

Chapter 1

Mathematical Prerequisites

Digital holography is a discipline that associates the techniques of traditional optical holography with current computational methods [GOO 67, HUA 71, KRO 72, LYO 04, SCH 05]. In the framework of the scalar theory of diffraction [BOR 99, GOO 72, GOO 05], digital holography tackles, based on diffraction formulae, the propagation of a light wave in an optical system, the study of interference between coherent light waves, and the reconstruction of surface waves diffracted by objects of various natures. In this context, the propagation of a light wave can be considered as the transformation of a two-dimensional signal by a linear system—the optical system. Various representations of the scalar amplitude of a light wave carrying information use special mathematical functions; the transformation of a light wave across a linear system uses a fundamental mathematical tool: two-dimensional Fourier analysis. The digital treatment of optical information leads us to treat the problems of sampling and discretization, under the restriction given by Shannon's theorem. Thus, the mathematical prerequisites for a good understanding of this book concern the frequently used mathematical functions, the two-dimensional Fourier transform, and the notions of the sampling theorem [GOO 72].

1.1. Frequently used special functions

Many mathematical functions that we will present in this section are frequently used in this book. To understand their properties, we give a brief account of their physical meaning.

1.1.1. The “rectangle” function

The one-dimensional rectangle function is defined by:

$$\text{rect}(x) = \begin{cases} 1 & (|x| \leq 1/2) \\ 0 & \text{otherwise} \end{cases} \quad [1.1]$$

This function is represented in Figure 1.1.

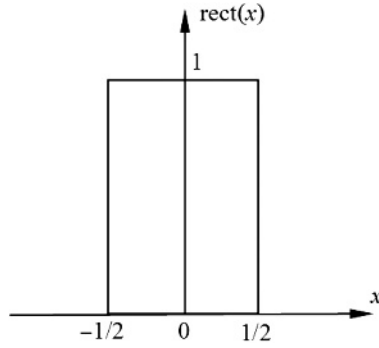


Figure 1.1. *Rectangle function*

Depending on the nature of the variable x , the rectangle function has various meanings. For example, if x is a spatial variable (a spatial coordinate in millimeters), we can use the function to represent the transmittance from a slit pierced in an opaque screen. In this book, we generally use the two-dimensional rectangle function that is obtained by the product of two one-dimensional functions. As an example, the following function is very useful:

$$f(x,y) = \text{rect}\left(\frac{x-x_0}{a}\right) \text{rect}\left(\frac{y-y_0}{b}\right) \quad [1.2]$$

This function is shown in Figure 1.2. It allows us to simply represent the transmittance from an aperture of a rectangular shape, centered on the point with coordinates (x_0, y_0) and of lengths a and b along the x - and y -axes, respectively. This binary function is very useful for considering the amplitude of an optical wave limited to a rectangular region, by eliminating the values outside the zone of interest.

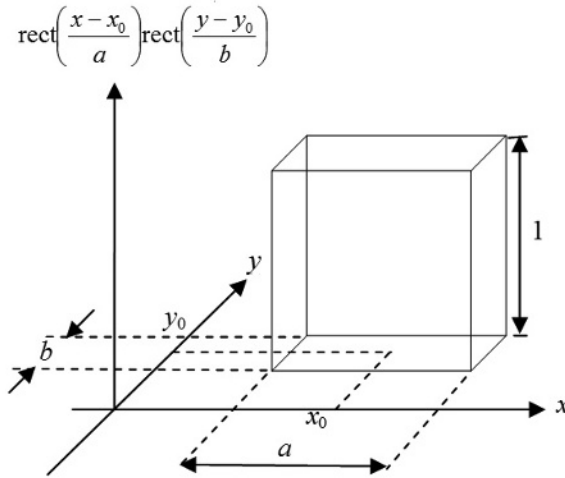


Figure 1.2. Two-dimensional rectangle function centered on (x_0, y_0)

1.1.2. The “sinc” function

The one-dimensional sinc function is defined by:

$$\text{sinc}(x) = \frac{\sin \pi x}{\pi x} \quad [1.3]$$

Its curve is presented in Figure 1.3.

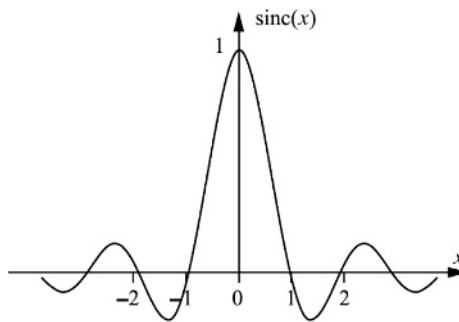


Figure 1.3. The sinc function

Also, the two-dimensional sinc function is formed by the product of two functions of independent variables:

$$\text{sinc}(x, y) = \text{sinc}(x)\text{sinc}(y) \quad [1.4]$$

Let us consider two positive values a and b ; Figure 1.4 shows the curve of the function $\text{sinc}^2(x/a, y/b)$. In Chapter 2, we will see that such a function represents the intensity distribution of Fraunhofer diffraction from a rectangular aperture illuminated by a coherent wave.

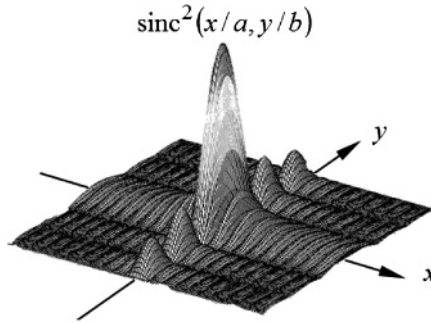


Figure 1.4. Two-dimensional sinc function

1.1.3. The “sign” function

The one-dimensional sign function is defined as:

$$\text{sgn}(x) = \begin{cases} 1 & x > 0 \\ 0 & x = 0 \\ -1 & x < 0 \end{cases} \quad [1.5]$$

The curve of this function is given in Figure 1.5.

If a function is multiplied by the function $\text{sgn}(x-a)$, for $a < 0$, the sign of the function will be inverted. If a coherent optical field is multiplied by this function, the resulting change corresponds to a phase shift of π . We can also form a two-dimensional sign function by taking the product of two one-dimensional functions.

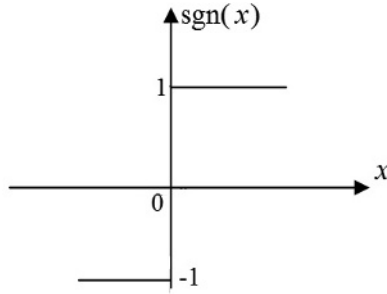


Figure 1.5. *The sign function*

1.1.4. The “triangle” function

The triangle function is defined as:

$$\Lambda(x) = \begin{cases} 1 - |x|, & |x| \leq 1 \\ 0, & \text{otherwise} \end{cases} \quad [1.6]$$

The curve of this function is given in Figure 1.6. Later, we will see that the Fourier transform of the function $\Lambda(x)$ is $\text{sinc}^2(f_x)$ (with the f_x coordinate corresponding to the *spatial* frequency). This function will be very useful in the Fourier analysis of optical diffracting functions (e.g. diffraction grating). As noted earlier, we can form a two-dimensional triangle function by taking the product of two one-dimensional functions.

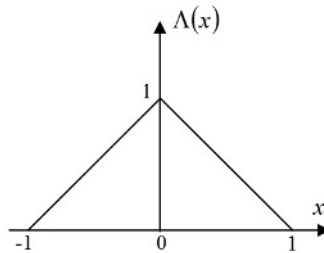


Figure 1.6. *Triangle function*

1.1.5. The “disk” function

In practice, an optical system is generally constructed with lenses whose mounts (cylinders) are circular in form. Their pupils are therefore circular and the disk

function is often used to model the diffraction of circular elements (iris diaphragms, mounts, etc.). The definition of this function, in polar and Cartesian coordinates, is:

$$\text{circ}(r) = \text{circ}\left(\sqrt{x^2 + y^2}\right) = \begin{cases} 1 & r = \sqrt{x^2 + y^2} \leq 1 \\ 0 & r = \sqrt{x^2 + y^2} > 1 \end{cases} \quad [1.7]$$

The surface of the disk function is given in Figure 1.7.

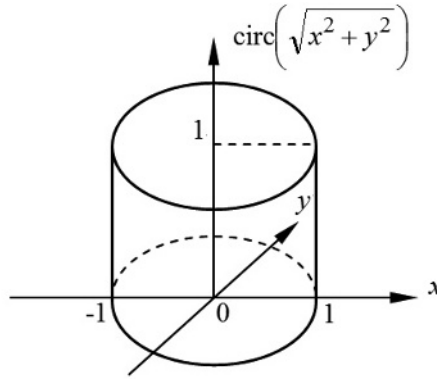


Figure 1.7. *The disk function*

1.1.6. The Dirac δ function

1.1.6.1. Definition

In the field of optical treatment of information, the Dirac δ distribution (henceforth called the δ “function”) in two dimensions is very widely used. Strictly speaking, δ is a distribution but for convenience we will hereafter call it a function. According to the Huygens–Fresnel principle of the propagation of light, a wave front can be considered as the sum of spherical “secondary” sources [BOR 99, GOO 72, GOO 05]. The two-dimensional δ function is often used to individually describe point sources. The fundamental property of the δ function is that, as for an infinitely narrow pulse of infinite height, the sum $\int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \delta(x, y) dx dy$ is equivalent to one (x and y being Cartesian coordinates). The δ function can be defined by various mathematical expressions, one of which is presented here.

Let us consider a series of the function $f_N(x) = N \text{ rect}(Nx)$ ($N = 1, 2, 3, \dots$). Figure 1.8 shows the curves corresponding to the number $N = 1, 2, 4$. It is evident

that the greater the value of N , the narrower the non-zero zone of the function. It is not difficult to imagine that if N tends to infinity, the value of the function $f_N(x) = N \text{rect}(Nx)$ will be infinite as well. On the other hand, the surface enclosed by the curve of the function and the x -axis stays unchanged, and equals one. Thus, by using the rectangular function, the one-dimensional δ function can also be defined as:

$$\delta(x) = \lim_{N \rightarrow \infty} N \text{rect}(Nx) \quad [1.8]$$

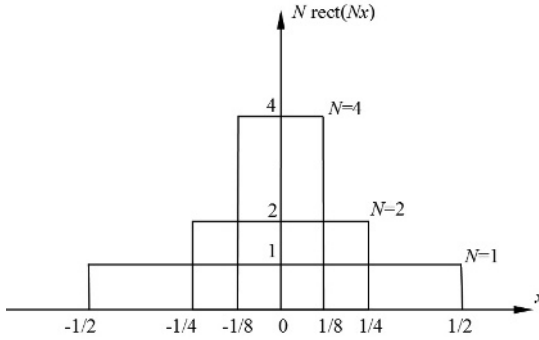


Figure 1.8. Graph of $f_N(x)$ for $N=1, 2, 4$

Evidently, we can also define the two-dimensional δ function as:

$$\delta(x, y) = \lim_{N \rightarrow \infty} N^2 \text{rect}(Nx) \text{rect}(Ny) \quad [1.9]$$

To facilitate the use of the δ function, we give some equivalent definitions:

$$\delta(x, y) = \lim_{N \rightarrow \infty} N^2 \exp\left[-N^2 \pi (x^2 + y^2)\right] \quad [1.10]$$

$$\delta(x, y) = \lim_{N \rightarrow \infty} N^2 \text{sinc}(Nx) \text{sinc}(Ny) \quad [1.11]$$

$$\delta(x, y) = \lim_{N \rightarrow \infty} \frac{N^2}{\pi} \text{circ}\left(N \sqrt{x^2 + y^2}\right) \quad [1.12]$$

$$\delta(x, y) = \lim_{N \rightarrow \infty} N \frac{J_1\left(2\pi N \sqrt{x^2 + y^2}\right)}{\sqrt{x^2 + y^2}} \quad [1.13]$$

In the last expression, J_1 is a first-order Bessel function of the first kind. Depending on the problem being studied, these definitions can be more or less appropriate and we can also choose which definition to apply in each case.

1.1.6.2. *Fundamental properties*

We will now consider some of the mathematical properties of the δ function. These properties will be used frequently in this book.

1.1.6.2.1. Contraction–dilation of coordinates

If a is any constant, we have:

$$\delta(ax) = \frac{1}{|a|} \delta(x) \quad [1.14]$$

1.1.6.2.2. Product

If the function $\varphi(x)$ is continuous at the point x_0 , we have:

$$\varphi(x) \delta(x - x_0) = \varphi(x_0) \delta(x - x_0) \quad [1.15]$$

1.1.6.2.3. Convolution

Let us consider the convolution of two functions δ and φ :

$$\delta(x) * \varphi(x) = \int_{-\infty}^{\infty} \delta(x_0) \varphi(x - x_0) dx_0 \quad [1.16]$$

Then we have:

$$\delta(x) * \varphi(x) = \varphi(x) * \delta(x) = \varphi(x) \quad [1.17]$$

The δ function is the unity of the convolution product.

1.1.6.2.4. Translation

The property of translation of the δ function is often used for theoretical analyses and proofs. Here we present this property and the corresponding proof. If $\varphi(x)$ is continuous at the point x_0 , then we have:

$$\int_{-\infty}^{\infty} \delta(x - x_0) \varphi(x) dx = \varphi(x_0) \quad [1.18]$$

PROOF.— Let $x - x_0 = x'$, on the left of the previous expression we can write:

$$\int_{-\infty}^{\infty} \delta(x) \varphi(x + x_0) dx = \int_{-\infty}^{-\varepsilon} \delta(x) \varphi(x + x_0) dx + \int_{-\varepsilon}^{+\varepsilon} \delta(x) \varphi(x + x_0) dx + \int_{+\varepsilon}^{\infty} \delta(x) \varphi(x + x_0) dx \quad [1.19]$$

If $\varepsilon \rightarrow 0$, the first and third terms on the right will be zero, therefore:

$$\begin{aligned} \int_{-\infty}^{\infty} \delta(x - x_0) \varphi(x) dx &= \lim_{\varepsilon \rightarrow 0} \int_{-\varepsilon}^{+\varepsilon} \delta(x) \varphi(x + x_0) dx \\ &= \varphi(x_0) \int_{-\varepsilon}^{+\varepsilon} \delta(x) dx = \varphi(x_0) \end{aligned} \quad [1.20]$$

In the same way, we can show that the two-dimensional δ function possesses the same property of translation.

$$\int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \delta(x - x_0, y - y_0) \varphi(x, y) dx dy = \varphi(x_0, y_0) \quad [1.21]$$

1.1.7. The “comb” function

The comb function is a periodic series of δ functions. It is frequently used to model the sample of continuous functions. The definition of the one-dimensional comb function is:

$$\text{comb}(x) = \sum_{n=-\infty}^{\infty} \delta(x - n) \quad (n = 1, 2, 3, \dots) \quad [1.22]$$

Figure 1.9 shows the curves of $\delta(x)$ and $\text{comb}(x)$. The two-dimensional comb function can be defined by the product of two one-dimensional comb functions:

$$\text{comb}(x, y) = \sum_{n=-\infty}^{\infty} \delta(x - n) \sum_{m=-\infty}^{\infty} \delta(y - m) \quad (n, m = 1, 2, 3, \dots) \quad [1.23]$$

Since the comb function is a periodic series of δ functions, it has analogous properties and is used in numerous analyses of optical signals.

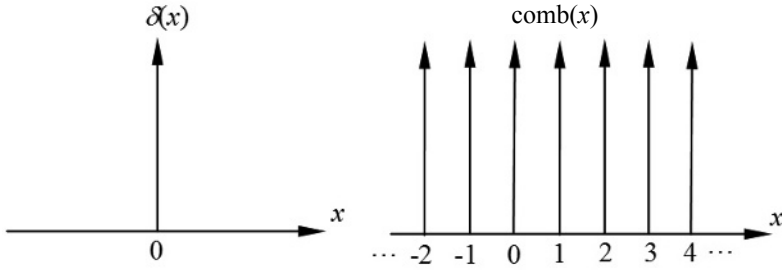


Figure 1.9. The $\delta(x)$ and $\text{comb}(x)$ functions

1.2. Two-dimensional Fourier transform

The Fourier transform is a very useful mathematical tool for the study of both linear and nonlinear phenomena. As the propagation of the optical field can be considered as a process of linear transformation of the “object” field to the “image” field, we are immediately interested in the two-dimensional Fourier transform [BOR 99, GOO 72].

1.2.1. Definition and existence conditions

The Fourier transform of a complex function $g(x, y)$ of two independent variables, which we write here as $F\{g(x, y)\}$, is defined as ($j = \sqrt{-1}$):

$$F\{g(x, y)\} = G(f_x, f_y) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} g(x, y) \exp[-j 2\pi(f_x x + f_y y)] dx dy \quad [1.24]$$

Thus defined, the transform is itself a complex-valued function of the two independent variables $G(f_x, f_y)$, called the spectral function, or spectrum, of the original function $g(x, y)$. The two variables f_x and f_y are considered, without loss of generality, as frequencies. In optics, (x, y) are spatial variables and (f_x, f_y) are *spatial frequencies* (mm^{-1}). Similarly, the inverse Fourier transform of the function $G(f_x, f_y)$, which we write as $F^{-1}\{G(f_x, f_y)\}$, is defined as:

$$F^{-1}\{G(f_x, f_y)\} = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} G(f_x, f_y) \exp[j 2\pi(f_x x + f_y y)] df_x df_y \quad [1.25]$$

We note that the direct and inverse transformations are completely analogous mathematical operations. They differ only by the sign of the exponent in the double integral. However, for some functions, these two integrals cannot exist in a mathematical sense. Therefore, we will briefly discuss the conditions of their existence. Among the various conditions, we concern ourselves with the following:

- $g(x, y)$ must be absolutely integrable in the xy -plane;
- $g(x, y)$ must have a finite number of discontinuities and a finite number of maxima and minima in any rectangle of finite area;
- $g(x, y)$ cannot have any infinite discontinuities.

In general, one of these three conditions can be ignored if we can guarantee strict adherence to the other conditions, but this is beyond the scope of discussions in this book.

For the representation of real physical waves by ideal mathematical functions, in the analysis of tools, one or more of the existing conditions presented above may be more or less unsatisfied [GOO 72]. However, as Bracelet [BRA 65] remarked, “the physical possibility is a sufficient condition of validity to justify the existence of a transformation”. Furthermore, the functions of interest to us are included in the scope of Fourier analysis, and it is evidently necessary to generalize definition [1.24] somewhat. Thus, it is possible to find a transformation that has meaning for functions that do not strictly satisfy the existing conditions, provided that these functions can be defined as the limit of a sequence of transformable functions. In transforming each term of this sequence, we generate a new sequence whose limit is called the generalized Fourier transform of the original function. These generalized transforms can be handled in the same way as the ordinary transforms, and the distinction between the two is often ignored. For a more detailed discussion of this generalization, the reader may refer to the work of Lighthill [LIG 60].

To simplify the study of Fourier analysis, including this generalization, Table 1.1 shows the Fourier transforms of some functions expressed in Cartesian coordinates.

1.2.2. Theorems related to the Fourier transform

We now present some important mathematical theorems followed by a brief account of their physical meaning [GOO 72]. The theorems mentioned below will be used frequently as they constitute fundamental tools for the use of Fourier transforms; they allow us to simplify the calculation of solutions to problems in Fourier analysis.

Original function	Fourier transform
$f(x, y) * g(x, y)$	$\tilde{f}(f_x, f_y) \times \tilde{g}(f_x, f_y)$
$f(x, y) \exp[2j\pi(xf_0 + yf_0)]$	$\tilde{f}(f_x - f_0, f_y - f_0)$
1	$\delta(f_x, f_y)$
$\delta(x, y)$	1
$\delta(x - x_0, y - y_0)$	$\exp[-j2\pi(f_x x_0 + f_y y_0)]$
$\text{rect}(x) \text{rect}(y)$	$\text{sinc}(f_x) \text{sinc}(f_y)$
$\Lambda(x) \Lambda(y)$	$\text{sinc}^2(f_x) \text{sinc}^2(f_y)$
$\text{sgn}(x) \text{sgn}(y)$	$\frac{1}{j\pi f_x} \times \frac{1}{j\pi f_y}$
$\exp[-\pi(x^2 + y^2)]$	$\exp[-\pi(f_x^2 + f_y^2)]$
$\exp[-j2\pi(ax + by)]$	$\delta(f_x - a, f_y - b)$
$\text{circ}(\sqrt{x^2 + y^2})$	$J_1(2\pi\sqrt{f_x^2 + f_y^2}) / \sqrt{f_x^2 + f_y^2}$
$\cos(2\pi f_0 x)$	$\frac{1}{2}[\delta(f_x - f_0) + \delta(f_x + f_0)]$
$\frac{1}{2}[\delta(x - x_0) + \delta(x + x_0)]$	$\cos(2\pi f_x x_0)$
$\sin(2\pi f_0 x)$	$\frac{1}{2j}[\delta(f_x - f_0) - \delta(f_x + f_0)]$
$\frac{j}{2}[\delta(x - x_0) - \delta(x + x_0)]$	$\sin(2\pi f_x x_0)$
$\text{comb}(x) \text{comb}(y)$	$\text{comb}(f_x) \text{comb}(f_y)$
$\exp[j\pi(a^2 x^2 + b^2 y^2)]$	$\frac{j}{ ab } \exp[-j\pi(f_x^2 / a^2 + f_y^2 / b^2)]$
$\exp[-j\pi(a^2 x^2 + b^2 y^2)]$	$-\frac{j}{ ab } \exp[j\pi(f_x^2 / a^2 + f_y^2 / b^2)]$
$\exp(-(a x + b y))$	$\frac{4ab}{(a^2 + 4\pi^2 f_x^2)(b^2 + 4\pi^2 f_y^2)}$

Table 1.1. Fourier transforms of some functions expressed in Cartesian coordinates

1.2.2.1. Linearity

The transform of the sum of two functions is simply the sum of their respective transforms:

$$F\{\alpha g(x, y) + \beta h(x, y)\} = \alpha F\{g(x, y)\} + \beta F\{h(x, y)\} \quad [1.26]$$

Where α and β are complex constants.

1.2.2.2. Similarity

If $F\{g(x, y)\} = G(f_x, f_y)$, and a and b are two real constants (different from 0), then:

$$F\{g(ax, by)\} = \frac{1}{|ab|} G\left(\frac{f_x}{a}, \frac{f_y}{b}\right) \quad [1.27]$$

This theorem is also known as the “contraction/dilation” theorem. It means that a “dilation” of the coordinates of the spatial domain (x, y) is expressed as a “contraction” of the coordinates in the frequency domain (f_x, f_y) and by a change in the amplitude and the width of the spectrum.

1.2.2.3. Translation

If $F\{g(x, y)\} = G(f_x, f_y)$, then:

$$F\{g(x-a, y-b)\} = G(f_x, f_y) \exp[-j2\pi(f_x a + f_y b)] \quad [1.28]$$

The translation of a function in the spatial domain introduces a linear phase variation in the frequency domain.

1.2.2.4. Parseval's theorem

If $F\{g(x, y)\} = G(f_x, f_y)$, then:

$$\int_{-\infty}^{\infty} \int_{-\infty}^{\infty} |g(x, y)|^2 dx dy = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} |G(f_x, f_y)|^2 df_x df_y \quad [1.29]$$

This theorem is generally interpreted as an expression of the conservation of energy between the spatial domain and the spatial frequency domain.

1.2.2.5. *The convolution theorem*

If $F\{g(x, y)\} = G(f_x, f_y)$ and $F\{h(x, y)\} = H(f_x, f_y)$, then:

$$F\left\{\int_{-\infty}^{\infty} \int_{-\infty}^{\infty} g(\xi, \eta) h(x - \xi, y - \eta) d\xi d\eta\right\} = G(f_x, f_y) H(f_x, f_y) \quad [1.30]$$

The Fourier transform of the convolution of two functions in the spatial domain is equivalent to the multiplication of their respective transformations. We will see in Chapter 3 that the Fourier transform can be calculated by the Fast Fourier Transform (FFT). This theorem offers the opportunity to calculate a convolution using FFT algorithms.

1.2.2.6. *The autocorrelation theorem*

If $F\{g(x, y)\} = G(f_x, f_y)$, then:

$$F\left\{\int_{-\infty}^{\infty} \int_{-\infty}^{\infty} g(\xi, \eta) g^*(\xi - x, \eta - y) d\xi d\eta\right\} = |G(f_x, f_y)|^2 \quad [1.31]$$

$$F\left\{|g(\xi, \eta)|^2\right\} = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} G(\xi, \eta) G^*(\xi + f_x, \eta + f_y) d\xi d\eta \quad [1.32]$$

This theorem can be considered as a particular case of the convolution theorem.

1.2.2.7. *The duality theorem*

Let us consider two functions f and g linked by the following integral development:

$$g(\alpha, \beta) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f(u, v) \exp[-j2\pi(u\alpha + v\beta)] du dv \quad [1.33]$$

We pose $(u, v) = (x, y)$ and $(\alpha, \beta) = (f_x, f_y)$, then:

$$F\{f(x, y)\} = g(f_x, f_y) \quad [1.34]$$

We now pose $(\alpha, \beta) = (x, y)$ and $(u, v) = (f_x, f_y)$, then:

$$g(x, y) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f(f_x, f_y) \exp[-j 2\pi(xf_x + yf_y)] du dv \quad [1.35]$$

being equally:

$$g(x, y) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f(-f_x, -f_y) \exp[+j 2\pi(xf_x + yf_y)] du dv \quad [1.36]$$

giving

$$g(x, y) = F^{-1} \{ f(-f_x, -f_y) \}, \quad [1.37]$$

and by applying the Fourier transform operator to the left and right,

$$F \{ g(x, y) \} = f(-f_x, -f_y). \quad [1.38]$$

Hence the property of duality of the Fourier transforms:

if

$$F \{ f(x, y) \} = g(f_x, f_y) \quad [1.39]$$

then

$$F \{ g(x, y) \} = f(-f_x, -f_y) \quad [1.40]$$

This property is very useful for determining Fourier transforms as it means that a pair of functions where one is the transform of the other generate a second pair of functions where one is the transform of the other. For example, if we consider the function $\text{rect}(x)\text{rect}(y)$, whose Fourier transform is $\sin c(f_x)\sin c(f_y)$ (see Table 1.1), then we can easily deduce that the Fourier transform of $\sin c(x)\sin c(y)$ is $\text{rect}(-f_x)\text{rect}(-f_y) = \text{rect}(f_x)\text{rect}(f_y)$ by the parity of the rect function.

1.2.3. Fourier transforms in polar coordinates

For a two-dimensional function with circular symmetry, it is more convenient to use polar coordinates. We consider a plane described by rectangular (x, y) and polar

(r, θ) coordinates and the corresponding spectral coordinates are (f_x, f_y) and (ρ, φ) , respectively. We then have:

$$\begin{cases} x = r \cos \theta \\ y = r \sin \theta \end{cases} \quad [1.41]$$

$$\begin{cases} f_x = \rho \cos \varphi \\ f_y = \rho \sin \varphi \end{cases} \quad [1.42]$$

Let $f(x, y)$ be an original function with spectral function $F(f_x, f_y)$. We can rewrite these as functions of polar coordinates:

$$g(r, \theta) = f(r \cos \theta, r \sin \theta) \quad [1.43]$$

$$G(\rho, \varphi) = F(\rho \cos \varphi, \rho \sin \varphi) \quad [1.44]$$

By substituting these two relations into [1.23] and [1.25], we obtain direct and inverse Fourier transforms, respectively, in polar coordinates:

$$G(\rho, \varphi) = \int_0^{2\pi} \int_0^{+\infty} r g(r, \theta) \exp[-j2\pi r \rho \cos(\theta - \varphi)] dr d\theta \quad [1.45]$$

$$g(r, \theta) = \int_0^{2\pi} \int_0^{+\infty} \rho G(\rho, \varphi) \exp[j2\pi r \rho \cos(\theta - \varphi)] d\rho d\varphi \quad [1.46]$$

Most optical systems are circularly symmetric, and in this case the function $f(r, \theta)$ depends only on the variable r . We, therefore, have $g(r, \theta) = g_R(r)$. We substitute this relation into [1.45] and, using the identity of the Bessel function:

$$J_0(a) = \frac{1}{2\pi} \int_0^{2\pi} \exp[-ja \cos(\theta - \varphi)] d\theta \quad [1.47]$$

we can deduce the Fourier transform of $g_R(r)$ in polar coordinates:

$$G(\rho, \theta) = G_R(\rho) = 2\pi \int_0^{+\infty} r g_R(r) J_0(2\pi r \rho) dr \quad [1.48]$$

where $J_0(a)$ is a zero-order Bessel function of the first kind. Thus, the Fourier transform of a circularly symmetric function is itself circularly symmetric and the expression [1.48] is called a Fourier–Bessel transform or Hankel transform of zero

order. In the same way, by substituting $G_R(\rho) = G(\rho, \varphi)$ into [1.46], we determine the expression of inverse Fourier transform in polar coordinates:

$$g(r) = 2\pi \int_0^{+\infty} \rho G_R(\rho) J_0(2\pi r \rho) d\rho \quad [1.49]$$

We note that the mathematical forms of the direct and inverse transformations are the same.

1.3. Linear systems

An optical system allows the transformation of an input signal into an output signal. The device situated between the two planes (“input” and “output”) perpendicular to the direction of propagation will be henceforth called an “optical system”. An optical system may have linear or nonlinear properties. In most cases, considering the system to be linear as a first approximation, we are able to obtain sufficiently precise representations of the observed phenomena. Here we will consider only linear systems.

1.3.1. Definition

From a mathematical point of view, a linear system corresponds to a transformation operation. We conveniently represent such a system by an operator $L\{\}$, at whose output the two-dimensional function $f(x, y)$ becomes a new function $p(x', y')$. This is expressed as:

$$p(x', y') = L\{f(x, y)\} \quad [1.50]$$

$f(x, y)$ and $p(x', y')$ are called the input function and the output function of the system, respectively.

Let us consider some input functions $f_1(x, y), f_2(x, y), \dots, f_n(x, y)$ and some output functions $p_1(x', y'), p_2(x', y'), \dots, p_n(x', y')$. We then have:

$$\begin{aligned} p_1(x', y') &= L\{f_1(x, y)\} \\ p_2(x', y') &= L\{f_2(x, y)\} \\ &\vdots \\ p_n(x', y') &= L\{f_n(x, y)\} \end{aligned} \quad [1.51]$$

Assuming a_1, a_2, \dots, a_n to be complex constants, if the set of a system's input and output functions satisfy:

$$\begin{aligned} p(x', y') &= L\{f_1(x, y) + f_2(x, y) + \dots + f_n(x, y)\} \\ &= L\{f_1(x, y)\} + L\{f_2(x, y)\} + \dots + L\{f_n(x, y)\} \\ &= p_1(x', y') + p_2(x', y') + \dots + p_n(x', y') \end{aligned} \quad [1.52]$$

and

$$\begin{aligned} p(x', y') &= L\{a_1 f_1(x, y) + a_2 f_2(x, y) + \dots + a_n f_n(x, y)\} \\ &= a_1 L\{f_1(x, y)\} + a_2 L\{f_2(x, y)\} + \dots + a_n L\{f_n(x, y)\} \\ &= a_1 p_1(x', y') + a_2 p_2(x', y') + \dots + a_n p_n(x', y') \end{aligned} \quad [1.53]$$

then this system can be considered linear.

The linear approach presents a considerable advantage: it allows us to express the response of a system to any input function in the form of a response to “elementary” functions into which the input has been decomposed. In conclusion, if we can decompose, by a simple method, the input function into “elementary” functions for which the response of the system is well known, we will obtain the output function by the sum of these responses.

1.3.2. Impulse response and superposition integrals

Using the translation property of the two-dimensional δ function, we can express a function $f(x, y)$ describing a light wave in the input plane as:

$$f(x, y) = \iint_{\infty} f(x_0, y_0) \delta(x - x_0, y - y_0) dx_0 dy_0 \quad [1.54]$$

The physical meaning of this expression is that the distribution of the input optical signal $f(x, y)$ can be considered as the linear combination of δ functions weighted by the value $f(x_0, y_0)$ and shifted with respect to each other, the elementary functions of the decomposition being precisely these δ functions. Since the system is linear, its response to the input signal $f(x, y)$ is determined by:

$$p(x, y) = L\left\{ \iint_{\infty} f(x_0, y_0) \delta(x - x_0, y - y_0) dx_0 dy_0 \right\} \quad [1.55]$$

We notice that the number $f(x_0, y_0)$ is a simple weighting factor applied to the elementary function $\delta(x-x_0, y-y_0)$. For any point with coordinates (x_0, y_0) , $f(x_0, y_0)$ is constant. According to its property of linearity, the operator $L\{\}$ can move inside the summation (integral) sign, giving:

$$p(x, y) = \iint_{\infty} f(x_0, y_0) L\{ \delta(x-x_0, y-y_0) \} dx_0 dy_0 \quad [1.56]$$

If we consider $h(x, y; x_0, y_0)$ as the response of the system at the point (x, y) of the output space, when the input is a δ function situated at the point (x_0, y_0) , we have:

$$h(x, y; x_0, y_0) = L\{ \delta(x-x_0, y-y_0) \} \quad [1.57]$$

The function h is called the impulse response of the system. The magnitude of the input and output of the system can then be related by the following equation:

$$p(x, y) = \iint_{\infty} f(x_0, y_0) h(x, y; x_0, y_0) dx_0 dy_0 \quad [1.58]$$

This fundamental expression goes by the name of the “superposition integral”.

To completely determine the output signal, we note that we must know the responses to local impulses at every possible point of the input plane. In general, the determination of the impulse responses is very complex. However, we will see in the following section that for an important subclass of linear systems called invariant linear systems, which are invariant in the space, we can determine the impulse responses in a simple way. In most cases, an optical system can be approximated by a space-invariant linear system.

1.3.3. Definition of a two-dimensional linear shift-invariant system

Two-dimensional invariant linear systems are an important subclass of linear systems. If the impulse response of a system $h(x, y; x_0, y_0)$ depends only on the distances $(x-x_0)$ and $(y-y_0)$, this system is considered as a linear shift-invariant system, that is:

$$h(x, y; x_0, y_0) = h(x-x_0, y-y_0) \quad [1.59]$$

Thus, the optical system is invariant in space if the output signal of a point in the input plane changes position only but not shape when the source point moves around

the input plane. For a shift-invariant optical system, the superposition integral can be rewritten as:

$$p(x, y) = \iint_{-\infty}^{+\infty} f(x_0, y_0) h(x - x_0, y - y_0) dx_0 dy_0 = f(x, y) * h(x, y) \quad [1.60]$$

This relation corresponds to the two-dimensional convolution of the input function with the impulse response of the system. Consequently, if an optical system is a linear shift-invariant system, on the condition that we are able to determine the impulse response of a point in the input plane (which is often considered on the axis of the system), whatever the optical input signal $f(x_0, y_0)$, the output signal $p(x, y)$ can be determined using expression [1.60].

1.3.4. Transfer functions

By taking the Fourier transform of both sides of [1.60] and using the convolution theorem, we obtain:

$$P(f_x, f_y) = F(f_x, f_y) H(f_x, f_y) \quad [1.61]$$

with:

$$F(f_x, f_y) = \int \int_{-\infty}^{+\infty} f(x, y) \exp[-j2\pi(f_x x + f_y y)] dx dy \quad [1.61a]$$

$$P(f_x, f_y) = \int \int_{-\infty}^{+\infty} p(x, y) \exp[-j2\pi(f_x x + f_y y)] dx dy \quad [1.61b]$$

$$H(f_x, f_y) = \int \int_{-\infty}^{+\infty} h(x, y) \exp[-j2\pi(f_x x + f_y y)] dx dy \quad [1.61c]$$

Expression [1.61] shows that the spectral function of the output signal is the product of the spectrum of the input signal with the function $H(f_x, f_y)$. This product is the frequency response to one of the elementary functions of the input signal and the function $H(f_x, f_y)$ is called the transfer function of the system. The transfer function of the system is determined by the Fourier transform of the impulse response [1.61c]. The output signal can be determined by the inverse Fourier transform of the spectrum of the output signal, that is:

$$p(x, y) = F^{-1} \{ P(f_x, f_y) H(f_x, f_y) \} \quad [1.62]$$

The spectrum of the output signal $P(f_x, f_y)$ can be calculated by the Fourier transform [1.61b]. If the transfer function of the system can be determined, the output signal can be obtained by [1.61b].

1.4. The sampling theorem

It is often convenient to represent a continuous function $g(x, y)$ by a table of sampled values taken at a discrete set of points in the xy -plane. Current numerical methods allow the presentation, storage, and propagation of almost all information of a physical nature. It is intuitive that if the samples of the continuous function $g(x, y)$ are taken at points sufficiently close together, the given samples are able to reliably represent the original function using a simple interpolation. However, for a given function, the question is to know the maximum sampling interval that we must respect. The answer is less evident. Yet, for a particular class of functions known as “bandwidth-limited functions”, the reconstruction can be carried out exactly, on the condition that the interval between two samples is not larger than a certain limit. A bandwidth-limited function is such that its Fourier transform is only non-zero on a finite region of the frequency space. The sampling theorem was initially proven by Whittaker [WHI 15] and was later revisited by Shannon [SHA 49] during his studies on information theory. This principle, which allows us to determine the maximum sampling interval, is called the Shannon–Whittaker sampling theorem.

The following section states the two-dimensional sampling theorem and refers to the work of Goodman [GOO 72].

1.4.1. *Sampling a continuous function*

Let us consider a set of samples of the function $g(x, y)$, taken over a rectangular mesh. The sampled function $g_s(x, y)$ is defined as:

$$g_s(x, y) = \text{comb}\left(\frac{x}{X}\right) \text{comb}\left(\frac{y}{Y}\right) g(x, y) \quad [1.63]$$

This function therefore consists of a set of δ functions separated by intervals of length X along the x -axis and of length Y along the y -axis, as shown in Figure 1.10, whose amplitude is the value of the function $g(x, y)$ at the point being considered.

The volume enclosed by the δ function representation in the space and the xy -plane is proportional to the value of $g(x, y)$ at each point of the sampling mesh.

Applying the convolution theorem, we obtain the spectrum $G_s(f_x, f_y)$ of $g_s(x, y)$ by convoluting the transform of $(x/X)\text{comb}(y/Y)$ with the transform of $g(x, y)$, that is:

$$G_s(f_x, f_y) = F \left\{ \text{comb} \left(\frac{x}{X} \right) \text{comb} \left(\frac{y}{Y} \right) \right\} * G(f_x, f_y) \quad [1.64]$$

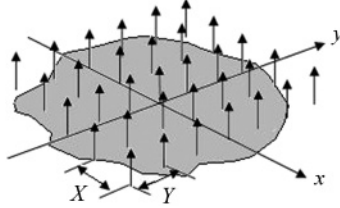


Figure 1.10. Two-dimensional sampling

Since:

$$\begin{aligned} F \left\{ \text{comb} \left(\frac{x}{X} \right) \text{comb} \left(\frac{y}{Y} \right) \right\} &= XY \text{comb}(Xf_x) \text{comb}(Yf_y) \\ &= \sum_{n=-\infty}^{\infty} \sum_{m=-\infty}^{\infty} \delta \left(f_x - \frac{n}{X} \right) \delta \left(f_y - \frac{m}{Y} \right) \end{aligned} \quad [1.65]$$

then we have:

$$G_s(f_x, f_y) = \sum_{n=-\infty}^{\infty} \sum_{m=-\infty}^{\infty} G \left(f_x - \frac{n}{X}, f_y - \frac{m}{Y} \right) \quad [1.66]$$

The spectrum of $g_s(x, y)$ can therefore be simply deduced by considering the spectrum of $g(x, y)$ localized at each point with coordinates $(n/X, m/Y)$ in the $f_x f_y$ -plane, as shown in Figure 1.11.

Since we assumed that the function $g(x, y)$ had a spectrum of limited scope, its spectrum $G(f_x, f_y)$ is only non-zero in the corresponding frequency space domain. If X and Y are sufficiently small, in other words, if the samples are taken on points that are sufficiently close to each other, the intervals $1/X$ and $1/Y$ between the various spectral regions will be large enough to ascertain that the neighboring regions do not overlap. To determine the maximum interval between two sampled points, let us suppose $2B_X$ and $2B_Y$ to be the dimensions following the respective directions of the

f_x - and f_y -axes of the smallest rectangle containing the whole spectral domain of $g(x, y)$. As shown in Figure 1.11, if the following two inequalities:

$$\begin{aligned} X &\leq \frac{1}{2B_x} \\ Y &\leq \frac{1}{2B_y} \end{aligned} \quad [1.67]$$

are satisfied, the different terms of the spectrum [1.66] of the sampled function are separated by the distances $1/X$ and $1/Y$ in the f_x and f_y directions, respectively. The maximum dimensions of the mesh of the sample network, which allow an exact restoration of the original function, are therefore $1/2B_x$ and $1/2B_y$. Having determined the maximum allowed distances between samples, we now study how to obtain the spectrum of $g(x, y)$ by a filter function, and how to reconstruct the original function $g(x, y)$.

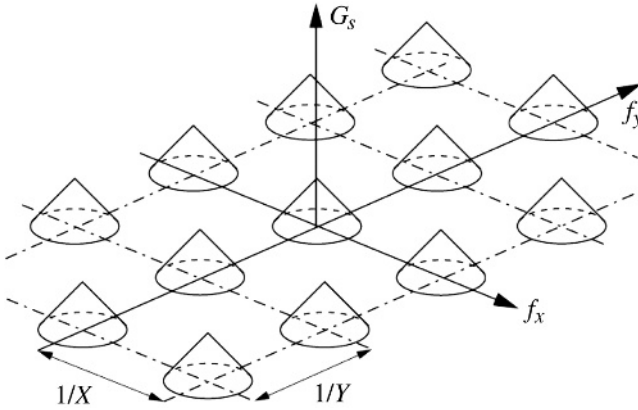


Figure 1.11. *Spectrum of the sampled function*

1.4.2. Reconstruction of the original function

Following Figure 1.11, we consider a two-dimensional rectangular function with sides $2B_x$ and $2B_y$ along the f_x - and f_y -axes, respectively. The filter function is:

$$H(f_x, f_y) = \text{rect}\left(\frac{f_x}{2B_x}\right) \cdot \text{rect}\left(\frac{f_y}{2B_y}\right) \quad [1.68]$$

We note that $G(f_x, f_y)$ is obtained from $G_s(f_x, f_y)$ since:

$$G(f_x, f_y) \equiv G_s(f_x, f_y) H(f_x, f_y) \quad [1.69]$$

This means that, if the sampled function $g_s(x, y)$ is considered as the input signal of a system, the function $g(x, y)$ will be considered as the output signal. Thus, $H(f_x, f_y)$ is the transfer function of the system. In this case, the identity [1.64] translates into the spatial domain by:

$$g_s(x, y) * h(x, y) \equiv g(x, y) \quad [1.70]$$

where

$$\begin{aligned} g_s(x, y) &= \text{comb}\left(\frac{x}{X}\right) \text{comb}\left(\frac{y}{Y}\right) g(x, y) \\ &= XY \sum_{n=-\infty}^{\infty} \sum_{m=-\infty}^{\infty} g(nX, mY) \delta(x - nX, y - mY) \end{aligned} \quad [1.71]$$

and $h(x, y)$ is the impulse response of the filter, which is written as:

$$h(x, y) = F^{-1} \left\{ \text{rect}\left(\frac{f_x}{2B_x}\right) \text{rect}\left(\frac{f_y}{2B_y}\right) \right\} = 4B_x B_y \text{sinc}(2B_x x) \text{sinc}(2B_y y) \quad [1.72]$$

Consequently:

$$\begin{aligned} g(x, y) &= 4B_x B_y XY \sum_{n=-\infty}^{\infty} \sum_{m=-\infty}^{\infty} \\ &g(nX, mY) \text{sinc}\left[2B_x(x - nX)\right] \text{sinc}\left[2B_y(y - mY)\right] \end{aligned} \quad [1.73]$$

Finally, when we choose the maximum allowed values $1/2B_x$ and $1/2B_y$ for the sampling intervals X and Y , the identity becomes:

$$g(x, y) = \sum_{n=-\infty}^{\infty} \sum_{m=-\infty}^{\infty} g\left(\frac{n}{2B_x}, \frac{m}{2B_y}\right) \text{sinc}\left(2B_x\left(x - \frac{n}{2B_x}\right)\right) \text{sinc}\left(2B_y\left(y - \frac{m}{2B_y}\right)\right) \quad [1.74]$$

Expression [1.74] represents a fundamental result that we will henceforth call the Whittaker–Shannon sampling theorem. It states that the exact reconstruction of a bandwidth-limited function can be carried out from the sampled values of the

function, taken after a suitable rectangular mesh sampling. The reconstruction is carried out by interpolating each sample point by an interpolation function constituted by the product of two sinc functions.

Note that this result is not the only possible sampling theorem. We chose two rather arbitrary sampling frames in the course of this study; with different assumptions we would have obtained a different sampling theorem. We first arbitrarily chose a rectangular sampling frame. Also, we chose the particular transfer function given by [1.68]. By making different choices we would establish other, equally valid theorems. For more detail, readers may refer to the articles by [BRA 56], [PET 62], and [LIN 59].

1.4.3. *Space-bandwidth product*

For a bandwidth-limited function $g(x, y)$, which is mainly non-zero in a region of the xy -plane bounded by $-L_X \leq x \leq L_X$ and $-L_Y \leq y \leq L_Y$ and whose maximum sampling intervals along the f_x - and f_y -axes are $1/2B_X$ and $1/2B_Y$, respectively, to thus satisfy the sampling theorem, the minimum value of the number of sampling points able to represent the function $g(x, y)$ is therefore:

$$N = (2L_X \times 2L_Y)(2B_X \times 2B_Y) = 16L_X L_Y B_X B_Y \quad [1.75]$$

This relation is called the *space-bandwidth* product of the function $g(x, y)$ [GOO 05] and expresses the value of the product of the space and frequency surfaces in which the function $g(x, y)$ and its spectrum $F\{g(x, y)\}$ are bounded. As a result, for a two-dimensional bandwidth-limited function, the space-bandwidth product determines the minimum number of degrees of freedom, N , that correctly represent it. When $g(x, y)$ is real, its number of degrees of freedom is N , since the samples are real; if $g(x, y)$ is a complex function, its number of degrees of freedom becomes $2N$ as each sample must be represented by two real values. Given the theorems of similarity and translation relating to the Fourier transform, the dilation of the coordinates and the translation of the function in the spatial or spectral domain do not affect the space-bandwidth product of the considered function. This means that, for a given function, the number of degrees of freedom is constant. This number can therefore be considered as a significant piece of information expressing the complexity of the function, and as a criterion that allows us to verify whether there is any loss of information during the sampling process.

