$ is zero. Then the variance of $f(x_j)$ is $E(|f(x_j)|^2)$. The Fourier transform of $f(x_j) \delta(x - x_j)$ is

$$F_j(\gamma) = f(x_j) e^{i\gamma(x-x_j)},$$

so the signals received at the points $(s, D, 0)$ are values of

$$F(\gamma) = \sum_{j=1}^J f(x_j) e^{i\gamma(x-x_j)}.$$

Correlating $F(\gamma)$ and $F(\alpha)$, we obtain

$$E(F(\gamma) \overline{F(\alpha)}) = \sum_{j=1}^J E(|f(x_j)|^2) e^{i(\gamma-\alpha)x_j}.$$

Notice that the correlation is a function, not of γ and α separately, but of their difference. This prompts the following definition of the *autocorrelation function* of $f(x)$:

$$R(\tau) = E(F(\gamma) \overline{F(\gamma - \tau)}) = \sum_{j=1}^J E(|f(x_j)|^2) e^{i\tau x_j}.$$

The inverse Fourier transform of the function $R(\tau)$ is then

$$r(x) = \sum_{j=1}^J E(|f(x_j)|^2) \delta(x - x_j).$$

This suggests that we can estimate the values $E(|f(x_j)|^2)$ by using our measured data to obtain estimates of the cross-correlation function $R(\tau)$. We shall return to this issue later.

4.2 Random Variables as Models

When we use mathematical tools, such as differential equations, probability, or systems of linear equations, to describe a real-world situation, we say that we are employing a *mathematical model*. Such models must be sufficiently sophisticated to capture the essential features of the situation, while remaining computationally manageable. In this chapter we are interested in one particular type of mathematical model, the *random variable*.

Imagine that you are holding a baseball four feet off the ground. If you drop it, it will land on the ground directly below where you held it. The height of the ball at any time during the fall is described by the function $h(t)$ satisfying the ordinary differential equation $h''(t) = -32\frac{\text{ft}}{\text{sec}^2}$. Solving this differential equation with the initial conditions $h(0) = 4 \text{ ft}$, $h'(0) = 0\frac{\text{ft}}{\text{sec}}$, we find that $h(t) = 4 - 16t^2$. Solving $h(T) = 0$ for T we find the elapsed time T until impact is $T = 0.5 \text{ sec.}$ The velocity of the ball at impact is $h'(T) = -32T = -16\frac{\text{ft}}{\text{sec}}$.

Now imagine that, instead of a baseball, you are holding a feather. The feather and the baseball are both subject to the same laws of gravity, but now other aspects of the situation, which we could safely ignore in the case of the baseball, become important in the case of the feather. Like the baseball, the feather is subjected to air resistance and to whatever fluctuations in air currents may be present during its fall. Unlike the baseball, however, the effects of the air matter to the flight of the feather; in fact, they become the dominant factors. When we designed our differential-equation model for the falling baseball we performed no experiments to help us understand its behavior. We simply ignored all other aspects of the situation, and included only gravity in our mathematical model. Even the modeling of gravity was slightly simplified, in that we assumed a constant gravitational acceleration, even though Newton's Laws tell us that it increases as we approach the center of the earth. When we drop the ball and find that our model is accurate we feel no need to change it. When we drop the feather we discover immediately that a new model is needed; but what?

The first thing we observe is that the feather falls in a manner that is impossible to predict with accuracy. Dropping it once again, we notice that it behaves differently this time, landing in a different place and, perhaps, taking longer to land. How are we to model such a situation, in which repeated experiments produce different results? Can we say nothing useful about what will happen when we drop the feather the next time?

As we continue to drop the feather, we notice that, while the feather usually does not fall directly beneath the point of release, it does not fall too far away. Suppose we draw a grid of horizontal and vertical lines on the ground, dividing the ground into a pattern of squares of equal area. Now we repeatedly drop the feather and record the proportion of times the feather is (mainly) contained within each square; we also record the

elapsed time. As we are about to drop the feather the next time, we may well assume that the outcome will be consistent with the behavior we have observed during the previous drops. While we cannot say for certain where the feather will fall, nor what the elapsed time will be, we feel comfortable making a *probabilistic statement* about the likelihood that the feather will land in any given square and about the elapsed time.

The squares into which the feather may land are finite, or, if we insist on creating an infinite grid, discretely infinite, while the elapsed time can be any positive real number. Let us number the squares as $n = 1, 2, 3, \dots$ and let p_n be the proportion of drops that resulted in the feather landing mainly in square n . Then $p_n \geq 0$ and $\sum_{n=1}^{\infty} p_n = 1$. The sequence $p = \{p_n | n = 1, 2, \dots\}$ is then a *discrete probability sequence* (dps), or a *probability sequence*, or a *discrete probability*. Now let N be the number of the square that will contain the feather on the next drop. All we can say about N is that, according to our model, the probability that N will equal n is p_n . We call N a *discrete random variable* with *probability sequence* p .

It is difficult to be more precise about what probability really means. When we say that the probability is p_n that the feather will land in square n on the next drop, where does that probability reside? Do we believe that the numbers p_n are *in the feather* somehow? Do these numbers simply describe our own ignorance, so are *in our heads*? Are they a combination of the two, in our heads as a result of our having experienced what the feather did previously? Perhaps it is best simply to view probability as a type of mathematical model that we choose to adopt in certain situations.

Now let T be the elapsed time for the next feather to hit the ground. What can we say about T ? Based on our prior experience, we are willing to say that, for any interval $[a, b]$ within $(0, \infty)$, the probability that T will take on a value within $[a, b]$ is the proportion of prior drops in which the elapsed time was between a and b . Then T is a *continuous random variable*, in that the values it may take on (in theory, at least) lie in a continuum. To help us calculate the probabilities associated with T we use our prior experience to specify a function $f_T(t)$, called the *probability density function* (pdf) of T , having the property that the probability that T will lie between a and b can be calculated as $\int_a^b f_T(t) dt$. Such $f_T(t)$ will have the properties $f_T(t) \geq 0$ for all positive t and $\int_0^{\infty} f_T(t) dt = 1$.

In the case of the falling feather we had to perform experiments to determine appropriate ps p and pdf $f_T(t)$. In practice, we often describe our random variables using a ps or pdf from a well-studied parametric family of such mathematical objects. Popular examples of such ps and pdf, such as Poisson probabilities and Gaussian pdf, are discussed early in most courses in probability theory.

It is simplest to discuss the main points of random signal processing within the context of discrete signals, so we return there now.

4.3 Discrete Random Signal Processing

Previously, we have encountered specific discrete functions, such as δ_k , u , e_ω , whose values at each integer n are given by an exact formula. In signal processing we must also concern ourselves with discrete functions whose values are not given by such formulas, but rather, seem to obey only probabilistic laws. We shall need such discrete functions to model noise. For example, imagine that, at each time n , a fair coin is tossed and $x(n) = 1$ if the coin shows heads, $x(n) = -1$ if the coin shows tails. We cannot determine the value of $x(n)$ from any formula; we must simply toss the coins. Given any discrete function x with values $x(n)$ that are either 1 or -1 , we cannot say if x was generated by such a coin-flipping manner. In fact, any such x could have been the result of coin flips. All we can say is how likely it is that a particular x was so generated. For example, if $x(n) = 1$ for n even and $x(n) = -1$ for n odd, we feel, intuitively, that it is highly unlikely that such an x came from random coin tossing. What bothers us, of course, is that the values $x(n)$ seem so predictable; randomness seems to require some degree of unpredictability. If we were given two such sequences, the first being the one described above, with 1 and -1 alternating, and the second exhibiting no obvious pattern, and asked to select the one generated by independent random coin tossing, we would clearly choose the second one. There is a subtle point here, however. When we say that we are “given an infinite sequence” what do we really mean? Because the issue here is not the infinite nature of the sequences, let us reformulate the discussion in terms of finite vectors of length, say, 100, with entries 1 or -1 . If we are shown a print-out of two such vectors, the first with alternating 1 and -1 , and the second vector exhibiting no obvious pattern, we would immediately say that it was the second one that was generated by the coin-flipping procedure, even though the two vectors are equally likely to have been so generated. The point is that we associate randomness with the absence of a pattern, more than with probability. When there is a pattern, the vector can be described in a way that is significantly shorter than simply listing its entries. Indeed, it has been suggested that a vector is random if it cannot be described in a manner shorter than simply listing its members.

4.3.1 The Simplest Random Sequence

We say that a sequence $x = \{x(n)\}$ is a *random sequence* or a *discrete random process* if $x(n)$ is a random variable for each integer n . A simple, yet remarkably useful, example is the random-coin-flip sequence, which we shall denote by $c = \{c(n)\}$. In this model a coin is flipped for each n and $c(n) = 1$ if the coin comes up heads, with $c(n) = -1$ if the coin comes up tails. It will be convenient to allow for the coin to be *biased*, that is,

for the probabilities of heads and tails to be unequal. We denote by p the probability that heads occurs and $1 - p$ the probability of tails; the coin is called *unbiased* or *fair* if $p = 1/2$. To find the *expected value* of $c(n)$, written $E(c(n))$, we multiply each possible value of $c(n)$ by its probability and sum; that is,

$$E(c(n)) = (+1)p + (-1)(1 - p) = 2p - 1.$$

If the coin is fair then $E(c(n)) = 0$. The variance of the random variable $c(n)$, measuring its tendency to deviate from its expected value, is $\text{var}(c(n)) = E([c(n) - E(c(n))]^2)$. We have

$$\text{var}(c(n)) = [+1 - (2p - 1)]^2 p + [-1 - (2p - 1)]^2 (1 - p) = 4p - 4p^2.$$

If the coin is fair then $\text{var}(c(n)) = 1$. It is important to note that we do not change the coin at any time during the generation of the random sequence c ; in particular, the p does not depend on n .

The random-coin-flip sequence c is the simplest example of a discrete random process or a random discrete function. It is important to remember that a random discrete function is not any one particular discrete function, but rather a probabilistic model chosen to allow us to talk about the probabilities associated with the values of the $x(n)$. In the next section we shall use this discrete random process to generate a wide class of discrete random processes, obtained by viewing $c = c(n)$ as the input into a linear, shift-invariant (LSI) filter.

4.4 Random Discrete Functions or Discrete Random Processes

A linear, shift-invariant (LSI) operator T with impulse response function $h = \{h(k)\}$ operates on any input sequence $x = \{x(n)\}$ to produce the output sequence $y = \{y(n)\}$ according to the convolution formula

$$y(n) = \sum_{k=-\infty}^{\infty} h(k)x(n-k) = \sum_{k=-\infty}^{\infty} x(k)h(n-k). \quad (4.1)$$

We learn more about the system that T represents when we select as input sinusoids at fixed frequencies. Let ω be a fixed frequency in the interval $[-\pi, \pi)$ and let $x = e_\omega$, so that $x(n) = e^{in\omega}$ for each integer n . Then Equation (4.1) shows us that the output is

$$y(n) = H(e^{i\omega})x(n),$$

where

$$H(e^{i\omega}) = \sum_{k=-\infty}^{\infty} h(k)e^{-ik\omega}. \quad (4.2)$$

This function of ω is called the *frequency-response function* of the system. We can learn even more about the system by selecting as input the sequence $x(n) = z^n$, where z is an arbitrary complex number. Then Equation (4.1) gives the output as

$$y(n) = H(z)x(n),$$

where

$$H(z) = \sum_{k=-\infty}^{\infty} h(k)z^{-k}. \quad (4.3)$$

Note that if we select $z = e^{i\omega}$ then $H(z) = H(e^{i\omega})$ as given by Equation (4.2). The function $H(z)$ of the complex variable z is the z -transform of the sequence h and also the *transfer function* of the system determined by h .

Now we take this approach one step further. Let us select as our input $x = \{x(n)\}$ the random-coin-flip sequence $c = \{c(n)\}$, with $p = 0.5$. It is important to note that such an x is not one specific discrete function, but a random model for such functions. The output $y = \{y(n)\}$ is again a random sequence, with

$$y(n) = \sum_{k=-\infty}^{\infty} h(k)c(n-k). \quad (4.4)$$

Clearly, in order for the infinite sum to converge we would need to place restrictions on the sequence h ; if $h(k)$ is zero except for finitely many values of k then we have no problem. We shall put off discussion of convergence issues and focus on statistical properties of the output random sequence y .

Let u and v be (possibly complex-valued) random variables with expected values $E(u)$ and $E(v)$, respectively. The covariance between u and v is defined to be

$$\text{cov}(u, v) = E([u - E(u)]\overline{[v - E(v)]}),$$

and the cross-correlation between u and v is

$$\text{corr}(u, v) = E(u\bar{v}).$$

It is easily shown that $\text{cov}(u, v) = \text{corr}(u, v) - E(u)\overline{E(v)}$. When $u = v$ we get $\text{cov}(u, u) = \text{var}(u)$ and $\text{corr}(u, u) = E(|u|^2)$. If $E(u) = E(v) = 0$ then $\text{cov}(u, v) = \text{corr}(u, v)$.

To illustrate, let $u = c(n)$ and $v = c(n-m)$. Then, since the coin is fair, $E(c(n)) = E(c(n-m)) = 0$ and

$$\text{cov}(c(n), c(n-m)) = \text{corr}(c(n), c(n-m)) = E(c(n)\overline{c(n-m)}).$$

Because the $c(n)$ are independent, $E(c(n)\overline{c(n-m)}) = 0$ for m not equal to 0, and $E(|c(n)|^2) = \text{var}(c(n)) = 1$. Therefore

$$\text{cov}(c(n), c(n-m)) = \text{corr}(c(n), c(n-m)) = 0, \text{ for } m \neq 0,$$

and

$$\text{cov}(c(n), c(n)) = \text{corr}(c(n), c(n)) = 1.$$

Returning now to the output sequence $y = \{y(n)\}$ we compute the correlation $\text{corr}(y(n), y(n-m)) = E(y(n)\overline{y(n-m)})$. Using the convolution formula Equation (4.4) we find that

$$\text{corr}(y(n), y(n-m)) = \sum_{k=-\infty}^{\infty} \sum_{j=-\infty}^{\infty} h(k)\overline{h(j)}\text{corr}(c(n-k), c(n-m-j)).$$

Since

$$\text{corr}(c(n-k), c(n-m-j)) = 0, \text{ for } k \neq m+j,$$

we have

$$\text{corr}(y(n), y(n-m)) = \sum_{k=-\infty}^{\infty} h(k)\overline{h(k-m)}. \quad (4.5)$$

The expression of the right side of Equation (4.5) is the definition of the *autocorrelation* of the sequence h , denoted $\rho_h(m)$; that is,

$$\rho_h(m) = \sum_{k=-\infty}^{\infty} h(k)\overline{h(k-m)}. \quad (4.6)$$

It is important to note that the expected value of $y(n)$ is

$$E(y(n)) = \sum_{k=-\infty}^{\infty} h(k)E(c(n-k)) = 0$$

and the correlation $\text{corr}(y(n), y(n-m))$ depends only on m ; neither quantity depends on n and the sequence y is therefore called *weak-sense stationary*. Let's consider an example.

Take $h(0) = h(1) = 0.5$ and $h(k) = 0$ otherwise. Then the system is the two-point moving-average, with

$$y(n) = 0.5x(n) + 0.5x(n-1).$$

With $x(n) = c(n)$ we have

$$y(n) = 0.5c(n) + 0.5c(n-1).$$

In the case of the random-coin-flip sequence c each $c(n)$ is unrelated to any other $c(m)$; the coin flips are independent. This is no longer the case for the $y(n)$; one effect of the filter h is to introduce correlation into the output. To illustrate, since $y(0)$ and $y(1)$ both depend, to some degree, on the value $c(0)$, they are related. Using Equation (4.6) we have

$$\rho_h(0) = h(0)h(0) + h(1)h(1) = 0.25 + 0.25 = 0.5,$$

$$\rho_h(-1) = h(0)h(1) = 0.25,$$

$$\rho_h(+1) = h(1)h(0) = 0.25,$$

and

$$\rho_h(m) = 0, \text{ otherwise.}$$

So we see that $y(n)$ and $y(n - m)$ are related, for $m = -1, 0, +1$, but not otherwise.

4.5 Correlation Functions and Power Spectra

As we have seen, any nonrandom sequence $h = \{h(k)\}$ has its autocorrelation function defined, for each integer m , by

$$\rho_h(m) = \sum_{k=-\infty}^{\infty} h(k)\overline{h(k-m)}.$$

For a random sequence $y(n)$ that is wide-sense stationary, its correlation function is defined to be

$$\rho_y(m) = E(y(n)\overline{y(n-m)}).$$

The *power spectrum* of h is defined for ω in $[-\pi, \pi]$ by

$$S_h(\omega) = \sum_{m=-\infty}^{\infty} \rho_h(m)e^{-im\omega}.$$

It is easy to see that

$$S_h(\omega) = |H(e^{i\omega})|^2,$$

so that $S_h(\omega) \geq 0$. The power spectrum of the random sequence $y = \{y(n)\}$ is defined as

$$S_y(\omega) = \sum_{m=-\infty}^{\infty} \rho_y(m)e^{-im\omega}.$$

Although it is not immediately obvious, we also have $S_y(\omega) \geq 0$. One way to see this is to consider

$$Y(e^{i\omega}) = \sum_{n=-\infty}^{\infty} y(n)e^{-in\omega}$$

and to calculate

$$E(|Y(e^{i\omega})|^2) = \sum_{m=-\infty}^{\infty} E(y(n)\overline{y(n-m)})e^{-im\omega} = S_y(\omega).$$

Given any power spectrum $S_y(\omega)$ we can construct $H(e^{i\omega})$ by selecting an arbitrary phase angle θ and letting

$$H(e^{i\omega}) = \sqrt{S_y(\omega)}e^{i\theta}.$$

We then obtain the nonrandom sequence h associated with $H(e^{i\omega})$ using

$$h(n) = \int_{-\pi}^{\pi} H(e^{i\omega})e^{in\omega}d\omega/2\pi.$$

It follows that $\rho_h(m) = \rho_y(m)$ for each m and $S_h(\omega) = S_y(\omega)$ for each ω .

What we have discovered is that, when the input to the system is the random-coin-flip sequence c , the output sequence y has a correlation function $\rho_y(m)$ that is equal to the autocorrelation of the sequence h . As we just saw, for any weak-sense stationary random sequence y with expected value $E(y(n))$ constant and correlation function $\text{corr}(y(n), y(n-m))$ independent of n , there is a LSI system h with $\rho_h(m) = \rho_y(m)$ for each m . Therefore, any weak-sense stationary random sequence y can be viewed as the output of an LSI system, when the input is the random-coin-flip sequence $c = \{c(n)\}$.

4.6 Random Sinusoidal Sequences

If $A = |A|e^{i\theta}$, with amplitude $|A|$ a positive-valued random variable and phase angle θ a random variable taking values in the interval $[-\pi, \pi]$ then A is a complex-valued random variable. For a fixed frequency ω_0 we define a random sinusoidal sequence $s = \{s(n)\}$ by $s(n) = Ae^{in\omega_0}$. We assume that θ has the uniform distribution over $[-\pi, \pi]$ so that the expected value of $s(n)$ is zero. The correlation function for s is

$$\rho_s(m) = E(s(n)\overline{s(n-m)}) = E(|A|^2)e^{im\omega_0}$$

and the power spectrum of s is

$$S_s(\omega) = E(|A|^2) \sum_{m=-\infty}^{\infty} e^{im(\omega_0-\omega)},$$

so that, by Equation (3.7), we have

$$S_s(\omega) = E(|A|^2)\delta(\omega - \omega_0).$$

We generalize this example to the case of multiple independent sinusoids. Suppose that, for $j = 1, \dots, J$, we have fixed frequencies ω_j and independent complex-valued random variables A_j . We let our random sequence be defined by

$$s(n) = \sum_{j=1}^J A_j e^{in\omega_j}.$$

Then the correlation function for x is

$$\rho_s(m) = \sum_{j=1}^J E(|A_j|^2) e^{im\omega_j}$$

and the power spectrum for s is

$$S_s(\omega) = \sum_{j=1}^J E(|A_j|^2) \delta(\omega - \omega_j).$$

A commonly used model in signal processing is that of independent sinusoids in additive noise.

Let $q = \{q(n)\}$ be an arbitrary weak-sense stationary discrete random sequence, with correlation function $\rho_q(m)$ and power spectrum $S_q(\omega)$. We say that q is white noise if $\rho_q(m) = 0$ for m not equal to zero, or, equivalently, if the power spectrum $S_q(\omega)$ is constant over the interval $[-\pi, \pi]$. The *independent sinusoids in additive noise* model is a random sequence of the form

$$x(n) = \sum_{j=1}^J A_j e^{in\omega_j} + q(n).$$

The *signal power* is defined to be $\rho_s(0)$, which is the sum of the $E(|A_j|^2)$, while the noise power is $\rho_q(0)$. The *signal-to-noise ratio* (SNR) is the ratio of signal power to noise power.

It is often the case that the SNR is quite low and it is desirable to process the x to enhance this ratio. The data we have is typically finitely many values of $x(n)$, say for $n = 1, 2, \dots, N$. One way to process the data is to estimate $\rho_x(m)$ for some small number of integers m around zero, using, for example, the *lag products* estimate

$$\hat{\rho}_x(m) = \frac{1}{N-m} \sum_{n=1}^{N-m} x(n) \overline{x(n-m)},$$

for $m = 0, 1, \dots, M < N$ and $\hat{\rho}_x(-m) = \overline{\hat{\rho}_x(m)}$. Because $\rho_q(m) = 0$ for m not equal to zero, we will have $\hat{\rho}_x(m)$ approximating $\rho_s(m)$ for nonzero values of m , thereby reducing the effect of the noise.

The additive noise is said to be *correlated* or *non-white* if it is not the case that $\rho_x(m) = 0$ for all nonzero m . In this case the noise power spectrum is not constant, and so may be concentrated in certain regions of the interval $[-\pi, \pi]$.

4.7 Spread-Spectrum Communication

In this section we return to the random-coin-flip model, this time allowing the coin to be biased, that is, p need not be 0.5. Let $s = \{s(n)\}$ be a random sequence, such as $s(n) = Ae^{in\omega_0}$, with $E(s(n)) = \mu$ and correlation function $\rho_s(m)$. Define a second random sequence x by

$$x(n) = s(n)c(n).$$

The random sequence x is generated from the random signal s by randomly changing its signs. We can show that

$$E(x(n)) = \mu(2p - 1)$$

and, for m not equal to zero,

$$\rho_x(m) = \rho_s(m)(2p - 1)^2,$$

with $\rho_x(0) = \rho_s(0) + 4p(1 - p)\mu^2$. Therefore, if $p = 1$ or $p = 0$ we get $\rho_x(m) = \rho_s(m)$ for all m , but for $p = 0.5$ we get $\rho_x(m) = 0$ for m not equal to zero. If the coin is unbiased, then the random sign changes convert the original signal s into white noise. Generally, we have

$$S_x(\omega) = (2p - 1)^2 S_s(\omega) + (1 - (2p - 1)^2)(\mu^2 + \rho_s(0)),$$

which says that the power spectrum of x is a combination of the signal power spectrum and a white-noise power spectrum, approaching the white-noise power spectrum as p approaches 0.5. If the original signal power spectrum is concentrated within a small interval, then the effect of the random sign changes is to spread that spectrum. Once we know what the sequence c is we can recapture the original signal from $s(n) = x(n)c(n)$. The use of such a spread spectrum permits the sending of multiple narrow-band signals, without confusion, as well as protecting against any narrow-band additive interference.

4.8 Stochastic Difference Equations

The ordinary first-order differential equation $y'(t) + ay(t) = f(t)$, with initial condition $y(0) = 0$ has for its solution $y(t) = e^{-at} \int_0^t e^{as} f(s) ds$.

One way to look at such differential equations is to consider $f(t)$ to be the input to a system having $y(t)$ as its output. The system determines which terms will occur on the left side of the differential equation. In many applications the input $f(t)$ is viewed as random noise and the output is then a continuous-time random process. Here we want to consider the discrete analog of such differential equations.

We replace the first derivative with the first difference, $y(n+1) - y(n)$ and we replace the input with the random-coin-flip sequence $c = \{c(n)\}$, to obtain the random difference equation

$$y(n+1) - y(n) + ay(n) = c(n). \quad (4.7)$$

With $b = 1 - a$ and $0 < b < 1$ we have

$$y(n+1) - by(n) = c(n). \quad (4.8)$$

The solution is $y = \{y(n)\}$ given by

$$y(n) = b^n \sum_{k=-\infty}^n b^{-k} c(k). \quad (4.9)$$

Comparing this with the solution of the differential equation, we see that the term b^n plays the role of $e^{-at} = (e^{-a})^t$, so that $b = 1 - a$ is substituting for e^{-a} . The infinite sum replaces the infinite integral, with $b^{-k}c(k)$ replacing the integrand $e^{as}f(s)$.

The solution sequence y given by Equation (4.9) is a weak-sense stationary random sequence and its correlation function is

$$\rho_y(m) = b^{|m|}/(1 - b^2).$$

Since

$$b^n \sum_{k=-\infty}^n b^{-k} = 1 - b$$

the random sequence $(1 - b)^{-1}y(n)$ is an infinite *moving-average* random sequence formed from the random sequence c .

We can derive the solution in Equation (4.9) using z-transforms. The expression $y(n) - by(n-1)$ can be viewed as the output of a LSI system with $h(0) = 1$ and $h(1) = -b$. Then $H(z) = 1 - bz^{-1} = (z - b)/z$ and the inverse $H(z)^{-1} = z/(z - b)$ describes the inverse system. Since

$$H(z)^{-1} = z/(z - b) = 1/(1 - bz^{-1}) = 1 + bz^{-1} + b^2z^{-2} + \dots$$

the inverse system applied to input $c = \{c(n)\}$ is

$$y(n) = 1c(n) + bc(n-1) + b^2c(n-2) + \dots = b^n \sum_{k=-\infty}^n b^{-k} c(k).$$

Bibliography

- [1] Bracewell, R.C. (1979) Image Reconstruction in Radio Astronomy, in [54], pp. 81–104.
- [2] Browne, J. and A. DePierro, A. (1996) “A row-action alternative to the EM algorithm for maximizing likelihoods in emission tomography.” *IEEE Trans. Med. Imag.* **15**, pp. 687–699.
- [3] Byrne, C., and Fitzgerald, R. (1982) Reconstruction from partial information with applications to tomography, *SIAM J. Appl. Math.*, **42(4)**, 933–940.
- [4] Byrne, C., Fitzgerald, R., Fiddy, M., Hall, T., and Darling, A. (1983) Image restoration and resolution enhancement, *J. Optical Soc. America*, **73**, 1481–1487.
- [5] Byrne, C., Levine, B.M., and Dainty, J.C. (1984) “Stable estimation of the probability density function of intensity from photon frequency counts.” *JOSA Communications* **1(11)**, pp. 1132–1135.
- [6] Byrne, C., and Fiddy, M. (1988) Images as power spectra; reconstruction as Wiener filter approximation, *Inverse Problems*, **4**, 399–409.
- [7] Byrne, C. (1993) “Iterative image reconstruction algorithms based on cross-entropy minimization.” *IEEE Transactions on Image Processing* **IP-2**, pp. 96–103.
- [8] Byrne, C. (1995) “Erratum and addendum to ‘Iterative image reconstruction algorithms based on cross-entropy minimization’.” *IEEE Transactions on Image Processing* **IP-4**, pp. 225–226.
- [9] Byrne, C. (1996) “Iterative reconstruction algorithms based on cross-entropy minimization.” in *Image Models (and their Speech Model Cousins)*, S.E. Levinson and L. Shepp, editors, IMA Volumes in Mathematics and its Applications, Volume 80, pp. 1–11. New York: Springer-Verlag.

- [10] Byrne, C. (1996) “Block-iterative methods for image reconstruction from projections.” *IEEE Transactions on Image Processing* **IP-5**, pp. 792–794.
- [11] Byrne, C. (1997) “Convergent block-iterative algorithms for image reconstruction from inconsistent data.” *IEEE Transactions on Image Processing* **IP-6**, pp. 1296–1304.
- [12] Byrne, C. (1998) “Accelerating the EMLL algorithm and related iterative algorithms by rescaled block-iterative (RBI) methods.” *IEEE Transactions on Image Processing* **IP-7**, pp. 100–109.
- [13] Byrne, C. (1999) “Iterative projection onto convex sets using multiple Bregman distances.” *Inverse Problems* **15**, pp. 1295–1313.
- [14] Byrne, C. (2000) “Block-iterative interior point optimization methods for image reconstruction from limited data.” *Inverse Problems* **16**, pp. 1405–1419.
- [15] Byrne, C. (2001) “Bregman-Legendre multidistance projection algorithms for convex feasibility and optimization.” in *Inherently Parallel Algorithms in Feasibility and Optimization and their Applications*, Butnariu, D., Censor, Y., and Reich, S., editors, pp. 87–100. Amsterdam: Elsevier Publ.,
- [16] Byrne, C. (2001) “Likelihood maximization for list-mode emission tomographic image reconstruction.” *IEEE Transactions on Medical Imaging* **20(10)**, pp. 1084–1092.
- [17] Byrne, C. (2002) “Iterative oblique projection onto convex sets and the split feasibility problem.” *Inverse Problems* **18**, pp. 441–453.
- [18] Byrne, C. (2004) “A unified treatment of some iterative algorithms in signal processing and image reconstruction.” *Inverse Problems* **20**, pp. 103–120.
- [19] Byrne, C. (2005) Choosing parameters in block-iterative or ordered-subset reconstruction algorithms, *IEEE Transactions on Image Processing*, **14 (3)**, pp. 321–327.
- [20] Byrne, C. (2005) “Signal Processing: A Mathematical Approach” , AK Peters, Publ., Wellesley, MA.
- [21] Byrne, C. and Censor, Y. (2001) “Proximity function minimization using multiple Bregman projections, with applications to split feasibility and Kullback-Leibler distance minimization.” *Annals of Operations Research* **105**, pp. 77–98.

- [22] Byrne, C. and Fiddy, M. (1987) "Estimation of continuous object distributions from Fourier magnitude measurements." *JOSA A* **4**, pp. 412–417.
- [23] Byrne, C. and Fitzgerald, R. (1979) "A unifying model for spectrum estimation." in *Proceedings of the RADC Workshop on Spectrum Estimation- October 1979*, Griffiss AFB, Rome, NY.
- [24] Byrne, C. and Fitzgerald, R. (1984) "Spectral estimators that extend the maximum entropy and maximum likelihood methods." *SIAM J. Applied Math.* **44(2)**, pp. 425–442.
- [25] Censor, Y. (1981) "Row-action methods for huge and sparse systems and their applications." *SIAM Review*, **23**: 444–464.
- [26] Censor, Y., Eggermont, P.P.B., and Gordon, D. (1983) "Strong underrelaxation in Kaczmarz's method for inconsistent systems." *Numerische Mathematik* **41**, pp. 83–92.
- [27] Censor, Y. and Segman, J. (1987) "On block-iterative maximization." *J. of Information and Optimization Sciences* **8**, pp. 275–291.
- [28] Censor, Y. and Zenios, S.A. (1997) *Parallel Optimization: Theory, Algorithms and Applications*. New York: Oxford University Press.
- [29] Chang, J.-H., Anderson, J.M.M., and Votaw, J.R. (2004) "Regularized image reconstruction algorithms for positron emission tomography." *IEEE Transactions on Medical Imaging* **23(9)**, pp. 1165–1175.
- [30] Cimmino, G. (1938) "Calcolo approssimato per soluzioni die sistemi di equazioni lineari." *La Ricerca Scientifica XVI, Series II, Anno IX* **1**, pp. 326–333.
- [31] Combettes, P. (1993) "The foundations of set theoretic estimation." *Proceedings of the IEEE* **81 (2)**, pp. 182–208.
- [32] Cooley, J. and Tukey, J. (1965) "An algorithm for the machine calculation of complex Fourier series." *Math. Comp.*, **19**, pp. 297–301.
- [33] Csiszár, I. and Tusnády, G. (1984) "Information geometry and alternating minimization procedures." *Statistics and Decisions* **Supp. 1**, pp. 205–237.
- [34] Csiszár, I. (1989) "A geometric interpretation of Darroch and Ratcliff's generalized iterative scaling." *The Annals of Statistics* **17 (3)**, pp. 1409–1413.

- [35] Csiszár, I. (1991) “Why least squares and maximum entropy? An axiomatic approach to inference for linear inverse problems.” *The Annals of Statistics* **19** (4), pp. 2032–2066.
- [36] Dainty, J. C. and Fiddy, M. (1984) “The essential role of prior knowledge in phase retrieval.” *Optica Acta* **31**, pp. 325–330.
- [37] Darroch, J. and Ratcliff, D. (1972) “Generalized iterative scaling for log-linear models.” *Annals of Mathematical Statistics* **43**, pp. 1470–1480.
- [38] Dempster, A.P., Laird, N.M. and Rubin, D.B. (1977) “Maximum likelihood from incomplete data via the EM algorithm.” *Journal of the Royal Statistical Society, Series B* **37**, pp. 1–38.
- [39] De Pierro, A. (1995) “A modified expectation maximization algorithm for penalized likelihood estimation in emission tomography.” *IEEE Transactions on Medical Imaging* **14**, pp. 132–137.
- [40] De Pierro, A. and Iusem, A. (1990) “On the asymptotic behavior of some alternate smoothing series expansion iterative methods.” *Linear Algebra and its Applications* **130**, pp. 3–24.
- [41] Dhanantwari, A., Stergiopoulos, S., and Iakovidis, I. (2001) “Correcting organ motion artifacts in x-ray CT medical imaging systems by adaptive processing. I. Theory.” *Med. Phys.* **28**(8), pp. 1562–1576.
- [42] Duda, R., Hart, P., and Stork, D. (2001) *Pattern Classification*, Wiley.
- [43] Eggermont, P.P.B., Herman, G.T., and Lent, A. (1981) “Iterative algorithms for large partitioned linear systems, with applications to image reconstruction.” *Linear Algebra and its Applications* **40**, pp. 37–67.
- [44] Fessler, J., Ficaro, E., Clinthorne, N., and Lange, K. (1997) Grouped-coordinate ascent algorithms for penalized-likelihood transmission image reconstruction, *IEEE Transactions on Medical Imaging*, **16** (2), pp. 166–175.
- [45] Fiddy, M. (1983) “The phase retrieval problem.” in *Inverse Optics*, SPIE Proceedings 413 (A.J. Devaney, editor), pp. 176–181.
- [46] Fienup, J. (1979) “Space object imaging through the turbulent atmosphere.” *Optical Engineering* **18**, pp. 529–534.
- [47] Fienup, J. (1987) “Reconstruction of a complex-valued object from the modulus of its Fourier transform using a support constraint.” *Journal of the Optical Society of America A* **4**(1), pp. 118–123.

- [48] Frieden, B. R. (1982) *Probability, Statistical Optics and Data Testing*. Berlin: Springer-Verlag.
- [49] Geman, S., and Geman, D. (1984) "Stochastic relaxation, Gibbs distributions and the Bayesian restoration of images." *IEEE Transactions on Pattern Analysis and Machine Intelligence* **PAMI-6**, pp. 721–741.
- [50] Gifford, H., King, M., de Vries, D., and Soares, E. (2000) "Channelized Hotelling and human observer correlation for lesion detection in hepatic SPECT imaging" *Journal of Nuclear Medicine* **41(3)**, pp. 514–521.
- [51] Gordon, R., Bender, R., and Herman, G.T. (1970) "Algebraic reconstruction techniques (ART) for three-dimensional electron microscopy and x-ray photography." *J. Theoret. Biol.* **29**, pp. 471–481.
- [52] Green, P. (1990) "Bayesian reconstructions from emission tomography data using a modified EM algorithm." *IEEE Transactions on Medical Imaging* **9**, pp. 84–93.
- [53] Hebert, T. and Leahy, R. (1989) "A generalized EM algorithm for 3-D Bayesian reconstruction from Poisson data using Gibbs priors." *IEEE Transactions on Medical Imaging* **8**, pp. 194–202.
- [54] Herman, G.T. (ed.) (1979) "Image Reconstruction from Projections", *Topics in Applied Physics, Vol. 32*, Springer-Verlag, Berlin.
- [55] Herman, G.T., and Natterer, F. (eds.) "Mathematical Aspects of Computerized Tomography" , *Lecture Notes in Medical Informatics, Vol. 8*, Springer-Verlag, Berlin.
- [56] Herman, G.T., Censor, Y., Gordon, D., and Lewitt, R. (1985) Comment (on the paper [87]), *Journal of the American Statistical Association* **80**, pp. 22–25.
- [57] Herman, G. T. and Meyer, L. (1993) "Algebraic reconstruction techniques can be made computationally efficient." *IEEE Transactions on Medical Imaging* **12**, pp. 600–609.
- [58] Holte, S., Schmidlin, P., Linden, A., Rosenqvist, G. and Eriksson, L. (1990) "Iterative image reconstruction for positron emission tomography: a study of convergence and quantitation problems." *IEEE Transactions on Nuclear Science* **37**, pp. 629–635.
- [59] Hudson, H.M. and Larkin, R.S. (1994) "Accelerated image reconstruction using ordered subsets of projection data." *IEEE Transactions on Medical Imaging* **13**, pp. 601–609.

- [60] Huesman, R., Klein, G., Moses, W., Qi, J., Ruetter, B., and Vi-rador, P. (2000) “List-mode maximum likelihood reconstruction ap-plied to positron emission mammography (PEM) with irregular sam-pling.” *IEEE Transactions on Medical Imaging* **19** (5), pp. 532–537.
- [61] Hutton, B., Kyme, A., Lau, Y., Skerrett, D., and Fulton, R. (2002) “A hybrid 3-D reconstruction/registration algorithm for correction of head motion in emission tomography.” *IEEE Transactions on Nuclear Science* **49** (1), pp. 188–194.
- [62] Kaczmarz, S. (1937) “Angenäherte Auflösung von Systemen linearer Gleichungen.” *Bulletin de l’Academie Polonaise des Sciences et Lettres* **A35**, pp. 355–357.
- [63] Kak, A., and Slaney, M. (2001) “Principles of Computerized Tomo-graphic Imaging”, SIAM, Philadelphia, PA.
- [64] King, M., Glick, S., Pretorius, H., Wells, G., Gifford, H., Narayanan, M., and Farncombe, T. (2004) Attenuation, Scatter, and Spatial Res-olution Compensation in SPECT, in [88], pp. 473–498.
- [65] Kullback, S. and Leibler, R. (1951) “On information and suffi-ciency.” *Annals of Mathematical Statistics* **22**, pp. 79–86.
- [66] Landweber, L. (1951) “An iterative formula for Fredholm integral equations of the first kind.” *Amer. J. of Math.* **73**, pp. 615–624.
- [67] Lange, K. and Carson, R. (1984) “EM reconstruction algorithms for emission and transmission tomography.” *Journal of Computer Assisted Tomography* **8**, pp. 306–316.
- [68] Lange, K., Bahn, M. and Little, R. (1987) “A theoretical study of some maximum likelihood algorithms for emission and transmission tomography.” *IEEE Trans. Med. Imag.* **MI-6(2)**, pp. 106–114.
- [69] Leahy, R., Hebert, T., and Lee, R. (1989) “Applications of Markov random field models in medical imaging.” in *Proceedings of the Confer-ence on Information Processing in Medical Imaging* Lawrence-Berkeley Laboratory, Berkeley, CA.
- [70] Leahy, R. and Byrne, C. (2000) “Guest editorial: Recent development in iterative image reconstruction for PET and SPECT.” *IEEE Trans. Med. Imag.* **19**, pp. 257–260.
- [71] Levitan, E. and Herman, G. (1987) “A maximum *a posteriori* proba-bility expectation maximization algorithm for image reconstruction in emission tomography.” *IEEE Transactions on Medical Imaging* **6**, pp. 185–192.

- [72] Liao, C.-W., Fiddy, M., and Byrne, C. (1997) "Imaging from the zero locations of far-field intensity data." *Journal of the Optical Society of America -A* **14** (12), pp. 3155–3161.
- [73] McLachlan, G.J. and Krishnan, T. (1997) *The EM Algorithm and Extensions*. New York: John Wiley and Sons, Inc.
- [74] Meidunas, E. (2001) *Re-scaled Block Iterative Expectation Maximization Maximum Likelihood (RBI-EMML) Abundance Estimation and Sub-pixel Material Identification in Hyperspectral Imagery*, MS thesis, Department of Electrical Engineering, University of Massachusetts Lowell.
- [75] Narayanan, M., Byrne, C. and King, M. (2001) "An interior point iterative maximum-likelihood reconstruction algorithm incorporating upper and lower bounds with application to SPECT transmission imaging." *IEEE Transactions on Medical Imaging* **TMI-20** (4), pp. 342–353.
- [76] Natterer, F. (1986) *Mathematics of Computed Tomography*. New York: John Wiley and Sons, Inc.
- [77] Natterer, F., and Wübbeling, F. (2001) *Mathematical Methods in Image Reconstruction*. Philadelphia, PA: SIAM Publ.
- [78] Parra, L. and Barrett, H. (1998) "List-mode likelihood: EM algorithm and image quality estimation demonstrated on 2-D PET." *IEEE Transactions on Medical Imaging* **17**, pp. 228–235.
- [79] Peters, T. (1981) Resolution improvement to CT systems using aperture-function correction, in [55], pp. 241–251.
- [80] P. Pretorius, M. King, T-S Pan, D. deVries, S. Glick, and C. Byrne (1998) Reducing the influence of the partial volume effect on SPECT activity quantitation with 3D modelling of spatial resolution in iterative reconstruction, *Phys.Med. Biol.* **43**, pp. 407–420.
- [81] Rockmore, A., and Macovski, A. (1976) A maximum likelihood approach to emission image reconstruction from projections, *IEEE Transactions on Nuclear Science*, **NS-23**, pp. 1428–1432.
- [82] Schmidlin, P. (1972) "Iterative separation of sections in tomographic scintigrams." *Nucl. Med.* **15**(1).
- [83] Shepp, L., and Vardi, Y. (1982) Maximum likelihood reconstruction for emission tomography, *IEEE Transactions on Medical Imaging*, **MI-1**, pp. 113–122.

- [84] Soares, E., Byrne, C., Glick, S., Appledorn, R., and King, M. (1993) Implementation and evaluation of an analytic solution to the photon attenuation and nonstationary resolution reconstruction problem in SPECT, *IEEE Transactions on Nuclear Science*, **40** (4), pp. 1231–1237.
- [85] Tanabe, K. (1971) “Projection method for solving a singular system of linear equations and its applications.” *Numer. Math.* **17**, pp. 203–214.
- [86] Twomey, S. (1996) *Introduction to the Mathematics of Inversion in Remote Sensing and Indirect Measurement*. New York: Dover Publ.
- [87] Vardi, Y., Shepp, L.A. and Kaufman, L. (1985) “A statistical model for positron emission tomography.” *Journal of the American Statistical Association* **80**, pp. 8–20.
- [88] Wernick, M. and Aarsvold, J., editors (2004) *Emission Tomography: The Fundamentals of PET and SPECT*. San Diego: Elsevier Academic Press.
- [89] Wright, G.A. (1997) “Magnetic Resonance Imaging” *IEEE Signal Processing Magazine*, **14** (1), pp. 56–66.

Index

- z -transform, 45
- array aperture, 11
- attenuated Radon transform, 30
- autocorrelation function, 64
- backprojection, 34
- best linear unbiased estimator, 86
- BLUE, 86
- broadband signal, 3
- central slice theorem, 33
- channelized Hotelling observer, 90
- classification, 85
- complex amplitude, 106
- complex exponential function, 105
- complex sinusoid, 105
- convolution, 9
- convolution filter, 9
- cross-entropy, 80
- detection, 85
- DFT, 87
- Dirac delta, 8
- discrete Fourier transform, 46, 87
- discrete-time Fourier transform, 46
- discrimination, 85
- emission tomography, 27
- estimation, 85
- expectation maximization (EM) algorithm, 79
- exponential Radon transform, 30
- farfield assumption, 4
- filtered backprojection, 34
- Fisher linear discriminant, 93
- fixed point, 82
- Fourier coefficients, 46
- Fourier Inversion Formula, 18
- Fourier inversion formula, 7
- Fourier transform, 3, 6
- Fourier-series expansion, 46
- Fourier-transform pair, 7
- frequency, 105
- frequency-domain extrapolation, 10
- frequency-response function, 9
- Hilbert transform, 36
- Hotelling linear discriminant, 90
- Hotelling observer, 90
- identification, 85
- interior-point methods, 113
- Kuhn-Karush-Tucker condition, 81
- Kullback-Leibler distance, 80
- line of response, 28
- linear sensor array, 11
- maximum a posteriori, 83
- modulation transfer function, 9
- narrowband signal, 3, 4
- Nyquist spacing, 15
- optical transfer function, 9
- oversampled data, 58
- partial volume effect, 31
- PET, 27

phase problem, 7
point-spread function, 9
Poisson, 32
positron emission tomography, 27

Radon transform, 33

sampling, 14
scatter, 30
Shannon Sampling Theorem, 15
sifting property, 8
single photon emission tomography, 27
sinusoids, 105
SPECT, 27
spill-over, 31
steepest descent algorithm, 113
synthetic-aperture radar, 12
system transfer function, 9

uniform line array, 14, 15