

Linear algebra, signal processing, and wavelets.
A unified approach.

Øyvind Ryan

January 6, 2014

Contents

1 Fourier analysis and applications to sound processing	1
1 Sound and Fourier series	3
1.1 Loudness: Sound pressure and decibels	4
1.2 The pitch of a sound	6
1.3 Fourier series: Basic concepts	11
1.4 Examples of Fourier series	16
1.5 Complex Fourier series	22
1.6 Rate of convergence for Fourier series	28
1.7 Some properties of Fourier series	32
1.8 Operations on sound: filters	35
1.9 The MP3 standard	37
1.10 Summary	40
2 Digital sound and Discrete Fourier analysis	41
2.1 Digital sound	41
2.2 Simple operations on digital sound	43
2.3 Discrete Fourier analysis: Basic concepts	49
2.4 The Discrete Fourier Transform	51
2.5 Connection between the DFT and Fourier series	58
2.6 Applications of the DFT	60
2.7 Reconstruction of a function from its samples. The sampling theorem.	66
2.8 The Fast Fourier Transform (FFT)	69
2.9 Summary	81
3 Operations on digital sound: digital filters	83
3.1 Matrix representations of filters	83
3.2 Formal definition of filters and the vector frequency response	87
3.3 The continuous frequency response and properties	94
3.4 Assembling the filter matrix and compact notation	101
3.5 Some examples of filters	104
3.6 Time-invariance of filters	119

3.7	More general filters	120
3.8	Implementation of filters	122
3.9	Summary	124
4	Symmetric filters and the DCT	125
4.1	Symmetric vectors and the DCT	126
4.2	Improvements from using the DCT to interpolate functions and approximate analog filters	139
4.3	Efficient implementations of the DCT	143
4.4	Summary	151
II	Wavelets and applications to image processing	153
5	Motivation for wavelets and some simple examples	155
5.1	Why wavelets?	156
5.2	A wavelet based on piecewise constant functions	157
5.3	m -level Discrete Wavelet Transform	170
5.4	A wavelet based on piecewise linear functions	179
5.5	Alternative wavelet based on piecewise linear functions	187
5.6	Multiresolution analysis: A generalization	193
5.7	Summary	197
6	The filter representation of wavelets	199
6.1	The filters of a wavelet transformation	200
6.2	Properties of the filter bank transforms of a wavelet	209
6.3	Filter-based algorithm for the DWT and the IDWT	220
6.4	A generalization of the filter representation, and its use in audio coding	225
6.5	Summary	234
7	Constructing interesting wavelets	237
7.1	From filters to scaling functions and mother wavelets	238
7.2	Vanishing moments	245
7.3	Characterization of wavelets w.r.t. number of vanishing moments	249
7.4	A design strategy suitable for lossless compression	256
7.5	A design strategy suitable for lossy compression	259
7.6	Orthonormal wavelets	261
7.7	Summary	263
8	The polyphase representation and wavelets	265
8.1	The polyphase representation	266
8.2	Motivation for the lifting factorization	268
8.3	The lifting factorization	270
8.4	The proof of Theorem 6.14	277
8.5	Cosine-modulated filter banks and the MP3 standard	279
8.6	Summary	290

9 Digital images	291
9.1 What is an image?	292
9.2 Some simple operations on images	295
9.3 Filter-based operations on images	303
9.4 Change of coordinates for images	318
9.5 Summary	324
10 Using tensor products to apply wavelets to images	327
10.1 Tensor product of function spaces	327
10.2 Tensor product of function spaces in a wavelet setting	330
10.3 Experiments with images using wavelets	337
10.4 An application to the FBI standard for compression of fingerprint images	348
10.5 Summary	353
A Basic Linear Algebra	355
A.1 Matrices	355
A.2 Vector spaces	355
A.3 Inner products and orthogonality	356
A.4 Coordinates and change of coordinates	356
A.5 Eigenvectors and eigenvalues	357
A.6 Diagonalization	357
B Signal processing and linear algebra: a translation guide	359
B.1 Complex numbers	359
B.2 Functions	359
B.3 Vectors	360
B.4 Inner products and orthogonality	360
B.5 Matrices and filters	360
B.6 Convolution	361
B.7 Polyphase factorizations and lifting	361
B.8 Transforms in general	362
B.9 Perfect reconstruction systems	362
B.10 Z-transform and frequency response	362
Bibliography	365
Nomenclature	367
Mathematics index	372
Index for MATLAB commands	376

Part I

Fourier analysis and applications to sound processing

Sound and Fourier series

A major part of the information we receive and perceive every day is in the form of audio. Most sounds are transferred directly from the source to our ears, like when we have a face to face conversation with someone or listen to the sounds in a forest or a street. However, a considerable part of the sounds are generated by loudspeakers in various kinds of audio machines like cell phones, digital audio players, home cinemas, radios, television sets and so on. The sounds produced by these machines are either generated from information stored inside, or electromagnetic waves are picked up by an antenna, processed, and then converted to sound. It is this kind of sound we are going to study in this chapter. The sound that is stored inside the machines or picked up by the antennas is usually represented as *digital sound*. This has certain limitations, but at the same time makes it very easy to manipulate and process the sound on a computer.

What we perceive as sound corresponds to the physical phenomenon of slight variations in air pressure near our ears. Larger variations mean louder sounds, while faster variations correspond to sounds with a higher pitch. The air pressure varies continuously with time, but at a given point in time it has a precise value. This means that sound can be considered to be a mathematical function.

Observation 1.1. A sound can be represented by a mathematical function, with time as the free variable. When a function represents a sound, it is often referred to as a *continuous sound*.

In the following we will briefly discuss the basic properties of sound: first the significance of the size of the variations, and then how many variations there are per second, the *frequency* of the sound. We also consider the important fact that any reasonable sound may be considered to be built from very simple basis sounds. Since a sound may be viewed as a function, the mathematical equivalent of this is that any decent function may be constructed from very simple basis functions. Fourier-analysis is the theoretical study of this, and in the last part of this chapter we establish the framework for this study, and analyze this on some examples for sound.

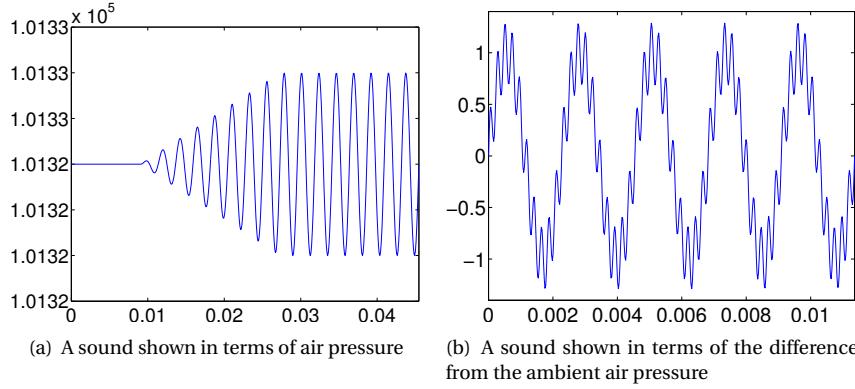


Figure 1.1: Two examples of audio signals.

1.1 Loudness: Sound pressure and decibels

An example of a simple sound is shown in Figure 1.1(a) where the oscillations in air pressure are plotted against time. We observe that the initial air pressure has the value 101 325 (we will shortly return to what unit is used here), and then the pressure starts to vary more and more until it oscillates regularly between the values 101 323 and 101 327. In the area where the air pressure is constant, no sound will be heard, but as the variations increase in size, the sound becomes louder and louder until about time $t = 0.6$ where the size of the oscillations becomes constant. The following summarises some basic facts about air pressure.

Fact 1.2 (Air pressure). Air pressure is measured by the SI-unit Pa (Pascal) which is equivalent to N/m^2 (force / area). In other words, 1 Pa corresponds to the force exerted on an area of $1\ m^2$ by the air column above this area. The normal air pressure at sea level is 101 325 Pa.

Fact 1.2 explains the values on the vertical axis in Figure 1.1(a): The sound was recorded at the normal air pressure of 101 325 Pa. Once the sound started, the pressure started to vary both below and above this value, and after a short transient phase the pressure varied steadily between 101 324 Pa and 101 326 Pa, which corresponds to variations of size 1 Pa about the fixed value. Everyday sounds typically correspond to variations in air pressure of about 0.00002–2 Pa, while a jet engine may cause variations as large as 200 Pa. Short exposure to variations of about 20 Pa may in fact lead to hearing damage. The volcanic eruption at Krakatoa, Indonesia, in 1883, produced a sound wave with variations as large as almost 100 000 Pa, and the explosion could be heard 5000 km away.

When discussing sound, one is usually only interested in the variations in air pressure, so the ambient air pressure is subtracted from the measurement. This corresponds to subtracting 101 325 from the values on the vertical axis in Figure 1.1(a). In Figure 1.1(b) the subtraction has been performed for another sound, and we see

that the sound has a slow, cos-like, variation in air pressure, with some smaller and faster variations imposed on this. This combination of several kinds of systematic oscillations in air pressure is typical for general sounds. The size of the oscillations is directly related to the loudness of the sound. We have seen that for audible sounds the variations may range from 0.00002 Pa all the way up to 100 000 Pa. This is such a wide range that it is common to measure the loudness of a sound on a logarithmic scale. Often air pressure is normalized so that it lies between -1 and 1 : The value 0 then represents the ambient air pressure, while -1 and 1 represent the lowest and highest representable air pressure, respectively. The following fact box summarises the previous discussion of what a sound is, and introduces the logarithmic decibel scale.

Fact 1.3 (Sound pressure and decibels). The physical origin of sound is variations in air pressure near the ear. The *sound pressure* of a sound is obtained by subtracting the average air pressure over a suitable time interval from the measured air pressure within the time interval. A square of this difference is then averaged over time, and the sound pressure is the square root of this average.

It is common to relate a given sound pressure to the smallest sound pressure that can be perceived, as a level on a decibel scale,

$$L_p = 10 \log_{10} \left(\frac{p^2}{p_{\text{ref}}^2} \right) = 20 \log_{10} \left(\frac{p}{p_{\text{ref}}} \right).$$

Here p is the measured sound pressure while p_{ref} is the sound pressure of a just perceivable sound, usually considered to be 0.00002 Pa.

The square of the sound pressure appears in the definition of L_p since this represents the *power* of the sound which is relevant for what we perceive as loudness.

The sounds in Figure 1.1 are synthetic in that they were constructed from mathematical formulas (see Exercises 2.2.1 and 2.2.2). The sounds in Figure 1.2 on the other hand show the variation in air pressure when there is no mathematical formula involved, such as is the case for a song. In (a) there are so many oscillations that it is impossible to see the details, but if we zoom in as in (c) we can see that there is a continuous function behind all the ink. It is important to realise that in reality the air pressure varies more than this, even over the short time period in (c). However, the measuring equipment may not be able to pick up those variations, and it is also doubtful whether we would be able to perceive such rapid variations.

Exercises for Section 1.1

- Compute the loudness of the Krakatoa explosion on the decibel scale, assuming that the variation in air pressure peaked at 100 000 Pa.

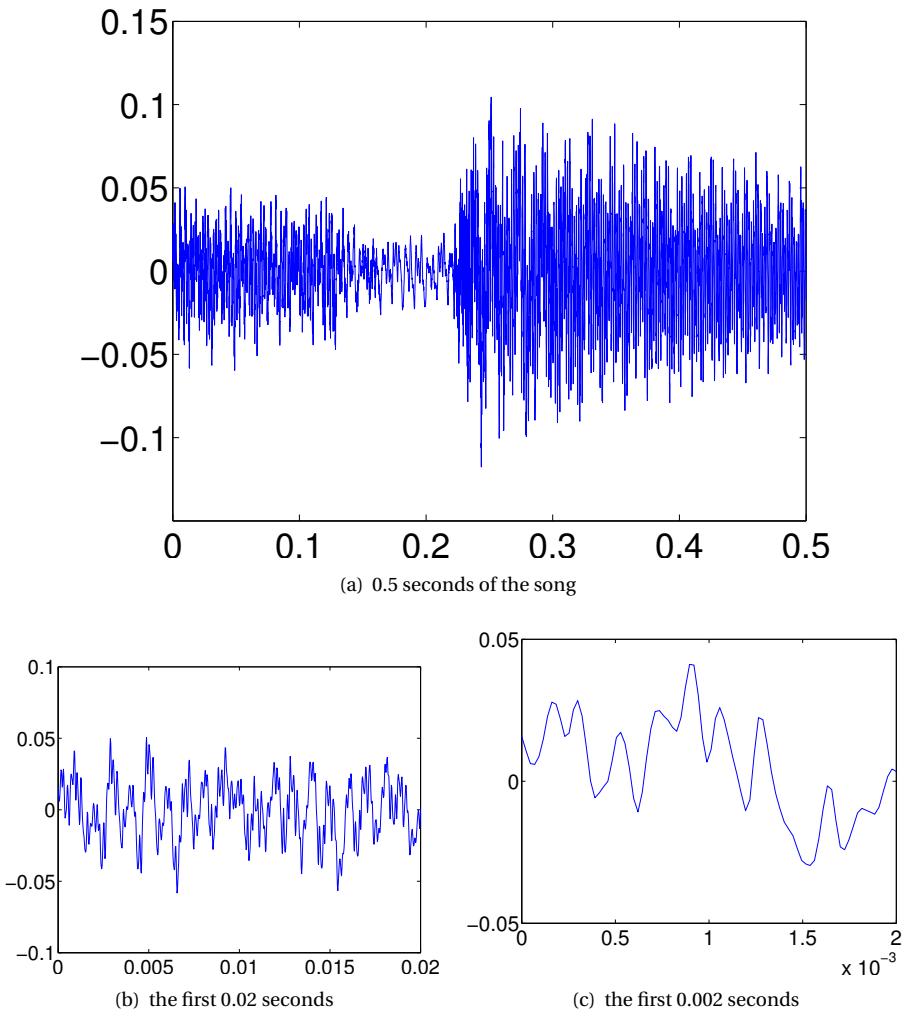


Figure 1.2: Variations in air pressure during parts of a song.

1.2 The pitch of a sound

Besides the size of the variations in air pressure, a sound has another important characteristic, namely the frequency (speed) of the variations. For most sounds the frequency of the variations varies with time, but if we are to perceive variations in air pressure as sound, they must fall within a certain range.

Fact 1.4. For a human with good hearing to perceive variations in air pressure as sound, the number of variations per second must be in the range 20–20 000.

To make these concepts more precise, we first recall what it means for a function to be periodic.

Definition 1.5. A real function f is said to be periodic with period T if

$$f(t + T) = f(t)$$

for all real numbers t .

Note that all the values of a periodic function f with period T are known if $f(t)$ is known for all t in the interval $[0, T)$. The prototypes of periodic functions are the trigonometric ones, and particularly $\sin t$ and $\cos t$ are of interest to us. Since $\sin(t + 2\pi) = \sin t$, we see that the period of $\sin t$ is 2π and the same is true for $\cos t$.

There is a simple way to change the period of a periodic function, namely by multiplying the argument by a constant.

Observation 1.6 (Frequency). If v is an integer, the function $f(t) = \sin(2\pi vt)$ is periodic with period $T = 1/v$. When t varies in the interval $[0, 1]$, this function covers a total of v periods. This is expressed by saying that f has *frequency* v .

Figure 1.3 illustrates observation 1.6. The function in (a) is the plain $\sin t$ which covers one period when t varies in the interval $[0, 2\pi]$. By multiplying the argument by 2π , the period is squeezed into the interval $[0, 1]$ so the function $\sin(2\pi t)$ has frequency $v = 1$. Then, by also multiplying the argument by 2, we push two whole periods into the interval $[0, 1]$, so the function $\sin(2\pi 2t)$ has frequency $v = 2$. In (d) the argument has been multiplied by 5 — hence the frequency is 5 and there are five whole periods in the interval $[0, 1]$. Note that any function on the form $\sin(2\pi vt + a)$ has frequency v , regardless of the value of a .

Since sound can be modelled by functions, it is reasonable to say that a sound with frequency v is a trigonometric function with frequency v .

Definition 1.7. The function $\sin(2\pi vt)$ represents what we will call a *pure tone* with frequency v . Frequency is measured in Hz (Herz) which is the same as s^{-1} (the time t is measured in seconds).

A pure tone with frequency 440 Hz sounds like this, and a pure tone with frequency 1500 Hz sounds like this. In Section 2.1 we will explain how we generated these sounds so that they could be played on a computer.

Any sound may be considered to be a function. In the next section we will explain why any reasonable function may be written as a sum of simple sin- and cos-functions with integer frequencies. When this is translated into properties of sound, we obtain an important principle.

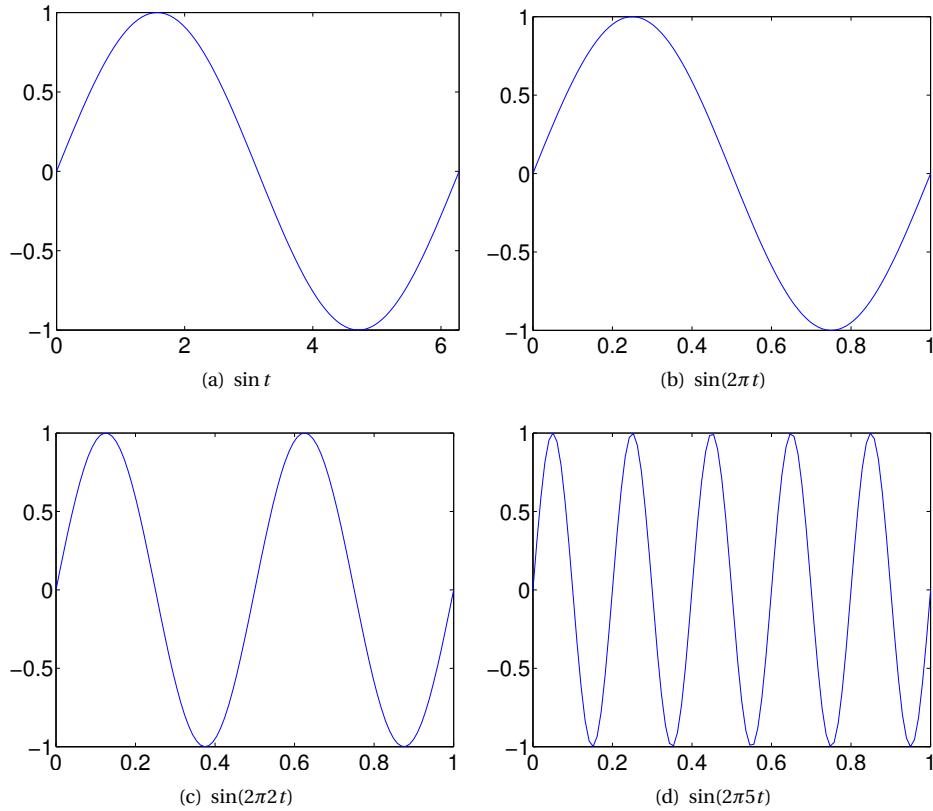


Figure 1.3: Versions of sin with different frequencies.

Observation 1.8 (Decomposition of sound into pure tones). Any sound f is a sum of pure tones at different frequencies. The amount of each frequency required to form f is the frequency content of f . Any sound can be reconstructed from its frequency content.

The most basic consequence of observation 1.8 is that it gives us an understanding of how any sound can be built from the simple building blocks of pure tones. This also means that we can store a sound f by storing its frequency content, as an alternative to storing f itself. This also gives us a possibility for lossy compression of digital sound: It turns out that, in a typical audio signal, most information is found in the lower frequencies, and some frequencies will be almost completely absent. This can be exploited for compression if we change the frequencies with small contribution a little bit and set them to 0, and then store the signal by only storing the nonzero part of the frequency content. When the sound is to be played back, we first convert the adjusted values to the adjusted frequency content back to a normal

function representation with an inverse mapping.

Fact 1.9 (Basic idea behind audio compression). Suppose an audio signal f is given. To compress f , perform the following steps:

1. Rewrite the signal f in a new format where frequency information becomes accessible.
2. Remove those frequencies that only contribute marginally to human perception of the sound.
3. Store the resulting sound by coding the adjusted frequency content with some lossless coding method.

This lossy compression strategy is essentially what is used in practice by commercial audio formats. The difference is that commercial software does everything in a more sophisticated way and thereby gets better compression rates. We will return to this in later chapters.

We will see later that Observation 1.8 also is the basis for many operations on sound. The same observation also makes it possible to explain more precisely what it means that we only perceive sounds with a frequency in the range 20–20000 Hz:

Fact 1.10. Humans can only perceive variations in air pressure as sound if the Fourier series of the sound signal contains at least one sufficiently large term with frequency in the range 20–20 000 Hz.

With appropriate software it is easy to generate a sound from a mathematical function; we can 'play' the function. If we play a function like $\sin(2\pi 440t)$, we hear a pleasant sound with a very distinct pitch, as expected. There are, however, many other ways in which a function can oscillate regularly. The function in Figure 1.1(b) for example, definitely oscillates 2 times every second, but it does not have frequency 2 Hz since it is not a pure tone. This sound is also not that pleasant to listen to. We will consider two more important examples of this, which are very different from smooth, trigonometric functions.

Example 1.11. We define the *square wave* of period T as the function which repeats with period T , and is 1 on the first half of each period, and -1 on the second half. This means that we can define it as the function

$$f_s(t) = \begin{cases} 1, & \text{if } 0 \leq t < T/2; \\ -1, & \text{if } T/2 \leq t < T. \end{cases} \quad (1.1)$$

In Figure 1.4(a) we have plotted the square wave when $T = 1/440$. This period is chosen so that it corresponds to the pure tone we already have listened to, and you can listen to this square wave here. In Exercise 2.1.4 you will learn how to generate this sound. We hear a sound with the same pitch as $\sin(2\pi 440t)$, but note that the square wave is less pleasant to listen to: There seems to be some sharp corners in

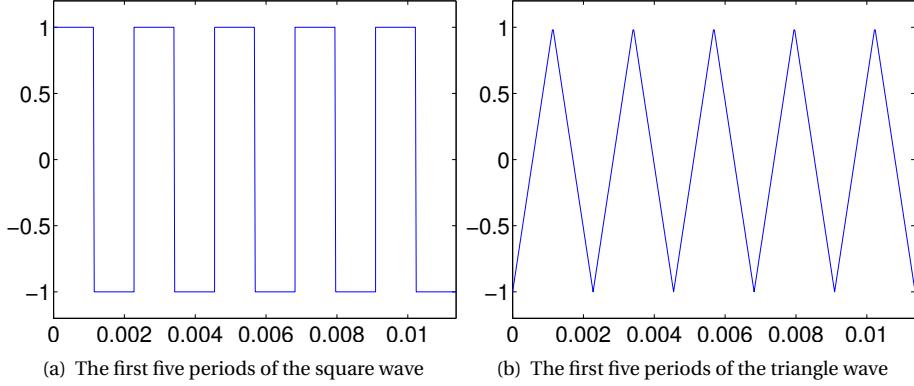


Figure 1.4: The square wave and the triangle wave, two functions with regular oscillations, but which are not simple, trigonometric functions.

the sound, translating into a rather shrieking, piercing sound. We will later explain this by the fact that the square wave can be viewed as a sum of many frequencies, and that all the different frequencies pollute the sound so that it is not pleasant to listen to.



Example 1.12. We define the *triangle wave* of period T as the function which repeats with period T , and increases linearly from -1 to 1 on the first half of each period, and decreases linearly from 1 to -1 on the second half of each period. This means that we can define it as the function

$$f_t(t) = \begin{cases} 4t/T - 1, & \text{if } 0 \leq t < T/2; \\ 3 - 4t/T, & \text{if } T/2 \leq t < T. \end{cases} \quad (1.2)$$

In Figure 1.4(b) we have plotted the triangle wave when $T = 1/440$. Again, this same choice of period gives us an audible sound, and you can listen to the triangle wave here. Again you will note that the triangle wave has the same pitch as $\sin(2\pi 440t)$, and is less pleasant to listen to than this pure tone. However, one can argue that it is somewhat more pleasant to listen to than a square wave. This will also be explained in terms of pollution with other frequencies later. ♣

In Section 1.3 we will begin to peek behind the curtains as to why these waves sound so different, even though we recognize them as having the exact same pitch.

Exercises for Section 1.2

1. Consider a sum of two pure tones, $f(t) = A_1 \sin(2\pi\nu_1 t) + A_2 \sin(2\pi\nu_2 t)$. For which values of A_1, A_2, ν_1, ν_2 is f periodic? What is the period of f when it is periodic?

1.3 Fourier series: Basic concepts

In Section 1.2 we identified audio signals with functions and discussed informally the idea of decomposing a sound into basis sounds (pure sounds) to make its frequency content available. In this chapter we will make this kind of decomposition more precise by discussing how a given function can be expressed in terms of the basic trigonometric functions. This is similar to Taylor series where functions are approximated by combinations of polynomials. But it is also different from Taylor series because we use trigonometric series rather than power series, and the approximations are computed in a very different way. The theory of approximation of functions with trigonometric functions is generally referred to as *Fourier analysis*. This is a central tool in practical fields like image- and signal processing, but it is also an important field of research within pure mathematics.

In the start of this chapter we had no constraints on the function f . Although Fourier analysis can be performed for very general functions, it turns out that it takes its simplest form when we assume that the function is periodic. Periodic functions are fully known when we know their values on a period $[0, T]$. In this case we will see that we can carry out the Fourier analysis in finite dimensional vector spaces of functions. This makes linear algebra a very useful tool in Fourier analysis: Many of the tools from your linear algebra course will be useful, in a situation that at first may seem far from matrices and vectors.

The basic idea of Fourier series is to approximate a given function by a combination of simple cos and sin functions. This means that we have to address at least three questions:

1. How general do we allow the given function to be?
2. What exactly are the combinations of cos and sin that we use for the approximations?
3. How do we determine the approximation?

Each of these questions will be answered in this section. Since we restrict to periodic functions, we will without much loss of generality assume that the functions are defined on $[0, T]$, where T is some positive number. Mostly we will also assume that f is continuous, but the theory can also be extended to functions which are only Riemann-integrable, and more precisely, to square integrable functions.

Definition 1.13 (Continuous and square-integrable functions). The set of continuous, real functions defined on an interval $[0, T]$ is denoted $C[0, T]$.

A real function f defined on $[0, T]$ is said to be square integrable if f^2 is Riemann-integrable, i.e., if the Riemann integral of f^2 on $[0, T]$ exists,

$$\int_0^T f(t)^2 dt < \infty.$$

The set of all square integrable functions on $[0, T]$ is denoted $L^2[0, T]$.

The sets of continuous and square-integrable functions can be equipped with an inner-product, a generalisation of the so-called dot-product for vectors.

Theorem 1.14. Both $L^2[0, T]$ and $C[0, T]$ are vector spaces. Moreover, if the two functions f and g lie in $L^2[0, T]$ (or in $C[0, T]$), then the product fg is also in $L^2[0, T]$ (or in $C[0, T]$). Moreover, both spaces are inner product spaces¹, with inner product² defined by

$$\langle f, g \rangle = \frac{1}{T} \int_0^T f(t)g(t) dt, \quad (1.3)$$

and associated norm

$$\|f\| = \sqrt{\frac{1}{T} \int_0^T f(t)^2 dt}. \quad (1.4)$$

The mysterious factor $1/T$ is included so that the constant function $f(t) = 1$ has norm 1, i.e., its role is as a normalizing factor.

Definition 1.13 and Theorem 1.14 answer the first question above, namely how general we allow our functions to be. Theorem 1.14 also gives an indication of how we are going to determine approximations—we are going to use inner products. We recall from linear algebra that the projection of a function f onto a subspace W with respect to an inner product $\langle \cdot, \cdot \rangle$ is the function $g \in W$ which minimizes $\|f - g\|$, also called *the error* in the approximation³. This projection is therefore also called a best approximation of f from W and is characterised by the fact that the function $f - g$, also called the *error function*, should be orthogonal to the subspace W , i.e. we should have

$$\langle f - g, h \rangle = 0, \quad \text{for all } h \in W.$$

More precisely, if $\phi = \{\phi_i\}_{i=1}^m$ is an orthogonal basis for W , then the best approximation g is given by

$$g = \sum_{i=1}^m \frac{\langle f, \phi_i \rangle}{\langle \phi_i, \phi_i \rangle} \phi_i. \quad (1.5)$$

The error $\|f - g\|$ is often referred to as the *least square error*.

We have now answered the second of our primary questions. What is left is a description of the subspace W of trigonometric functions. This space is spanned by the pure tones we discussed in Section 1.2.

Definition 1.15 (Fourier series). Let $V_{N,T}$ be the subspace of $C[0, T]$ spanned by the set of functions given by

$$\begin{aligned} \mathcal{D}_{N,T} = \{1, \cos(2\pi t/T), \cos(2\pi 2t/T), \dots, \cos(2\pi Nt/T), \\ \sin(2\pi t/T), \sin(2\pi 2t/T), \dots, \sin(2\pi Nt/T)\}. \end{aligned} \quad (1.6)$$

³See Section 6.3 in [20] for a review of projections and least squares approximations.

The space $V_{N,T}$ is called the N 'th order Fourier space. The N th-order Fourier series approximation of f , denoted f_N , is defined as the best approximation of f from $V_{N,T}$ with respect to the inner product defined by (1.3).

The space $V_{N,T}$ can be thought of as the space spanned by the pure tones of frequencies $1/T, 2/T, \dots, N/T$, and the Fourier series can be thought of as linear combination of all these pure tones. From our discussion in Section 1.2, we should expect that if N is sufficiently large, $V_{N,T}$ can be used to approximate most sounds in real life. The approximation f_N of a sound f from a space $V_{N,T}$ can also serve as a compressed version if many of the coefficients can be set to 0 without the error becoming too big.

Note that all the functions in the set $\mathcal{D}_{N,T}$ are periodic with period T , but most have an even shorter period. More precisely, $\cos(2\pi nt/T)$ has period T/n , and frequency n/T . In general, the term *fundamental frequency* is used to denote the lowest frequency of a given periodic function.

Definition 1.15 characterises the Fourier series. The next lemma gives precise expressions for the coefficients.

Theorem 1.16. The set $\mathcal{D}_{N,T}$ is an orthogonal basis for $V_{N,T}$. In particular, the dimension of $V_{N,T}$ is $2N+1$, and if f is a function in $L^2[0, T]$, we denote by a_0, \dots, a_N and b_1, \dots, b_N the coordinates of f_N in the basis $\mathcal{D}_{N,T}$, i.e.

$$f_N(t) = a_0 + \sum_{n=1}^N (a_n \cos(2\pi nt/T) + b_n \sin(2\pi nt/T)). \quad (1.7)$$

The a_0, \dots, a_N and b_1, \dots, b_N are called the (real) Fourier coefficients of f , and they are given by

$$a_0 = \langle f, 1 \rangle = \frac{1}{T} \int_0^T f(t) dt, \quad (1.8)$$

$$a_n = 2 \langle f, \cos(2\pi nt/T) \rangle = \frac{2}{T} \int_0^T f(t) \cos(2\pi nt/T) dt \quad \text{for } n \geq 1, \quad (1.9)$$

$$b_n = 2 \langle f, \sin(2\pi nt/T) \rangle = \frac{2}{T} \int_0^T f(t) \sin(2\pi nt/T) dt \quad \text{for } n \geq 1. \quad (1.10)$$

Proof: To prove orthogonality, assume first that $m \neq n$. We compute the inner

product

$$\begin{aligned}
& \langle \cos(2\pi mt/T), \cos(2\pi nt/T) \rangle \\
&= \frac{1}{T} \int_0^T \cos(2\pi mt/T) \cos(2\pi nt/T) dt \\
&= \frac{1}{2T} \int_0^T (\cos(2\pi mt/T + 2\pi nt/T) + \cos(2\pi mt/T - 2\pi nt/T)) \\
&= \frac{1}{2T} \left[\frac{T}{2\pi(m+n)} \sin(2\pi(m+n)t/T) + \frac{T}{2\pi(m-n)} \sin(2\pi(m-n)t/T) \right]_0^T \\
&= 0.
\end{aligned}$$

Here we have added the two identities $\cos(x \pm y) = \cos x \cos y \mp \sin x \sin y$ together to obtain an expression for $\cos(2\pi mt/T) \cos(2\pi nt/T) dt$ in terms of $\cos(2\pi mt/T + 2\pi nt/T)$ and $\cos(2\pi mt/T - 2\pi nt/T)$. By testing all other combinations of sin and cos also, we obtain the orthogonality of all functions in $\mathcal{D}_{N,T}$ in the same way.

We find the expressions for the Fourier coefficients from the general formula (1.5). We first need to compute the following inner products of the basis functions,

$$\begin{aligned}
\langle \cos(2\pi mt/T), \cos(2\pi mt/T) \rangle &= \frac{1}{2} \\
\langle \sin(2\pi mt/T), \sin(2\pi mt/T) \rangle &= \frac{1}{2} \\
\langle 1, 1 \rangle &= 1,
\end{aligned}$$

which are easily derived in the same way as above. The orthogonal decomposition theorem (1.5) now gives

$$\begin{aligned}
f_N(t) &= \frac{\langle f, 1 \rangle}{\langle 1, 1 \rangle} 1 + \sum_{n=1}^N \frac{\langle f, \cos(2\pi nt/T) \rangle}{\langle \cos(2\pi nt/T), \cos(2\pi nt/T) \rangle} \cos(2\pi nt/T) \\
&\quad + \sum_{n=1}^N \frac{\langle f, \sin(2\pi nt/T) \rangle}{\langle \sin(2\pi nt/T), \sin(2\pi nt/T) \rangle} \sin(2\pi nt/T) \\
&= \frac{\frac{1}{T} \int_0^T f(t) dt}{1} + \sum_{n=1}^N \frac{\frac{1}{T} \int_0^T f(t) \cos(2\pi nt/T) dt}{\frac{1}{2}} \cos(2\pi nt/T) \\
&\quad + \sum_{n=1}^N \frac{\frac{1}{T} \int_0^T f(t) \sin(2\pi nt/T) dt}{\frac{1}{2}} \sin(2\pi nt/T) \\
&= \frac{1}{T} \int_0^T f(t) dt + \sum_{n=1}^N \left(\frac{2}{T} \int_0^T f(t) \cos(2\pi nt/T) dt \right) \cos(2\pi nt/T) \\
&\quad + \sum_{n=1}^N \left(\frac{2}{T} \int_0^T f(t) \sin(2\pi nt/T) dt \right) \sin(2\pi nt/T).
\end{aligned}$$

Equation (1.8)-(1.10) now follow by comparison with Equation (1.7). ■

Since f is a function in time, and the a_n, b_n represent contributions from different frequencies, the Fourier series can be thought of as a change of coordinates,

from what we vaguely can call the *time domain*, to what we can call the *frequency domain* (or *Fourier domain*). We will call the basis $\mathcal{D}_{N,T}$ the *N'th order Fourier basis* for $V_{N,T}$. We note that $\mathcal{D}_{N,T}$ is not an orthonormal basis; it is only orthogonal.

In the signal processing literature, Equation (1.7) is known as *the synthesis equation*, since the original function f is synthesized as a sum of trigonometric functions. Similarly, equations (1.8)-(1.10) are called *analysis equations*.

A major topic in harmonic analysis is to state conditions on f which guarantees the convergence of its Fourier series. We will not discuss this in detail here, since it turns out that, by choosing N large enough, any reasonable periodic function can be approximated arbitrarily well by its N th-order Fourier series approximation. More precisely, we have the following result for the convergence of the Fourier series, stated without proof.

Theorem 1.17 (Convergence of Fourier series). Suppose that f is periodic with period T , and that

1. f has a finite set of discontinuities in each period.
2. f contains a finite set of maxima and minima in each period.
3. $\int_0^T |f(t)| dt < \infty$.

Then we have that $\lim_{N \rightarrow \infty} f_N(t) = f(t)$ for all t , except at those points t where f is not continuous.

The conditions in Theorem 1.17 are called the *Dirichlet conditions* for the convergence of the Fourier series. They are just one example of conditions that ensure the convergence of the Fourier series. There also exist much more general conditions that secure convergence. These can require deep mathematical theory in order to prove, depending on the generality.

An illustration of Theorem 1.17 is shown in Figure 1.5 where the cubic polynomial $f(x) = -\frac{1}{3}x^3 + \frac{1}{2}x^2 - \frac{3}{16}x + 1$ is approximated by a 9th order Fourier series. The trigonometric approximation is periodic with period 1 so the approximation becomes poor at the ends of the interval since the cubic polynomial is not periodic. The approximation is plotted on a larger interval in Figure 1.5(b), where its periodicity is clearly visible.

What you should have learnt in this section

The inner product which we use for function spaces. Definition of the Fourier spaces, and the orthogonality of the Fourier basis. Fourier series approximations as best approximations. Formulas for the Fourier coefficients.

Exercises for Section 1.3

1. Find a function f which is Riemann-integrable on $[0, T]$, and so that $\int_0^T f(t)^2 dt$ is infinite.

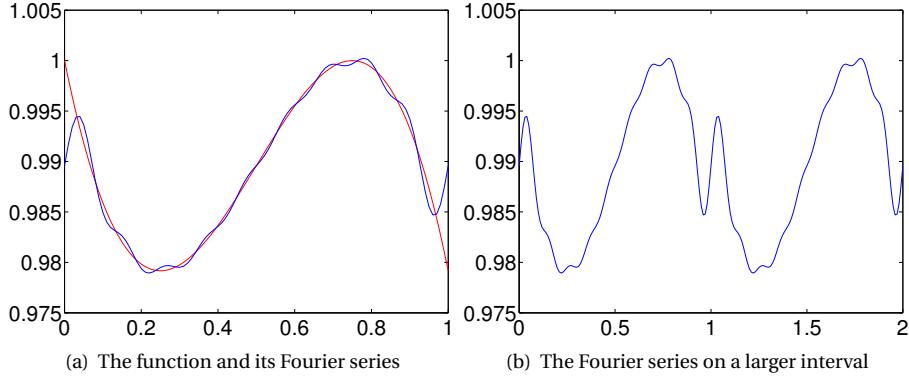


Figure 1.5: The cubic polynomial $f(x) = -\frac{1}{3}x^3 + \frac{1}{2}x^2 - \frac{3}{16}x + 1$ on the interval $[0, 1]$, together with its Fourier series approximation from $V_{9,1}$.

2. Given the two Fourier spaces $V_{N_1, T_1}, V_{N_2, T_2}$. Find necessary and sufficient conditions in order for $V_{N_1, T_1} \subset V_{N_2, T_2}$.

1.4 Examples of Fourier series

Let us compute the Fourier series of some interesting functions.

Example 1.18. Let us compute the Fourier coefficients of the square wave, as defined by Equation (1.1) in Example 1.11. If we first use Equation (1.8) we obtain

$$a_0 = \frac{1}{T} \int_0^T f_s(t) dt = \frac{1}{T} \int_0^{T/2} dt - \frac{1}{T} \int_{T/2}^T dt = 0.$$

Using Equation (1.9) we get

$$\begin{aligned} a_n &= \frac{2}{T} \int_0^T f_s(t) \cos(2\pi nt/T) dt \\ &= \frac{2}{T} \int_0^{T/2} \cos(2\pi nt/T) dt - \frac{2}{T} \int_{T/2}^T \cos(2\pi nt/T) dt \\ &= \frac{2}{T} \left[\frac{T}{2\pi n} \sin(2\pi nt/T) \right]_0^{T/2} - \frac{2}{T} \left[\frac{T}{2\pi n} \sin(2\pi nt/T) \right]_{T/2}^T \\ &= \frac{2}{T} \frac{T}{2\pi n} ((\sin(n\pi) - \sin 0) - (\sin(2n\pi) - \sin(n\pi))) = 0. \end{aligned}$$

Finally, using Equation (1.10) we obtain

$$\begin{aligned}
b_n &= \frac{2}{T} \int_0^T f_s(t) \sin(2\pi nt/T) dt \\
&= \frac{2}{T} \int_0^{T/2} \sin(2\pi nt/T) dt - \frac{2}{T} \int_{T/2}^T \sin(2\pi nt/T) dt \\
&= \frac{2}{T} \left[-\frac{T}{2\pi n} \cos(2\pi nt/T) \right]_0^{T/2} + \frac{2}{T} \left[\frac{T}{2\pi n} \cos(2\pi nt/T) \right]_{T/2}^T \\
&= \frac{2}{T} \frac{T}{2\pi n} ((-\cos(n\pi) + \cos 0) + (\cos(2n\pi) - \cos(n\pi))) \\
&= \frac{2(1 - \cos(n\pi))}{n\pi} \\
&= \begin{cases} 0, & \text{if } n \text{ is even;} \\ 4/(n\pi), & \text{if } n \text{ is odd.} \end{cases}
\end{aligned}$$

In other words, only the b_n -coefficients with n odd in the Fourier series are nonzero. The Fourier series of the square wave is thus

$$\frac{4}{\pi} \sin(2\pi t/T) + \frac{4}{3\pi} \sin(2\pi 3t/T) + \frac{4}{5\pi} \sin(2\pi 5t/T) + \frac{4}{7\pi} \sin(2\pi 7t/T) + \dots \quad (1.11)$$

With $N = 20$, there are 10 trigonometric terms in this sum. The corresponding Fourier series can be plotted on the same interval with the following code.

```
t=0:(1/fs):3;
y=zeros(1,length(t));
for n=1:2:19
    y = y + (4/(n*pi))*sin(2*pi*n*t/T);
end
plot(t,y)
```

In Figure 1.6(a) we have plotted the Fourier series of the square wave when $T = 1/440$, and when $N = 20$. In Figure 1.6(b) we have also plotted the values of the first 100 Fourier coefficients b_n , to see that they actually converge to zero. This is clearly necessary in order for the Fourier series to converge.

Even though f oscillates regularly between -1 and 1 with period T , the discontinuities mean that it is far from the simple $\sin(2\pi t/T)$ which corresponds to a pure tone of frequency $1/T$. From Figure 1.6(b) we see that the dominant coefficient in the Fourier series is b_1 , which tells us how much there is of the pure tone $\sin(2\pi t/T)$ in the square wave. This is not surprising since the square wave oscillates T times every second as well, but the additional nonzero coefficients pollute the pure sound. As we include more and more of these coefficients, we gradually approach the square wave, as shown for $N = 20$.

There is a connection between how fast the Fourier coefficients go to zero, and how we perceive the sound. A pure sine sound has only one nonzero coefficient, while the square wave Fourier coefficients decrease as $1/n$, making the sound less

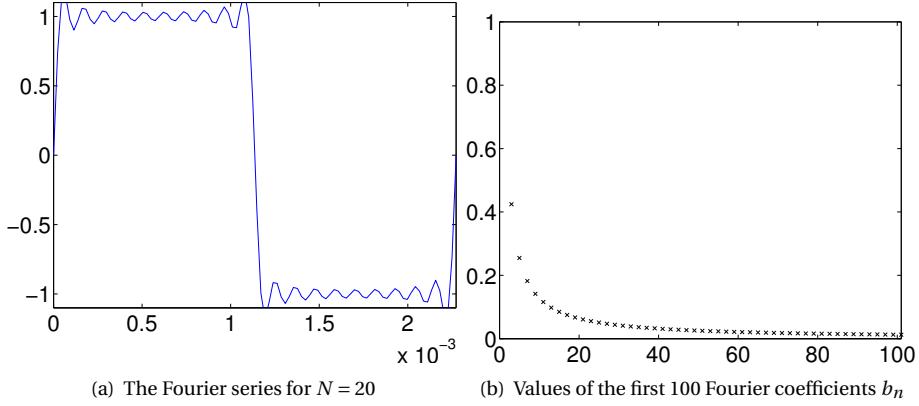


Figure 1.6: The Fourier series of the square wave of Example 1.18

pleasant. This explains what we heard when we listened to the sound in Example 1.11. Also, it explains why we heard the same pitch as the pure tone, since the first frequency in the Fourier series has the same frequency as the pure tone we listened to, and since this had the highest value.

Let us listen to the Fourier series approximations of the square wave. For $N = 1$ and with $T = 1/440$ as above, it sounds like this. This sounds exactly like the pure sound with frequency 440Hz, as noted above. For $N = 5$ the Fourier series approximation sounds like this, and for $N = 9$ it sounds like this. Indeed, these sounds are more like the square wave itself, and as we increase N we can hear how the introduction of more frequencies gradually pollutes the sound more and more. In Exercise 2.1.5 you will be asked to write a program which verifies this. ♣

Example 1.19. Let us also compute the Fourier coefficients of the triangle wave, as defined by Equation (1.2) in Example 1.12. We now have

$$a_0 = \frac{1}{T} \int_0^{T/2} \frac{4}{T} \left(t - \frac{T}{4} \right) dt + \frac{1}{T} \int_{T/2}^T \frac{4}{T} \left(\frac{3T}{4} - t \right) dt.$$

Instead of computing this directly, it is quicker to see geometrically that the graph of f_t has as much area above as below the x -axis, so that this integral must be zero. Similarly, since f_t is symmetric about the midpoint $T/2$, and $\sin(2\pi nt/T)$ is antisymmetric about $T/2$, we have that $f_t(t)\sin(2\pi nt/T)$ also is antisymmetric about $T/2$, so that

$$\int_0^{T/2} f_t(t)\sin(2\pi nt/T) dt = - \int_{T/2}^T f_t(t)\sin(2\pi nt/T) dt.$$

This means that, for $n \geq 1$,

$$b_n = \frac{2}{T} \int_0^{T/2} f_t(t)\sin(2\pi nt/T) dt + \frac{2}{T} \int_{T/2}^T f_t(t)\sin(2\pi nt/T) dt = 0.$$

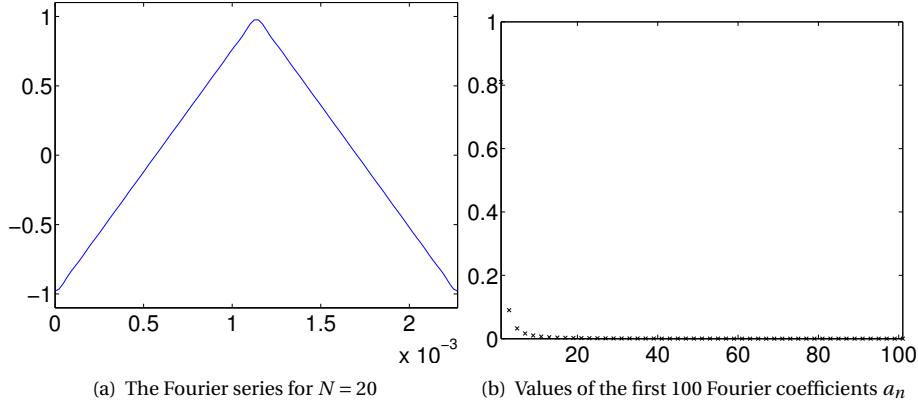


Figure 1.7: The Fourier series of the triangle wave of Example 1.19

For the final coefficients, since both f and $\cos(2\pi nt/T)$ are symmetric about $T/2$, we get for $n \geq 1$,

$$\begin{aligned}
 a_n &= \frac{2}{T} \int_0^{T/2} f_t(t) \cos(2\pi nt/T) dt + \frac{2}{T} \int_{T/2}^T f_t(t) \cos(2\pi nt/T) dt \\
 &= \frac{4}{T} \int_0^{T/2} f_t(t) \cos(2\pi nt/T) dt = \frac{4}{T} \int_0^{T/2} \frac{4}{T} \left(t - \frac{T}{4} \right) \cos(2\pi nt/T) dt \\
 &= \frac{16}{T^2} \int_0^{T/2} t \cos(2\pi nt/T) dt - \frac{4}{T} \int_0^{T/2} \cos(2\pi nt/T) dt \\
 &= \frac{4}{n^2 \pi^2} (\cos(n\pi) - 1) \\
 &= \begin{cases} 0, & \text{if } n \text{ is even;} \\ -8/(n^2 \pi^2), & \text{if } n \text{ is odd.} \end{cases}
 \end{aligned}$$

where we have dropped the final tedious calculations (use integration by parts). From this it is clear that the Fourier series of the triangle wave is

$$-\frac{8}{\pi^2} \cos(2\pi t/T) - \frac{8}{3^2 \pi^2} \cos(2\pi 3t/T) - \frac{8}{5^2 \pi^2} \cos(2\pi 5t/T) - \frac{8}{7^2 \pi^2} \cos(2\pi 7t/T) + \dots \quad (1.12)$$

In Figure 1.7 we have repeated the plots used for the square wave, for the triangle wave. As before, we have used $T = 1/440$. The figure clearly shows that the Fourier series coefficients decay much faster.

Let us also listen to different Fourier series approximations of the triangle wave. For $N = 1$ and with $T = 1/440$ as above, it sounds like this. Again, this sounds exactly like the pure sound with frequency 440Hz. For $N = 5$ the Fourier series approximation sounds like this, and for $N = 9$ it sounds like this. Again these sounds are more like the triangle wave itself, and as we increase N we can hear that the introduction of more frequencies pollutes the sound. However, since the triangle wave Fourier

coefficients decrease as $1/n^2$ instead of $1/n$ as for the square wave, the sound is, although unpleasant due to pollution by many frequencies, not as unpleasant as the square wave. Also, it converges faster to the triangle wave itself, as also can be heard. In Exercise 2.1.5 you will be asked to write a program which verifies this. ♣

There is an important lesson to be learnt from the previous examples: Even if the signal is nice and periodic, it may not have a nice representation in terms of trigonometric functions. Thus, trigonometric functions may not be the best bases to use for expressing other functions. Unfortunately, many more such cases can be found, as the next example shows.

Example 1.20. Let us consider a periodic function which is 1 on $[0, T_0]$, but 0 is on $[T_0, T]$. This is a signal with short duration when T_0 is small compared to T . We compute that $y_0 = T_0/T$, and

$$a_n = \frac{2}{T} \int_0^{T_0} \cos(2\pi nt/T) dt = \frac{1}{\pi n} [\sin(2\pi nt/T)]_0^{T_0} = \frac{\sin(2\pi nT_0/T)}{\pi n}$$

for $n \geq 1$. Similar computations hold for b_n . We see that $|a_n|$ is of the order $1/(\pi n)$, and that infinitely many n contribute. This function may be thought of as a simple building block, corresponding to a small time segment. However, we see that it is not a simple building block in terms of trigonometric functions. This time segment building block may be useful for restricting a function to smaller time segments, and later on we will see that it still can be useful. ♣

1.4.1 Fourier series for symmetric and antisymmetric functions

In Example 1.18 we saw that the Fourier coefficients b_n vanished, resulting in a sine-series for the Fourier series of the square wave. Similarly, in Example 1.19 we saw that a_n vanished, resulting in a cosine-series for the triangle wave. This is not a coincidence, and is captured by the following result, since the square wave was defined so that it was antisymmetric about 0, and the triangle wave so that it was symmetric about 0.

Theorem 1.21 (Symmetry and antisymmetry). If f is antisymmetric about 0 (that is, if $f(-t) = -f(t)$ for all t), then $a_n = 0$, so the Fourier series is actually a sine-series. If f is symmetric about 0 (which means that $f(-t) = f(t)$ for all t), then $b_n = 0$, so the Fourier series is actually a cosine-series.

Proof: Note first that we can write

$$a_n = \frac{2}{T} \int_{-T/2}^{T/2} f(t) \cos(2\pi nt/T) dt \quad b_n = \frac{2}{T} \int_{-T/2}^{T/2} f(t) \sin(2\pi nt/T) dt,$$

i.e. we can change the integration bounds from $[0, T]$ to $[-T/2, T/2]$. This follows from the fact that all $f(t)$, $\cos(2\pi nt/T)$ and $\sin(2\pi nt/T)$ are periodic with period T .

Suppose first that f is symmetric. We obtain

$$\begin{aligned} b_n &= \frac{2}{T} \int_{-T/2}^{T/2} f(t) \sin(2\pi nt/T) dt \\ &= \frac{2}{T} \int_{-T/2}^0 f(t) \sin(2\pi nt/T) dt + \frac{2}{T} \int_0^{T/2} f(t) \sin(2\pi nt/T) dt \\ &= \frac{2}{T} \int_{-T/2}^0 f(t) \sin(2\pi nt/T) dt - \frac{2}{T} \int_0^{-T/2} f(-t) \sin(-2\pi nt/T) dt \\ &= \frac{2}{T} \int_{-T/2}^0 f(t) \sin(2\pi nt/T) dt - \frac{2}{T} \int_{-T/2}^0 f(t) \sin(2\pi nt/T) dt = 0. \end{aligned}$$

where we have made the substitution $u = -t$, and used that sin is antisymmetric. The case when f is antisymmetric can be proved in the same way, and is left as an exercise. \blacksquare

In fact, the connection between symmetric and antisymmetric functions, and sine- and cosine series can be made even stronger by observing the following:

1. Any cosine series $a_0 + \sum_{n=1}^N a_n \cos(2\pi nt/T)$ is a symmetric function.
2. Any sine series $\sum_{n=1}^N b_n \sin(2\pi nt/T)$ is an antisymmetric function.
3. Any periodic function can be written as a sum of a symmetric- and an anti-symmetric function by writing

$$f(t) = \frac{f(t) + f(-t)}{2} + \frac{f(t) - f(-t)}{2}. \quad (1.13)$$

4. If $f_N(t) = a_0 + \sum_{n=1}^N (a_n \cos(2\pi nt/T) + b_n \sin(2\pi nt/T))$, then

$$\begin{aligned} \frac{f_N(t) + f_N(-t)}{2} &= a_0 + \sum_{n=1}^N a_n \cos(2\pi nt/T) \\ \frac{f_N(t) - f_N(-t)}{2} &= \sum_{n=1}^N b_n \sin(2\pi nt/T). \end{aligned}$$

What you should have learnt in this section

Using Matlab to plot Fourier series. For symmetric /antisymmetric functions, Fourier series are actually cosine/sine series.

Exercises for Section 1.4

1. Prove the second part of Theorem 1.21, i.e. show that if f is antisymmetric about 0 (i.e. $f(-t) = -f(t)$ for all t), then $a_n = 0$, i.e. the Fourier series is actually a sine-series.
2. Find the Fourier series coefficients of the periodic functions with period T defined by being $f(t) = t$, $f(t) = t^2$, and $f(t) = t^3$, on $[0, T]$.

3. Write down difference equations for finding the Fourier coefficients of $f(t) = t^{k+1}$ from those of $f(t) = t^k$, and write a program which uses this recursion. Use the program to verify what you computed in Exercise 2.
4. Use the previous exercise to find the Fourier series for $f(x) = -\frac{1}{3}x^3 + \frac{1}{2}x^2 - \frac{3}{16}x + 1$ on the interval $[0, 1]$. Plot the 9th order Fourier series for this function. You should obtain the plots from Figure 1.5.

1.5 Complex Fourier series

In Section 1.3 we saw how a function can be expanded in a series of sines and cosines. These functions are related to the complex exponential function via Eulers formula

$$e^{ix} = \cos x + i \sin x$$

where i is the imaginary unit with the property that $i^2 = -1$. Because the algebraic properties of the exponential function are much simpler than those of cos and sin, it is often an advantage to work with complex numbers, even though the given setting is real numbers. This is definitely the case in Fourier analysis. More precisely, we will make the substitutions

$$\cos(2\pi nt/T) = \frac{1}{2} (e^{2\pi int/T} + e^{-2\pi int/T}) \quad (1.14)$$

$$\sin(2\pi nt/T) = \frac{1}{2i} (e^{2\pi int/T} - e^{-2\pi int/T}) \quad (1.15)$$

in Definition 1.15. From these identities it is clear that the set of complex exponential functions $e^{2\pi int/T}$ also is a basis of periodic functions (with the same period) for $V_{N,T}$. We may therefore reformulate Definition 1.15 as follows:

Definition 1.22 (Complex Fourier basis). We define the set of functions

$$\mathcal{F}_{N,T} = \{e^{-2\pi i Nt/T}, e^{-2\pi i (N-1)t/T}, \dots, e^{-2\pi i t/T}, \dots, e^{2\pi i Nt/T}\}, \quad (1.16)$$

$$1, e^{2\pi i t/T}, \dots, e^{2\pi i (N-1)t/T}, e^{2\pi i Nt/T}\}, \quad (1.17)$$

and call this the order N complex Fourier basis for $V_{N,T}$.

The function $e^{2\pi int/T}$ is also called a pure tone with frequency n/T , just as sines and cosines are. We would like to show that these functions also are orthogonal. To show this, we need to say more on the inner product we have defined by Equation (1.3). A weakness with this definition is that we have assumed real functions f and g , so that this can not be used for the complex exponential functions $e^{2\pi int/T}$. For general complex functions we will extend the definition of the inner product as follows:

$$\langle f, g \rangle = \frac{1}{T} \int_0^T f \bar{g} dt. \quad (1.18)$$

The associated norm now becomes

$$\|f\| = \sqrt{\frac{1}{T} \int_0^T |f(t)|^2 dt}. \quad (1.19)$$

The motivation behind Equation 1.18, where we have conjugated the second function, lies in the definition of an *inner product for vector spaces over complex numbers*. From before we are used to vector spaces over real numbers, but vector spaces over complex numbers are defined through the same set of axioms as for real vector spaces, only replacing real numbers with complex numbers. For complex vector spaces, the axioms defining an inner product are the same as for real vector spaces, except for that the axiom

$$\langle f, g \rangle = \langle g, f \rangle \quad (1.20)$$

is replaced with the axiom

$$\langle f, g \rangle = \overline{\langle g, f \rangle}, \quad (1.21)$$

i.e. a conjugation occurs when we switch the order of the functions. This new axiom can be used to prove the property $\langle f, cg \rangle = \bar{c}\langle f, g \rangle$, which is a somewhat different property from what we know for real inner product spaces. This follows by writing

$$\langle f, cg \rangle = \overline{\langle cg, f \rangle} = \overline{c\langle g, f \rangle} = \bar{c}\overline{\langle g, f \rangle} = \bar{c}\langle f, g \rangle.$$

Clearly the inner product given by (1.18) satisfies Axiom (1.21). With this definition it is quite easy to see that the functions $e^{2\pi int/T}$ are orthonormal. Using the orthogonal decomposition theorem we can therefore write

$$\begin{aligned} f_N(t) &= \sum_{n=-N}^N \frac{\langle f, e^{2\pi int/T} \rangle}{\langle e^{2\pi int/T}, e^{2\pi int/T} \rangle} e^{2\pi int/T} = \sum_{n=-N}^N \langle f, e^{2\pi int/T} \rangle e^{2\pi int/T} \\ &= \sum_{n=-N}^N \left(\frac{1}{T} \int_0^T f(t) e^{-2\pi int/T} dt \right) e^{2\pi int/T}. \end{aligned}$$

We summarize this in the following theorem, which is a version of Theorem 1.16 which uses the complex Fourier basis:

Theorem 1.23. We denote by $y_{-N}, \dots, y_0, \dots, y_N$ the coordinates of f_N in the basis $\mathcal{F}_{N,T}$, i.e.

$$f_N(t) = \sum_{n=-N}^N y_n e^{2\pi int/T}. \quad (1.22)$$

The y_n are called the complex Fourier coefficients of f , and they are given by

$$y_n = \langle f, e^{2\pi int/T} \rangle = \frac{1}{T} \int_0^T f(t) e^{-2\pi int/T} dt. \quad (1.23)$$

Let us consider some examples where we compute complex Fourier series.

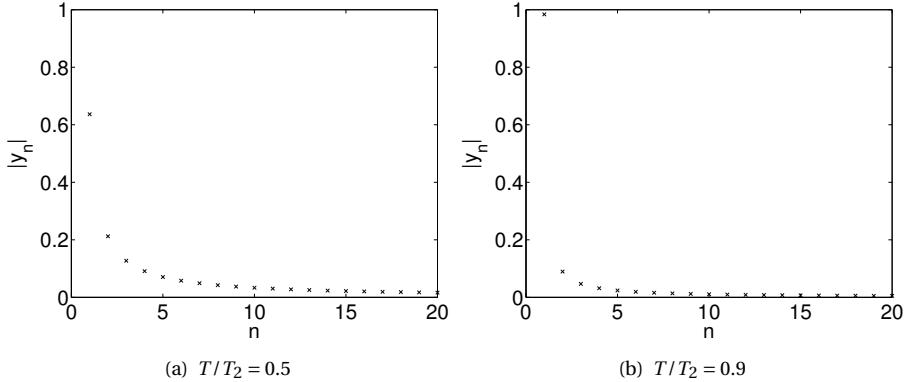


Figure 1.8: Plot of $|y_n|$ when $f(t) = e^{2\pi i t/T_2}$, and $T_2 > T$

Example 1.24. Let us consider the pure sound $f(t) = e^{2\pi i t/T_2}$ with period T_2 , but let us consider it only on the interval $[0, T]$ instead, where $T < T_2$. Note that this f is not periodic, since we only consider the part $[0, T]$ of the period $[0, T_2]$. The Fourier coefficients are

$$\begin{aligned} y_n &= \frac{1}{T} \int_0^T e^{2\pi i t/T_2} e^{-2\pi i nt/T} dt = \frac{1}{2\pi i T(1/T_2 - n/T)} \left[e^{2\pi i t(1/T_2 - n/T)} \right]_0^T \\ &= \frac{1}{2\pi i (T/T_2 - n)} \left(e^{2\pi i T/T_2} - 1 \right). \end{aligned}$$

Here it is only the term $1/(T/T_2 - n)$ which depends on n , so that y_n can only be large when n is close to T/T_2 . In Figure 1.8 we have plotted $|y_n|$ for two different combinations of T, T_2 . In both examples it is seen that many Fourier coefficients contribute, but this is more visible when $T/T_2 = 0.5$. When $T/T_2 = 0.9$, most contribution is seen to be in the y_1 -coefficient. This sounds reasonable, since f then is closest to the pure tone $f(t) = e^{2\pi i t/T}$ of frequency $1/T$ (which in turn has $y_1 = 1$ and all other $y_n = 0$).

3

Apart from computing complex Fourier series, there is an important lesson to be learnt from the previous example: In order for a periodic function to be approximated by other periodic functions, their period must somehow match. Let us consider another example as well.

Example 1.25. What often is the case is that a sound changes in content over time. Assume that it is equal to a pure tone of frequency n_1/T on $[0, T/2)$, and equal to a pure tone of frequency n_2/T on $[T/2, T)$, i.e.

$$f(t) = \begin{cases} e^{2\pi i n_1 t/T} & \text{on } [0, T_2] \\ e^{2\pi i n_2 t/T} & \text{on } [T_2, T) \end{cases}.$$

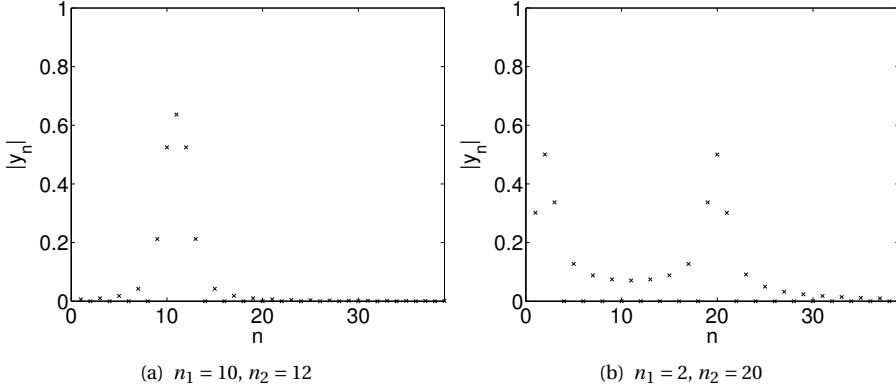


Figure 1.9: Plot of $|y_n|$ when we have two different pure tones at the different parts of a period.

When $n \neq n_1, n_2$ we have that

$$\begin{aligned} y_n &= \frac{1}{T} \left(\int_0^{T/2} e^{2\pi i n_1 t/T} e^{-2\pi i n t/T} dt + \int_{T/2}^T e^{2\pi i n_2 t/T} e^{-2\pi i n t/T} dt \right) \\ &= \frac{1}{T} \left(\left[\frac{T}{2\pi i(n_1 - n)} e^{2\pi i(n_1 - n)t/T} \right]_0^{T/2} + \left[\frac{T}{2\pi i(n_2 - n)} e^{2\pi i(n_2 - n)t/T} \right]_{T/2}^T \right) \\ &= \frac{e^{\pi i(n_1 - n)} - 1}{2\pi i(n_1 - n)} + \frac{1 - e^{\pi i(n_2 - n)}}{2\pi i(n_2 - n)}. \end{aligned}$$

Let us restrict to the case when n_1 and n_2 are both even. We see that

$$y_n = \begin{cases} \frac{1}{2} + \frac{1}{\pi i(n_2 - n_1)} & n = n_1, n_2 \\ 0 & n \text{ even, } n \neq n_1, n_2 \\ \frac{n_1 - n_2}{\pi i(n_1 - n)(n_2 - n)} & n \text{ odd} \end{cases}$$

Here we have computed the cases $n = n_1$ and $n = n_2$ as above. In Figure 1.9 we have plotted $|y_n|$ for two different combinations of n_1, n_2 . We see from the figure that, when n_1, n_2 are close, the Fourier coefficients are close to those of a pure tone with $n \approx n_1, n_2$, but that also other frequencies contribute. When n_1, n_2 are further apart, we see that the Fourier coefficients are like the sum of the two base frequencies, but that other frequencies contribute also here. ♣

There is an important lesson to be learnt from this as well: We should be aware of changes in a sound over time, and it may not be smart to use a frequency representation over a large interval when we know that there are simpler frequency representations on the smaller intervals. The following example shows that, in some cases it is not necessary to compute the Fourier integrals at all, in order to compute the Fourier series.

Example 1.26. Let us compute the complex Fourier series of the function $f(t) = \cos^3(2\pi t/T)$, where T is the period of f . We can write

$$\begin{aligned}\cos^3(2\pi t/T) &= \left(\frac{1}{2}(e^{2\pi it/T} + e^{-2\pi it/T})\right)^3 \\ &= \frac{1}{8}(e^{2\pi i3t/T} + 3e^{2\pi it/T} + 3e^{-2\pi it/T} + e^{-2\pi i3t/T}) \\ &= \frac{1}{8}e^{2\pi i3t/T} + \frac{3}{8}e^{2\pi it/T} + \frac{3}{8}e^{-2\pi it/T} + \frac{1}{8}e^{-2\pi i3t/T}.\end{aligned}$$

From this we see that the complex Fourier series is given by $y_1 = y_{-1} = \frac{3}{8}$, and that $y_3 = y_{-3} = \frac{1}{8}$. In other words, it was not necessary to compute the Fourier integrals in this case, and we see that the function lies in $V_{3,T}$, i.e. there are finitely many terms in the Fourier series. In general, if the function is some trigonometric function, we can often use trigonometric identities to find an expression for the Fourier series. ♣

If we reorder the real and complex Fourier bases so that the two functions $\{\cos(2\pi nt/T), \sin(2\pi nt/T)\}$ and $\{e^{2\pi int/T}, e^{-2\pi int/T}\}$ have the same index in the bases, equations (1.14)-(1.15) give us that the change of coordinates matrix⁴ from $\mathcal{D}_{N,T}$ to $\mathcal{F}_{N,T}$, denoted $P_{\mathcal{F}_{N,T} \leftarrow \mathcal{D}_{N,T}}$, is represented by repeating the matrix

$$\frac{1}{2} \begin{pmatrix} 1 & 1/i \\ 1 & -1/i \end{pmatrix}$$

along the diagonal (with an additional 1 for the constant function 1). In other words, since a_n, b_n are coefficients relative to the real basis and y_n, y_{-n} the corresponding coefficients relative to the complex basis, we have for $n > 0$,

$$\begin{pmatrix} y_n \\ y_{-n} \end{pmatrix} = \frac{1}{2} \begin{pmatrix} 1 & 1/i \\ 1 & -1/i \end{pmatrix} \begin{pmatrix} a_n \\ b_n \end{pmatrix}.$$

This can be summarized by the following theorem:

Theorem 1.27 (Change of coefficients between real and complex Fourier bases).

The complex Fourier coefficients y_n and the real Fourier coefficients a_n, b_n of a function f are related by

$$\begin{aligned}y_0 &= a_0, \\ y_n &= \frac{1}{2}(a_n - ib_n), \\ y_{-n} &= \frac{1}{2}(a_n + ib_n),\end{aligned}$$

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

⁴See Section 4.7 in [20], to review the mathematics behind change of coordinates.

Combining with Theorem 1.21, Theorem 1.27 can help us state properties of complex Fourier coefficients for symmetric- and antisymmetric functions. We look into this in Exercise 8.

Due to the somewhat nicer formulas for the complex Fourier coefficients when compared to the real Fourier coefficients, we will write most Fourier series in complex form in the following.

What you should have learnt in this section

The complex Fourier basis and its orthonormality.

Exercises for Section 1.5

1. Show that the complex functions $e^{2\pi int/T}$ are orthonormal.
2. Compute the complex Fourier series of the function $f(t) = \sin^2(2\pi t/T)$.
3. Repeat Exercise 1.3.2, computing the complex Fourier series instead of the real Fourier series.
4. In this exercise we will find a connection with certain Fourier series and the rows in Pascal's triangle.
 - a. Show that both $\cos^n(t)$ and $\sin^n(t)$ are in $V_{N,2\pi}$ for $1 \leq n \leq N$.
 - b. Write down the N 'th order complex Fourier series for $f_1(t) = \cos t$, $f_2(t) = \cos^2 t$, og $f_3(t) = \cos^3 t$.
 - c. In (b) you should be able to see a connection between the Fourier coefficients and the three first rows in Pascal's triangle. Formulate and prove a general relationship between row n in Pascal's triangle and the Fourier coefficients of $f_n(t) = \cos^n t$.
5. Compute the complex Fourier coefficients of the square wave using Equation 1.23, i.e. repeat the calculations from Example 1.18 for the complex case. Use Theorem 1.27 to verify your result.
6. Repeat Exercise 5 for the triangle wave.
7. Use Equation (1.23) to compute the complex Fourier coefficients of the periodic functions with period T defined by, respectively, $f(t) = t$, $f(t) = t^2$, and $f(t) = t^3$, on $[0, T]$. Use Theorem 1.27 to verify your calculations from Exercise 1.4.2.
8. In this exercise we will prove a version of Theorem 1.21 for complex Fourier coefficients.
 - a. If f is symmetric about 0, show that y_n is real, and that $y_{-n} = y_n$.
 - b. If f is antisymmetric about 0, show that the y_n are purely imaginary, $y_0 = 0$, and that $y_{-n} = -y_n$.
 - c. Show that $\sum_{n=-N}^N y_n e^{2\pi int/T}$ is symmetric when $y_{-n} = y_n$ for all n , and rewrite it as a cosine-series.
 - d. Show that $\sum_{n=-N}^N y_n e^{2\pi int/T}$ is antisymmetric when $y_0 = 0$ and $y_{-n} = -y_n$ for all n , and rewrite it as a sine-series.

1.6 Rate of convergence for Fourier series

We have earlier mentioned criteria which guarantee that the Fourier series converges. Another important topic is the rate of convergence, given that it actually converges. If the series converges quickly, we may only need a few terms in the Fourier series to obtain a reasonable approximation. We have already seen examples which illustrate different convergence rates: The square wave seemed to have very slow convergence rate near the discontinuities, while the triangle wave did not seem to have the same problem.

Before discussing results concerning convergence rates we consider a simple lemma which will turn out to be useful.

Lemma 1.28. Assume that f is differentiable. Then $(f_N)'(t) = (f'_N)(t)$. In other words, the derivative of the Fourier series equals the Fourier series of the derivative.

Proof: We first compute

$$\begin{aligned} \langle f, e^{2\pi i n t/T} \rangle &= \frac{1}{T} \int_0^T f(t) e^{-2\pi i n t/T} dt \\ &= \frac{1}{T} \left(\left[-\frac{T}{2\pi i n} f(t) e^{-2\pi i n t/T} \right]_0^T + \frac{T}{2\pi i n} \int_0^T f'(t) e^{-2\pi i n t/T} dt \right) \\ &= \frac{T}{2\pi i n} \frac{1}{T} \int_0^T f'(t) e^{-2\pi i n t/T} dt = \frac{T}{2\pi i n} \langle f', e^{2\pi i n t/T} \rangle. \end{aligned}$$

where we used integration by parts, and that $-\frac{T}{2\pi i n} f(t) e^{-2\pi i n t/T}$ are periodic with period T . It follows that $\langle f, e^{2\pi i n t/T} \rangle = \frac{T}{2\pi i n} \langle f', e^{2\pi i n t/T} \rangle$. From this we get that

$$\begin{aligned} (f_N)'(t) &= \left(\sum_{n=-N}^N \langle f, e^{2\pi i n t/T} \rangle e^{2\pi i n t/T} \right)' = \frac{2\pi i n}{T} \sum_{n=-N}^N \langle f, e^{2\pi i n t/T} \rangle e^{2\pi i n t/T} \\ &= \sum_{n=-N}^N \langle f', e^{2\pi i n t/T} \rangle e^{2\pi i n t/T} = (f'_N)(t). \end{aligned}$$

where we substituted the connection between the inner products we just found. ■

Example 1.29. The connection between the Fourier series of the function and its derivative can be used to simplify the computation of Fourier series for new functions. Let us see how we can use this to compute the Fourier series of the triangle wave, which was quite a tedious job in Example 1.19. However, the relationship $f'_t(t) = \frac{4}{T} f_s(t)$ is straightforward to see from the plots of the square wave f_s and the triangle wave f_t . From this relationship and from Equation (1.11) for the Fourier series of the square wave it follows that

$$((f_t)')_N(t) = \frac{4}{T} \left(\frac{4}{\pi} \sin(2\pi t/T) + \frac{4}{3\pi} \sin(2\pi 3t/T) + \frac{4}{5\pi} \sin(2\pi 5t/T) + \dots \right).$$

If we integrate this we obtain

$$(f_t)_N(t) = -\frac{8}{\pi^2} \left(\cos(2\pi t/T) + \frac{1}{3^2} \cos(2\pi 3t/T) + \frac{1}{5^2} \cos(2\pi 5t/T) + \dots \right) + C.$$

What remains is to find the integration constant C . This is simplest found if we set $t = T/4$, since then all cosine terms are 0. Clearly then $C = 0$, and we arrive at the same expression as in Equation (1.12) for the Fourier series of the triangle wave. This approach clearly had less computations involved. There is a minor point here which we have not addressed: the triangle wave is not differentiable at two points, as required by Lemma 1.28. It is, however, not too difficult to see that this result still holds in cases where we have a finite number of nondifferentiable points only. ♣

We get the following corollary to Lemma 1.28:

Corollary 1.30. If the complex Fourier coefficients of f are y_n and f is differentiable, then the Fourier coefficients of $f'(t)$ are $\frac{2\pi i n}{T} y_n$.

If we turn this around, we note that the Fourier coefficients of $f(t)$ are $T/(2\pi i n)$ times those of $f'(t)$. If f is s times differentiable, we can repeat this argument to show that the Fourier coefficients of $f(t)$ are $(T/(2\pi i n))^s$ times those of $f^{(s)}(t)$. In other words, the Fourier coefficients of a function which is many times differentiable decay to zero very fast.

Observation 1.31. The Fourier series converges quickly when the function is many times differentiable.

An illustration is found in examples 1.18 and 1.19, where we saw that the Fourier series coefficients for the triangle wave converged more quickly to zero than those of the square wave. This is explained by the fact that the square wave is discontinuous, while the triangle wave is continuous with a discontinuous first derivative. Also, the functions considered in examples 1.24 and 1.25 are not continuous, which partially explain why we there saw contributions from many frequencies.

The requirement of continuity in order to obtain quickly converging Fourier series may seem like a small problem. However, often the function is not defined on the whole real line: it is often only defined on the interval $[0, T]$. If we extend this to a periodic function on the whole real line, by repeating one period as shown in Figure 1.10(a), there is no reason why the new function should be continuous at the boundaries $0, T, 2T$ etc., even though the function we started with may be continuous on $[0, T]$. This would require that $f(0) = \lim_{t \rightarrow T} f(t)$. If this does not hold, the function may not be well approximated with trigonometric functions, due to a slowly convergence Fourier series. We can therefore ask ourselves the following question:

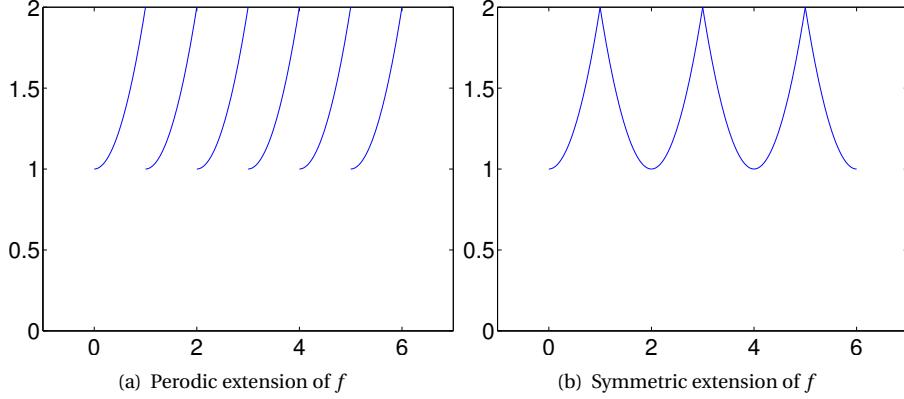


Figure 1.10: Two different extensions of f to a periodic function on the whole real line

Idea 1.32. Assume that f is continuous on $[0, T]$. Can we construct another periodic function which agrees with f on $[0, T]$, and which is both continuous and periodic (maybe with period different from T)?

If this is possible the Fourier series of the new function could produce better approximations for f . It turns out that the following extension strategy does the job:

Definition 1.33 (Symmetric extension of a function). Let f be a function defined on $[0, T]$. By the *symmetric extension* of f , denoted \check{f} , we mean the function defined on $[0, 2T]$ by

$$\check{f}(t) = \begin{cases} f(t), & \text{if } 0 \leq t \leq T; \\ f(2T - t), & \text{if } T < t \leq 2T. \end{cases}$$

Clearly the following holds:

Theorem 1.34. If f is continuous on $[0, T]$, then \check{f} is continuous on $[0, 2T]$, and $\check{f}(0) = \check{f}(2T)$. If we extend \check{f} to a periodic function on the whole real line (which we also will denote by \check{f}), this function is continuous, agrees with f on $[0, T]$, and is a symmetric function.

This also means that the Fourier series of \check{f} is a cosine series, so that it is determined by the cosine-coefficients a_n . The symmetric extension of f is shown in Figure 1.10(b). \check{f} is symmetric since, for $0 \leq t \leq T$,

$$\check{f}(-t) = \check{f}(2T - t) = f(2T - (2T - t)) = f(t) = \check{f}(t).$$

In summary, we now have two possibilities for approximating a function f defined only on $[0, T]$, where the latter addresses a shortcoming of the first:

1. By the Fourier series of f
2. By the Fourier series of \check{f} restricted to $[0, T]$ (which actually is a cosine-series)

Example 1.35. Let f be the function with period T defined by $f(t) = 2t/T - 1$ for $0 \leq t < T$. In each period the function increases linearly from -1 to 1 . Because f is discontinuous at the boundaries, we would expect the Fourier series to converge slowly. Since the function is antisymmetric, the coefficients a_n are zero, and we compute b_n as

$$\begin{aligned} b_n &= \frac{2}{T} \int_0^T \frac{2}{T} \left(t - \frac{T}{2} \right) \sin(2\pi nt/T) dt = \frac{4}{T^2} \int_0^T \left(t - \frac{T}{2} \right) \sin(2\pi nt/T) dt \\ &= \frac{4}{T^2} \int_0^T t \sin(2\pi nt/T) dt - \frac{2}{T} \int_0^T \sin(2\pi nt/T) dt \\ &= -\frac{2}{\pi n}, \end{aligned}$$

so that the Fourier series is

$$-\frac{2}{\pi} \sin(2\pi t/T) - \frac{2}{2\pi} \sin(2\pi 2t/T) - \frac{2}{3\pi} \sin(2\pi 3t/T) - \frac{2}{4\pi} \sin(2\pi 4t/T) - \dots,$$

which indeed converges slowly to 0. Let us now instead consider the symmetric extension of f . Clearly this is the triangle wave with period $2T$, and the Fourier series of this was

$$\begin{aligned} &-\frac{8}{\pi^2} \cos(2\pi t/(2T)) - \frac{8}{3^2\pi^2} \cos(2\pi 3t/(2T)) - \frac{8}{5^2\pi^2} \cos(2\pi 5t/(2T)) \\ &- \frac{8}{7^2\pi^2} \cos(2\pi 7t/(2T)) + \dots. \end{aligned}$$

Comparing the two series, we see that the coefficient at frequency n/T in the first series has value $-2/(n\pi)$, while in the second series it has value

$$-\frac{8}{(2n)^2\pi^2} = -\frac{2}{n^2\pi^2}.$$

The second series clearly converges faster than the first.

If we use $T = 1/880$, the symmetric extension will be the triangle wave of Example 1.19. Its Fourier series for $N = 10$ is shown in Figure 1.7(b) and the Fourier series for $N = 20$ is shown in Figure 1.11. The value $N = 10$ is used since this corresponds to the same frequencies as the previous figure for $N = 20$. It is clear from the plot that the Fourier series for f itself is not a very good approximation. However, we cannot differentiate between the Fourier series and the function itself for the triangle wave.



What you should have learnt in this section

The convergence rate of a Fourier series depends on the regularity of the function. How this motivates the symmetric extension of a function.

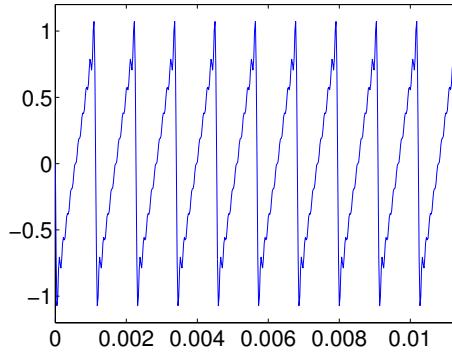


Figure 1.11: The Fourier series for $N = 10$ for the function in Example 1.35

Exercises for Section 1.6

1. Show that the complex Fourier coefficients y_n of f , and the cosine-coefficients a_n of \check{f} are related by $a_{2n} = y_n + y_{-n}$. This result is not enough to obtain the entire Fourier series of \check{f} , but at least it gives us half of it.

1.7 Some properties of Fourier series

We continue by establishing some important properties of Fourier series, in particular the Fourier coefficients for some important functions. In these lists, we will use the notation $f \rightarrow y_n$ to indicate that y_n is the n 'th Fourier coefficient of $f(t)$.

Theorem 1.36 (Fourier series pairs). The functions 1 , $e^{2\pi i nt/T}$, and $\chi_{-a,a}$ have the Fourier coefficients

$$\begin{aligned} 1 \rightarrow \mathbf{e}_0 &= (1, 0, 0, 0, \dots) \\ e^{2\pi i nt/T} \rightarrow \mathbf{e}_n &= (0, 0, \dots, 1, 0, 0, \dots) \\ \chi_{-a,a} \rightarrow &\frac{\sin(2\pi n a/T)}{\pi n}. \end{aligned}$$

The 1 in \mathbf{e}_n is at position n and the function $\chi_{-a,a}$ is the characteristic function of the interval $[-a, a]$, defined by

$$\chi_{-a,a}(t) = \begin{cases} 1, & \text{if } t \in [-a, a]; \\ 0, & \text{otherwise.} \end{cases}$$

The first two pairs are easily verified, so the proofs are omitted. The case for $\chi_{-a,a}$ is very similar to the square wave, but easier to prove, and therefore also omitted.

Theorem 1.37 (Fourier series properties). The mapping $f \rightarrow y_n$ is linear: if $f \rightarrow x_n, g \rightarrow y_n$, then

$$af + bg \rightarrow ax_n + by_n$$

For all n . Moreover, if f is real and periodic with period T , the following properties hold:

1. $y_n = \overline{y_{-n}}$ for all n .
2. If $g(t) = f(-t)$ and $f \rightarrow y_n$, then $g \rightarrow \overline{y_n}$. In particular,
 - (a) if $f(t) = f(-t)$ (i.e. f is symmetric), then all y_n are real, so that b_n are zero and the Fourier series is a cosine series.
 - (b) if $f(t) = -f(-t)$ (i.e. f is antisymmetric), then all y_n are purely imaginary, so that the a_n are zero and the Fourier series is a sine series.
3. If $g(t) = f(t - d)$ (i.e. g is the function f delayed by d) and $f \rightarrow y_n$, then $g \rightarrow e^{-2\pi i n d / T} y_n$.
4. If $g(t) = e^{2\pi i d t / T} f(t)$ with d an integer, and $f \rightarrow y_n$, then $g \rightarrow y_{n-d}$.
5. Let d be a number. If $f \rightarrow y_n$, then $f(d + t) = f(d - t)$ for all t if and only if the argument of y_n is $-2\pi n d / T$ for all n .

Proof: The proof of linearity is left to the reader. Property 1 follows immediately by writing

$$\begin{aligned} y_n &= \frac{1}{T} \int_0^T f(t) e^{-2\pi i n t / T} dt = \overline{\frac{1}{T} \int_0^T f(t) e^{2\pi i n t / T} dt} \\ &= \overline{\frac{1}{T} \int_0^T f(t) e^{-2\pi i (-n)t / T} dt} = \overline{y_{-n}}. \end{aligned}$$

Also, if $g(t) = f(-t)$, we have that

$$\begin{aligned} \frac{1}{T} \int_0^T g(t) e^{-2\pi i n t / T} dt &= \frac{1}{T} \int_0^T f(-t) e^{-2\pi i n t / T} dt = -\frac{1}{T} \int_0^{-T} f(t) e^{2\pi i n t / T} dt \\ &= \frac{1}{T} \int_0^T f(t) e^{2\pi i n t / T} dt = \overline{y_n}. \end{aligned}$$

Property 2 follows from this, since the remaining statements here were established in Theorems 1.21, 1.27, and Exercise 1.5.8. To prove property 3, we observe that the

Fourier coefficients of $g(t) = f(t - d)$ are

$$\begin{aligned} \frac{1}{T} \int_0^T g(t) e^{-2\pi i n t / T} dt &= \frac{1}{T} \int_0^T f(t - d) e^{-2\pi i n t / T} dt \\ &= \frac{1}{T} \int_0^T f(t) e^{-2\pi i n (t+d) / T} dt \\ &= e^{-2\pi i n d / T} \frac{1}{T} \int_0^T f(t) e^{-2\pi i n t / T} dt = e^{-2\pi i n d / T} y_n. \end{aligned}$$

For property 4 we observe that the Fourier coefficients of $g(t) = e^{2\pi i d t / T} f(t)$ are

$$\begin{aligned} \frac{1}{T} \int_0^T g(t) e^{-2\pi i n t / T} dt &= \frac{1}{T} \int_0^T e^{2\pi i d t / T} f(t) e^{-2\pi i n t / T} dt \\ &= \frac{1}{T} \int_0^T f(t) e^{-2\pi i (n-d)t / T} dt = y_{n-d}. \end{aligned}$$

If $f(d+t) = f(d-t)$ for all t , we define the function $g(t) = f(t+d)$ which is symmetric about 0, so that it has real Fourier coefficients. But then the Fourier coefficients of $f(t) = g(t-d)$ are $e^{-2\pi i n d / T}$ times the (real) Fourier coefficients of g by property 3. It follows that y_n , the Fourier coefficients of f , has argument $-2\pi n d / T$. The proof in the other direction follows by noting that any function where the Fourier coefficients are real must be symmetric about 0, once the Fourier series is known to converge. This proves property 5. \blacksquare

Let us analyze these properties, to see that they match the notion we already have for frequencies and sound. We will say that two sounds “essentially are the same” if the absolute values of each Fourier coefficient are equal. Note that this does not mean that the sounds sound the same, it merely says that the contributions at different frequencies are comparable.

The first property says that the positive and negative frequencies in a (real) sound essentially are the same. The second says that, when we play a sound backwards, the frequency content is essentially the same. This is certainly the case for all pure sounds. The third property says that, if we delay a sound, the frequency content also is essentially the same. This also matches our intuition on sound, since we think of the frequency representation as something which is time-independent. The fourth property says that, if we multiply a sound with a pure tone, the frequency representation is shifted (delayed), according to the value of the frequency. This is something we see in early models for the transmission of audio, where an audio signal is transmitted after having been multiplied with what is called a *carrier wave*. You can think of the carrier signal as a pure tone. The result is a signal where the frequencies have been shifted with the frequency of the carrier wave. The point of shifting the frequency of the transmitted signal is to make it use a frequency range in which one knows that other signals do not interfere. The last property looks a bit mysterious. We will not have use for this property before the next chapter.

From Theorem 1.37 we also see that there exist several cases of duality between a function and its Fourier series:

1. Delaying a function corresponds to multiplying the Fourier coefficients with a complex exponential. Vice versa, multiplying a function with a complex exponential corresponds to delaying the Fourier coefficients.
2. Symmetry/antisymmetry for a function corresponds to the Fourier coefficients being real/purely imaginary. Vice versa, a function which is real has Fourier coefficients which are conjugate symmetric.

Actually, one can show that these dualities are even stronger if we had considered Fourier series of complex functions instead of real functions. We will not go into this.

What you should have learnt in this section

Some simple Fourier series pairs. Certain properties of Fourier series, for instance how delay of a function or multiplication with a complex exponential affect the Fourier coefficients.

Exercises for Section 1.7

1. Define the function f with period T on $[-T/2, T/2]$ by

$$f(t) = \begin{cases} 1, & \text{if } -T/4 \leq t < T/4; \\ -1, & \text{if } |T/4| \leq t < |T/2|. \end{cases}$$

f is just the square wave, shifted with $T/4$. Compute the Fourier coefficients of f directly, and use Property 3 in Theorem 1.37 to verify your result.

2. Find a function f which has the complex Fourier series

$$\sum_{n \text{ odd}} \frac{4}{\pi(n+4)} e^{2\pi i nt/T}.$$

Hint: Attempt to use one of the properties in Theorem 1.37 on the Fourier series of the square wave.

1.8 Operations on sound: filters

It is easy to see how we can use Fourier coefficients to analyse or improve sound: Noise in a sound often corresponds to the presence of some high frequencies with large coefficients, and by removing these, we remove the noise. For example, we could set all the coefficients except the first one to zero. This would change the unpleasant square wave to the pure tone $\sin(2\pi 440t)$, which we started our experiments with. Doing so is an example of an important operation on sound called *filtering*.

Definition 1.38 (Analog filters). An operation on sound is called an *analog filter* if it preserves the different frequencies in the sound. In other words, s is an analog filter if, for any sound $f = \sum_v c(v)e^{2\pi i v t}$, the output $s(f)$ is a sound which can be written on the form

$$s(f) = s\left(\sum_v c(v)e^{2\pi i v t}\right) = \sum_v c(v)\lambda_s(v)e^{2\pi i v t},$$

where $\lambda_s(v)$ is a function describing how s treats the different frequencies. $\lambda_s(v)$ uniquely determines s , and is also called the *frequency response* of s .

The following is clear:

Theorem 1.39. The following hold for an analog filter s :

1. When f is periodic with period T , $s(f)$ is also periodic with period T .
2. When $s(f)$ we have that $(s(f))_N = s(f_N)$, i.e. s maps the N 'th order Fourier series of f to the N 'th order Fourier series of $s(f)$.
3. Any pure tone is an eigenvector of s .

The analog filters we will look at have the following form:

Theorem 1.40. Assume that $g \in L^1(\mathbb{R})$. The operation

$$f(t) \rightarrow h(t) = \int_{-\infty}^{\infty} g(s)f(t-s)ds. \quad (1.24)$$

is an analog filter. Analog filters which can be expressed like this are also called *convolutions*. Also

1. When $f \in L^2(\mathbb{R})$, then $h \in L^2(\mathbb{R})$.
2. The frequency response of the filter is $\lambda_s(v) = \int_{-\infty}^{\infty} g(s)e^{-2\pi i vs}ds$

The function g is also called a *convolution kernel*. We also write s_g for the analog filter with convolution kernel g .

The name convolution kernel comes from the fact that filtering operations are also called convolution operations in the literature. In the analog filters we will look at later, the convolution kernel will always have finite support, so that the filter takes the form $f(t) \rightarrow h(t) = \int_a^b g(s)f(t-s)ds$ for some constants a, b . Also note that the integral above may not exist, so that one needs to put some restrictions on the functions, such that $f \in L^2(\mathbb{R})$. Note also that all analog filters may not be expressed as convolutions.

Proof: We compute

$$s(e^{2\pi i \nu t}) = \int_{-\infty}^{\infty} g(s) e^{2\pi i \nu(t-s)} ds = \int_{-\infty}^{\infty} g(s) e^{-2\pi i \nu s} ds e^{2\pi i \nu t} = \lambda_s(f) e^{2\pi i \nu t},$$

which shows that s is a filter with the stated frequency response. That $h \in L^2(\mathbb{R})$, when $f \in L^2(\mathbb{R})$ follows from Minkowski's inequality for integrals [12]. \blacksquare

The function g is arbitrary, so that this strategy leads to a wide class of analog filters. We may ask the question of whether the general analog filter always has this form. We will not go further into this, although one can find partially affirmative answers to this question.

We also need to say something about the connection between filters and symmetric functions. We saw that the symmetric extension of a function took the form of a cosine-series, and that this converged faster to the symmetric extension than the Fourier series did to the function. If a filter preserves cosine-series it will also preserve symmetric extensions, and therefore also map fast-converging Fourier series to fast-converging Fourier series. The following result will be useful in this respect:

Theorem 1.41. If the frequency response of a filter satisfies $\lambda_s(\nu) = \lambda_s(-\nu)$ for all frequencies ν , then the filter preserves cosine series and sine series.

Proof: We have that

$$\begin{aligned} s(\cos(2\pi n t / T)) &= s\left(\frac{1}{2}(e^{2\pi i n t / T} + e^{-2\pi i n t / T})\right) \\ &= \frac{1}{2}\lambda_s(n/T)e^{2\pi i n t / T} + \frac{1}{2}\lambda_s(-n/T)e^{-2\pi i n t / T} \\ &= \lambda_s(n/T)\left(\frac{1}{2}(e^{2\pi i n t / T} + e^{-2\pi i n t / T})\right) = \lambda_s(n/T)\cos(2\pi n t / T). \end{aligned}$$

This means that s preserves cosine-series. A similar computation holds for sine-series holds as well. \blacksquare

An analog filter where $\lambda_s(\nu) = \lambda_s(-\nu)$ is also called a *symmetric filter*. As an example, consider the analog filter $s(f_1) = \int_{-a}^a g(s)f_1(t-s)ds$ where g is symmetric around 0 and supported on $[-a, a]$. s is a symmetric filter since

$$\lambda_s(\nu) = \int_{-a}^a g(s)e^{-2\pi i \nu s} ds = \int_{-a}^a g(s)e^{2\pi i \nu s} ds = \lambda_s(-\nu).$$

Filters are much used in practice, but the way we have defined them here makes them not very useful for computation. We will handle the problem of making filters suitable for computation in Chapter 3.

1.9 The MP3 standard

Digital audio first became commonly available when the CD was introduced in the early 1980s. As the storage capacity and processing speeds of computers increased,

it became possible to transfer audio files to computers and both play and manipulate the data, in ways such as in the previous section. However, audio was represented by a large amount of data and an obvious challenge was how to reduce the storage requirements. Lossless coding techniques like Huffman and Lempel-Ziv coding were known and with these kinds of techniques the file size could be reduced to about half of that required by the CD format. However, by allowing the data to be altered a little bit it turned out that it was possible to reduce the file size down to about ten percent of the CD format, without much loss in quality. The MP3 audio format takes advantage of this.

MP3, or more precisely *MPEG-1 Audio Layer 3*, is part of an audio-visual standard called MPEG. MPEG has evolved over the years, from MPEG-1 to MPEG-2, and then to MPEG-4. The data on a DVD disc can be stored with either MPEG-1 or MPEG-2, while the data on a bluray-disc can be stored with either MPEG-2 or MPEG-4. MP3 was developed by Philips, CCETT (Centre commun d'études de télévision et télécommunications), IRT (Institut für Rundfunktechnik) and Fraunhofer Society, and became an international standard in 1991. Virtually all audio software and music players support this format. MP3 is just a sound format. It leaves a substantial amount of freedom in the encoder, so that different encoders can exploit properties of sound in various ways, in order to alter the sound in removing inaudible components therein. As a consequence there are many different MP3 encoders available, of varying quality. In particular, an encoder which works well for higher bit rates (high quality sound) may not work so well for lower bit rates.

With MP3, the sound is split into *frequency bands*, each band corresponding to a particular frequency range. In the simplest model, 32 frequency bands are used. A frequency analysis of the sound, based on what is called a *psycho-acoustic model*, is the basis for further transformation of these bands. The psycho-acoustic model computes the significance of each band for the human perception of the sound. When we hear a sound, there is a mechanical stimulation of the ear drum, and the amount of stimulus is directly related to the size of the sample values of the digital sound. The movement of the ear drum is then converted to electric impulses that travel to the brain where they are perceived as sound. The perception process uses a transformation of the sound so that a steady oscillation in air pressure is perceived as a sound with a fixed frequency. In this process certain kinds of perturbations of the sound are hardly noticed by the brain, and this is exploited in lossy audio compression.

More precisely, when the psycho-acoustic model is applied to the frequency content resulting from our frequency analysis, *scale factors* and *masking thresholds* are assigned for each band. The computed masking thresholds have to do with a phenomenon called *masking*. A simple example of this is that a loud sound will make a simultaneous low sound inaudible. For compression this means that if certain frequencies of a signal are very prominent, most of the other frequencies can be removed, even when they are quite large. If the sounds are below the masking threshold, it is simply omitted by the encoder, since the model says that the sound should be inaudible.

Masking effects are just one example of what is called psycho-acoustic effects, and all such effects can be taken into account in a psycho-acoustic model. Another

obvious such effect regards computing the scale factors: the human auditory system can only perceive frequencies in the range 20 Hz – 20 000 Hz. An obvious way to do compression is therefore to remove frequencies outside this range, although there are indications that these frequencies may influence the listening experience inaudibly. The computed scaling factors tell the encoder about the precision to be used for each frequency band: If the model decides that one band is very important for our perception of the sound, it assigns a big scale factor to it, so that more effort is put into encoding it by the encoder (i.e. it uses more bits to encode this band).

Using appropriate scale factors and masking thresholds provide compression, since bits used to encode the sound are spent on parts important for our perception. Developing a useful psycho-acoustic model requires detailed knowledge of human perception of sound. Different MP3 encoders use different such models, so they may produce very different results, worse or better.

The information remaining after frequency analysis and using a psycho-acoustic model is coded efficiently with (a variant of) Huffman coding. MP3 supports bit rates from 32 to 320 kb/s and the sampling rates 32, 44.1, and 48 kHz. The format also supports variable bit rates (the bit rate varies in different parts of the file). An MP3 encoder also stores metadata about the sound, such as the title of the audio piece, album and artist name and other relevant data.

MP3 too has evolved in the same way as MPEG, from MP1 to MP2, and to MP3, each one more sophisticated than the other, providing better compression. MP3 is not the latest development of audio coding in the MPEG family: AAC (Advanced Audio Coding) is presented as the successor of MP3 by its principal developer, Fraunhofer Society, and can achieve better quality than MP3 at the same bit rate, particularly for bit rates below 192 kb/s. AAC became well known in April 2003 when Apple introduced this format (at 128 kb/s) as the standard format for their iTunes Music Store and iPod music players. AAC is also supported by many other music players, including the most popular mobile phones.

The technologies behind AAC and MP3 are very similar. AAC supports more sample rates (from 8 kHz to 96 kHz) and up to 48 channels. AAC uses the same transformation as MP3, but AAC processes 1 024 samples at a time. AAC also uses much more sophisticated processing of frequencies above 16 kHz and has a number of other enhancements over MP3. AAC, as MP3, uses Huffman coding for efficient coding of the transformed values. Tests seem quite conclusive that AAC is better than MP3 for low bit rates (typically below 192 kb/s), but for higher rates it is not so easy to differentiate between the two formats. As for MP3 (and the other formats mentioned here), the quality of an AAC file depends crucially on the quality of the encoding program.

There are a number of variants of AAC, in particular AAC Low Delay (AAC-LD). This format was designed for use in two-way communication over a network, for example the internet. For this kind of application, the encoding (and decoding) must be fast to avoid delays (a delay of at most 20 ms can be tolerated).

1.10 Summary

We discussed the basic question of what is sound is, and concluded that sound could be modeled as a sum of frequency components. If the function was periodic we could define its Fourier series, which can be thought of as an approximation scheme for periodic functions using finite-dimensional spaces of trigonometric functions. We established the basic properties of Fourier series, and some duality relationships between the function and its Fourier series. We have also computed the Fourier series of the square wave and the triangle wave, and we saw that we could speed up the convergence of the Fourier series by instead considering the symmetric extension of the function.

We also discussed the MP3 standard for compression of sound, and its relation to a psychoacoustic model which describes how the human auditory system perceives sound. There exist a wide variety of documents on this standard. In [24], an overview is given, which, although written in a signal processing friendly language and representing most relevant theory such as for the psychoacoustic model, does not dig into all the details.

we also defined analog filters, which were operations which operate on continuous sound, without any assumption on periodicity. In signal processing literature onde defines the *Continuous-time Fourier transform*, or CTFT. We will not use this concept in this book. We have instead disguised this concept as the frequency response of an analog filter. To be more precise: in the literature, the CTFT of g is nothing but the frequency response of an analog filter with g as convolution kernel.

Chapter 2

Digital sound and Discrete Fourier analysis

In Chapter 1 we saw how a periodic function can be decomposed into a linear combination of sines and cosines, or equivalently, a linear combination of complex exponential functions. This kind of decomposition is, however, not very convenient from a computational point of view. First of all, the coefficients are given by integrals that in most cases cannot be evaluated exactly, so some kind of numerical integration technique needs to be applied. Secondly, functions are defined for all time instances. On computers and various kinds of media players, however, the sound is *digital*, meaning that it is represented by a large number of function values, and not by a function defined for all time instances.

In this chapter our starting point is simply a vector which represents the sound values (rather than the function $f(t)$), and as before we would like to decompose this in terms of linear combinations of vectors built from complex exponentials. As before it turns out that this is simplest when we assume that the values in the vector repeat periodically. Then a vector of finite dimension can be used to represent all sound values, and operations and the computation of (discrete) Fourier series simply amounts to multiplying the vector by a matrix, and there are efficient algorithms for doing this. It turns out that these algorithms can also be used for computing good approximations to the Fourier series in Chapter 1.

2.1 Digital sound

We start by defining what a digital sound is and by establishing some notation and terminology.

Definition 2.1 (Digital sound). A digital sound is a sequence $\mathbf{x} = \{x_i\}_{i=0}^{N-1}$ that corresponds to measurements of the air pressure of a sound f , recorded at a fixed

rate of f_s (the sampling frequency or sampling rate) measurements per second, i.e.,

$$x_k = f(k/f_s), \quad \text{for } k = 0, 1, \dots, N.$$

The measurements are often referred to as samples. The time between successive measurements is called the sampling period and is usually denoted T_s . The length of the vector is usually assumed to be N , and it is indexed from 0 to $N - 1$. If the sound is in stereo there will be two arrays \mathbf{x}_1 and \mathbf{x}_2 , one for each channel. Measuring the sound is also referred to as sampling the sound, or analog to digital (AD) conversion.

Note that this indexing convention for vectors is not standard in mathematics. Note also that Matlab indexes vectors from 1, so algorithms given here must be adjusted accordingly.

In most cases, a digital sound is sampled from an analog (continuous) audio signal. This is usually done with a technique called Pulse Code Modulation (PCM). The audio signal is sampled at regular intervals and the sampled values stored in a suitable number format. Both the sampling frequency, and the accuracy and number format used for storing the samples, may vary for different kinds of audio, and both influence the quality of the resulting sound. For simplicity the quality is often measured by the number of bits per second, i.e., the product of the sampling rate and the number of bits (binary digits) used to store each sample. This is also referred to as the *bit rate*. For the computer to be able to play a digital sound, samples must be stored in a file or in memory on a computer. To do this efficiently, digital sound formats are used. A couple of them are described in the examples below.

Example 2.2. In the classical CD-format the audio signal is sampled 44 100 times per second and the samples stored as 16-bit integers. This works well for music with a reasonably uniform dynamic range, but is problematic when the range varies. Suppose for example that a piece of music has a very loud passage. In this passage the samples will typically make use of almost the full range of integer values, from $-2^{15} - 1$ to 2^{15} . When the music enters a more quiet passage the sample values will necessarily become much smaller and perhaps only vary in the range -1000 to 1000 , say. Since $2^{10} = 1024$ this means that in the quiet passage the music would only be represented with 10-bit samples. This problem can be avoided by using a floating-point format instead, but very few audio formats appear to do this.

The bit rate for CD-quality stereo sound is $44100 \times 2 \times 16 \text{ bits/s} = 1411.2 \text{ kb/s}$. This quality measure is particularly popular for lossy audio formats where the uncompressed audio usually is the same (CD-quality). However, it should be remembered that even two audio files in the same file format and with the same bit rate may be of very different quality because the encoding programs may be of different quality. ♣

Example 2.3. For telephony it is common to sample the sound 8000 times per second and represent each sample value as a 13-bit integer. These integers are then converted to a kind of 8-bit floating-point format with a 4-bit significand. Telephony therefore generates a bit rate of 64 000 bits per second, i.e. 64 kb/s. ♣

Newer formats with higher quality are available. Music is distributed in various formats on DVDs (DVD-video, DVD-audio, Super Audio CD) with sampling rates up to 192 000 and up to 24 bits per sample. These formats also support surround sound (up to seven channels in contrast to the two stereo channels on a CD). In the following we will assume all sound to be digital. Later we will return to how we reconstruct audible sound from digital sound.

2.2 Simple operations on digital sound

Simple operations and computations with digital sound can be done in any programming environment. Let us take a look at how these. From Definition 2.1, digital sound is just an array of sample values $\mathbf{x} = (x_i)_{i=0}^{N-1}$, together with the sample rate f_s . Performing operations on the sound therefore amounts to doing the appropriate computations with the sample values and the sample rate. The most basic operation we can perform on a sound is simply playing it.

2.2.1 Playing a sound

You may already have listened to pure tones, square waves and triangle waves in the last section. The corresponding sound files were generated in a way we will describe shortly, placed in a directory available on the internet, and linked to from these notes. A program on your computer was able to play these files when you clicked on them. Let us take a closer look at the different steps here. You will need these steps in Exercise 3, where you will be asked to implement a Matlab-function which plays a pure sound with a given frequency on your computer.

First we need to know how we can obtain the samples of a pure tone. The following code does this for a pure tone with frequency f , over a period of 3 seconds, and with sampling rate f_s .

```
t=0:(1/fs):3;
sd=sin(2*pi*f*t);
```

Matlab code will be displayed in this way throughout these notes. In order to listen to these sound samples, we need to take a look at the functions built into Matlab for playing sound. We have the two functions

```
playblocking(playerobj)
playblocking(playerobj, [start stop])
```

These simply play the audio segment encapsulated by the object `playerobj`. `playblocking` means that the method playing the sound will block until it has finished playing. We will have use for this functionality later on, since we may play sounds in successive order. With the first function the entire audio segment is played. With the second function the playback starts at sample `start`, and ends at sample `stop`. These functions are just software interfaces to the sound card in your computer. It basically sends the array of sound samples and sample rate to the sound card, which uses

some method for reconstructing the sound to an analog sound signal. This analog signal is then sent to the loudspeakers and we hear the sound.

Fact 2.4. The basic command in a programming environment that handles sound takes as input an array of sound samples x and a sample rate s , and plays the corresponding sound through the computer's loudspeakers.

The mysterious `playerobj` object above can be obtained from the sound samples (represented by a vector S) and the sampling rate (fs) by the function:

```
playerobj=audioplayer(S,fs)
```

The sound samples can have different data types. We will always assume that they are of type `double`. Matlab requires that they have values between -1 and 1 (i.e. these represent the range of numbers which can be played through the sound card of the computer). Also, S can actually be a matrix: Each column in the matrix represents a sound channel. Sounds we generate on our own from a mathematical function (as for the pure tone above) will typically have only one channel, so that S has only one column. If S originates from a stereo sound file, it will have two columns.

You can create S on your own, and set the sampling rate to whatever value you like. However, we can also fill in the sound samples from a sound file. To do this from a file in the wav-format named `filename`, simply write

```
[S,fs]=wavread(filename)
```

The wav-format format was developed by Microsoft and IBM, and is one of the most common file formats for CD-quality audio. It uses a 32-bit integer to specify the file size at the beginning of the file, which means that a WAV-file cannot be larger than 4 GB. In addition to filling in the sound samples in the vector S , this function also returns the sampling rate fs used in the file. The function

```
wavwrite(S,fs,filename)
```

can similarly be used to write the data stored in the vector S to the wav-file by the name `filename`. In the following we will both fill in the vector S on our own by using values from mathematical functions, as well as from a file. As an example of the first, we can listen to and write to a file the pure tone of frequency 440Hz, which we listened to in Section 1.2, with the help of the following code:

```
antsec=3;
fs=40000;
t=linspace(0,antsec,fs*antsec);
S=sin(2*pi*440*t);
playerobj=audioplayer(S,fs);
playblocking(playerobj);
wavwrite(S,fs,'puretone440.wav');
```

The code creates a pure tone which lasts for three seconds (if you want the tone to last longer, you can change the value of the variable `antsec`). We also tell the computer that there are 40000 samples per second. This value is not coincidental, and we will return to this. In fact, the sound file for the pure tone embedded into this document was created in this way! In the same way we can listen to the square wave with the help of the following code:

```
antsec=3;
fs=44100;
samplesperperiod=round(fs/440);
oneperiod=[ones(1,round(samplesperperiod/2)) ...
    -ones(1,round(samplesperperiod/2))];
allsamples=zeros(1,antsec*440*length(oneperiod));
for k=1:(antsec*440)
    allsamples(((k-1)*length(oneperiod)+1):k*length(oneperiod))=oneperiod;
end
playerobj=audioplayer(allsamples,fs);
playblocking(playerobj);
```

The code creates 440 copies of the square wave per second by first computing the number of samples needed for one period when it is known that we should have a total of 40000 samples per second, and then constructing the samples needed for one period. In the same fashion we can listen to the triangle wave simply by replacing the code for generating the samples for one period with the following:

```
oneperiod=[linspace(-1,1,round(samplesperperiod/2)) ...
    linspace(1,-1,round(samplesperperiod/2))];
```

Instead of using the formula for the triangle wave, directly, we have used the function `linspace`.

As an example of how to fill in the sound samples from a file, the code

```
[S fs] = wavread('castanets.wav');
```

reads the file `castanets.wav`, and stores the sound samples in the matrix `S`. In this case there are two sound channels, so there are two columns in `S`. To work with sound from only one channel, we extract the second channel as follows:

```
x=S(:,2);
```

`wavread` returns sound samples with floating point precision. If we have made any changes to the sound samples, we need to secure that they are between -1 and 1 before we play them. If the sound samples are stored in `x`, this can be achieved as follows:

```
x = x / max(abs(x));
```

`x` can now be played just as the signals we constructed from mathematical formulas above.

It may be that some other environment than Matlab gives you the `play` functionality on your computer. Even if no environment on your computer supports such `play`-functionality at all, you may still be able to play the result of your computations if there is support for saving the sound in some standard format like mp3. The resulting file can then be played by the standard audio player on your computer.

Example 2.5 (Changing the sample rate). We can easily play back a sound with a different sample rate than the standard one. If we in the code above instead wrote `fs=80000`, the sound card will assume that the time distance between neighbouring samples is half the time distance in the original. The result is that the sound takes half as long, and the frequency of all tones is doubled. For voices the result is a characteristic Donald Duck-like sound.

Conversely, the sound can be played with half the sample rate by setting `fs=20000`. Then the length of the sound is doubled and all frequencies are halved. This results in low pitch, roaring voices.

Fact 2.6. A digital sound can be played back with a double or half sample rate by replacing

```
playerobj=audioplayer(S,fs);
```

with

```
playerobj=audioplayer(S,2*fs);
```

and

```
playerobj=audioplayer(S,fs/2);
```

respectively.

The sample file `castanets.wav` played at double sampling rate sounds like this, while it sounds like this when it is played with half the sampling rate. ♣

Example 2.7 (Playing the sound backwards). At times a popular game has been to play music backwards to try and find secret messages. In the old days of analog music on vinyl this was not so easy, but with digital sound it is quite simple; we just need to reverse the samples. To do this we just loop through the array and put the last samples first.

Fact 2.8. Let $x = (x_i)_{i=0}^{N-1}$ be the samples of a digital sound. Then the samples $y = (y_i)_{i=0}^{N-1}$ of the reverse sound are given by

$$y_i = x_{N-i-1}, \text{ for } i = 0, 1, \dots, N-1.$$

When we reverse the sound samples with Matlab, we have to reverse the elements in both sound channels. This can be performed as follows

```

sz=size(S,1);
newS=[S(sz:(-1):1,1) S(sz:(-1):1,2)];

```

Performing this on our sample file you generate a sound which sounds like this. ♣

Example 2.9 (Adding noise). To remove noise from recorded sound can be very challenging, but adding noise is simple. There are many kinds of noise, but one kind is easily obtained by adding random numbers to the samples of a sound.

Fact 2.10. Let x be the samples of a digital sound of length N . A new sound y with noise added can be obtained by adding a random number to each sample,

```

y=x+c*(2*rand(1,N)-1);

```

where `rand` is a Matlab function that returns random numbers in the interval $[0, 1]$, and c is a constant (usually smaller than 1) that dampens the noise. The effect of writing $(2*rand(1,N)-1)$ is that random numbers between -1 and 1 are returned instead of random numbers between 0 and 1 .

Adding noise in this way will produce a general hissing noise similar to the noise you hear on the radio when the reception is bad. As before you should add noise to both channels. Note also that the sound samples may be outside $[-1, 1]$ after adding noise, so that you should scale the samples before writing them to file. The factor c is important, if it is too large, the noise will simply drown the signal y : `castanets.wav` with noise added with $c = 0.4$ sounds like this, while with $c = 0.1$ it sounds like this.



In addition to the operations listed above, the most important operations on digital sound are *digital filters*. These are given a separate treatment in Chapter 3.

What you should have learnt in this section

Matlab operations for reading, writing, and listening to sound. Construct sounds from mathematical formulas. Changing the sample rate and adding noise.

Exercises for Section 2.2

- Define the following sound signal

$$f(t) = \begin{cases} 0 & 0 \leq t \leq 4/440 \\ 2\frac{440t-4}{8} \sin(2\pi 440t) & 4/440 \leq t \leq 12/440 \\ 2\sin(2\pi 440t) & 12/440 \leq t \leq 20/440 \end{cases}$$

This corresponds to the sound plotted in Figure 1.1(a), where the sound is un audible in the beginning, and increases linearly in loudness over time with a given frequency until maximum loudness is achieved. Write a Matlab program which generates this sound, and listen to it.

- 2.** Find two constant a and b so that the function $f(t) = a\sin(2\pi 440t) + b\sin(2\pi 4400t)$ resembles the plot from Figure 1.1(b) as closely as possible. Generate the samples of this sound, and listen to it with Matlab.
- 3.** Let us write some code so that we can experiment with different pure sounds

- a.** Write a function

```
function playpuresound(f)
```

which generates the samples over a period of 3 seconds for a pure tone with frequency f , with sampling frequency $f_s = 2.5f$ (we will explain this value later).

- b.** Use the function `playpuresound` to listen to pure sounds of frequency 440Hz and 1500Hz, and verify that they are the same as the sounds you already have listened to in this section.
- c.** How high frequencies are you able to hear with the function `playpuresound`? How low frequencies are you able to hear?

- 4.** Write functions

```
function playsquare(T)
function playtriangle(T)
```

which plays the square wave of Example 1.11 and the triangle wave of Example 1.12, respectively, where T is given by the parameter. In your code, let the samples of the waves be taken at a frequency of 40000 samples per second. Verify that you generate the same sounds as you played in these examples when you set $T = \frac{1}{440}$.

- 5.** Let us write programs so that we can listen to the Fourier approximations of the square wave and the triangle wave.

- a.** Write functions

```
function playsquaretrunk(T,N)
function playtrianglertrunk(T,N)
```

which plays the order N Fourier approximation of the square wave and the triangle wave, respectively, for three seconds. Verify that you can generate the sounds you played in examples 1.18 and 1.19.

- b.** For these Fourier approximations, how high must you choose N for them to be indistinguishable from the square/triangle waves themselves? Also describe how the characteristics of the sound changes when n increases.

- 6.** In this exercise we will experiment as in the first examples of this section.

- a.** Write a function

```
function playdifferentfs()
```

which plays the sound samples of `castanets.wav` with the same sample rate as the original file, then with twice the sample rate, and then half the sample rate. You should start with reading the file into a matrix (as explained in this section). Are the sounds the same as those you heard in Example 2.5?

- b.** Write a function

```
function playreverse()
```

which plays the sound samples of `castanets.wav` backwards. Is the sound the same as the one you heard in Example 2.7?

- c.** Write the new sound samples from b. to a new `wav`-file, as described above, and listen to it with your favourite mediaplayer.

7. In this exercise, we will experiment with adding noise to a signal.

- a.** Write a function

```
function playnoise(c)
```

which plays the sound samples of `castanets.wav` with noise added for the damping constant c as described above. Your code should add noise to both channels of the sound, and scale the sound samples so that they are between -1 and 1 .

- b.** With your program, generate the two sounds played in Example 2.9, and verify that they are the same as those you heard.

- c.** Listen to the sound samples with noise added for different values of c . For which range of c is the noise audible?

2.3 Discrete Fourier analysis: Basic concepts

In this section we will parallel the developments we did for Fourier series, assuming instead that vectors (rather than functions) are involved. As with Fourier series we will assume that the vector is periodic. This means that we can represent it with the values from only the first period. In the following we will only work with these values, but we will remind ourselves from time to time that the values actually come from a periodic vector. As for functions, we will denote the periodic vector as the periodic extension of the finite vector. To illustrate this, we have in Figure 2.1 shown a vector \mathbf{x} in (a), and the periodic extension of \mathbf{x} in (b). At the outset our vectors will have real components, but since we use complex exponentials we must be able to work with complex vectors also. We therefore first need to define the standard inner product and norm for complex vectors.

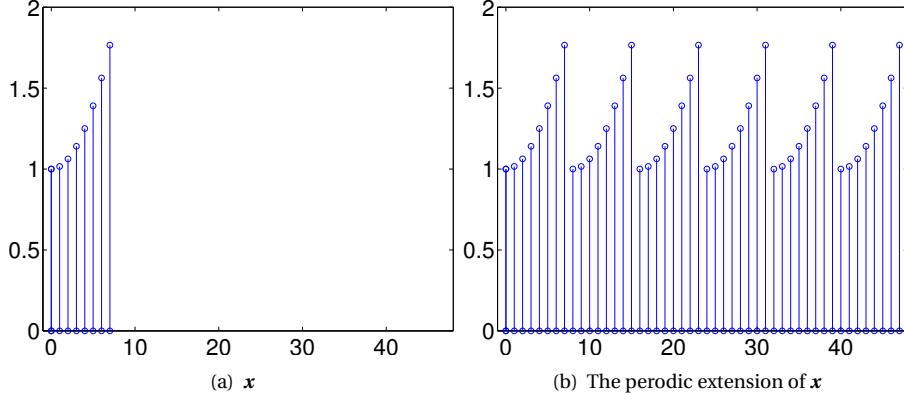


Figure 2.1: A vector and its periodic extension.

Definition 2.11. For complex vectors of length N the Euclidean inner product is given by

$$\langle \mathbf{x}, \mathbf{y} \rangle = \sum_{k=0}^{N-1} x_k \overline{y_k}. \quad (2.1)$$

The associated norm is

$$\|\mathbf{x}\| = \sqrt{\sum_{k=0}^{N-1} |x_k|^2}. \quad (2.2)$$

In the previous chapter we saw that, using a Fourier series, a function with period T could be approximated by linear combinations of the functions (the pure tones) $\{e^{2\pi i n t / T}\}_{n=0}^N$. This can be generalised to vectors (digital sounds), but then the pure tones must of course also be vectors.

Definition 2.12 (Discrete Fourier analysis). In Discrete Fourier analysis, a vector $\mathbf{x} = (x_0, \dots, x_{N-1})$ is represented as a linear combination of the N vectors

$$\boldsymbol{\phi}_n = \frac{1}{\sqrt{N}} \left(1, e^{2\pi i n / N}, e^{2\pi i 2n / N}, \dots, e^{2\pi i kn / N}, \dots, e^{2\pi i n(N-1) / N} \right).$$

These vectors are called the normalised complex exponentials, or the pure digital tones of order N . The whole collection $\mathcal{F}_N = \{\boldsymbol{\phi}_n\}_{n=0}^{N-1}$ is called the N -point Fourier basis.

The following lemma shows that the vectors in the Fourier basis are orthonormal, so they do indeed form a basis.

Lemma 2.13. The normalised complex exponentials $\{\phi_n\}_{n=0}^{N-1}$ of order N form an orthonormal basis in \mathbb{R}^N .

Proof: Let n_1 and n_2 be two distinct integers in the range $[0, N - 1]$. The inner product of ϕ_{n_1} and ϕ_{n_2} is then given by

$$\begin{aligned}\langle \phi_{n_1}, \phi_{n_2} \rangle &= \frac{1}{N} \langle e^{2\pi i n_1 k/N}, e^{2\pi i n_2 k/N} \rangle \\ &= \frac{1}{N} \sum_{k=0}^{N-1} e^{2\pi i n_1 k/N} e^{-2\pi i n_2 k/N} \\ &= \frac{1}{N} \sum_{k=0}^{N-1} e^{2\pi i (n_1 - n_2) k/N} \\ &= \frac{1}{N} \frac{1 - e^{2\pi i (n_1 - n_2)}}{1 - e^{2\pi i (n_1 - n_2)/N}} \\ &= 0.\end{aligned}$$

In particular, this orthogonality means that the the complex exponentials form a basis. Clearly also $\langle \phi_n, \phi_n \rangle = 1$, so that the N -point Fourier basis is in fact an orthonormal basis. ■

Note that the normalising factor $\frac{1}{\sqrt{N}}$ was not present for pure tones in the previous chapter. Also, the normalising factor $\frac{1}{T}$ from the last chapter is not part of the definition of the inner product in this chapter. These are small differences which have to do with slightly different notation for functions and vectors, and which will not cause confusion in what follows.

2.4 The Discrete Fourier Transform

The focus in Discrete Fourier analysis is to change coordinates from the standard basis to the Fourier basis, performing some operations on this “Fourier representation”, and then change coordinates back to the standard basis. Such operations are of crucial importance, and in this section we study some of their basic properties. We start with the following definition.

Definition 2.14 (Discrete Fourier Transform). The change of coordinates from the standard basis of \mathbb{R}^N to the Fourier basis \mathcal{F}_N is called the discrete Fourier transform (or DFT). The $N \times N$ matrix F_N that represents this change of basis is called the (N -point) Fourier matrix. If \mathbf{x} is a vector in \mathbb{R}^N , its coordinates $\mathbf{y} = (y_0, y_1, \dots, y_{N-1})$ relative to the Fourier basis are called the DFT coefficients of \mathbf{x} (in other words $\mathbf{y} = F_N \mathbf{x}$). The DFT of \mathbf{x} is sometimes written as $\hat{\mathbf{x}}$.

We will normally write \mathbf{x} for the given vector in \mathbb{R}^N , and \mathbf{y} for the DFT of this vector. In applied fields, the Fourier basis vectors are also called *synthesis vectors*,

since they can be used to “synthesize” the vector \mathbf{x} , with weights provided by the DFT coefficients $\mathbf{y} = (y_n)_{n=0}^{N-1}$. To be more precise, we have that the change of coordinates performed by the DFT can be written as

$$\mathbf{x} = y_0 \boldsymbol{\phi}_0 + y_1 \boldsymbol{\phi}_1 + \cdots + y_{N-1} \boldsymbol{\phi}_{N-1} = (\boldsymbol{\phi}_0 \quad \boldsymbol{\phi}_1 \cdots \boldsymbol{\phi}_{N-1}) \mathbf{y} = F_N^{-1} \mathbf{y}, \quad (2.3)$$

where we have used the inverse of the defining relation $\mathbf{y} = F_N \mathbf{x}$, and that the $\boldsymbol{\phi}_n$ are the columns in F_N^{-1} (this follows from the fact that F_N^{-1} is the change of coordinates matrix from the Fourier basis to the standard basis, and the Fourier basis vectors are clearly the columns in this matrix). Equation (2.3) is also called the synthesis equation.

Example 2.15. Let \mathbf{x} be the vector of length N defined by $x_k = \cos(2\pi 5k/N)$, and \mathbf{y} the vector of length N defined by $y_k = \sin(2\pi 7k/N)$. Let us see how we can compute $F_N(2\mathbf{x} + 3\mathbf{y})$. By the definition of the DFT as a change of coordinates, $F_N(\boldsymbol{\phi}_n) = \mathbf{e}_n$. We therefore get

$$\begin{aligned} F_N(2\mathbf{x} + 3\mathbf{y}) &= F_N(2 \cos(2\pi 5 \cdot / N) + 3 \sin(2\pi 7 \cdot / N)) \\ &= F_N\left(2 \frac{1}{2}(e^{2\pi i 5 \cdot / N} + e^{-2\pi i 5 \cdot / N}) + 3 \frac{1}{2i}(e^{2\pi i 7 \cdot / N} - e^{-2\pi i 7 \cdot / N})\right) \\ &= F_N\left(\sqrt{N}\boldsymbol{\phi}_5 + \sqrt{N}\boldsymbol{\phi}_{N-5} - \frac{3i}{2}\sqrt{N}(\boldsymbol{\phi}_7 - \boldsymbol{\phi}_{N-7})\right) \\ &= \sqrt{N}(F_N(\boldsymbol{\phi}_5) + F_N(\boldsymbol{\phi}_{N-5}) - \frac{3i}{2}F_N\boldsymbol{\phi}_7 + \frac{3i}{2}F_N\boldsymbol{\phi}_{N-7}) \\ &= \sqrt{N}\mathbf{e}_5 + \sqrt{N}\mathbf{e}_{N-5} - \frac{3i}{2}\sqrt{N}\mathbf{e}_7 + \frac{3i}{2}\sqrt{N}\mathbf{e}_{N-7}. \end{aligned}$$

♣

Let us also find the matrix F_N itself. From Lemma 2.13 we know that the columns of F_N^{-1} are orthonormal. If the matrix was real, it would have been called orthogonal, and the inverse matrix could have been obtained by transposing. F_N^{-1} is complex, however, and it is easy to see that the conjugation present in the definition of the inner product (2.1), implies that the inverse of F_N can be obtained if we also conjugate, in addition to transpose, i.e. $(F_N)^{-1} = (\overline{F_N})^T$. We call $(\overline{A})^T$ the *conjugate transpose* of A , and write A^H for it. We thus have that $(F_N)^{-1} = (F_N)^H$. Matrices which fulfill this are called *unitary*, which thus is the parallel to orthogonal matrices in the world of complex matrices.

Theorem 2.16. The Fourier matrix F_N is the unitary $N \times N$ -matrix with entries given by

$$(F_N)_{nk} = \frac{1}{\sqrt{N}} e^{-2\pi i nk/N},$$

for $0 \leq n, k \leq N-1$.

Note that in the signal processing literature, it is not common to include the normalizing factor $1/\sqrt{N}$ in the definition of the DFT. From our more mathematical point of view this is useful since it makes the Fourier matrix unitary.

Let us now consider the change of coordinate from the Fourier basis back to the standard basis. This operation is also very common, so it also deserves its own definition.

Definition 2.17 (IDFT). If $\mathbf{y} \in \mathbb{R}^N$ the vector $\mathbf{x} = (F_N)^H \mathbf{y}$ is referred to as the inverse discrete Fourier transform or (IDFT) of \mathbf{y} .

That \mathbf{y} is the DFT of \mathbf{x} and \mathbf{x} is the IDFT of \mathbf{y} can also be expressed in component form

$$x_k = \frac{1}{\sqrt{N}} \sum_{n=0}^{N-1} y_n e^{2\pi i n k / N}, \quad (2.4)$$

$$y_n = \frac{1}{\sqrt{N}} \sum_{k=0}^{N-1} x_k e^{-2\pi i n k / N}. \quad (2.5)$$

In applied fields such as signal processing, it is more common to state the DFT and IDFT in these component forms, rather than in the matrix forms $\mathbf{x} = (F_N)^H \mathbf{y}$ and $\mathbf{y} = F_N \mathbf{x}$.

Let us now see how these formulas work out in practice by considering some examples.

Example 2.18 (DFT on a square wave). Let us attempt to apply the DFT to a signal \mathbf{x} which is 1 on indices close to 0, and 0 elsewhere. Assume that

$$x_{-L} = \dots = x_{-1} = x_0 = x_1 = \dots = x_L = 1,$$

while all other values are 0. This is similar to a square wave, with some modifications: First of all we assume symmetry around 0, while the square wave of Example 1.11 assumes antisymmetry around 0. Secondly the values of the square wave are now 0 and 1, contrary to -1 and 1 before. Finally, we have a different proportion of where the two values are assumed. Nevertheless, we will also refer to the current digital sound as a square wave.

Since indices with the DFT are between 0 and $N-1$, and since \mathbf{x} is assumed to have period N , the indices $[-L, L]$ where our signal is 1 translates to the indices $[0, L]$ and $[N-L, N-1]$ (i.e., it is 1 on the first and last parts of the vector). Elsewhere our signal is zero. Since $\sum_{k=N-L}^{N-1} e^{-2\pi i n k / N} = \sum_{k=-L}^{-1} e^{-2\pi i n k / N}$ (since $e^{-2\pi i n k / N}$ is periodic with

period N), the DFT of \mathbf{x} is

$$\begin{aligned}
y_n &= \frac{1}{\sqrt{N}} \sum_{k=0}^L e^{-2\pi i n k / N} + \frac{1}{\sqrt{N}} \sum_{k=N-L}^{N-1} e^{-2\pi i n k / N} \\
&= \frac{1}{\sqrt{N}} \sum_{k=0}^L e^{-2\pi i n k / N} + \frac{1}{\sqrt{N}} \sum_{k=-L}^{-1} e^{-2\pi i n k / N} \\
&= \frac{1}{\sqrt{N}} \sum_{k=-L}^L e^{-2\pi i n k / N} \\
&= \frac{1}{\sqrt{N}} e^{2\pi i n L / N} \frac{1 - e^{-2\pi i n (2L+1) / N}}{1 - e^{-2\pi i n / N}} \\
&= \frac{1}{\sqrt{N}} e^{2\pi i n L / N} e^{-\pi i n (2L+1) / N} e^{\pi i n / N} \frac{e^{\pi i n (2L+1) / N} - e^{-\pi i n (2L+1) / N}}{e^{\pi i n / N} - e^{-\pi i n / N}} \\
&= \frac{1}{\sqrt{N}} \frac{\sin(\pi n (2L+1) / N)}{\sin(\pi n / N)}.
\end{aligned}$$

This computation does in fact also give us the IDFT of the same vector, since the IDFT just requires a change of sign in all the exponents. From this example we see that, in order to represent \mathbf{x} in terms of frequency components, all components are actually needed. The situation would have been easier if only a few frequencies were needed. ♣

Example 2.19. In most cases it is difficult to compute a DFT by hand, due to the entries $e^{-2\pi i n k / N}$ in the matrices, which typically can not be represented exactly. The DFT is therefore usually calculated on a computer only. However, in the case $N = 4$ the calculations are quite simple. In this case the Fourier matrix takes the form

$$F_4 = \frac{1}{2} \begin{pmatrix} 1 & 1 & 1 & 1 \\ 1 & -i & -1 & i \\ 1 & -1 & 1 & -1 \\ 1 & i & -1 & -i \end{pmatrix}.$$

We now can compute the DFT of a vector like $(1, 2, 3, 4)^T$ simply as

$$F_4 \begin{pmatrix} 1 \\ 2 \\ 3 \\ 4 \end{pmatrix} = \frac{1}{2} \begin{pmatrix} 1+2+3+4 \\ 1-2i-3+4i \\ 1-2+3-4 \\ 1+2i-3-4i \end{pmatrix} = \begin{pmatrix} 5 \\ -1+i \\ -1 \\ -1-i \end{pmatrix}.$$

♣

Example 2.20 (Direct implementation of the DFT). Matlab supports complex arithmetic, so the DFT can be implemented very simply and directly by the code

```

function y=DFTImpl(x)
    N=length(x);
    FN=zeros(N);
    for n=1:N

```

```

FN(n,:)=exp(-2*pi*i*(n-1)*(0:(N-1))/N)/sqrt(N);
end
y=FN*x;

```

Note that n has been replaced by $n - 1$ in this code since n runs from 1 to N (array indices must start at 1 in Matlab).

A direct implementation of the IDFT, which we could call `IDFTImpl` can be done similarly. Multiplying a full $N \times N$ matrix by a vector requires roughly N^2 arithmetic operations. The DFT algorithm above will therefore take a long time when N becomes moderately large, particularly in Matlab. It turns out that a much more efficient algorithm exists for computing the DFT, which we will study at the end of this chapter. Matlab also has a built-in implementation of the DFT which uses such an efficient algorithm. ♣

The DFT has properties which are very similar to those of Fourier series, as they were listed in Theorem 1.37. The following theorem sums this up:

Theorem 2.21 (DFT properties). Let \mathbf{x} be a real vector of length N . The DFT has the following properties:

1. $(\hat{\mathbf{x}})_{N-n} = \overline{(\hat{\mathbf{x}})_n}$ for $0 \leq n \leq N - 1$.
2. If \mathbf{z} is the vector with the components of \mathbf{x} reversed so that $z_k = x_{N-k}$ for $0 \leq k \leq N - 1$, then $\hat{\mathbf{z}} = \overline{\hat{\mathbf{x}}}$. In particular,
 - (a) if $x_k = x_{N-k}$ for all n (so \mathbf{x} is symmetric), then $\hat{\mathbf{x}}$ is a real vector.
 - (b) if $x_k = -x_{N-k}$ for all k (so \mathbf{x} is antisymmetric), then $\hat{\mathbf{x}}$ is a purely imaginary vector.
3. If d is an integer and \mathbf{z} is the vector with components $z_k = x_{k-d}$ (the vector \mathbf{x} with its elements delayed by d), then $(\hat{\mathbf{z}})_n = e^{-2\pi i dn/N} (\hat{\mathbf{x}})_n$.
4. If d is an integer and \mathbf{z} is the vector with components $z_k = e^{2\pi i dk/N} x_k$, then $(\hat{\mathbf{z}})_n = (\hat{\mathbf{x}})_{n-d}$.

Proof: The methods used in the proof are very similar to those used in the proof of Theorem 1.37. From the definition of the DFT we have

$$\begin{aligned}
 (\hat{\mathbf{x}})_{N-n} &= \frac{1}{\sqrt{N}} \sum_{k=0}^{N-1} e^{-2\pi i k(N-n)/N} x_k = \frac{1}{\sqrt{N}} \sum_{k=0}^{N-1} e^{2\pi i kn/N} x_k \\
 &= \frac{1}{\sqrt{N}} \sum_{k=0}^{N-1} \overline{e^{-2\pi i kn/N} x_k} = \overline{(\hat{\mathbf{x}})_n}
 \end{aligned}$$

which proves property 1. To prove property 2, we write

$$\begin{aligned}
(\widehat{\mathbf{z}})_n &= \frac{1}{\sqrt{N}} \sum_{k=0}^{N-1} z_k e^{-2\pi i k n / N} = \frac{1}{\sqrt{N}} \sum_{k=0}^{N-1} x_{N-k} e^{-2\pi i k n / N} \\
&= \frac{1}{\sqrt{N}} \sum_{u=1}^N x_u e^{-2\pi i (N-u) n / N} = \frac{1}{\sqrt{N}} \sum_{u=0}^{N-1} x_u e^{2\pi i u n / N} \\
&= \frac{1}{\sqrt{N}} \sum_{u=0}^{N-1} x_u e^{-2\pi i u n / N} = \overline{(\widehat{\mathbf{x}})_n}.
\end{aligned}$$

If \mathbf{x} is symmetric it follows that $\mathbf{z} = \mathbf{x}$, so that $(\widehat{\mathbf{x}})_n = \overline{(\widehat{\mathbf{x}})_n}$. Therefore \mathbf{x} must be real. The case of antisymmetry follows similarly.

To prove property 3 we observe that

$$\begin{aligned}
(\widehat{\mathbf{z}})_n &= \frac{1}{\sqrt{N}} \sum_{k=0}^{N-1} x_{k-d} e^{-2\pi i k n / N} = \frac{1}{\sqrt{N}} \sum_{k=0}^{N-1} x_k e^{-2\pi i (k+d) n / N} \\
&= e^{-2\pi i d n / N} \frac{1}{\sqrt{N}} \sum_{k=0}^{N-1} x_k e^{-2\pi i k n / N} = e^{-2\pi i d n / N} (\widehat{\mathbf{x}})_n.
\end{aligned}$$

For the proof of property 4 we note that the DFT of \mathbf{z} is

$$(\widehat{\mathbf{z}})_n = \frac{1}{\sqrt{N}} \sum_{k=0}^{N-1} e^{2\pi i d k / N} x_n e^{-2\pi i k n / N} = \frac{1}{\sqrt{N}} \sum_{k=0}^{N-1} x_n e^{-2\pi i (n-d) k / N} = (\widehat{\mathbf{x}})_{n-d}.$$

This completes the proof. ■

These properties have similar interpretations as the ones listed in Theorem 1.37 for Fourier series. Property 1 says that we need to store only about one half of the DFT coefficients, since the remaining coefficients can be obtained by conjugation. In particular, when N is even, we only need to store $y_0, y_1, \dots, y_{N/2}$. This also means that, if we plot the (absolute value) of the DFT of a real vector, we will see a symmetry around the index $n = N/2$. The theorem generalizes the properties from Theorem 1.37, except for the last property where the signal had a point of symmetry. We will delay the generalization of this property to later.

Example 2.22. To see how we can use the fourth property of Theorem 2.21, consider a vector $\mathbf{x} = (x_0, x_1, x_2, x_3, x_4, x_5, x_6, x_7)$ with length $N = 8$, and assume that \mathbf{x} is so that $F_8(\mathbf{x}) = (1, 2, 3, 4, 5, 6, 7, 8)$. Consider the vector \mathbf{z} with components $z_k = e^{2\pi i 2k/8} x_k$. Let us compute $F_8(\mathbf{z})$. Since multiplication of \mathbf{x} with $e^{2\pi i kd/N}$ delays the output $\mathbf{y} = F_N(\mathbf{x})$ with d elements, setting $d = 2$, the $F_8(\mathbf{z})$ can be obtained by delaying $F_8(\mathbf{x})$ by two elements, so that $F_8(\mathbf{z}) = (7, 8, 1, 2, 3, 4, 5, 6)$. It is straightforward to compute this directly also:

$$\begin{aligned}
(F_N \mathbf{z})_n &= \sum_{k=0}^{N-1} z_k e^{-2\pi i k n / N} = \sum_{k=0}^{N-1} e^{2\pi i 2k / N} x_k e^{-2\pi i k n / N} \\
&= \sum_{k=0}^{N-1} x_k e^{-2\pi i k (n-2) / N} = (F_N(\mathbf{x}))_{n-2}.
\end{aligned}$$



What you should have learnt in this section

The definition of the Fourier basis and its orthonormality. The definition of the Discrete Fourier Transform, its inverse, and its unitarity. How to apply the DFT to a sum of sinusoids. Properties of the DFT, such as how it treats delayed vectors, and multiplication with a complex exponential.

Exercises for Section 2.4

1. Compute the 4 point DFT of the vector $(2, 3, 4, 5)^T$.
2. As in Example 2.19, state the exact cartesian form of the Fourier matrix for the cases $N = 6$, $N = 8$, and $N = 12$.
3. We have a real vector \mathbf{x} with length N , and define the vector \mathbf{z} by delaying all elements in \mathbf{x} with 5 cyclically, i.e. $z_5 = x_0$, $z_6 = x_1, \dots, z_N = x_{N-5}$, and $z_0 = x_{N-4}, \dots, z_4 = x_{N-1}$. If $|(\mathcal{F}_N \mathbf{x})_n| = 2$, what is then $|(\mathcal{F}_N \mathbf{z})_n|$? Justify the answer.
4. Given a real vector \mathbf{x} of length 8 where $(\mathcal{F}_8(\mathbf{x}))_2 = 2 - i$, what is $(\mathcal{F}_8(\mathbf{x}))_6$?
5. Let \mathbf{x} be the vector of length N where $x_k = \cos^2(2\pi k/N)$. What is then $\mathcal{F}_N(\mathbf{x})$?
6. Let \mathbf{x} be the vector with entries $x_k = c^k$. Show that the DFT of \mathbf{x} is given by the vector with components

$$y_n = \frac{1}{\sqrt{N}} \frac{1 - c^N}{1 - ce^{-2\pi i n/N}}$$

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

7. If \mathbf{x} is complex, Write the DFT in terms of the DFT on real sequences. Hint: Split into real and imaginary parts, and use linearity of the DFT.
8. As in Example 2.20, write a function

```
function x=IDFTImpl(y)
```

which computes the IDFT.

9. Assume that N is even.

- a. Show that, if $x_{k+N/2} = x_k$ for all $0 \leq k < N/2$, then $y_n = 0$ when n is odd.
- b. Show that, if $x_{k+N/2} = -x_k$ for all $0 \leq k < N/2$, then $y_n = 0$ when n is even.
- c. Show also the converse statements in a. and b..
- d. Also show the following:
 1. $x_n = 0$ for all odd n if and only if $y_{k+N/2} = y_k$ for all $0 \leq k < N/2$.
 2. $x_n = 0$ for all even n if and only if $y_{k+N/2} = -y_k$ for all $0 \leq k < N/2$.

10. Let $\mathbf{x}_1, \mathbf{x}_2$ be real vectors, and set $\mathbf{x} = \mathbf{x}_1 + i\mathbf{x}_2$. Use Theorem 2.21 to show that

$$\begin{aligned} (\mathcal{F}_N(\mathbf{x}_1))_k &= \frac{1}{2} \left((\mathcal{F}_N(\mathbf{x}))_k + \overline{(\mathcal{F}_N(\mathbf{x}))_{N-k}} \right) \\ (\mathcal{F}_N(\mathbf{x}_2))_k &= \frac{1}{2i} \left((\mathcal{F}_N(\mathbf{x}))_k - \overline{(\mathcal{F}_N(\mathbf{x}))_{N-k}} \right) \end{aligned}$$

This shows that, sometimes, one DFT can be used to compute two different DFT's.

2.5 Connection between the DFT and Fourier series

So far we have focused on the DFT as a tool to rewrite a vector in terms of the Fourier basis vectors. In practice, the given vector \mathbf{x} will often be sampled from some real data given by a function $f(t)$. We may then compare the frequency content of the vector \mathbf{x} and the frequency content of f , and ask how these are related: What is the relationship between the Fourier coefficients of f and the DFT-coefficients of \mathbf{x} ?

In order to study this, assume for simplicity that f lies in a Fourier space $V_{M,T}$ for some M . This means that f equals its Fourier approximation f_M ,

$$f(t) = f_M(t) = \sum_{n=-M}^M z_n e^{2\pi i n t / T}, \quad (2.6)$$

where

$$z_n = \frac{1}{T} \int_0^T f(t) e^{-2\pi i n t / T} dt.$$

We here have changed our notation for the Fourier coefficients from z_n to y_n , in order not to confuse them with the DFT coefficients. We recall that in order to represent the frequency n/T fully, we need the corresponding exponentials with both positive and negative arguments, i.e., both $e^{2\pi i n t / T}$ and $e^{-2\pi i n t / T}$.

Fact 2.23. Suppose f is given by its Fourier series (2.6). Then the total frequency content for the frequency n/T is given by the two coefficients z_n and z_{-n} .

The following connection between the Fourier coefficients of f and the DFT of the samples of f also states how the DFT can be used to compute a Fourier series.

Proposition 2.24 (Relation between Fourier coefficients and DFT coefficients).

Let

$$f(t) = \sum_{n=-M}^M z_n e^{2\pi i n t / T},$$

be a function in $V_{M,T}$, and let $N = 2M + 1$. Suppose that \mathbf{x} are N uniform samples from f over one period, i.e.

$$x_k = f(kT/N), \quad \text{for } k = 0, 1, \dots, N-1.$$

and let \mathbf{y} be the DFT of \mathbf{x} . Then $\mathbf{z} = (y_{M+1}, \dots, y_{2M}, y_0, \dots, y_M) / \sqrt{N}$. In particular, the total contribution to f from frequency n/T , where n is an integer in the range $0 \leq n \leq M$, is given by y_n and y_{N-n} .

Proof. The vector \mathbf{x} can be expressed in terms of its DFT \mathbf{y} as

$$x_k = \frac{1}{\sqrt{N}} \sum_{n=0}^{N-1} y_n e^{2\pi i n k / N}. \quad (2.7)$$

If we evaluate f at the sample points we have

$$f(kT/N) = \sum_{n=-M}^M z_n e^{2\pi i n k / N}, \quad (2.8)$$

and a comparison now gives

$$\sum_{n=-M}^M z_n e^{2\pi i n k / N} = \frac{1}{\sqrt{N}} \sum_{n=0}^{N-1} y_n e^{2\pi i n k / N} \quad \text{for } k = 0, 1, \dots, N-1.$$

Exploiting the fact that both \mathbf{y} and the complex exponentials are periodic with period N , we can rewrite this as

$$\sum_{n=-M}^M z_n e^{2\pi i n k / N} = \frac{1}{\sqrt{N}} \sum_{n=-M}^{-1} y_{n+2M+1} e^{2\pi i n k / N} + \frac{1}{\sqrt{N}} \sum_{n=0}^M y_n e^{2\pi i n k / N}.$$

This system of N equations can be written on the form

$$G\mathbf{z} = G(y_{M+1}, \dots, y_{2M}, y_0, \dots, y_M) / \sqrt{N},$$

where G is the $N \times N$ -matrix with entries $\frac{1}{\sqrt{N}} e^{2\pi i n k / N}$, $-M \leq n < M$. It is straightforward to show that the matrix G is invertible (we can argue in the same way as we did when we showed that the Fourier basis was orthonormal), so that

$$\mathbf{z} = (y_{M+1}, \dots, y_{2M}, y_0, \dots, y_M) / \sqrt{N}.$$

This proves the result. ■

The proof above bases itself on taking $N = 2M + 1$ samples. However, if we take more samples (including the old samples), say $N = 2K + 1$ samples with $K > M$, it is clear that we will get the Fourier coefficients in $V_{K,T}$ (since $V_{M,T} \subset V_{K,T}$), where the new coefficients are 0. Theorem 2.24 enables us to find (unique) trigonometric functions which interpolate (pass through) a set of data points. We have in elementary calculus courses seen how to determine a polynomial of degree $N - 1$ that interpolates a set of N data points - such polynomials are called interpolating polynomials. Instead of expressing the interpolating polynomial as a Fourier series, it is often expressed with the DFT, as follows:

Corollary 2.25 (Interpolation with the Fourier basis). Let f be a function defined on the interval $[0, T]$, and let \mathbf{x} be the samples

$$x_k = f(kT/N) \quad \text{for } k = 0, 1, \dots, N-1.$$

There is exactly one linear combination $g(t)$ on the form

$$g(t) = \frac{1}{\sqrt{N}} \sum_{n=0}^{N-1} y_n e^{2\pi i n t / T} \quad (2.9)$$

which satisfies the conditions

$$g(kT/N) = f(kT/N), \quad k = 0, 1, \dots, N-1$$

and it is determined by $\mathbf{y} = F_N \mathbf{x}$.

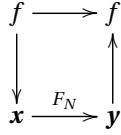


Figure 2.2: How we can interpolate f from $V_{M,T}$ with help of the DFT. The left vertical arrow represents sampling. The right vertical arrow represents interpolation, i.e. computing Equation (2.9).

Proof: Insert $t = 0, t = T/N, t = 2T/N, \dots, t = (N-1)T/N$ in Equation 2.9 to arrive at the equations

$$f(kT/N) = \frac{1}{\sqrt{N}} \sum_{n=0}^{N-1} y_n e^{2\pi i n k / N} \quad 0 \leq k \leq N-1.$$

This gives us an equation system for finding the y_n with the invertible Fourier matrix as coefficient matrix, and the result follows. ■

Corollary 2.25 also enables us to approximate functions in a simple way. To elaborate on this, recall that the Fourier series approximation f_N is a best approximation to f from $V_{M,T}$. We usually can't compute f_N exactly, however, since this requires us to compute the Fourier integrals. We could instead form the samples \mathbf{x} of f , and apply Corollary 2.25 to these. If N is high, f_N is a good approximation to f , so that the samples of f_N is a good approximation to \mathbf{x} . By continuity of the DFT, it follows that $\mathbf{y} = F_N \mathbf{x}$ is a good approximation to the DFT of the samples of f_N , so that $\tilde{f}(t) = \sum_{n=0}^{N-1} y_n e^{2\pi i n t / T}$ is a good approximation to f_N , and therefore also to f . We have illustrated this in Figure 2.2. The new function \tilde{f} has the same values as f in the sample points, but note that it may not be equal to f_N . The interpolating function \tilde{f} is thus another kind of approximation to f , and very useful, since it can be computed without evaluating the Fourier integrals, contrary to the Fourier series:

Idea 2.26. The function \tilde{f} resulting from sampling, taking the DFT, and interpolation, as shown in Figure 2.2, also gives an approximation to f . \tilde{f} is a worse approximation to f than f_N (since f_N is a best approximation), but it is much more useful since it avoids evaluation of the Fourier integrals, and depends only on the samples.

2.6 Applications of the DFT

The DFT can be extremely useful for operations on sound. In this section we will concentrate on two such operations. The first one is adjustment of frequencies in sound, the second is compression of sound. To elaborate on this we first need to establish a connection between the frequency and DFT index. A pure tone of frequency

v has the form $e^{2\pi i v t}$. When we sample this with frequency f_s , we collect values at the time instances $t = kT_s = k/f_s$, so that we get the digital pure sound $e^{2\pi i v k/f_s}$. This corresponds to DFT index n (i.e. the term $e^{2\pi i kn/N}$) if n satisfies $v/f_s = n/N$. Here $N = T f_s$ is the total number of samples in one period.

Observation 2.27 (Connection between DFT index and frequency). Assume that we take N samples of a sound with sampling frequency f_s . Then DFT index n corresponds the frequency $v = nf_s/N$.

Let us in a similar way make an observation on what we should interpret as high and low frequency contributions. Low frequency contribution in a Fourier series is the contribution from

$$e^{-2\pi i Lt/T}, \dots, e^{-2\pi i t/T}, 1, e^{2\pi i t/T}, \dots, e^{2\pi i Lt/T}$$

i.e. $\sum_{n=-L}^L z_n e^{2\pi i nt/T}$ (L represents a limit we set for what to be interpreted as low frequencies). This means that low frequencies correspond to indices n so that $-L \leq n \leq L$. However, since DFT coefficients have indices between 0 and $N - 1$, and $e^{-2\pi i kn/N} = e^{2\pi i k(N-n)/N}$, low frequencies correspond to indices n in $[0, L] \cup [N - L, N - 1]$. This also means that the highest frequencies correspond to DFT indices near $N/2$:

Observation 2.28 (DFT indices for high and low frequencies). When y is the DFT of x , the low frequencies in x correspond to the indices in y near 0 and N . The high frequencies in x correspond to the indices in y near $N/2$.

We will use this observation in the following example, when we use the DFT to distinguish between high and low frequencies in a sound.

Example 2.29 (Using the DFT to adjust frequencies in sound). Since the DFT coefficients represent the contribution in a sound at given frequencies, we can listen to the different frequencies of a sound by adjusting the DFT coefficients. Let us first see how we can listen to the lower frequencies only. As explained, these correspond to DFT-indices n in $[0, L] \cup [N - L, N - 1]$. In Matlab these have indices from 1 to $L + 1$, and from $N - L + 1$ to N . The remaining frequencies, i.e. the higher frequencies which we want to eliminate, thus have Matlab-indices between $L + 2$ and $N - L$. We can now perform a DFT, eliminate these high frequencies, and perform an inverse DFT, to recover the sound signal where these frequencies have been eliminated. With the help of the DFT implementation from Example 2.20, all this can be achieved with the following code:

```
y=DFTImpl(x);
y((L+2):(N-L))=zeros(N-(2*L+1),1);
newx=IDFTImpl(y);
```

To test this in practice, we also need to obtain the actual sound samples. If we use our sample file `castanets.wav`, you will see that the code runs very slowly. In fact

it seems to never complete. The reason is that `DFTImp1` attempts to construct a matrix F_N with as many rows and columns as there are sound samples in the file, and there are just too many samples, so that F_N grows too big, and matrix multiplication with it gets too time-consuming. We will shortly see much better strategies for applying the DFT to a sound file, but for now we will simply attempt instead to split the sound file into smaller blocks, each of size $N = 32$, and perform the code above on each block. This is less time-consuming, since big matrices are avoided. You will be spared the details for actually splitting the sound file into blocks: you can find the function `playDFTlower(L)` which performs this splitting, sets the relevant frequency components to 0, and plays the resulting sound samples. If you try this for $L = 7$ (i.e. we keep only 15 of the DFT coefficients) the result sounds like this. You can hear the disturbance in the sound, but we have not lost that much even if more than half the DFT coefficients are dropped. If we instead try $L = 3$ the result will sound like this. The quality is much poorer now. However we can still recognize the song, and this suggests that most of the frequency information is contained in the lower frequencies.

Similarly we can listen to high frequencies by including only DFT coefficients with index close to $\frac{N}{2}$. The function `playDFThigher(L)` sets all DFT coefficients to zero, except for those with indices $\frac{N}{2}-L, \dots, \frac{N}{2}, \dots, \frac{N}{2}+L$. Let us verify that there is less information in the higher frequencies by trying the same values for L as above for this function. For $L = 7$ (i.e. we keep only the middle 15 DFT coefficients) the result sounds like this, for $L = 3$ the result sounds like this. Both sounds are quite unrecognizable, confirming that most information is contained in the lower frequencies.

♣

Note that there may be a problem in the previous example: when we restrict to the values in a given block, we actually look at a different signal. The new signal repeats the values in the block in periods, while the old signal consists of one much bigger block. What are the differences in the frequency representations of the two signals?

Assume that the entire sound has length M . The frequency representation of this is computed as an M -point DFT (the signal is actually repeated with period M), and we write the sound samples as a sum of frequencies: $x_k = \sum_{n=0}^{M-1} y_n e^{2\pi i kn/M}$. Let us consider the effect of restricting to a block for each of the contributing pure tones $\frac{1}{\sqrt{N}} e^{2\pi i k n_0 / M}$, $0 \leq n_0 \leq M-1$ (where we have scaled with $\frac{1}{\sqrt{N}}$). When we restrict this to a block of size N , we get the signal $\{e^{2\pi i k n_0 / M}\}_{k=0}^{N-1}$. Depending on n_0 , this may not be a Fourier basis vector! Its N -point DFT gives us its frequency representation, and the absolute value of this is

$$\begin{aligned} |y_n| &= \left| \frac{1}{N} \sum_{k=0}^{N-1} e^{2\pi i k n_0 / M} e^{-2\pi i k n / N} \right| = \left| \frac{1}{N} \sum_{k=0}^{N-1} e^{2\pi i k (n_0 / M - n / N)} \right| \\ &= \left| \frac{1}{N} \frac{1 - e^{2\pi i N (n_0 / M - n / N)}}{1 - e^{2\pi i (n_0 / M - n / N)}} \right| = \frac{1}{N} \left| \frac{\sin(\pi N (n_0 / M - n / N))}{\sin(\pi (n_0 / M - n / N))} \right|. \end{aligned} \quad (2.10)$$

If $n_0 = M/N$, this gives $y_{n_0} = 1$, and $y_n = 0$ when $n \neq n_0$. Thus, splitting the signal into blocks gives another pure tone when $n_0 = M/N$. When n_0 is different from

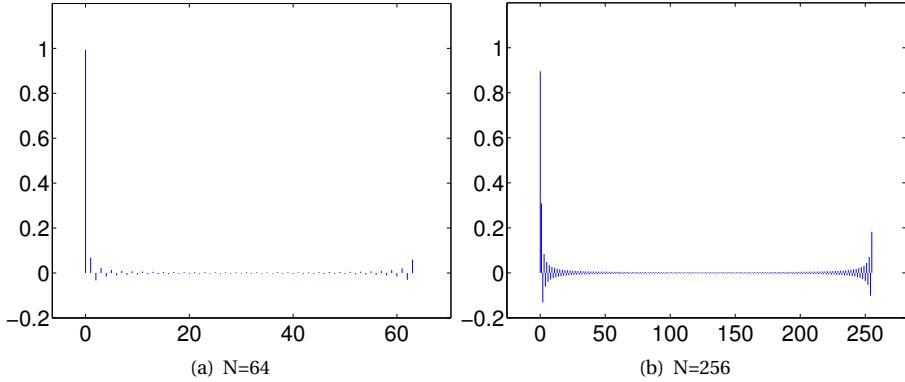


Figure 2.3: The frequency representation obtained when restricting to a block of size N of the signal.

M/N the situation is different. Let us set $M = 1000$, $n_0 = 1$, and experiment with different values of N . Figure 2.3 shows the y_n values for different values of N . We see that the frequency representation is now very different, and that many frequencies contribute. The explanation is that the pure tone is not a pure tone when $N = 64$ and $N = 256$, since at this scale such frequencies are too high to be represented exactly. The closest pure tone in frequency is $n = 0$, and we see that this has the biggest contribution, but other frequencies also contribute. The other frequencies contribute much more when $N = 256$, as can be seen from the peak in the closest frequency $n = 0$. In conclusion, when we split into blocks, the frequency representation may change in an undesirable way. This is a common problem in signal processing theory, that one in practice needs to restrict to smaller segments of samples, but that this restriction may have undesired effects.

Another problem when we restrict to a shorter periodic signal is that we may obtain discontinuities at the boundaries between the new periods, even if there were no discontinuities in the original signal. And, as we know from the square wave, discontinuities introduce undesired frequencies. We have already mentioned that symmetric extensions may be used to remedy this.

The MP3 standard also applies a DFT to the sound data. In its simplest form it applies a 512 point DFT. There are some differences to how this is done when compared to Example 2.29, however. In our example we split the sound into disjoint blocks, and applied a DFT to each of them. The MP3 standard actually splits the sound into blocks which overlap, as this creates a more continuous frequency representation. Another difference is that the MP3 standard applies a *window* to the sound samples, and the effect of this is that the new signal has a frequency representation which is closer to the original one, when compared to the signal obtained by using the block values unchanged as above. We will go into details on this in Section 3.3.1.

Example 2.30. We can achieve compression of a sound by setting the Fourier coefficients which are small to zero. The idea is that frequencies with small values at the

corresponding frequency indices contribute little to our perception of the sound, so that they can be discarded. As a result we obtain a sound with less frequency components, which is thus more suitable for compression. To test this in practice, we first need to set a threshold, which decides what frequencies should be kept or not. The following code then sets frequencies below the threshold to zero:

```
y=DFTImpl(x);
y=(abs(y)>=threshold).*y;
newx=IDFTImpl(y);
```

In this code 1 represents a value of true in the logical expression which is evaluated, 0 represents false. The value is 1 if and only if the absolute value of the corresponding element is greater than or equal to threshold. As in the previous example, we can apply this code to small blocks of the signal at a time, and listen to the result by playing it. We have implemented a function `playDFTthreshold(threshold)` which splits our sample audio file into blocks of the same size as above, applies the code above with the given threshold, and plays the result. The code also writes to the display how large percentage of the DFT indices were set to 0. If you run this function with `threshold` equal to 0.02, the result sounds like this, and the function says that about 74.1% of the DFT indices were set to zero. You can clearly hear the disturbance in the sound, but we have not lost that much. If we instead try `threshold` equal to 0.1, the result will sound like this, and the function says that about 93.5% of the DFT indices were set to zero. The quality is much poorer now, even if we still can recognize the song. This suggests that most of the frequency information is contained in frequencies with the highest values. In figure 2.4 we have illustrated this principle for compression. The samples of the sound are shown in (a) and (b). In (c) all values of the DFT with absolute value smaller than 0.02 have been set to zero. The sound is then reconstructed with the IDFT, and the result shown in (d). The two signals in (a) and (d) visually look almost the same even though the signal in (d) can be represented with less than 10 % of the information present in (a).

Note that using a neglection threshold in this way is too simple in practice: The neglection threshold in general should depend on the frequency, since the human auditory system is more sensitive to certain frequency information. ♣

Example 2.31. The previous example is a rather simple procedure to obtain compression. The disadvantage with it is that it only affects small frequencies. A more neutral way to obtain compression is to let each DFT index occupy a certain number of bits. This provides us with compression if this number of bits is less than what actually is used to represent the sound. This is closer to what modern audio standards do. Let us say that we would like to use only n bits (keeping the bit used for the sign out of this). We can then use the following code to quantize the DFT-coefficients in one block:

```
y=DFTImpl(x);
y=y*2^n;
y=round(y);
y=y/2^n;
```

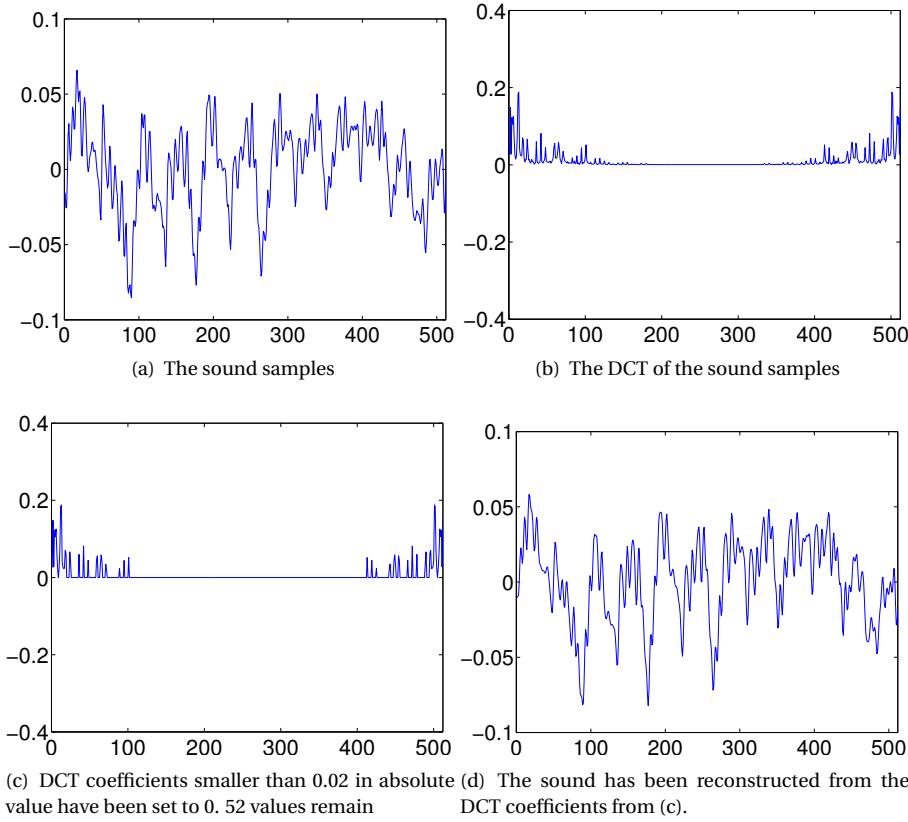


Figure 2.4: Experimenting with the DFT on a small part (512 sound samples) of a song.

```
newx=IDFTImpl(y);
```

If we replace the code from the previous examples, we hear similar effect to the sound. This quantization procedure is also too simple, however, since since the human auditory system is more sensitive to certain frequency information, and should thus allocate a higher number of bits for such frequencies. Modern audio standards take this into account, but we will not go into details on this. ♣

What you should have learnt in this section

Translation between DFT index and frequency. In particular DFT indices for high and low frequencies. How one can use the DFT to adjust frequencies in sound.

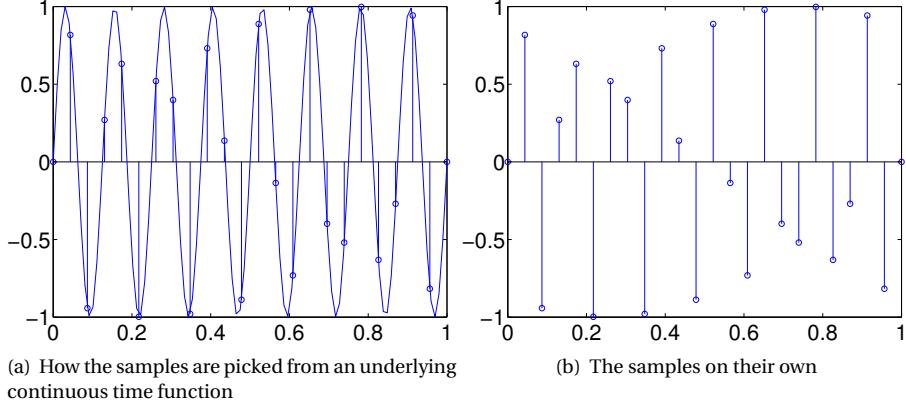


Figure 2.5: An example of sampling.

Exercises for Section 2.6

1. Suppose that DFTImpl of IDFTImpl are implementations of the DFT and the IDFT, respectively. Explain what the code below does, line by line:

```
[S,fs]=wavread('castanets.wav');
S=S(1:2^(17),1);
y=DFTImpl(S);
y((2^17/4+1):(3*2^17/4))=zeros(2^16,1);
newS=IDFTImpl(y);
newS=newS/max(abs(newS));
playerobj=audioplayer(newS,fs);
playblocking(playerobj)
```

Comment in particular why we perform the command on the third line from bottom. What changes in the sound do you expect to hear?

2. In the code from the previous exercise it turns out that $f_s = 44100\text{Hz}$. Which frequencies in the sound file will be changed by line 4 in the code?

2.7 Reconstruction of a function from its samples. The sampling theorem.

The second interesting facet to Theorem 2.24 has to do with when exact reconstruction of a function from its samples is possible. An example of sampling a function is illustrated in Figure 2.5. From (b) it is clear that information is lost when we discard everything but the sample values. There may however be an exception to this, if we assume that the function satisfies some property. Assume for instance that f is equal to a finite Fourier series. This means that f can be written on the form (2.6), so that the highest frequency in the signal is bounded by M/T . Our analysis prior to

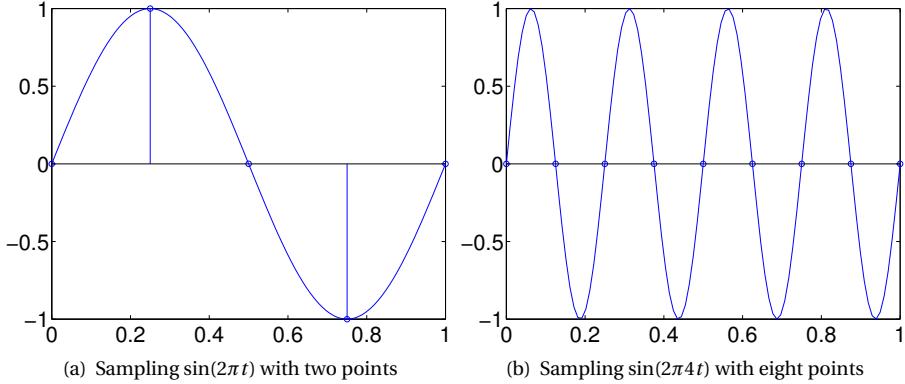


Figure 2.6: Sampling two different pure tones.

Theorem 2.24 states that all such functions can be reconstructed exactly from their samples, as long as the number of samples is $N \geq 2M + 1$, taken uniformly over a period. Moreover, the DFT is central in the reconstruction formula. Dividing by T we get $\frac{N}{T} \geq \frac{2M+1}{T}$, which states that the sampling frequency ($f_s = N/T$ is the number of samples per second) should be bigger than two times the highest frequency (M/T). In Figure 2.6 we try to get some intuition on this by considering some pure tones. In (a) we consider one period of $\sin(2\pi t)$. Since this is in $V_{M,T} = V_{1,2\pi}$, reconstruction should be possible if we have $N \geq 2M + 1 = 3$ samples. Four sample points, as seen in (a), should thus be enough to reconstruct the function. In (b) we have plotted $\sin(2\pi 4t)$, with 8 sample points taken uniformly from the first period. Here $M = 4$, so that we require $2M + 1 = 9$ sample point from the first period. In theory, these sample points may therefore not be enough to reconstruct a unique function in $V_{M,T}$ passing through these points. This is seen to be the case, since clearly the zero function is another function from $V_{M,T}$ besides $\sin(2\pi 4t)$, which passes through these points.

Let us restate the reconstruction of f , so that it only uses the samples. The reconstruction formula was

$$f(t) = \frac{1}{\sqrt{N}} \sum_{n=-M}^M z_n e^{2\pi i n t / T},$$

where the z_n can be found from the equations

$$f(kT_s) = \frac{1}{\sqrt{N}} \sum_{n=-M}^M z_n e^{2\pi i n k / N} \quad -M \leq k \leq M,$$

where we have substituted $N = T/T_s$ (deduced from $T = NT_s$ with T_s being the sampling period). From this equation we see that

$$z_n = \frac{1}{\sqrt{N}} \sum_{k=-M}^M f(kT_s) e^{-2\pi i n k / N},$$

and inserting this in the reconstruction formula we get

$$\begin{aligned}
f(t) &= \frac{1}{N} \sum_{n=-M}^M \sum_{k=-M}^M f(kT_s) e^{-2\pi i n k / N} e^{2\pi i n t / T} \\
&= \sum_{k=-M}^M \frac{1}{N} \left(\sum_{n=-M}^M f(kT_s) e^{2\pi i n(t/T - k/N)} \right) \\
&= \sum_{k=-M}^M \frac{1}{N} e^{-2\pi i M(t/T - k/N)} \frac{1 - e^{2\pi i N(t/T - k/N)}}{1 - e^{2\pi i (t/T - k/N)}} f(kT_s) \\
&= \sum_{k=-M}^M \frac{1}{N} \frac{\sin(\pi(t - kT_s)/T_s)}{\sin(\pi(t - kT_s)/T)} f(kT_s)
\end{aligned}$$

Let us summarize our findings as follows:

Theorem 2.32. (Sampling theorem and the ideal interpolation formula for periodic functions) Let f be a periodic function with period T , and assume that f has no frequencies higher than v Hz. Then f can be reconstructed exactly from its samples $f(-MT_s), \dots, f(MT_s)$ (where T_s is the sampling period, $N = \frac{T}{T_s}$ is the number of samples per period, and $M = 2N + 1$) when the sampling rate $f_s = \frac{1}{T_s}$ is bigger than $2v$. Moreover, the reconstruction can be performed through the formula

$$f(t) = \sum_{k=-M}^M f(kT_s) \frac{1}{N} \frac{\sin(\pi(t - kT_s)/T_s)}{\sin(\pi(t - kT_s)/T)}. \quad (2.11)$$

Formula (2.12) is also called the ideal interpolation formula for periodic functions. Such formulas, where one reconstructs a function based on a weighted sum of the sample values, are more generally called *interpolation formulas*. The function $\frac{1}{N} \frac{\sin(\pi(t - kT_s)/T_s)}{\sin(\pi(t - kT_s)/T)}$ is also called an interpolation kernel. Note that f itself may not be equal to a finite Fourier series, and reconstruction is in general not possible then. The ideal interpolation formula can in such cases still be used, but the result we obtain may be different from $f(t)$.

It turns out that the interpolation formula above can be rewritten without the dependence on T and N , i.e. so that the interpolation formula is valid for all numbers of samples. This formula is what is usually listed in the literature:

Theorem 2.33. (Sampling theorem and the ideal interpolation formula for periodic functions, general version) Assume that f is periodic with period T , and has no frequencies higher than v Hz. Then f can be reconstructed exactly from its samples $\dots, f(-2T_s), f(-T_s), f(0), f(T_s), f(2T_s), \dots$ when T is a multiple of T_s , and when the sampling rate is bigger than $2v$. Moreover, the reconstruction can be performed through the formula

$$f(t) = \sum_{k=-M}^M f(kT_s) \frac{\sin(\pi(t - kT_s)/T_s)}{\pi(t - kT_s)/T_s}. \quad (2.12)$$

Proof: Note first that f can also be viewed as a function with period sT for any integer $s > 1$. Writing sT for T , and sN for N in the previous interpolation formula, we get

$$\begin{aligned} f(t) &= \sum_{k=-sM}^{sM} f(kT_s) \frac{1}{sN} \frac{\sin(\pi(t - kT_s)/T_s)}{\sin(\pi(t - kT_s)/(sT))} \\ &= \sum_{k=-rM}^{rM} f(kT_s) \frac{1}{sN} \frac{\sin(\pi(t - kT_s)/T_s)}{\sin(\pi(t - kT_s)/(sT))} \\ &\quad + \sum_{rM < |k| \leq sM} f(kT_s) \frac{1}{sN} \frac{\sin(\pi(t - kT_s)/T_s)}{\sin(\pi(t - kT_s)/(sT))}, \end{aligned}$$

where we have split the summation further (r is a number smaller than s). Note that $\lim_{s \rightarrow \infty} sN \sin(\pi(t - kT_s)/(sT)) = \pi(t - kT_s)/T_s$. Thus if we let $s \rightarrow \infty$ while keeping r fixed, the first sum above converges to

$$\sum_{k=-rM}^{rM} f(kT_s) \frac{\sin(\pi(t - kT_s)/T_s)}{\pi(t - kT_s)/T_s}.$$

The second sum can be written as

$$\sum_{k=-M}^M f(kT_s) \sum_{r < |u| \leq s} \frac{1}{sN} \frac{\sin(\pi(t - (k + uM)T_s)/T_s)}{\sin(\pi(t - (k + uM)T_s)/(sT))}.$$

The numerator here turns the inner sum into an alternating series, and the denominator is increasing when r is chosen big enough at the start. This means that the sum converges, and is bounded by the first term

$$\begin{aligned} &\frac{1}{sN \sin(\pi(t - (k + rM)T_s)/(sT))} \\ &\leq \frac{2}{sN \pi(t - (k + rM)T_s)/(sT)} = \frac{2T_s}{\pi(t - (k + rM)T_s)}. \end{aligned}$$

For a given t , this can be made a small as we like, if we choose r big enough. ■

In the literature one actually shows that this formula is valid also for functions which are not periodic.

2.8 The Fast Fourier Transform (FFT)

The main application of the DFT is as a tool to compute frequency information in large datasets. It is therefore important that these operations can be performed by efficient algorithms. Straightforward implementation from the definition is not efficient if the data sets are large. However, it turns out that the underlying matrices may be factored in a way that leads to much more efficient algorithms, and this is the topic of the present section, where we discuss the most widely used implementation of the DFT, usually referred to as the Fast Fourier Transform (FFT). Before we state the FFT algorithm, let us first comment on two assumptions we will make.

First of all, we will assume that N , the length of the vector that is to be transformed by the DFT, is a power of 2. We will see that this assumption leads to nice, recursive algorithms for the DFT, where each step divides N by 2. In many settings this power of 2 assumption can be done. As an example, in compression of sound, one restricts processing to a certain block of the sound data, since the entire sound is too big to be processed in one piece. One then has a freedom to how big these blocks are made, and one often assumes a power of 2 for these block sizes for simplicity. At the end of this section we will explain how the more general FFT can be computed when N is not a power of 2.

Secondly, we will assume that the input vector \mathbf{x} to the FFT algorithm is real. This is the case for most practical purposes, since the vector \mathbf{x} typically are measurements from some sensors, sound samples, or pixel values in an image. The assumption of real input reduces the dimensionality of the problem to half of that required with complex input. This reduction of dimensionality can easily be recognized in the output \mathbf{y} of the DFT as well, since $y_{N-n} = \overline{y_n}$ when the input to the DFT is real. Clearly $y_{N/2}$ is real due to this, and since y_0 also is real, the output of the DFT of a real vector with N even is uniquely determined from

$$y_0 = \Re y_0, \Re y_1, \Im y_1, \Re y_2, \Im y_2, \dots, \Re y_{N/2-1}, \Im y_{N/2-1}, y_{N/2} = \Re y_{N/2}. \quad (2.13)$$

We thus can restrict to computing these values, and the output of the DFT need only include these values.

Now, let us state the FFT algorithm. This first algorithm will in fact also work for complex numbers, and for any N which is an even number.

Theorem 2.34 (FFT algorithm when N is even). Let $\mathbf{y} = F_N \mathbf{x}$ be the N -point DFT of \mathbf{x} with N an even number. For any integer n in the interval $[0, N/2 - 1]$ the DFT \mathbf{y} of \mathbf{x} is then given by

$$y_n = \frac{1}{\sqrt{2}} \left((F_{N/2} \mathbf{x}^{(e)})_n + e^{-2\pi i n/N} (F_{N/2} \mathbf{x}^{(o)})_n \right), \quad (2.14)$$

$$y_{N/2+n} = \frac{1}{\sqrt{2}} \left((F_{N/2} \mathbf{x}^{(e)})_n - e^{-2\pi i n/N} (F_{N/2} \mathbf{x}^{(o)})_n \right), \quad (2.15)$$

where $\mathbf{x}^{(e)}, \mathbf{x}^{(o)}$ are the sequences of length $N/2$ consisting of the even and odd samples of \mathbf{x} , respectively. In other words,

$$(\mathbf{x}^{(e)})_k = x_{2k} \text{ for } 0 \leq k \leq N/2 - 1,$$

$$(\mathbf{x}^{(o)})_k = x_{2k+1} \text{ for } 0 \leq k \leq N/2 - 1.$$

Put differently, the formulas (2.14)–(2.15) reduces the computation of an N -point DFT to two $N/2$ -point DFT's. It turns out that this can speed up computations considerably, but let us first check that these formulas are correct.

Proof: Suppose first that $0 \leq n \leq N/2 - 1$. We start by splitting the sum in the

expression for the DFT into even and odd indices,

$$\begin{aligned}
y_n &= \frac{1}{\sqrt{N}} \sum_{k=0}^{N-1} x_k e^{-2\pi i n k / N} \\
&= \frac{1}{\sqrt{N}} \sum_{k=0}^{N/2-1} x_{2k} e^{-2\pi i n 2k / N} + \frac{1}{\sqrt{N}} \sum_{k=0}^{N/2-1} x_{2k+1} e^{-2\pi i n (2k+1) / N} \\
&= \frac{1}{\sqrt{2}} \frac{1}{\sqrt{N/2}} \sum_{k=0}^{N/2-1} x_{2k} e^{-2\pi i n k / (N/2)} \\
&\quad + e^{-2\pi i n / N} \frac{1}{\sqrt{2}} \frac{1}{\sqrt{N/2}} \sum_{k=0}^{N/2-1} x_{2k+1} e^{-2\pi i n k / (N/2)} \\
&= \frac{1}{\sqrt{2}} (F_{N/2} \mathbf{x}^{(e)})_n + \frac{1}{\sqrt{2}} e^{-2\pi i n / N} (F_{N/2} \mathbf{x}^{(o)})_n,
\end{aligned}$$

where we have substituted $\mathbf{x}^{(e)}$ and $\mathbf{x}^{(o)}$ as in the text of the theorem, and recognized the $N/2$ -point DFT in two places. For the second half of the DFT coefficients, i.e. $\{y_{N/2+n}\}_{0 \leq n \leq N/2-1}$, we similarly have

$$\begin{aligned}
y_{N/2+n} &= \frac{1}{\sqrt{N}} \sum_{k=0}^{N-1} x_k e^{-2\pi i (N/2+n) k / N} = \frac{1}{\sqrt{N}} \sum_{k=0}^{N-1} x_k e^{-\pi i k} e^{-2\pi i n k / N} \\
&= \frac{1}{\sqrt{N}} \sum_{k=0}^{N/2-1} x_{2k} e^{-2\pi i n 2k / N} - \frac{1}{\sqrt{N}} \sum_{k=0}^{N/2-1} x_{2k+1} e^{-2\pi i n (2k+1) / N} \\
&= \frac{1}{\sqrt{2}} \frac{1}{\sqrt{N/2}} \sum_{k=0}^{N/2-1} x_{2k} e^{-2\pi i n k / (N/2)} \\
&\quad - e^{-2\pi i n / N} \frac{1}{\sqrt{2}} \frac{1}{\sqrt{N/2}} \sum_{k=0}^{N/2-1} x_{2k+1} e^{-2\pi i n k / (N/2)} \\
&= \frac{1}{\sqrt{2}} (F_{N/2} \mathbf{x}^{(e)})_n - \frac{1}{\sqrt{2}} e^{-2\pi i n / N} (F_{N/2} \mathbf{x}^{(o)})_n.
\end{aligned}$$

This concludes the proof. ■

In addition we need the formula

$$y_{N/2} = \frac{1}{\sqrt{N}} \sum_{k=0}^{N/2-1} ((\mathbf{x}^{(e)})_k - (\mathbf{x}^{(o)})_k) \tag{2.16}$$

to obtain coefficient $N/2$, since this is the only coefficient which can't be obtained from $y_0, y_1, \dots, y_{N/2-1}$ by this symmetry. Note also that the special expression for y_0 ,

$$y_0 = \frac{1}{\sqrt{N}} \sum_{n=0}^{N-1} x_k, \tag{2.17}$$

so that both y_0 and $y_{N/2}$ can be computed without complex arithmetic. It turns out that Theorem 2.34 can be interpreted as a matrix factorization. For this we need to define the concept of a block matrix.

Definition 2.35. Let m_0, \dots, m_{r-1} and n_0, \dots, n_{s-1} be integers, and let $A^{(i,j)}$ be an $m_i \times n_j$ -matrix for $i = 0, \dots, r-1$ and $j = 0, \dots, s-1$. The notation

$$A = \left(\begin{array}{c|c|c|c} A^{(0,0)} & A^{(0,1)} & \cdots & A^{(0,s-1)} \\ \hline A^{(1,0)} & A^{(1,1)} & \cdots & A^{(1,s-1)} \\ \hline \vdots & \vdots & \vdots & \vdots \\ \hline A^{(r-1,0)} & A^{(r-1,1)} & \cdots & A^{(r-1,s-1)} \end{array} \right)$$

denotes the $(m_0 + m_1 + \dots + m_{r-1}) \times (n_0 + n_1 + \dots + n_{s-1})$ -matrix where the matrix entries occur as in the $A^{(i,j)}$ matrices, in the way they are ordered, and with solid lines indicating borders between the blocks. When A is written in this way it is referred to as a block matrix.

We will express the Fourier matrix in factored form involving block matrices. The following observation is just a formal way to split a vector into its even and odd components.

Observation 2.36. Define the permutation matrix P_N by

$$\begin{aligned} (P_N)_{i,2i} &= 1, \quad \text{for } 0 \leq i \leq N/2 - 1; \\ (P_N)_{i,2i-N+1} &= 1, \quad \text{for } N/2 \leq i < N; \\ (P_N)_{i,j} &= 0, \quad \text{for all other } i \text{ and } j; \end{aligned}$$

and let \mathbf{x} be a column vector. The mapping $\mathbf{x} \rightarrow P\mathbf{x}$ permutes the components of \mathbf{x} so that the even components are placed first and the odd components last,

$$P_N \mathbf{x} = \begin{pmatrix} \mathbf{x}^{(e)} \\ \mathbf{x}^{(o)} \end{pmatrix},$$

with $\mathbf{x}^{(e)}, \mathbf{x}^{(o)}$ defined as in Theorem 2.34.

Let $D_{N/2}$ be the $(N/2) \times (N/2)$ -diagonal matrix with entries $(D_{N/2})_{n,n} = e^{-2\pi i n/N}$ for $n = 0, 1, \dots, N/2 - 1$. It is clear from Equation (2.14) that the first half of \mathbf{y} is then given by obtained as

$$\frac{1}{\sqrt{2}} \left(\begin{array}{c|c} F_{N/2} & D_{N/2} F_{N/2} \end{array} \right) P_N \mathbf{x},$$

and from Equation (2.15) that the second half of \mathbf{y} can be obtained as

$$\frac{1}{\sqrt{2}} \left(\begin{array}{c|c} F_{N/2} & -D_{N/2} F_{N/2} \end{array} \right) P_N \mathbf{x}.$$

From these two formulas we can derive the promised factorisation of the Fourier matrix.

Theorem 2.37 (DFT matrix factorization). The Fourier matrix may be factored as

$$F_N = \frac{1}{\sqrt{2}} \begin{pmatrix} F_{N/2} & D_{N/2}F_{N/2} \\ F_{N/2} & -D_{N/2}F_{N/2} \end{pmatrix} = \frac{1}{\sqrt{2}} \begin{pmatrix} I & D_{N/2} \\ I & -D_{N/2} \end{pmatrix} \begin{pmatrix} F_{N/2} & \mathbf{0} \\ \mathbf{0} & F_{N/2} \end{pmatrix} P_N. \quad (2.18)$$

This factorization in terms of block matrices is commonly referred to as the *FFT factorization* of the Fourier matrix. The matrix P_N above can be dropped if we assume that the even indices are placed first, before the odd indices. Since the even indices have last bit zero, the odd indices have last bit one, this is achieved if we bit-reverse the elements in \mathbf{x} . We will do this in what follows, so that the matrix P_N will not be listed, also in cases where the FFT factorization is continued, to factorize $F_{N/2}$, and so on.

We will shortly see why the FFT factorization saves us many arithmetic operations, but let us already now establish why the FFT factorization corresponds to a factorization into sparse matrices (sparse matrix factorization is the starting point for most efficient factorizations): If we applied the FFT matrix factorization again to Equation 2.18, we obtain

$$F_N = \frac{1}{2} \begin{pmatrix} I & D_{N/4} & \mathbf{0} & \mathbf{0} \\ I & -D_{N/4} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & I & D_{N/4} \\ \mathbf{0} & \mathbf{0} & I & -D_{N/4} \end{pmatrix} \begin{pmatrix} F_{N/4} & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & F_{N/4} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & F_{N/4} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & F_{N/4} \end{pmatrix} \quad (2.19)$$

Clearly, if this factorization is repeated, we obtain a product of sparse matrices.

In an implementation based on formula (2.18), the two leftmost blocks in the block matrix in (2.18) correspond to applying the $\frac{N}{2}$ -point DFT to the even samples. The two rightmost blocks correspond to applying the $N/2$ -point DFT to the odd samples, and multiplying the result with $D_{N/2}$. The results from these transforms are finally added together. By repeating the splitting we will eventually come to the case where $N = 1$. Then F_1 is just the scalar 1, so the DFT is the trivial assignment $y_0 = x_0$. The FFT can therefore be implemented with the following Matlab code:

```
function y = FFTImpl(x)
N = length(x);
if N == 1
    y = x(1);
else
    ye = FFTImpl(x(1:2:(N-1)));
    yo = FFTImpl(x(2:2:N));
    D=exp(-2*pi*1j*(0:(N/2-1))'/N);
    yo=yo.*D;
    y = [ ye + yo; ye - yo]/sqrt(2);
end
```

Note that this function is recursive; it calls itself. If this is your first encounter with a recursive program, it is worth running through the code manually for a given value of N , such as $N = 4$.

2.8.1 The Inverse Fast Fourier Transform (IFFT)

The IDFT is very similar to the DFT, and it is straightforward to obtain a similar expression as in Theorem 2.34 by following the same proof. We would, however, like to have an expression for the IDFT which assumes that the DFT was obtained from real data. This is the same as restricting the IFFT implementations where $y_{N-n} = \overline{y_n}$. The following result addresses this.

Theorem 2.38 (IFFT algorithm). When \mathbf{x} is a real vector and $\mathbf{y} = F_N \mathbf{x}$, the IDFT of \mathbf{y} can be computed as follows: Define $\mathbf{y}^{(0)}$ and $\mathbf{y}^{(1)}$ as the vectors of length $N/2$ with entries $y_n + \overline{y_{N/2-n}}$ and $e^{2\pi i n/N} (y_n - \overline{y_{N/2-n}})$, respectively. Then $\mathbf{y}^{(0)}$ and $\mathbf{y}^{(1)}$ are also DFT's of real vectors, and \mathbf{x} can be computed from the formulas

$$\begin{aligned}\mathbf{x}^{(e)} &= \frac{1}{\sqrt{2}} (F_{N/2})^H \mathbf{y}^{(0)} \\ \mathbf{x}^{(o)} &= \frac{1}{\sqrt{2}} (F_{N/2})^H \mathbf{y}^{(1)}.\end{aligned}$$

Proof: If we apply Equation (2.14) to the indices n and $N/2 - n$, and using that $(F_{N/2} \mathbf{z})_{N-n} = (F_{N/2} \mathbf{z})_n$ for any real vector \mathbf{z} , we obtain

$$\begin{aligned}y_n &= \frac{1}{\sqrt{2}} \left((F_{N/2} \mathbf{x}^{(e)})_n + e^{-2\pi i n/N} (F_{N/2} \mathbf{x}^{(o)})_n \right) \\ \overline{y_{N/2-n}} &= \frac{1}{\sqrt{2}} \left((F_{N/2} \mathbf{x}^{(e)})_{N-n} - e^{-2\pi i n/N} (F_{N/2} \mathbf{x}^{(o)})_{N-n} \right).\end{aligned}$$

Adding these we get that $\sqrt{2}(F_{N/2} \mathbf{x}^{(e)})_n = y_n + \overline{y_{N/2-n}}$, so that $\mathbf{x}^{(e)}$ can be obtained by defining $\mathbf{y}^{(0)}$ as the vector with elements $y_n + \overline{y_{N/2-n}}$, and then computing $\frac{1}{\sqrt{2}} (F_{N/2})^H \mathbf{y}^{(0)}$. Similarly, by subtracting the two equations we obtain $\sqrt{2}(e^{-2\pi i n/N} F_{N/2} \mathbf{x}^{(o)})_n = y_n - \overline{y_{N/2-n}}$, so that $\mathbf{x}^{(o)}$ can be obtained by defining $\mathbf{y}^{(1)}$ as the vector with elements $e^{2\pi i n/N} (y_n - \overline{y_{N/2-n}})$, and then computing $\frac{1}{\sqrt{2}} (F_{N/2})^H \mathbf{y}^{(1)}$. Finally it is straightforward to check that $\mathbf{y}_{N-n}^{(0)} = \overline{\mathbf{y}_n^{(0)}}$ and $\mathbf{y}_{N-n}^{(1)} = \overline{\mathbf{y}_n^{(1)}}$, so that both are DFT's of real vectors. ■

With this result it is straightforward to modify FFTImpl.m so that it performs the inverse DFT. With the assumption that \mathbf{x} is real we could have reimplemented our functions so that only $y_0, y_1, \dots, y_{N/2-1}$ are output from FFTImpl, and so that these are the only input to IFFTImpl. We have, however, avoided this in our implementations. Matlab has built-in functions for computing the DFT and the IDFT using the FFT algorithm. These functions are called `fft` and `ifft`. Note, however, that these functions do not use the normalizing factor $1/\sqrt{N}$ that we have adopted here.

The Matlab help pages give a short description of these algorithms. Note in particular that Matlab makes no assumption about the length of the vector. Matlab may however check if the length of the vector is 2^r , and in those cases a variant of the algorithm discussed here is used. In general, fast algorithms exist when the vector length N can be factored as a product of small integers.

Many audio and image formats make use of the FFT. To get optimal speed these algorithms typically split the signals into blocks of length 2^r with r some integer in the range 5–10 and utilise a suitable variant of the algorithms discussed above.

2.8.2 Reduction in the number of multiplications with the FFT

Before we continue we also need to explain why the FFT and IFFT factorizations lead to more efficient implementations than the direct DFT and IDFT implementations. We first need some terminology for how we count the number of operations of a given type in an algorithm. In particular we are interested in the limiting behaviour when N becomes large, which is the motivation for the following definition.

Definition 2.39 (Order of an algorithm). Let R_N be the number of operations of a given type (such as multiplication or addition) in an algorithm, where N describes the dimension of the data in the algorithm (such as the size of the matrix or length of the vector), and let f be a positive function. The algorithm is said to be of order N , also written $O(f(N))$, if the number of operations grows as $f(N)$ for large N , or more precisely, if

$$\lim_{N \rightarrow \infty} \frac{R_N}{f(N)} = 1.$$

We will also use this notation for functions, and say that a real function g is $O(f(x))$ if $\lim g(x)/f(x) = 0$ where the limit mostly will be taken as $x \rightarrow \mathbf{0}$ (this means that $g(x)$ is much smaller than $f(x)$ when x approaches the limit).

Let us see how we can use this terminology to describe the complexity of the FFT algorithm. For simplicity, we will only consider the number of multiplications. Let therefore M_N be the number of multiplications needed by the N -point FFT as defined by Theorem 2.34. It is clear from the algorithm that

$$M_N = 2M_{N/2} + N/2. \quad (2.20)$$

The factor 2 corresponds to the two matrix multiplications, while the term $N/2$ denotes the multiplications in the exponent of the exponentials that make up the matrix $D_{N/2}$ (or $E_{N/2}$). The exponent $2\pi i/N$ may be computed once and for all outside the loops, and has therefore not been counted. Also, we have not counted the multiplications with $1/\sqrt{2}$. The reason is that, the factor $1/\sqrt{N}$ is not included in most definitions of the DFT, and one can show that this difference leads to FFT algorithms where the factor $1/\sqrt{2}$ is not present (see Exercise 3).

Note that all multiplications performed by the FFT are complex. It is normal to count the number of real multiplications instead, since any multiplication of two complex numbers can be performed as four multiplications of real numbers (and

two additions), by writing the number in terms of its real and imaginary part, and multiplying them together. Therefore, if we instead define M_N to be the number of real multiplications required by the FFT, we obtain the alternative difference equation

$$M_N = 2M_{N/2} + 2N.$$

It is straightforward to show that the IFFT implementation requires the same count. In fact, a simple trick in the FFT and IFFT implementations can also be done so that

$$M_N = 2M_{N/2} + N, \quad (2.21)$$

i.e. so that we only need to perform half of the additional real multiplications at each stage, see Exercise 9. In the same exercise you will be asked to find the solution of this difference equation and show that the number of real multiplications required by the FFT algorithm is $O(N \log_2 N)$. In contrast, the direct implementation of the DFT requires N^2 complex multiplications, and thus $4N^2$ real multiplications. In other words, the FFT and IFFT significantly reduce the number of multiplications.

Theorem 2.40 (Number of operations in the FFT and IFFT algorithms). The N -point FFT and IFFT algorithms both require $O(N \log_2 N)$ real multiplications. In comparison, the number of real multiplications required by direct implementations of the N -point DFT and IDFT is $4N^2$.

In a similar way one can show that the number of additions required by the algorithm is also roughly $O(N \log_2 N)$. Another reason why the FFT is efficient is that, since the FFT splits the calculation of the DFT into computing two DFT's of half the size, the FFT is well suited for parallel computing: the two smaller FFT's can be performed independently of one another, for instance in two different computing cores on the same computer.

2.8.3 The FFT when $N = N_1 N_2$

Applying an FFT to a vector of length 2^n is by far the most common thing to do. It turns out, however, that the idea behind the algorithm easily carries over to the case when N is any composite number, i.e. when $N = N_1 N_2$. This makes the FFT useful also in settings where we have a dictated number of elements in \mathbf{x} , which is not an even number. The approach we will present in this section will help us as long as N is not a prime number. The case when N is a prime number needs other techniques.

So, assume that $N = N_1 N_2$. Any time-index k can be written uniquely on the form $N_1 k + p$, with $0 \leq k < N_2$, and $0 \leq p < N_1$. We will make the following definition.

Definition 2.41 (Polyphase components of a vector). Let $\mathbf{x} \in \mathbb{R}^{N_1 N_2}$. We denote by $\mathbf{x}^{(p)}$ the vector in \mathbb{R}^{N_2} with entries $(\mathbf{x}^{(p)})_k = x_{N_1 k + p}$. $\mathbf{x}^{(p)}$ is also called the p 'th *polyphase component* of \mathbf{x} .

The previous vectors $\mathbf{x}^{(e)}$ and $\mathbf{x}^{(o)}$ can be seen as special cases of polyphase components. Polyphase components will also be useful later (see Chapter 8). Using the polyphase notation, we can write

$$\begin{aligned} F_N \mathbf{x} &= \frac{1}{\sqrt{N_1 N_2}} \sum_{k=0}^{N-1} x_k e^{-2\pi i n k / N} = \frac{1}{\sqrt{N_1 N_2}} \sum_{p=0}^{N_1-1} \sum_{k=0}^{N_2-1} (\mathbf{x}^{(p)})_k e^{-2\pi i n (N_1 k + p) / N} \\ &= \frac{1}{\sqrt{N_1 N_2}} \sum_{p=0}^{N_1-1} e^{-2\pi i np / N} \sum_{k=0}^{N_2-1} (\mathbf{x}^{(p)})_k e^{-2\pi i nk / N_2} \end{aligned}$$

Similarly, any frequency index n can be written uniquely on the form $N_2 q + n$, with $0 \leq q < N_1$, and $0 \leq n < N_2$, so that the DFT can also be written as

$$\begin{aligned} &\frac{1}{\sqrt{N_1 N_2}} \sum_{p=0}^{N_1-1} e^{-2\pi i (N_2 q + n) p / N} \sum_{k=0}^{N_2-1} (\mathbf{x}^{(p)})_k e^{-2\pi i (N_2 q + n) k / N_2} \\ &= \frac{1}{\sqrt{N_1}} \sum_{p=0}^{N_1-1} e^{-2\pi i qp / N_1} e^{-2\pi i np / N} \frac{1}{\sqrt{N_2}} \sum_{k=0}^{N_2-1} (\mathbf{x}^{(p)})_k e^{-2\pi i nk / N_2}. \end{aligned}$$

Now, if X is the $N_2 \times N_1$ -matrix X where the p 'th column is $\mathbf{x}^{(p)}$, we recognize the inner sum $\frac{1}{\sqrt{N_2}} \sum_{k=0}^{N_2-1} (\mathbf{x}^{(p)})_k e^{-2\pi i nk / N_2}$ as matrix multiplication with F_{N_2} and X , so that this can be written as $(F_{N_2} X)_{n,p}$. The entire sum can thus be written as

$$\frac{1}{\sqrt{N_1}} \sum_{p=0}^{N_1-1} e^{-2\pi i qp / N_1} e^{-2\pi i np / N} (F_{N_2} X)_{n,p}.$$

Now, define Y as the matrix where X is multiplied componentwise with the matrix with (n, p) -component $e^{-2\pi i np / N}$. The entire sum can then be written as

$$\frac{1}{\sqrt{N_1}} \sum_{p=0}^{N_1-1} e^{-2\pi i qp / N_1} Y_{n,p} = (Y F_{N_1})_{n,q}$$

This means that the sum can be written as component (n, q) in the matrix $Y F_{N_1}$. Clearly $Y F_{N_1}$ is the matrix where the DFT is applied to all rows of Y . We have thus shown that component $N_2 q + n$ of $F_N \mathbf{x}$ equals $(Y F_{N_1})_{n,q}$. This means that $F_N \mathbf{x}$ can be obtained by stacking the columns of $Y F_{N_1}$ on top of one another. We can thus summarize our procedure as follows, which gives a recipe for splitting an FFT into smaller FFT's when N is not a prime number.

Theorem 2.42 (FFT algorithm when N is composite). When $N = N_1 N_2$, the FFT of a vector \mathbf{x} can be computed as follows

1. Form the $N_2 \times N_1$ -matrix X , where the p 'th column is $\mathbf{x}^{(p)}$.
2. Perform the DFT on all the columns in X , i.e. compute $F_{N_2} X$.
3. Multiply element (n, p) in the resulting matrix with $e^{-2\pi i np / N}$ (these are called *twiddle factors*), to obtain matrix Y .

4. Perform the DFT on all the rows in the resulting matrix, i.e. compute YF_{N_1} .
5. Form the vector where the columns of the resulting matrix are stacked on top of one another.

With this procedure a count of complex multiplication gives that

$$M_{N_1 N_2} = N_1 M_{N_2} + N_2 M_{N_1} + N_1 N_2,$$

so that $M_N = N_1 M_{N_2} + N_2 M_{N_1} + N$, and we can obtain similar estimates as before for the number of both real and complex multiplications with this.

From the algorithm one easily deduces how the IDFT can be computed also: All steps are invertible, and can be performed by IFFT or multiplication. We thus only need to perform the inverse steps in reverse order.

But what about the case when N is a prime number? Rader's algorithm [29] handles this case by expressing a DFT with N a prime number in terms of DFT's of length $N - 1$ (which is not a prime number). Our previous scenario can then be followed, but stops quickly again if $N - 1$ has prime factors of high order. Since there are some computational penalties in applying Rader's algorithm, it may thus be inefficient some cases. Winograd's FFT algorithm [36] extends Rader's algorithm to work for the case when $N = p^r$. This algorithm tends to reduce the number of multiplications, at the price of an increased number of additions. It is difficult to program and is rarely used in practice.

2.8.4 Applications of the FFT

The FFT has been stated as one the ten most important inventions of the 20'th century. With its invention, the Discrete Fourier Transform was within computational reach in many fields. Real time processing is one important application, such as processing of sound, image, and video.

Compression of sound is one of the many important applications of the FFT. The MP3 standard uses it to find the different frequency components in sound, and to match this information with a psychoacoustic model, in order to find how the different data should be compressed.

What you should have learnt in this section

How the FFT algorithm works by splitting into two FFT's of half the length. Simple FFT implementation. Reduction in the number of multiplications with the FFT.

Exercises for Section 2.8

1. In this exercise we will compare different methods for computing the DFT.

- a.** Write code which compares the execution times for an N -point DFT for the following three cases: Direct implementation of the DFT (as in Example 2.20), the FFT implementation used in this chapter, and the method `fft` which is built into matlab. Your code should take N 's input.
- b.** A big problem with the direct implementation of the DFT, besides that it is very slow, is that it runs out of memory for rather low N , since it allocates the full Fourier matrix. For what value of N does this occur (the answer will depend on the amount of memory available on your computer)?
- c.** It seems that Matlab's built-in FFT is still much faster than our FFT implementation. Try to explain what can be the cause of this.
- 2.**
- a.** Compute the DFT of the vectors $\mathbf{x}_1 = (1, 3, 5, 7)$, and of $\mathbf{x}_2 = (2, 4, 6, 8)$ (i.e. compute $F_4\mathbf{x}_1$ and $F_4\mathbf{x}_2$).
 - b.** Explain how you can compute the DFT of the vector $(1, 2, 3, 4, 5, 6, 7, 8)$ based on the computation from a. (you don't need to perform the actual computation). What are the benefits of this approach, and what is the algorithm called?
- 3.** As previously mentioned, in signal processing literature it is not common to include the normalizing factor $\frac{1}{\sqrt{N}}$ in the definition of the DFT, i.e. they instead define the Fourier matrix as the matrix with entries $e^{-2\pi i kn/N}$. Prove that, with this new definition of the DFT, the FFT algorithm instead takes the form
- $$y_n = (F_{N/2}\mathbf{x}^{(e)})_n + e^{-2\pi i n/N} (F_{N/2}\mathbf{x}^{(o)})_n$$
- $$y_{N/2+n} = (F_{N/2}\mathbf{x}^{(e)})_n - e^{-2\pi i n/N} (F_{N/2}\mathbf{x}^{(o)})_n$$
- We see here that the multiplication with $1/\sqrt{2}$ is not present in the algorithm anymore, which explains why these multiplications were not counted previously when we computed the number of multiplications required by an FFT.
- 4.** Based on the FFT-algorithm, explain why performing a DFT of even length N on real data is equivalent to perform a DFT of length $N/2$ on complex data. By some it has been argued that one can find improved FFT algorithms when one assumes that the data is real. This exercise thus speaks against the possibility of this.
- 5.** Assume that $N = N_1 N_2$, and that $\mathbf{x} \in \mathbb{R}^N$ satisfies $x_{k+rN_1} = x_k$ for all k, r , i.e. \mathbf{x} has period N_1 . Show that $y_n = 0$ for all n which are not multiplums of N_2 .
- 6.** Assume that $N = N_1 N_2$, and that $\mathbf{x}^{(p)} = \mathbf{0}$ fr $p \neq 0$. Show that the polyphase components $\mathbf{y}^{(p)}$ of $\mathbf{y} = F_N\mathbf{x}$ are constant vectors for all p .
- 7.** When we wrote down the difference equation $M_N = 2M_{N/2} + N$ for the number of multiplications in the FFT algorithm, you could argue that some multiplications were not counted. Some of these were explained in Exercise 3. Which other multiplications in the FFT algorithm were not counted when writing down this difference equation? Do you have a suggestion to why these multiplications were not counted?

8. Write down a difference equation for computing the number of real additions required by the FFT algorithm.

9. In this exercise we will first look at the small trick which can be used so that the number of multiplications needed by the FFT algorithm satisfies Equation (2.21), and then solve this difference equation to find the expression for the number of operations.

Assume that \mathbf{x} is a real signal. Equation (2.14) with indices n and $N/2 - n$ can be written

$$\begin{aligned} y_n &= \frac{1}{\sqrt{2}} \left((F_{N/2}\mathbf{x}^{(e)})_n + e^{-2\pi i n/N} (F_{N/2}\mathbf{x}^{(o)})_n \right) \\ y_{N/2-n} &= \frac{1}{\sqrt{2}} \left((F_{N/2}\mathbf{x}^{(e)})_{N/2-n} + e^{-2\pi i (N/2-n)/N} (F_{N/2}\mathbf{x}^{(o)})_{N/2-n} \right) \\ &= \frac{1}{\sqrt{2}} \left((F_{N/2}\mathbf{x}^{(e)})_{N/2-n} - e^{2\pi i n/N} \overline{(F_{N/2}\mathbf{x}^{(o)})_n} \right) \\ &= \frac{1}{\sqrt{2}} \left(\overline{(F_{N/2}\mathbf{x}^{(e)})_n} - e^{-2\pi i n/N} \overline{(F_{N/2}\mathbf{x}^{(o)})_n} \right). \end{aligned}$$

We see here that, if we have computed the terms in y_n (which needs an additional 4 real multiplications, since $e^{-2\pi i n/N}$ and $(F_{N/2}\mathbf{x}^{(o)})_n$ are complex), no further multiplications are needed in order to compute $y_{N/2-n}$, since its compression simply conjugates these terms before adding them.

a. Based on what we proved above, explain how we easily can change the implementation of FFT so that the number of multiplications required satisfies the difference equation $M_N = 2M_{N/2} + N$.

b. With $x_r = M_{2^r}$, show that x_r satisfies the difference equation $x_{r+1} - 2x_r = 2 \cdot 2^r$, and show that the general solution to this is $x_r = r2^r + C2^r$. Conclude that, with this simple trick, the FFT algorithm requires $O(N \log_2 N)$ multiplications.

c. Attempt to find a similar trick as above to slightly modify the IFFT algorithm so that it too requires only $O(N \log_2 N)$ multiplications.

Hint: In the vector $e^{2\pi i n/N} (y_n - \overline{y_{N/2-n}})$ prove that the second part of the vector can be obtained by conjugate symmetry of the first half.

d. Assume that we did not make the trick we started this exercise with, so that we instead have the difference equation $M_N = 2M_{N/2} + 2N$ for the number of real multiplications. Repeat what we did in b. for this new difference equation, to conclude that the FFT now requires $O(2N \log_2 N)$ real multiplications, i.e. that the trick lead to a 50% reduction in the number of real multiplications.

10 (The Split-radix algorithm). If we in Equation (2.18) split the rightmost $F_{N/2}$ by

using Equation (2.18) again, we obtain

$$F_N = \frac{1}{\sqrt{2}} \begin{pmatrix} F_{N/2} & D_{N/2} \frac{1}{\sqrt{2}} \begin{pmatrix} F_{N/4} & D_{N/4} F_{N/4} \\ F_{N/4} & -D_{N/4} F_{N/4} \end{pmatrix} \\ F_{N/2} & -D_{N/2} \frac{1}{\sqrt{2}} \begin{pmatrix} F_{N/4} & D_{N/4} F_{N/4} \\ F_{N/4} & -D_{N/4} F_{N/4} \end{pmatrix} \end{pmatrix}. \quad (2.22)$$

- a.** Let $E_{N/4}$ be the $(N/4) \times (N/4)$ diagonal matrix with $e^{-2\pi i n/N}$ on the diagonal. Show that $D_{N/2} = \begin{pmatrix} E_{N/4} & \mathbf{0} \\ \mathbf{0} & -iE_{N/4} \end{pmatrix}$.
- b.** Let $G_{N/4}$ be the diagonal matrix $E_{D/4} D_{N/4}$. Show that Equation (2.22) can be written as

$$F_N = \frac{1}{\sqrt{2}} \begin{pmatrix} F_{N/2} & \frac{1}{\sqrt{2}} \begin{pmatrix} E_{N/4} F_{N/4} & G_{N/4} F_{N/4} \\ -iE_{N/4} F_{N/4} & iG_{N/4} F_{N/4} \end{pmatrix} \\ F_{N/2} & \frac{1}{\sqrt{2}} \begin{pmatrix} -E_{N/4} F_{N/4} & -G_{N/4} F_{N/4} \\ iE_{N/4} F_{N/4} & -iG_{N/4} F_{N/4} \end{pmatrix} \end{pmatrix}.$$

- c.** Show that, if you have performed the first half of the multiplications in $E_{N/4} F_{N/4}$ and $G_{N/4} F_{N/4}$, there is no need to do the second half.
- d.** Show that, if x is real, this algorithm leads to a multiplication count which satisfies $M_N = M_{N/2} + 2M_{N/4} + N$. Show also that $M_N = O(\frac{2}{3} \log_2 N)$. This algorithm is also called the split-radix algorithm, and held until recently the record for the lowest count of arithmetic operations for any FFT algorithm.

2.9 Summary

We defined digital sound, and demonstrated how we could perform simple operations on digital sound such as adding noise, playing at different rates e.t.c.. Digital sound could be obtained by sampling the sounds from the previous chapter. We considered the analog of Fourier series for digital sound, which is called the Discrete Fourier Transform, and looked at its properties and its relation to Fourier series. We also saw that the sampling theorem guaranteed that there is no loss in considering the samples of a function, as long as the sampling rate is high enough compared to the highest frequency in the sound.

We obtained an implementation of the DFT, called the FFT, which is more efficient in terms of the number of arithmetic operations than a direct implementation of the DFT. The FFT has been cited as one of the ten most important algorithms of the 20'th century [4]. The original paper [6] by Cooley and Tukey dates back to 1965, and handles the case when N is composite. In the literature, one has been interested in the FFT algorithms where the number of (real) additions and multiplications (combined) is as low as possible. This number is also called the *flop count*. The presentation in this book thus differs from the literature in that we mostly count only the number of multiplications. The split radix algorithm [37, 10], which we reviewed

in Exercise 2.8. 10, held the record for the lowest flop count until quite recently [18]. It may seem strange that the total number of additions and multiplications are considered: Aren't multiplications more time-consuming than additions? When you consider how this is done mechanically, this is certainly the case: In fact, floating point multiplication can be considered as a combination of many floating point additions. Due to this, one can find many places in the literature where expressions are rewritten so that the multiplication count is reduced, at the cost of a higher addition count. Winograd's algorithm [36] is an example of this, where the number of additions is much higher than the number of multiplications. However, most modern CPU's have more complex hardware dedicated to computing multiplications, which can result in that one floating point multiplication can be performed in one cycle, just as one addition can. Another thing is that modern CPU's typically can perform many additions and multiplications in parallel, and the higher complexity in the multiplication hardware may result in that the CPU can run less multiplications in parallel, compared to additions. In other words, if we run test program on a computer, it may be difficult to detect any differences in performance between addition and multiplication, even though complex big-scale computing should in theory show some differences. There are also other important aspects of the FFT, besides the flop count. Another is memory use. It is possible to implement the FFT so that the output is computed into the same memory as the input, so that the FFT algorithm does not require extra memory besides the input buffer. Clearly, one should bit-reverse the input buffer in order to achieve this.

We have now defined two types of transforms to the frequency domain: Fourier series for continuous, periodic functions, and the DFT, for periodic vectors. In the literature there are in fact two other transforms also: The Continuous time Fourier transform (CTFT) for continuous functions which are not periodic, and the Discrete time Fourier transform (DTFT) for vectors which are not periodic [28]. In this book we will deliberately avoid these two transforms. The reason is that they assume that the signal to transform is of infinite duration. But, in practice, signals we analyze has a limited time scope.

The sampling theorem is also one of the most important results of the last century. It was discovered by Harry Nyquist and Claude Shannon [31], but also by others independently. One can show that the sampling theorem holds also for functions which are not periodic, as long as we have the same bound on the highest frequency. This is more common in the literature. In fact, the proof seen here where we restrict to periodic functions is not common. The advantage of the proof seen here is that we remain in a finite dimensional setting, and that we only need the DFT. More generally, proofs of the sampling theorem in the literature use the DTFT and the CTFT.

Chapter 3

Operations on digital sound: digital filters

In Section 1.8 we defined analog filters as operations on sound which preserved different frequencies. Such operations are important since they can change the frequency content in many ways. Analog filters can not be used computationally, however, since they are defined for all instances in time. As when we defined the DFT to make Fourier series computable, we would like to define *digital filters*, in order to make analog filters computable. It turns out that what we will define as digital filters can be computed by the following procedure:

$$z_n = \frac{1}{4}(x_{n-1} + 2x_n + x_{n+1}), \quad \text{for } n = 0, 1, \dots, N-1. \quad (3.1)$$

Here \mathbf{x} denotes the *input vector*, and \mathbf{z} the *output vector*. In other words, the output of a digital filter is constructed by combining several input elements linearly. The concrete filter defined by Equation (3.1) is called a *smoothing filter*. We will demonstrate that it smooths the variations in the sound, and this is where it gets its name from. We will start this chapter by looking at matrix representations for operations as given by Equation (3.1). Then we will formally define digital filters in terms of preservation of frequencies as we did for analog filters, and show that the formal definition is equivalent to such operations.

3.1 Matrix representations of filters

Let us consider Equation (3.1) in some more detail to get more intuition about filters. As before we assume that the input vector is periodic with period N , so that $x_{n+N} = x_n$. Our first observation is that the output vector \mathbf{z} is also periodic with period N since

$$z_{n+N} = \frac{1}{4}(x_{n+N-1} + 2x_{n+N} + x_{n+N+1}) = \frac{1}{4}(x_{n-1} + 2x_n + x_{n+1}) = z_n.$$

The filter is also clearly a linear transformation and may therefore be represented by an $N \times N$ matrix S that maps the vector $\mathbf{x} = (x_0, x_1, \dots, x_{N-1})$ to the vector $\mathbf{z} = (z_0, z_1, \dots, z_{N-1})$, i.e., we have $\mathbf{z} = S\mathbf{x}$. To find S , for $1 \leq n \leq N-2$ it is clear from Equation (3.1) that row n has the value $1/4$ in column $n-1$, the value $1/2$ in column n , and the value $1/4$ in column $n+1$. For row 0 we must be a bit more careful, since the index -1 is outside the legal range of the indices. This is where the periodicity helps us out so that

$$z_0 = \frac{1}{4}(x_{-1} + 2x_0 + x_1) = \frac{1}{4}(x_{N-1} + 2x_0 + x_1) = \frac{1}{4}(2x_0 + x_1 + x_{N-1}).$$

From this we see that row 0 has the value $1/4$ in columns 1 and $N-1$, and the value $1/2$ in column 0. In exactly the same way we can show that row $N-1$ has the entry $1/4$ in columns 0 and $N-2$, and the entry $1/2$ in column $N-1$. In summary, the matrix of the smoothing filter is given by

$$S = \frac{1}{4} \begin{pmatrix} 2 & 1 & 0 & 0 & \cdots & 0 & 0 & 0 & 1 \\ 1 & 2 & 1 & 0 & \cdots & 0 & 0 & 0 & 0 \\ 0 & 1 & 2 & 1 & \cdots & 0 & 0 & 0 & 0 \\ \vdots & \vdots \\ 0 & 0 & 0 & 0 & \cdots & 0 & 1 & 2 & 1 \\ 1 & 0 & 0 & 0 & \cdots & 0 & 0 & 1 & 2 \end{pmatrix}. \quad (3.2)$$

A matrix on this form is called a Toeplitz matrix. The general definition is as follows and may seem complicated, but is in fact quite straightforward:

Definition 3.1 (Toeplitz matrices). An $N \times N$ -matrix S is called a Toeplitz matrix if its elements are constant along each diagonal. More formally, $S_{k,l} = S_{k+s,l+s}$ for all nonnegative integers k , l , and s such that both $k+s$ and $l+s$ lie in the interval $[0, N-1]$. A Toeplitz matrix is said to be circulant if in addition

$$S_{(k+s) \bmod N, (l+s) \bmod N} = S_{k,l}$$

for all integers k , l in the interval $[0, N-1]$, and all s (Here \bmod denotes the remainder modulo N).

Toeplitz matrices are very popular in the literature and have many applications. A Toeplitz matrix is constant along each diagonal, while the additional property of being circulant means that each row and column of the matrix 'wraps over' at the edges. Clearly the matrix given by Equation (3.2) satisfies Definition 3.1 and is a circulant Toeplitz matrix. A Toeplitz matrix is uniquely identified by the values on its nonzero diagonals, and a circulant Toeplitz matrix is uniquely identified by the values on the main diagonal, and on the diagonals above (or under) it. While Toeplitz matrices here show up in the context of filters, they will also show up later in the context of wavelets.

Equation (3.1) leads us to the more general expression

$$z_n = \sum_k t_k x_{n-k}. \quad (3.3)$$

This general expression opens up for defining many types of operations. The values t_k will be called *filter coefficients*. The range of k is not specified, but is typically an interval around 0, since z_n usually is calculated by combining x_k 's with indices close to n . Both positive and negative indices are allowed. As an example, for formula (3.1) k ranges over $-1, 0$, and 1 , and we have that $t_{-1} = t_1 = 1/4$, and $t_0 = 1/2$. By following the same argument as above, the following is clear:

Proposition 3.2. Any operation defined by Equation (3.3) is a linear transformation which transforms a vector of period N to another of period N . It may therefore be represented by an $N \times N$ matrix S that maps the vector $\mathbf{x} = (x_0, x_1, \dots, x_{N-1})$ to the vector $\mathbf{z} = (z_0, z_1, \dots, z_{N-1})$, i.e., we have $\mathbf{z} = S\mathbf{x}$. Moreover, the matrix S is a circulant Toeplitz matrix, and the first column \mathbf{s} of this matrix is given by

$$s_k = \begin{cases} t_k, & \text{if } 0 \leq k < N/2; \\ t_{k-N} & \text{if } N/2 \leq k \leq N-1. \end{cases} \quad (3.4)$$

In other words, the first column of S can be obtained by placing the coefficients in (3.3) with positive indices at the beginning of \mathbf{s} , and the coefficients with negative indices at the end of \mathbf{s} .

This proposition will be useful for us, since it explains how to pass from the form (3.3), which is most common in practice, to the matrix form S .

Example 3.3. Let us apply Proposition 3.2 to the operation defined by formula (3.1):

1. for $k = 0$ Equation 3.4 gives $s_0 = t_0 = 1/2$.
2. For $k = 1$ Equation 3.4 gives $s_1 = t_1 = 1/4$.
3. For $k = N - 1$ Equation 3.4 gives $s_{N-1} = t_{-1} = 1/4$.

For all k different from $0, 1$, and $N - 1$, we have that $s_k = 0$. Clearly this gives the matrix in Equation (3.2). ♣

Proposition 3.2 is also useful when we have a circulant Toeplitz matrix S , and we want to find filter coefficients t_k so that $\mathbf{z} = S\mathbf{x}$ can be written on the form (3.3):

Example 3.4. Consider the matrix

$$S = \begin{pmatrix} 2 & 1 & 0 & 3 \\ 3 & 2 & 1 & 0 \\ 0 & 3 & 2 & 1 \\ 1 & 0 & 3 & 2 \end{pmatrix}.$$

This is a circulant Toeplitz matrix with $N = 4$, and we see that $s_0 = 2$, $s_1 = 3$, $s_2 = 0$, and $s_3 = 1$. The first equation in (3.4) gives that $t_0 = s_0 = 2$, and $t_1 = s_1 = 3$. The second equation in (3.4) gives that $t_{-2} = s_2 = 0$, and $t_{-1} = s_3 = 1$. By including only the t_k which are nonzero, the operation can be written as

$$z_n = t_{-1}x_{n-(-1)} + t_0x_n + t_1x_{n-1} + t_2x_{n-2} = x_{n+1} + 2x_0 + 3x_{n-1}.$$



What you should have learnt in this section

How to write down the circulant Toeplitz matrix from a digital filter expression, and vice versa. How to find the first column of this matrix (\mathbf{s}) from the filter coefficients (\mathbf{t}).

Exercises for Section 3.1

1. Assume that the filter S is defined by the formula

$$z_n = \frac{1}{4}x_{n+1} + \frac{1}{4}x_n + \frac{1}{4}x_{n-1} + \frac{1}{4}x_{n-2}.$$

Write down the filter coefficients t_k , and the matrix for S when $N = 8$.

2. Given the circulant Toeplitz matrix

$$S = \begin{pmatrix} 1 & 2 & 0 & 0 \\ 0 & 1 & 2 & 0 \\ 0 & 0 & 1 & 2 \\ 2 & 0 & 0 & 1 \end{pmatrix},$$

write down the filter coefficients t_k .

3. Assume that S is a circulant Toeplitz matrix so that only

$$S_{0,0}, \dots, S_{0,F} \text{ and } S_{0,N-E}, \dots, S_{0,N-1}$$

are nonzero on the first row, where E, F are given numbers. When implementing this filter on a computer we need to make sure that the vector indices lie in $[0, N-1]$. Show that $z_n = (\mathbf{S}\mathbf{x})_n$ can be split into the following different formulas, depending on n , to achieve this:

- a.** $0 \leq n < E$:

$$z_n = \sum_{k=0}^{n-1} S_{0,N+k-n}x_k + \sum_{k=n}^{F+n} S_{0,k-n}x_k + \sum_{k=N-1-E+n}^{N-1} S_{0,k-n}x_k. \quad (3.5)$$

- b.** $E \leq n < N - F$:

$$z_n = \sum_{k=n-E}^{n+F} S_{0,k-n}x_k. \quad (3.6)$$

- c.** $N - F \leq n < N$:

$$z_n = \sum_{k=0}^{n-(N-F)} S_{0,k-n}x_k + \sum_{k=n-E}^{n-1} S_{0,N+k-n}x_k + \sum_{k=n}^{N-1} S_{0,k-n}x_k. \quad (3.7)$$

From these three formulas we can write down a full implementation of the filter. This implementation is often more useful than writing down the entire matrix S , since we save computation when many of the matrix entries are zero.

3.2 Formal definition of filters and the vector frequency response

Let us now define digital filters formally, and establish their relationship to Toeplitz matrices. We have seen that a sound can be decomposed into different frequency components, and we would like to define filters as operations which adjust these frequency components in a predictable way. One such example is provided in Example 2.29, where we simply set some of the frequency components to 0. The natural starting point is to require for a filter that the output of a pure tone is a pure tone with the same frequency.

Definition 3.5 (Digital filters and vector frequency response). A linear transformation $S : \mathbb{R}^N \rightarrow \mathbb{R}^N$ is said to be a digital filter, or simply a filter, if it maps any Fourier vector in \mathbb{R}^N to a multiple of itself. In other words, for any integer n in the range $0 \leq n \leq N - 1$ there exists a value $\lambda_{S,n}$ so that

$$S(\phi_n) = \lambda_{S,n}\phi_n, \quad (3.8)$$

i.e., the N Fourier vectors are the eigenvectors of S . The vector of (eigen)values $\lambda_S = (\lambda_{S,n})_{n=0}^{N-1}$ is often referred to as the (vector) frequency response of S .

We will identify the linear transformation S with its matrix relative to the standard basis. Since the Fourier basis vectors are orthogonal vectors, S is clearly orthogonally diagonalizable. Since also the Fourier basis vectors are the columns in $(F_N)^H$, we have that

$$S = F_N^H D F_N \quad (3.9)$$

whenever S is a digital filter, where D has the frequency response (i.e. the eigenvalues) on the diagonal¹. In particular, if S_1 and S_2 are digital filters, we can write $S_1 = F_N^H D_1 F_N$ and $S_2 = F_N^H D_2 F_N$, so that

$$S_1 S_2 = F_N^H D_1 F_N F_N^H D_2 F_N = F_N^H D_1 D_2 F_N.$$

Since $D_1 D_2 = D_2 D_1$ for any diagonal matrices, we get the following corollary:

Corollary 3.6. The product of two digital filters is again a digital filter. Moreover, all digital filters commute, i.e. if S_1 and S_2 are digital filters, $S_1 S_2 = S_2 S_1$.

Clearly also $S_1 + S_2$ is a filter when S_1 and S_2 are. The set of all filters is thus a vector space, which also is closed under multiplication. Such a space is called an *algebra*. Since all filters commute, this algebra is also called a *commutative algebra*.

¹Recall that the orthogonal diagonalization of S takes the form $S = PDP^T$, where P contains as columns an orthonormal set of eigenvectors, and D is diagonal with the eigenvectors listed on the diagonal (see Section 7.1 in [20]).

There are several other equivalent characterizations of a digital filter as well. The next characterization helps us prove that the operations we started this chapter with actually are filters.

Theorem 3.7. A linear transformation S is a digital filter if and only if it is a circulant Toeplitz matrix.

Proof: That S is a filter is equivalent to the fact that $S = (F_N)^H D F_N$ for some diagonal matrix D . We observe that the entry at position (k, l) in S is given by

$$S_{k,l} = \frac{1}{N} \sum_{n=0}^{N-1} e^{2\pi i k n / N} \lambda_{S,n} e^{-2\pi i l n / N} = \frac{1}{N} \sum_{n=0}^{N-1} e^{2\pi i (k-l) n / N} \lambda_{S,n}.$$

Another entry on the same diagonal (shifted s rows and s columns) is

$$\begin{aligned} S_{(k+s) \bmod N, (l+s) \bmod N} &= \frac{1}{N} \sum_{n=0}^{N-1} e^{2\pi i ((k+s) \bmod N - (l+s) \bmod N) n / N} \lambda_{S,n} \\ &= \frac{1}{N} \sum_{n=0}^{N-1} e^{2\pi i (k-l) n / N} \lambda_{S,n} = S_{k,l}, \end{aligned}$$

which proves that S is a circulant Toeplitz matrix. Conversely, if S is a circulant Toeplitz matrix, it is an exercise to prove that $e^{2\pi i k n / N}$ are eigenvectors, so that S is a digital filter. ■

In particular, operations defined by (3.3) are digital filters, when restricted to vectors with period N . The following results enables us to compute the eigenvalues/frequency response easily, so that we do not need to form the characteristic polynomial and find its roots:

Theorem 3.8. Any digital filter is uniquely characterized by the values in the first column of its matrix. Moreover, if \mathbf{s} is the first column in S , the frequency response of S is given by

$$\boldsymbol{\lambda}_S = \sqrt{N} F_N \mathbf{s}. \quad (3.10)$$

Conversely, if we know the frequency response $\boldsymbol{\lambda}_S$, the first column \mathbf{s} of S is given by

$$\mathbf{s} = \frac{1}{\sqrt{N}} (F_N)^H \boldsymbol{\lambda}_S. \quad (3.11)$$

Proof: If we replace S by $(F_N)^H D F_N$ we find that

$$F_N \mathbf{s} = F_N S \begin{pmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix} = F_N F_N^H D F_N \begin{pmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix} = D F_N \begin{pmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix} = \frac{1}{\sqrt{N}} D \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix},$$

where we have used the fact that the first column in F_N has all entries equal to $1/\sqrt{N}$. But the diagonal matrix D has all the eigenvalues of S on its diagonal, and hence the last expression is the vector of eigenvalues λ_S , which proves (3.10). Equation (3.11) follows directly by applying the inverse DFT to (3.10). ■

The first column \mathbf{s} , which thus characterizes the filter, is also called the *impulse response*. This name stems from the fact that we can write $\mathbf{s} = S\mathbf{e}_0$, i.e. the vector \mathbf{s} is the output (often called response) to the vector \mathbf{e}_0 (often called an impulse). Equation (3.10) states that the frequency response can be written as

$$\lambda_{S,n} = \sum_{k=0}^{N-1} s_k e^{-2\pi i n k / N}, \quad \text{for } n = 0, 1, \dots, N-1, \quad (3.12)$$

where s_k are the components of \mathbf{s} .

Example 3.9. The identity matrix is a digital filter since $I = (F_N)^H I F_N$. Since \mathbf{e}_0 is the first column, it has impulse response $\mathbf{s} = \mathbf{e}_0$. Its frequency response has 1 in all components and therefore preserves all frequencies, as expected. ♣

Example 3.10. When only few of the coefficients s_k are nonzero, it is possible to obtain nice expressions for the frequency response. To see this, let us compute the frequency response of the filter defined from formula (3.1). We saw that the first column of the corresponding Toeplitz matrix satisfied $s_0 = 1/2$, and $s_{N-1} = s_1 = 1/4$. The frequency response is thus

$$\begin{aligned} \lambda_{S,n} &= \frac{1}{2} e^0 + \frac{1}{4} e^{-2\pi i n / N} + \frac{1}{4} e^{-2\pi i n(N-1) / N} \\ &= \frac{1}{2} e^0 + \frac{1}{4} e^{-2\pi i n / N} + \frac{1}{4} e^{2\pi i n / N} = \frac{1}{2} + \frac{1}{2} \cos(2\pi n / N). \end{aligned}$$



Equations (3.9), (3.10), and (3.11) are important relations between the matrix- and frequency representations of a filter. We see that the DFT is a crucial ingredient in these relations. A consequence is that, once you recognize a matrix as circulant Toeplitz, you do not need to make the tedious calculation of eigenvectors and eigenvalues which you are used to. Let us illustrate this with an example.

Example 3.11. Let us compute the eigenvalues and eigenvectors of the simple matrix

$$S = \begin{pmatrix} 4 & 1 \\ 1 & 4 \end{pmatrix}.$$

It is straightforward to compute the eigenvalues and eigenvectors of this matrix the way you learnt in your first course in linear algebra. However, this matrix is also a circulant Toeplitz matrix, so that we can use the results in this section to compute the eigenvalues and eigenvectors. Since here $N = 2$, we have that $e^{2\pi i n k / N} = e^{\pi i n k} = (-1)^{nk}$. This means that the Fourier basis vectors are $(1, 1)/\sqrt{2}$ and $(1, -1)/\sqrt{2}$, which also are the eigenvectors of S . The eigenvalues are the frequency response of S , which can be obtained as

$$\sqrt{N} F_N \mathbf{s} = \sqrt{2} \frac{1}{\sqrt{2}} \begin{pmatrix} 1 & 1 \\ 1 & -1 \end{pmatrix} \begin{pmatrix} 4 \\ 1 \end{pmatrix} = \begin{pmatrix} 5 \\ 3 \end{pmatrix}$$

The eigenvalues are thus 3 and 5. You could have obtained the same result with Matlab. Note that Matlab may not return the eigenvectors exactly as the Fourier basis vectors, since the eigenvectors are not unique (the multiple of an eigenvector is also an eigenvector). Matlab may here for instance switch the signs of the eigenvectors. We have no control over what Matlab actually chooses to do, since it uses some underlying numerical algorithm for computing eigenvectors which we can't influence.



In signal processing, the frequency content of a vector (i.e., its DFT) is also referred to as its spectrum. This may be somewhat confusing from a linear algebra perspective, because in this context the term spectrum is used to denote the eigenvalues of a matrix. But because of Theorem 3.8 this is not so confusing after all if we interpret the spectrum of a vector (in signal processing terms) as the spectrum of the corresponding digital filter (in linear algebra terms).

Certain vectors are easy to express in terms of the Fourier basis. This enables us to compute the output of such vectors from a digital filter easily, as the following example shows.

Example 3.12. Let us consider the filter S defined by $z_n = \frac{1}{6}(x_{n+2} + 4x_{n+1} + 6x_n + 4x_{n-1} + x_{n-2})$, and see how we can compute $S\mathbf{x}$ when

$$\mathbf{x} = (\cos(2\pi \cdot 0/N), \cos(2\pi \cdot 1/N), \dots, \cos(2\pi \cdot (N-1)/N)),$$

where N is the length of the vector. We note first that

$$\begin{aligned}\sqrt{N}\phi_5 &= \left(e^{2\pi i \cdot 0/N}, e^{2\pi i \cdot 1/N}, \dots, e^{2\pi i \cdot (N-1)/N}\right) \\ \sqrt{N}\phi_{N-5} &= \left(e^{-2\pi i \cdot 0/N}, e^{-2\pi i \cdot 1/N}, \dots, e^{-2\pi i \cdot (N-1)/N}\right),\end{aligned}$$

Since $e^{2\pi i \cdot 5k/N} + e^{-2\pi i \cdot 5k/N} = 2\cos(2\pi \cdot 5k/N)$, we get by adding the two vectors that $\mathbf{x} = \frac{1}{2}\sqrt{N}(\phi_5 + \phi_{N-5})$. Since the ϕ_n are eigenvectors, we have expressed \mathbf{x} as a sum of eigenvectors. The corresponding eigenvalues are given by the vector frequency response, so let us compute this. If $N = 8$, computing $S\mathbf{x}$ means to multiply with the 8×8 circulant Toeplitz matrix

$$\frac{1}{6} \begin{pmatrix} 6 & 4 & 1 & 0 & 0 & 0 & 1 & 4 \\ 4 & 6 & 4 & 1 & 0 & 0 & 0 & 1 \\ 1 & 4 & 6 & 4 & 1 & 0 & 0 & 0 \\ 0 & 1 & 4 & 6 & 4 & 1 & 0 & 0 \\ 0 & 0 & 1 & 4 & 6 & 4 & 1 & 0 \\ 0 & 0 & 0 & 1 & 4 & 6 & 4 & 1 \\ 1 & 0 & 0 & 0 & 1 & 4 & 6 & 4 \\ 4 & 1 & 0 & 0 & 0 & 1 & 4 & 6 \end{pmatrix}$$

We now see that

$$\begin{aligned}\lambda_{S,n} &= \frac{1}{6}(6 + 4e^{-2\pi i n/N} + e^{-2\pi i 2n/N} + e^{-2\pi i (N-2)n/N} + 4e^{-2\pi i (N-1)n/N}) \\ &= \frac{1}{6}(6 + 4e^{2\pi i n/N} + 4e^{-2\pi i n/N} + e^{2\pi i 2n/N} + e^{-2\pi i 2n/N}) \\ &= 1 + \frac{4}{3}\cos(2\pi n/N) + \frac{1}{3}\cos(4\pi n/N).\end{aligned}$$

The two values of this we need are

$$\begin{aligned}\lambda_{S,5} &= 1 + \frac{4}{3} \cos(2\pi 5/N) + \frac{1}{3} \cos(4\pi 5/N) \\ \lambda_{S,N-5} &= 1 + \frac{4}{3} \cos(2\pi(N-5)/N) + \frac{1}{3} \cos(4\pi(N-5)/N) \\ &= 1 + \frac{4}{3} \cos(2\pi 5/N) + \frac{1}{3} \cos(4\pi 5/N).\end{aligned}$$

Since these are equal, \mathbf{x} is a sum of eigenvectors with equal eigenvalues. This means that \mathbf{x} itself also is an eigenvector, with the same eigenvalue, so that

$$S\mathbf{x} = \left(1 + \frac{4}{3} \cos(2\pi 5/N) + \frac{1}{3} \cos(4\pi 5/N)\right) \mathbf{x}.$$

♣

3.2.1 Using digital filters to approximate analog filters

The formal definition of digital filters resembles that of analog filters, the difference being that the Fourier basis is now discrete. From this one may think that one can construct digital filters from analog filters. The following result clarifies this:

Theorem 3.13. Let s be an analog filter with frequency response $\lambda_s(f)$, and assume that $f \in V_{M,T}$ (so that also $s(f) \in V_{M,T}$). Let

$$\begin{aligned}\mathbf{x} &= (f(0 \cdot T/N), f(1 \cdot T/N), \dots, f((N-1)T/N)) \\ \mathbf{z} &= (s(f)(0 \cdot T/N), s(f)(1 \cdot T/N), \dots, s(f)((N-1)T/N))\end{aligned}$$

be vectors of $N = 2M + 1$ uniform samples from f and $s(f)$. Then the operation $S : \mathbf{x} \rightarrow \mathbf{z}$ (i.e. the operation which sends the samples of the input to the samples of the output) is well-defined on \mathbb{R}^N , and is an $N \times N$ -digital filter with frequency response $\lambda_{S,n} = \lambda_s(n/T)$.

Proof: With $N = 2M + 1$ we know that $f \in V_{M,T}$ is uniquely determined from \mathbf{x} . This means that $s(f)$ also is uniquely determined from \mathbf{x} , so that \mathbf{z} also is uniquely determined from \mathbf{x} . The operation $S : \mathbf{x} \rightarrow \mathbf{z}$ is therefore well-defined on \mathbb{R}^N .

Clearly also $s(e^{2\pi i n t / T}) = \lambda_s(n/T) e^{2\pi i n t / T}$. Since the samples of $e^{2\pi i n t / T}$ is the vector $e^{2\pi i k n / N}$, and the samples of $\lambda_s(n/T) e^{2\pi i n t / T}$ is $\lambda_s(n/T) e^{2\pi i k n / N}$, the vector $e^{2\pi i k n / N}$ is an eigenvector of S with eigenvalue $\lambda_s(n/T)$. Clearly then S is a digital filter with frequency response $\lambda_{S,n} = \lambda_s(n/T)$. ■

It is interesting that the digital frequency response above is obtained by sampling the analog frequency response. In this way we also see that it is easy to realize any digital filter as the restriction of an analog filter: any analog filter s will do where the frequency response has the values $\lambda_{S,n}$ at the points n/T . In the theorem it is essential that $f \in V_{M,T}$. There are many functions with the same samples, but where the samples of the output from the analog filter are different. When we restrict to $V_{M,T}$, however, the output samples are always determined from the input samples.

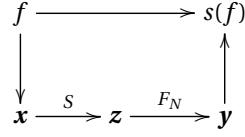


Figure 3.1: The connections between analog and digital filters, sampling and interpolation, provided by Theorem 3.13. The left vertical arrow represents sampling, the right vertical arrow represents interpolation.

Theorem 3.13 explains how digital filters can occur in practice. In the real world, a signal is modeled as a continuous function $f(t)$, and an operation on signals as an analog filter s . We can't compute the entire output $s(f)$ of the analog filter, but it is possible to apply the digital filter from Theorem 3.13 to the samples x of f . In general $f(t)$ may not lie in $V_{M,T}$, but we can denote by \tilde{f} the unique function in $V_{M,T}$ with the same samples as f (as in Section 2.5). By definition, Sx are the samples of $s(\tilde{f}) \in V_{M,T}$. $s(\tilde{f})$ can finally be found from these samples by using the procedure from Figure 2.2 for finding $s(\tilde{f})$. This procedure for finding $s(\tilde{f})$ is illustrated in Figure 3.1. Clearly, $s(\tilde{f})$ is an approximation to $s(f)$, since \tilde{f} is an approximation to f , and since s is continuous. Let us summarize this as follows:

Idea 3.14 (Approximating an analog filter). An analog filter s can be approximated through sampling, a digital filter, the DFT, and interpolation, as illustrated in Figure 3.1. S is the digital filter with frequency response $\lambda_{S,n} = \lambda_s(n/T)$. When $f \in V_{M,T}$, this approximation equals $s(f)$. When we increase the number of sample points/the size of the filter, the approximation becomes better. If there is a bound on the highest frequency in f , there exists an N so that when sampling of that size, the approximation equals $s(f)$.

Let us comment on why the last statements here are true. That the approximation equals $s(f)$ when $f \in V_{M,T}$ is obvious, since both f and $s(f) \in V_{M,T}$ are determined from their samples then. If there is a bound on the highest frequency in f , then f lies in $V_{M,T}$ for large enough M , so that we recover $s(f)$ as our approximation using $N = 2M + 1$. Finally, what happens when there is no bound on the highest frequency? We know that $s(f_N) = (s(f))_N$. Since f_N is a good approximation to f , the samples x of f are close to the samples of f_N . By continuity of the digital filter, $z = Sx$ will also be close to the samples of $(s(f))_N = s(f_N)$, so that (also by continuity) interpolating with z gives a good approximation to $(s(f))_N$, which is again a good approximation to $s(f)$. From this it follows that the digital filter is a better approximation when N is high.

What you should have learnt in this section

The formal definition of a digital filter in terms of having the Fourier vectors as eigenvectors. The definition of the vector frequency response in terms of the cor-

responding eigenvalues. Digital filters are circulant Toeplitz matrices. For filters, eigenvalues can be computed by taking the DFT of the first column \mathbf{s} , and there is no need to compute eigenvectors explicitly. How to apply a digital filter to a sum of sines or cosines, by splitting these into a sum of eigenvectors.

Exercises for Section 3.2

1. Let S be a circulant Toeplitz matrix. Show that the Fourier vectors $e^{2\pi i kn/N}$ are eigenvectors of S . This completes the proof of the converse part of Theorem 3.7.
2. Let S be a digital filter. Show that S is symmetric if and only if the frequency response satisfies $s_k = s_{N-k}$ for all k .
3. Consider the matrix
$$S = \begin{pmatrix} 4 & 1 & 3 & 1 \\ 1 & 4 & 1 & 3 \\ 3 & 1 & 4 & 1 \\ 1 & 3 & 1 & 4 \end{pmatrix}.$$
 - a. Compute the eigenvalues and eigenvectors of S using the results of this section. You should only need to perform one DFT in order to achieve this.
 - b. Verify the result from a. by computing the eigenvectors and eigenvalues the way you taught in your first course in linear algebra. This should be a much more tedious task.
 - c. Use Matlab to compute the eigenvectors and eigenvalues of S also. For some reason some of the eigenvectors seem to be different from the Fourier basis vectors, which you would expect from the theory in this section. Try to find an explanation for this.
4. Assume that S_1 and S_2 are two circulant Toeplitz matrices.
 - a. How can you express the eigenvalues of $S_1 + S_2$ in terms of the eigenvalues of S_1 and S_2 ?
 - b. How can you express the eigenvalues of $S_1 S_2$ in terms of the eigenvalues of S_1 and S_2 ?
 - c. If A and B are general matrices, can you find a formula which expresses the eigenvalues of $A + B$ and AB in terms of those of A and B ? If not, can you find a counterexample to what you found in a. and b.?
5. Consider the linear mapping S which keeps every second component in \mathbb{R}^N , i.e. $S(\mathbf{e}_{2k}) = \mathbf{e}_{2k}$, and $S(\mathbf{e}_{2k-1}) = \mathbf{0}$. Is S a digital filter?

3.3 The continuous frequency response and properties

If we make the substitution $\omega = 2\pi n/N$ in the formula for $\lambda_{S,n}$, we may interpret the frequency response as the values on a continuous function on $[0, 2\pi]$.

Theorem 3.15. The function $\lambda_S(\omega)$ defined on $[0, 2\pi]$ by

$$\lambda_S(\omega) = \sum_k t_k e^{-ik\omega}, \quad (3.13)$$

where t_k are the filter coefficients of S , satisfies

$$\lambda_{S,n} = \lambda_S(2\pi n/N) \text{ for } n = 0, 1, \dots, N-1$$

for any N . In other words, regardless of N , the vector frequency response lies on the curve λ_S .

Proof: For any N we have that

$$\begin{aligned} \lambda_{S,n} &= \sum_{k=0}^{N-1} s_k e^{-2\pi i nk/N} = \sum_{0 \leq k < N/2} s_k e^{-2\pi i nk/N} + \sum_{N/2 \leq k \leq N-1} s_k e^{-2\pi i nk/N} \\ &= \sum_{0 \leq k < N/2} t_k e^{-2\pi i nk/N} + \sum_{N/2 \leq k \leq N-1} t_{k-N} e^{-2\pi i nk/N} \\ &= \sum_{0 \leq k < N/2} t_k e^{-2\pi i nk/N} + \sum_{-N/2 \leq k \leq -1} t_k e^{-2\pi i n(k+N)/N} \\ &= \sum_{0 \leq k < N/2} t_k e^{-2\pi i nk/N} + \sum_{-N/2 \leq k \leq -1} t_k e^{-2\pi i nk/N} \\ &= \sum_{-N/2 \leq k < N/2} t_k e^{-2\pi i nk/N} = \lambda_S(\omega). \end{aligned}$$

where we have used Equation (3.4). ■

Both $\lambda_S(\omega)$ and $\lambda_{S,n}$ will be referred to as frequency responses in the following. To distinguish the two, while $\lambda_{S,n}$ is called the vector frequency response of S , $\lambda_S(\omega)$ is called the *continuous frequency response* of S . ω is called *angular frequency*.

The difference in the definition of the continuous- and the vector frequency response lies in that one uses the filter coefficients t_k , while the other uses the impulse response s_k . While these contain the same values, they are ordered differently. Had we used the impulse response to define the continuous frequency response, we would have needed to compute $\sum_{k=0}^{N-1} s_k e^{-\pi i \omega k}$, which does not converge when $N \rightarrow \infty$ (although it gives the right values at all points $\omega = 2\pi n/N$ for all N !). The filter coefficients avoid this convergence problem, however, since we assume that only t_k with $|k|$ small are nonzero. In other words, filter coefficients are used in the definition of the continuous frequency response so that we can find a continuous curve where we can find the vector frequency response values for all N .

The frequency response contains the important characteristics of a filter, since it says how it behaves for the different frequencies. When analyzing a filter, we therefore often plot the frequency response. Often we plot only the absolute value (or

the magnitude) of the frequency response, since this is what explains how each frequency is amplified or attenuated. Since λ_S is clearly periodic with period 2π , we may restrict angular frequency to the interval $[0, 2\pi)$. The conclusion in Observation 2.28 was that the low frequencies in a vector correspond to DFT indices close to 0 and $N - 1$, and high frequencies correspond to DFT indices close to $N/2$. This observation is easily translated to a statement about angular frequencies:

Observation 3.16. When plotting the frequency response on $[0, 2\pi)$, angular frequencies near 0 and 2π correspond to low frequencies, angular frequencies near π correspond to high frequencies

λ_S may also be viewed as a function defined on the interval $[-\pi, \pi]$. Plotting on $[-\pi, \pi]$ is often done in practice, since it makes clearer what corresponds to lower frequencies, and what corresponds to higher frequencies:

Observation 3.17. When plotting the frequency response on $[-\pi, \pi]$, angular frequencies near 0 correspond to low frequencies, angular frequencies near $\pm\pi$ correspond to high frequencies.

The following holds:

Theorem 3.18. Assume that s is an analog filter, and that we sample a periodic function at rate f_s over one period, and denote the corresponding digital filter by S . The analog and digital frequency responses are related by $\lambda_s(f) = \lambda_S(2\pi f f_s)$.

To see this, note first that S has frequency response $\lambda_{S,n} = \lambda_s(n/T) = \lambda_s(f)$, where $f = n/T$. We then rewrite $\lambda_{S,n} = \lambda_S(2\pi n/N) = \lambda_S(2\pi f T/N) = \lambda_S(2\pi f f_s)$.

Example 3.19. In Example 3.10 we computed the vector frequency response of the filter defined in formula (3.1). The filter coefficients are here $t_{-1} = 1/4$, $t_0 = 1/2$, and $t_1 = 1/4$. The continuous frequency response is thus

$$\lambda_S(\omega) = \frac{1}{4}e^{i\omega} + \frac{1}{2} + \frac{1}{4}e^{-i\omega} = \frac{1}{2} + \frac{1}{2}\cos\omega.$$

Clearly this matches the computation from Example 3.10. Figure 3.2 shows plots of this frequency response, plotted on the intervals $[0, 2\pi)$ and $[-\pi, \pi]$. Both the continuous frequency response and the vector frequency response for $N = 51$ are shown. Figure (b) shows clearly how the high frequencies are softened by the filter. ♣

Since the frequency response is essentially a DFT, it inherits several properties from Theorem 2.21. We will mostly use the continuous frequency response to express these properties.

Theorem 3.20. We have that

1. The continuous frequency response satisfies $\lambda_S(-\omega) = \overline{\lambda_S(\omega)}$.

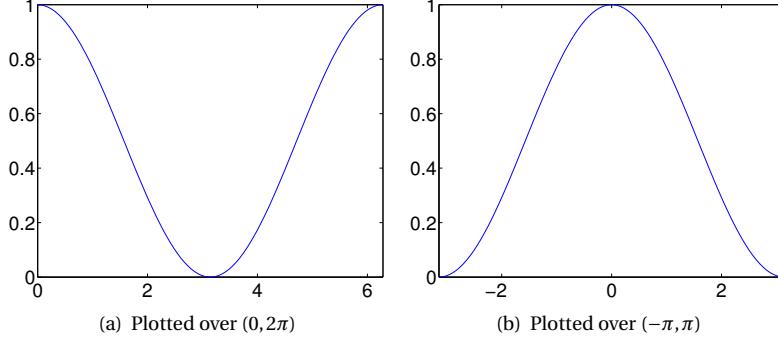


Figure 3.2: The (absolute value of the) frequency response of the moving average filter of Formula (3.1) from the beginning of this chapter.

2. If S is a digital filter, S^T is also a digital filter. Moreover, if the frequency response of S is $\lambda_S(\omega)$, then the frequency response of S^T is $\overline{\lambda_S(\omega)}$.
3. If S is symmetric, λ_S is real. Also, if S is antisymmetric (the element on the opposite side of the diagonal is the same, but with opposite sign), λ_S is purely imaginary.
4. A digital filter S is invertible if and only if $\lambda_{S,n} \neq 0$ for all n . In that case S^{-1} is also a digital filter, and $\lambda_{S^{-1},n} = 1/\lambda_{S,n}$.
5. If S_1 and S_2 are digital filters, then $S_1 S_2$ also is a digital filter, and $\lambda_{S_1 S_2}(\omega) = \lambda_{S_1}(\omega) \lambda_{S_2}(\omega)$.

Proof: Property 1. and 3. follow directly from Theorem 2.21. Transposing a matrix corresponds to reversing the first column of the matrix and thus also the filter coefficients. Due to this Property 2. also follows from Theorem 2.21. If $S = (F_N)^H D F_N$, and all $\lambda_{S,n} \neq 0$, we have that $S^{-1} = (F_N)^H D^{-1} F_N$, where D^{-1} is a digonal matrix with the values $1/\lambda_{S,n}$ on the diagonal. Clearly then S^{-1} is also a digital filter, and its frequency response is $\lambda_{S^{-1},n} = 1/\lambda_{S,n}$, which proves 4. The last property follows in the same was as we showed that filters commute:

$$S_1 S_2 = (F_N)^H D_1 F_N (F_N)^H D_2 F_N = (F_N)^H D_1 D_2 F_N.$$

The frequency response of $S_1 S_2$ is thus obtained by multiplying the frequency responses of S_1 and S_2 . ■

In particular the frequency response may not be real, although this was the case in the first example of this section. Theorem 3.20 applies also for the vector frequency response. Since the vector frequency response are the eigenvalues of the filter, the last property above says that, for filters, multiplication of matrices corre-

sponds to multiplication of eigenvalues. Clearly this is an important property which is shared with all other matrices which have the same eigenvectors.

Example 3.21. Assume that the filters S_1 and S_2 have the frequency responses $\lambda_{S_1}(\omega) = \cos(2\omega)$, $\lambda_{S_2}(\omega) = 1 + 3\cos\omega$. Let us see how we can use Theorem 3.20 to compute the filter coefficients and the matrix of the filter $S = S_1S_2$. We first notice that, since both frequency responses are real, all S_1 , S_2 , and $S = S_1S_2$ are symmetric. We rewrite the frequency responses as

$$\begin{aligned}\lambda_{S_1}(\omega) &= \frac{1}{2}(e^{2i\omega} + e^{-2i\omega}) = \frac{1}{2}e^{2i\omega} + \frac{1}{2}e^{-2i\omega} \\ \lambda_{S_2}(\omega) &= 1 + \frac{3}{2}(e^{i\omega} + e^{-i\omega}) = \frac{3}{2}e^{i\omega} + 1 + \frac{3}{2}e^{-i\omega}.\end{aligned}$$

We now get that

$$\begin{aligned}\lambda_{S_1S_2}(\omega) &= \lambda_{S_1}(\omega)\lambda_{S_2}(\omega) = \left(\frac{1}{2}e^{2i\omega} + \frac{1}{2}e^{-2i\omega}\right)\left(\frac{3}{2}e^{i\omega} + 1 + \frac{3}{2}e^{-i\omega}\right) \\ &= \frac{3}{4}e^{3i\omega} + \frac{1}{2}e^{2i\omega} + \frac{3}{4}e^{i\omega} + \frac{3}{4}e^{-i\omega} + \frac{1}{2}e^{-2i\omega} + \frac{3}{4}e^{-3i\omega}\end{aligned}$$

From this expression we see that the filter coefficients of S are $t_{\pm 1} = 3/4$, $t_{\pm 2} = 1/2$, $t_{\pm 3} = 3/4$. All other filter coefficients are 0. Using Theorem 3.2, we get that $s_1 = 3/4$, $s_2 = 1/2$, and $s_3 = 3/4$, while $s_{N-1} = 3/4$, $s_{N-2} = 1/2$, and $s_{N-3} = 3/4$ (all other s_k are 0). This gives us the matrix representation of S . ♣

3.3.1 Windowing operations

In this section we will take a look at a very important, and perhaps surprising, application of the continuous frequency response. Let us return to the computations from Example 2.29. There we saw that, when we restricted to a block of the signal, this affected the frequency representation. If we substitute with the angular frequencies $\omega = 2\pi n/N$ and $\omega_0 = 2\pi n_0/M$ in Equation (2.10), we get

$$y_n = \frac{1}{N} \sum_{k=0}^{N-1} e^{ik\omega_0} e^{-ik\omega} = \frac{1}{N} \sum_{k=0}^{N-1} e^{-ik(\omega - \omega_0)}$$

(here y_n were the DFT components of the sound after we had restricted to a block). This expression states that, when we restrict to a block of length N in the signal by discarding the other samples, a pure tone of angular frequency ω_0 suddenly gets a frequency contribution at angular frequency ω also, and the contribution is given by this formula. The expression is seen to be the same as the frequency response of the filter $\frac{1}{N}\{1, 1, \dots, 1\}$ (where 1 is repeated N times), evaluated at $\omega - \omega_0$. This filter is nothing but a (delayed) moving average filter. The frequency response of a moving average filter thus governs how the different frequencies pollute when we limit ourselves to a block of the signal. Since this frequency response has its peak at 0, angular frequencies ω close to ω_0 have biggest values, so that the pollution is mostly from frequencies close to ω_0 . But unfortunately, other frequencies also pollute.

One can also ask the question if there are better ways to restrict to a block of size N of the signal. We formulate the following idea.

Idea 3.22. Let (x_0, \dots, x_M) be a sound of length M . We would like to find values $\mathbf{w} = \{w_0, \dots, w_{N-1}\}$ so that the new sound $(w_0 x_0, \dots, w_{N-1} x_{N-1})$ of length N (i.e. where the samples are attenuated by the window samples, and where samples have been discarded) has a frequency representation as close to \mathbf{x} as possible. The vector \mathbf{w} is called a window of length N , and the new sound is called the windowed signal.

Above we encountered the window $\mathbf{w} = \{1, 1, \dots, 1\}$. This is called the rectangular window. To see how we can find a good window, note first that the DFT values in the windowed signal of length N is

$$y_n = \frac{1}{N} \sum_{k=0}^{N-1} w_k e^{ik\omega_0} e^{-ik\omega} = \frac{1}{N} \sum_{k=0}^{N-1} w_k e^{-ik(\omega-\omega_0)}.$$

This is the frequency response of $\frac{1}{N} \mathbf{w}$. In order to limit the pollution from other frequencies, we thus need to construct a window with a frequency response with smaller values than that of the rectangular window away from 0. Let us summarize our findings as follows:

Observation 3.23. Assume that we would like to construct a window of length N . It is desirable that the frequency response of the window has small values away from zero.

We will not go into techniques for how such frequency responses can be constructed, but only consider one example different from the rectangular window. We define the Hamming window by

$$w_n = 2(0.54 - 0.46 \cos(2\pi n / (N - 1))). \quad (3.14)$$

The frequency responses of the rectangular window and the Hamming window are compared in Figure 3.3 for $N = 32$. We see that the Hamming window has much smaller values away from 0, so that it is better suited as a window. However, the width of the “main lobe” (i.e. the main structure at the center), seems to be bigger. The window coefficients themselves are shown in Figure 3.4. It is seen that the frequency response of the Hamming window attenuates more and more as we get close to the boundaries. Many other windows are used in the literature. The concrete window from Exercise 5 is for instance used in the MP3 standard. It is applied to the sound, and after this an FFT is applied to the windowed sound in order to make a frequency analysis of that part of the sound. The effect of the window is that there is smaller loss in the frequency representation of the sound when we restrict to a block of sound samples. This is a very important part of the psychoacoustic model used in the MP3 encoder, since it has to make compression decisions based on the frequency information in the sound.

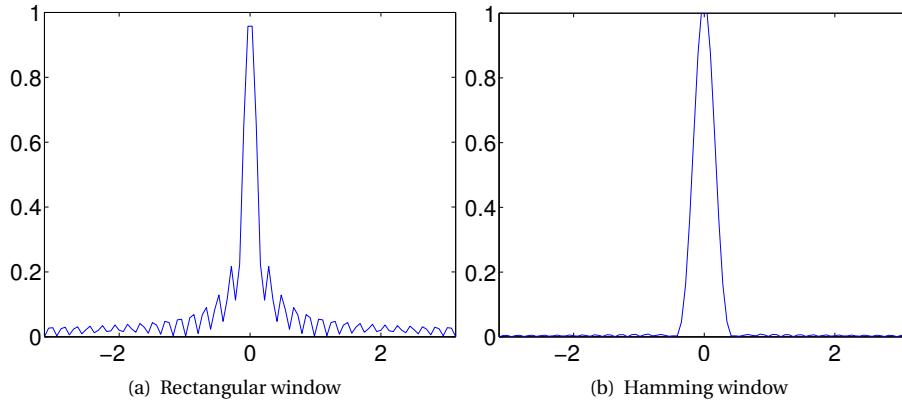


Figure 3.3: The frequency responses of the windows we have considered for restricting to a block of the signal.

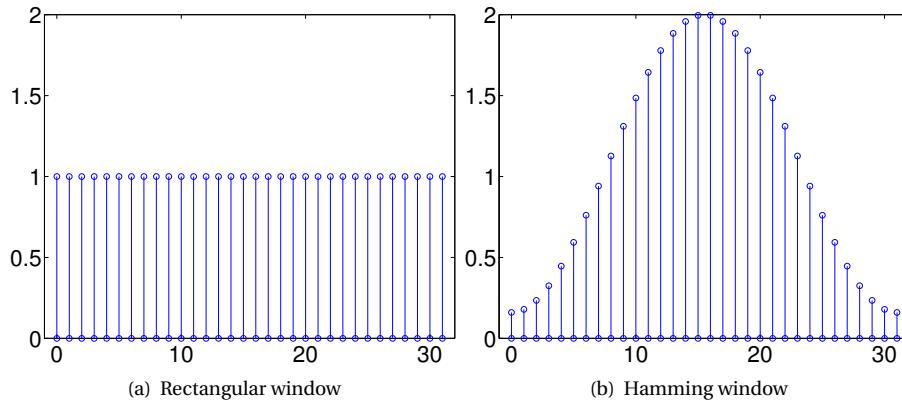


Figure 3.4: The window coefficients we have considered for restricting to a block of the signal.

What you should have learnt in this section

The definition of the continuous frequency response in terms of the filter coefficients t . Connection with the vector frequency response. Properties of the continuous frequency response. How to compute the frequency response of the product of two filters.

Exercises for Section 3.3

- 1.** Let again S be the filter defined by the equation

$$z_n = \frac{1}{4}x_{n+1} + \frac{1}{4}x_n + \frac{1}{4}x_{n-1} + \frac{1}{4}x_{n-2},$$

as in Exercise 3.1.1. Compute and plot (the magnitude of) $\lambda_S(\omega)$.

- 2.** A filter S is defined by the equation

$$z_n = \frac{1}{3}(x_n + 3x_{n-1} + 3x_{n-2} + x_{n-3}).$$

a. Compute and plot the (magnitude of the continuous) frequency response of the filter, i.e. $|\lambda_S(\omega)|$. Is the filter a lowpass filter or a highpass filter?

b. Find an expression for the vector frequency response $\lambda_{S,2}$. What is $S\mathbf{x}$ when \mathbf{x} is the vector of length N with components $e^{2\pi i 2k/N}$?

- 3.** A filter S_1 is defined by the equation

$$z_n = \frac{1}{16}(x_{n+2} + 4x_{n+1} + 6x_n + 4x_{n-1} + x_{n-2}).$$

a. Write down an 8×8 circulant Toeplitz matrix which corresponds to applying S_1 on a periodic signal with period $N = 8$.

b. Compute and plot (the continuous) frequency response of the filter. Is the filter a lowpass filter or a highpass filter?

c. Another filter S_2 has (continuous) frequency response $\lambda_{S_2}(\omega) = (e^{i\omega} + 2 + e^{-i\omega})/4$. Write down the filter coefficients for the filter S_1S_2 .

- 4.** Assume that the filters S_1 and S_2 have the frequency responses $\lambda_{S_1}(\omega) = 2 + 4\cos(\omega)$, $\lambda_{S_2}(\omega) = 3\sin(2\omega)$.

a. Compute and plot the frequency response of the filter S_1S_2 .

b. Write down the filter coefficients t_k and the impulse response \mathbf{s} for the filter S_1S_2 .

- 5.** The Hanning window is defined by $w_n = 1 - \cos(2\pi n/(N-1))$. Compute and plot the window coefficients and the continuous frequency response of this window for $N = 32$, and compare with the window coefficients and the frequency responses for the rectangular- and the Hamming window.

3.4 Assembling the filter matrix and compact notation

Let us return to how we first defined a filter in Equation (3.3):

$$z_n = \sum_k t_k x_{n-k}.$$

As mentioned, the range of k may not be specified. In some applications in signal processing there may in fact be infinitely many nonzero t_k . However, when \mathbf{x} is assumed to have period N , we may as well assume that the k 's range over an interval of length N (else many of the t_k 's can be added together to simplify the formula). Also, any such interval can be chosen. It is common to choose the interval so that it is centered around 0 as much as possible. For this, we can choose for instance $[|N/2| - N + 1, |N/2|]$. With this choice we can write Equation (3.3) as

$$z_n = \sum_{k=|N/2|-N+1}^{|N/2|} t_k x_{n-k}. \quad (3.15)$$

The index range in this sum is typically even smaller, since often much less than N of the t_k are nonzero (For Equation (3.1), there were only three nonzero t_k). In such cases one often uses a more compact notation for the filter:

Definition 3.24 (Compact notation for filters). Let k_{\min}, k_{\max} be the smallest and biggest index of a filter coefficient in Equation (3.15) so that $t_k \neq 0$ (if no such values exist, let $k_{\min} = k_{\max} = 0$), i.e.

$$z_n = \sum_{k=k_{\min}}^{k_{\max}} t_k x_{n-k}. \quad (3.16)$$

We will use the following compact notation for S :

$$S = \{t_{k_{\min}}, \dots, t_{-1}, \underline{t_0}, t_1, \dots, t_{k_{\max}}\}.$$

In other words, the entry with index 0 has been underlined, and only the nonzero t_k 's are listed. k_{\max} and k_{\min} are also called the start and end indices of S . By the length of S , denoted $l(S)$, we mean the number $k_{\max} - k_{\min}$.

One seldom writes out the matrix of a filter, but rather uses this compact notation.

Example 3.25. Using the compact notation for a filter, we would write $S = \{1/4, \underline{1/2}, 1/4\}$ for the filter given by formula (3.1)). For the filter

$$z_n = x_{n+1} + 2x_0 + 3x_{n-1}$$

from Example 3.4, we would write $S = \{1, \underline{2}, 3\}$. ♣

Applying a filter to a vector \mathbf{x} is also called taking the convolution of the two vectors \mathbf{t} and \mathbf{x} . Convolution is usually defined without the assumption that the

input vector is periodic, and without any assumption on the vector lengths (i.e. they may be sequences of infinite length):

Definition 3.26 (Convolution of vectors). By the *convolution* of two vectors \mathbf{x} and \mathbf{y} we mean the vector $\mathbf{x} * \mathbf{y}$ defined by

$$(\mathbf{x} * \mathbf{y})_n = \sum_k x_k y_{n-k}. \quad (3.17)$$

If both \mathbf{x} and \mathbf{y} have infinitely many nonzero entries, the sum is an infinite one, and may diverge.

The case where both vectors \mathbf{x} and \mathbf{y} have a finite number of nonzero elements deserves extra attention. Assume that x_0, \dots, x_{N-1} and y_0, \dots, y_{M-1} are the only nonzero elements. If $\mathbf{z} = \mathbf{x} * \mathbf{y}$, it is clear from Equation (3.17) that only the elements z_0, \dots, z_{M+N-2} can be nonzero. The convolution operation is therefore fully represented by the finite-dimensional operation from $\mathbb{R}^N \times \mathbb{R}^M \rightarrow \mathbb{R}^{M+N-1}$ defined by

$$(x_0, \dots, x_{N-1}) \times (y_0, \dots, y_{M-1}) \rightarrow (z_0, \dots, z_{M+N-2}). \quad (3.18)$$

Matlab has the built-in function `conv` for performing this operation. The `conv` function thus considers the convolution in terms of the (finite) nonzero parts of the vectors, without keeping track of the start and end indices. Exercise 7 explains how one may keep track of these indices. Put differently, Matlab's `conv`-operation pads our original values with zeros in both directions, instead of considering a periodic input vector. Due to its simplicity it is used much in practice, but it is not exactly the same as applying a digital filter. The following result states that the difference between convolution and filtering is only at the start and end of the vector.

Theorem 3.27. Let $S = \{t_{-E}, \dots, t_0, \dots, t_F\}$ where $E, F \geq 0$, and let $\mathbf{x} = (x_0, \dots, x_{N-1})$. We have that $\mathbf{t} * \mathbf{x}$ equals $S\mathbf{x}$ at the indices from F to $N - E - 1$.

Proof. We have that $(S\mathbf{x})_n = t_{-E}x_{n+E} + \dots + t_Fx_{n-F}$. Here the indices of \mathbf{x} lie between 0 and $N - 1$ if and only if $n + E \leq N - 1$ and $n - F \geq 0$, which happen if $F \leq n \leq N - E - 1$. We therefore do not access the added zeros outside $[0, N - 1]$, so that the result is equal. \blacksquare

We also have a very important connection between convolution and polynomials

Proposition 3.28. Assume that $p(x) = a_Nx^N + a_{N-1}x^{N-1} + \dots + a_1x + a_0$ and $q(x) = b_Mx^M + b_{M-1}x^{M-1} + \dots + b_1x + b_0$ are polynomials of degree N and M respectively. Then the coefficients of the polynomial pq can be obtained by computing `conv(a, b)`.

We can thus interpret a filter as a polynomial. In this setting, clearly the length $l(S)$ of the filter can be interpreted as the degree of the polynomial. Clearly also,

this polynomial is the frequency response, when we insert $e^{i\omega}$ for the variable. Also, applying two filters in succession is equivalent to applying the convolution of the filters, so that two filtering operations can be combined to one.

Since the number of nonzero filter coefficients is typically much less than N (the period of the input vector), the matrix S have many entries which are zero. Multiplication with such matrices requires less additions and multiplications than for other matrices: If S has k nonzero filter coefficients, S has Nk nonzero entries, so that kN multiplications and $(k-1)N$ additions are needed to compute Sx . This is much less than the N^2 multiplications and $(N-1)N$ additions needed in the general case. Perhaps more important is that we need not form the entire matrix, we can perform the matrix multiplication directly in a loop. Exercise 3 investigates this further. For large N we risk running into out of memory situations if we had to form the entire matrix.

What you should have learnt in this section

The compact filter notation for filters with a finite number of filter coefficients. The Matlab `conv`-operation for computing convolution, and thus apply a filter to a vector.

Exercises for Section 3.4

1. Compute and plot the continuous frequency response of the filter $S = \{1/4, 1/2, 1/4\}$. Where does the frequency response achieve its maximum and minimum value, and what are these values?
2. Plot the continuous frequency response of the filter $T = \{1/4, -1/2, 1/4\}$. Where does the frequency response achieve its maximum and minimum value, and what are these values? Can you write down a connection between this frequency response and that from Exercise 1?
3. Define the filter S by $S = \{1, 2, 3, 4, 5, 6\}$. Write down the matrix for S when $N = 8$. Plot (the magnitude of) $\lambda_S(\omega)$, and indicate the values $\lambda_{S,n}$ for $N = 8$ in this plot.
4. Given the circulant Toeplitz matrix

$$S = \frac{1}{5} \begin{pmatrix} 1 & 1 & 1 & \cdots & 1 \\ 1 & 1 & 1 & \cdots & 0 \\ 0 & 1 & 1 & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \cdots & 1 \\ 1 & 0 & 0 & \cdots & 1 \\ 1 & 1 & 0 & \cdots & 1 \\ 1 & 1 & 1 & \cdots & 1 \end{pmatrix}$$

Write down the filter coefficients of this matrix, and use the compact notation $\{t_{k_{min}}, \dots, t_{-1}, \underline{t_0}, t_1, \dots, t_{k_{max}}\}$. Compute and plot (the magnitude) of $\lambda_S(\omega)$.

5. Assume that $S = \{1, c, c^2, \dots, c^k\}$. Compute and plot $\lambda_S(\omega)$ when $k = 4$ and $k = 8$. How does the choice of k influence the frequency response? How does the choice of c influence the frequency response?

6. Compute the convolution of $\{1, 2, 1\}$ with itself. interpret the result in terms of two polynomials.

7. In this exercise we will find out how to keep track of the length and the start and end indices when we convolve two sequences.

a. Let \mathbf{x} be zero outside x_a, \dots, x_{a+N-1} , and \mathbf{y} be zero outside y_b, \dots, y_{b+M-1} . Show that $\mathbf{z} = \mathbf{x} * \mathbf{y}$ is zero outside $z_{a+b}, \dots, z_{a+b+M+N-2}$. Explain why this means that $l(\mathbf{x} * \mathbf{y}) = l(\mathbf{x}) + l(\mathbf{y})$ for general vectors.

b. Find expressions for the start- and end indices k_{min}, k_{max} for $\mathbf{x} * \mathbf{y}$, in terms of those of \mathbf{x} and \mathbf{y} .

8. Implement a Matlab function

```
function z=convimpl(x,y)
```

which from input vectors of dimension N and M returns an output vector of dimension $N + M - 1$, as dictated by Equation (3.18). Test your function together with the built-in Matlab conv-function to verify that the give the same result.

9. Show that if $S = \{t_0, \dots, t_F\}$ and $\mathbf{x} \in \mathbb{R}^N$, then $S \begin{pmatrix} \mathbf{x} \\ \mathbf{0}_F \end{pmatrix} = \mathbf{t} * \mathbf{x}$. Thus if we add zeros in a vector, filtering and convolution are the same.

3.5 Some examples of filters

We have now established the basic theory of filters, and it is time to study some specific examples.

Example 3.29 (Time delay filters). The simplest possible type of Toeplitz matrix is one where there is only one nonzero diagonal. Let us define the Toeplitz matrix E_d as the one which has first column equal to \mathbf{e}_d . In the notation above, and when $d > 0$, this filter can also be written as $S = \{0, \dots, 1\}$ where the 1 occurs at position d . We observe that

$$(E_d \mathbf{x})_n = \sum_{k=0}^{N-1} (E_d)_{n,k} x_k = \sum_{k=0}^{N-1} (E_d)_{(n-k) \bmod N, 0} x_k = x_{n-d},$$

since only when $(n - k) \bmod N = d$ do we have a contribution in the sum. It is thus clear that multiplication with E_d delays a vector by d samples. For this reason E_d is also called a *time delay filter*. The frequency response of the time delay filter is clearly the function $\lambda_S(\omega) = e^{-id\omega}$, which has magnitude 1. This filter therefore does not change the magnitude of the different frequencies. ♣

Example 3.30 (Adding echo). An echo is a copy of the sound that is delayed and softer than the original sound. The sample that comes t seconds before sample i has index $i - tf_s$ where f_s is the sampling rate. This also makes sense even if t is not an integer so we can use this to produce delays that are less than one second. The one complication with this is that the number tf_s may not be an integer. We can get round this by rounding it to the nearest integer. This corresponds to adjusting the echo slightly. The following holds:

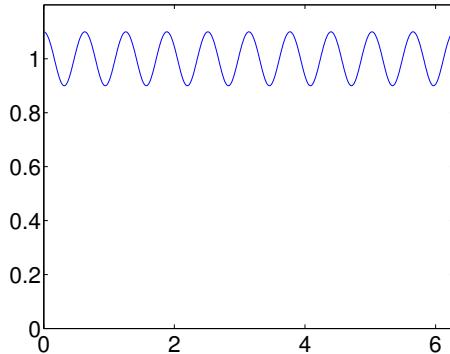


Figure 3.5: The frequency response of a filter which adds an echo with damping factor $c = 0.1$ and delay $d = 10$.

Fact 3.31. Let (x, s) be a digital sound. Then the sound y with samples given by

```
y=[x zeros(1,d)];  
y((d+1):N)=x((d+1):N)+c*x(1:(N-d));
```

will include an echo of the original sound. Here $d=\text{round}(ms)$ is the integer closest to $t f_s$, and c is a constant which is usually smaller than 1.

This is an example of a filtering operation where each output element is constructed from two input elements. As in the case of noise it is important to dampen the part that is added to the original sound, otherwise the echo will be too loud. Note also that the formula that creates the echo is not used at the beginning of the signal, since it is not audible until after d samples. Also, the echo is un audible if d is too small. You can listen to the sample file with echo added with $d = 10000$ and $c = 0.5$ here.

Using our compact filter notation, the filter which adds echo can be written as

$$S = \{1, 0, \dots, 0, c\},$$

where the damping factor c appears after the delay d . The frequency response of this is $\lambda_S(\omega) = 1 + ce^{-id\omega}$. This frequency response is not real, which means that the filter is not symmetric. In Figure 3.5 we have plotted the magnitude of this frequency response with $c = 0.1$ and $d = 10$. We see that the response varies between 0.9 and 1.1, so that the deviation from 1 is controlled by the damping factor c . Also, we see that the oscillation in the frequency response, as visible in the plot, is controlled by the delay d . ♣

Let us now take a look at some filters which adjust the bass and treble in sound. The fact that the filters are useful for these purposes will be clear when we plot the frequency response.

Example 3.32 (Reducing the treble with moving average filters). The treble in a sound is generated by the fast oscillations (high frequencies) in the signal. If we want to reduce the treble we have to adjust the sample values in a way that reduces those fast oscillations. A general way of reducing variations in a sequence of numbers is to replace one number by the average of itself and its neighbours, and this is easily done with a digital sound signal. If we let the new sound signal be $\mathbf{z} = (z_i)_{i=0}^{N-1}$ we can compute it as

```

z(1)=(x(N)+x(1)+x(2))/3;
for t=2: (N-1)
    z(t)=(x(t-1)+x(t)+x(t+1))/3;
end
z(N)=(x(N-1)+x(N)+x(1))/3;

```

In Example 3.30 we did not take into account that the signal is assumed periodic. Above this has been taken into account through the addition of the first and last lines (which correspond to the circulating part of the matrix). This filter is also called a *moving average filter*, and it can be written in compact form as

$$S = \left\{ \frac{1}{3}, \frac{1}{3}, \frac{1}{3} \right\}.$$

If we set $N = 4$, the corresponding circulant Toeplitz matrix for the filter is

$$S = \frac{1}{3} \begin{pmatrix} 1 & 1 & 0 & 1 \\ 1 & 1 & 1 & 0 \\ 0 & 1 & 1 & 1 \\ 1 & 0 & 1 & 1 \end{pmatrix}$$

The frequency response is $\lambda_S(\omega) = (e^{i\omega} + 1 + e^{-i\omega})/3 = (1 + 2\cos(\omega))/3$. More generally we can construct the moving average filter of $2L + 1$ elements, which is $S = \{1, \dots, \underline{1}, \dots, 1\}/(2L + 1)$, where there is symmetry around 0. In Example 2.18 we computed that the DFT of the vector $\{1, \dots, \underline{1}, \dots, 1\}$ was

$$\mathbf{x} = \frac{1}{\sqrt{N}} \frac{\sin(\pi n(2L + 1)/N)}{\sin(\pi n/N)}.$$

Due to this the frequency response of S is

$$\lambda_{S,n} = \frac{1}{2L + 1} \frac{\sin(\pi n(2L + 1)/N)}{\sin(\pi n/N)},$$

and thus

$$\lambda_S(\omega) = \frac{1}{2L + 1} \frac{\sin((2L + 1)\omega/2)}{\sin(\omega/2)}.$$

We clearly have

$$0 \leq \frac{1}{2L + 1} \frac{\sin((2L + 1)\omega/2)}{\sin(\omega/2)} \leq 1,$$

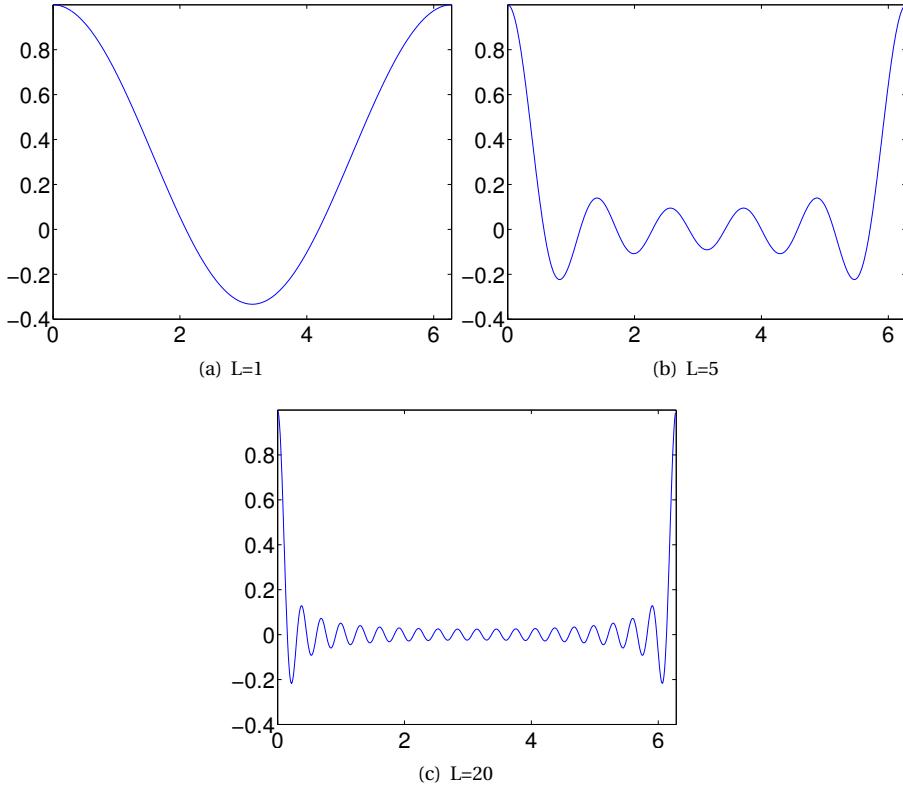


Figure 3.6: The frequency response of moving average filters of different length.

and this frequency response approaches 1 as $\omega \rightarrow 0$. The frequency response thus peaks at 0, and it is clear that this peak gets narrower and narrower as L increases, i.e. we use more and more samples in the averaging process. This appeals to our intuition that this kind of filters smooth the sound by keeping only lower frequencies. In Figure 3.6 we have plotted the frequency response for moving average filters with $L = 1$, $L = 5$, and $L = 20$. We see, unfortunately, that the frequency response is far from a filter which keeps some frequencies unaltered, while annihilating others (this is a desirable property): Although the filter distinguishes between high and low frequencies, it slightly changes the small frequencies. Moreover, the higher frequencies are not annihilated, even when we increase L to high values. ♣

In the previous example we mentioned filters which favour certain frequencies of interest, while annihilating the others. This is a desirable property for filters, so let us give names to such filters:

Definition 3.33. A filter S is called

1. a *lowpass filter* if $\lambda_S(\omega)$ is large when ω is close to 0, and $\lambda_S(\omega) \approx 0$ when ω

is close to π (i.e. S keeps low frequencies and annihilates high frequencies),

2. a *highpass filter* if $\lambda_S(\omega)$ is large when ω is close to π , and $\lambda_S(\omega) \approx 0$ when ω is close to 0 (i.e. S keeps high frequencies and annihilates low frequencies),
3. a *bandpass filter* if $\lambda_S(\omega)$ is large within some interval $[a, b] \subset [0, 2\pi]$, and $\lambda_S(\omega) \approx 0$ outside this interval.

This definition should be considered rather vague when it comes to what we mean by “ ω close to $0, \pi$ ”, and “ $\lambda_S(\omega)$ is large”: in practice, when we talk about low-pass and highpass filters, it may be that the frequency responses are still quite far from what is commonly referred to as *ideal lowpass or highpass filters*, where the frequency response only assumes the values 0 and 1 near 0 and π . The next example considers an ideal lowpass filter.

Example 3.34 (Ideal lowpass filters). By definition, the ideal lowpass filter keeps frequencies near 0 unchanged, and completely removes frequencies near π . We now have the theory in place in order to find the filter coefficients for such a filter: In Example 2.29 we implemented the ideal lowpass filter with the help of the DFT. Mathematically you can see that this code is equivalent to computing $(F_N)^H D F_N$ where D is the diagonal matrix with the entries $0, \dots, L$ and $N-L, \dots, N-1$ being 1, the rest being 0. Clearly this is a digital filter, with frequency response as stated. If the filter should keep the angular frequencies $|\omega| \leq \omega_c$ only, where ω_c describes the highest frequency we should keep, we should choose L so that $\omega_c = 2\pi L/N$. Again, in Example 2.18 we computed the DFT of this vector, and it followed from Theorem 2.21 that the IDFT of this vector equals its DFT. This means that we can find the filter coefficients by using Equation (3.11): Since the IDFT was $\frac{1}{\sqrt{N}} \frac{\sin(\pi k(2L+1)/N)}{\sin(\pi k/N)}$, the filter coefficients are

$$\frac{1}{\sqrt{N}} \frac{1}{\sqrt{N}} \frac{\sin(\pi k(2L+1)/N)}{\sin(\pi k/N)} = \frac{1}{N} \frac{\sin(\pi k(2L+1)/N)}{\sin(\pi k/N)}.$$

This means that the filter coefficients lie as N points uniformly spaced on the curve $\frac{1}{N} \frac{\sin(\omega(2L+1)/2)}{\sin(\omega/2)}$ between 0 and π . This curve has been encountered many other places in these notes. The filter which keeps only the frequency $\omega_c = 0$ has all filter coefficients being $\frac{1}{N}$ (set $L = 1$), and when we include all frequencies (set $L = N$) we get the filter where $x_0 = 1$ and all other filter coefficients are 0. When we are between these two cases, we get a filter where s_0 is the biggest coefficient, while the others decrease towards 0 along the curve we have computed. The bigger L and N are, the quicker they decrease to zero. All filter coefficients are usually nonzero for this filter, since this curve is zero only at certain points. This is unfortunate, since it means that the filter is time-consuming to compute. ♣

The two previous examples show an important duality between vectors which are 1 on some elements and 0 on others (also called window vectors), and the vector $\frac{1}{N} \frac{\sin(\pi k(2L+1)/N)}{\sin(\pi k/N)}$ (also called a sinc): filters of the one type correspond to frequency responses of the other type, and vice versa. The examples also show that, in some cases only the filter coefficients are known, while in other cases only the frequency

response is known. In any case we can deduce the one from the other, and both cases are important.

Filters are much more efficient when there are few nonzero filter coefficients. In this respect the second example displays a problem: in order to create filters with particularly nice properties (such as being an ideal lowpass filter), one may need to sacrifice computational complexity by increasing the number of nonzero filter coefficients. The trade-off between computational complexity and desirable filter properties is a very important issue in *filter design theory*.

Example 3.35. In order to decrease the computational complexity for the ideal low-pass filter in Example 3.34, one can for instance include only the first filter coefficients, i.e. $\{\frac{1}{N} \frac{\sin(\pi k(2L+1)/N)}{\sin(\pi k/N)}\}_{k=-N_0}^{N_0}$, ignoring the last ones. Hopefully this gives us a filter where the frequency response is not that different from the ideal lowpass filter. In Figure 3.7 we show the corresponding frequency responses. In the figure we have set $N = 128$, $L = 32$, so that the filter removes all frequencies $\omega > \pi/2$. N_0 has been chosen so that the given percentage of all coefficients are included. Clearly the figure shows that we should be careful when we omit filter coefficients: if we drop too many, the frequency response is far away from that of an ideal bandpass filter. In particular, we see that the new frequency response oscillates wildly near the discontinuity of the ideal lowpass filter. Such oscillations are called *Gibbs oscillations*.



Example 3.36 (Filters and the MP3 standard). We mentioned previously that the MP3 standard splits the sound into frequency bands. This splitting is actually performed by particular filters, which we will consider now.

In the example above, we saw that when we dropped the last filter coefficients in the ideal lowpass filter, there were some undesired effects in the frequency response of the resulting filter. Are there other and better approximations to the ideal lowpass filter which uses the same number of filter coefficients? This question is important, since the ear is sensitive to certain frequencies, and we would like to extract these frequencies for special processing, using as low computational complexity as possible. In the MP3-standard, such filters have been constructed. These filters are more advanced than the ones we have seen upto now. They have as many as 512 filter coefficients! We will not go into the details on how these filters are constructed, but only show how their frequency responses look. In Figure 3.8(a), the “prototype filter” which is used in the MP3 standard is shown. We see that this is very close to an ideal lowpass filter. Moreover, many of the undesirable effect from the previous example have been eliminated: The oscillations near the discontinuities are much smaller, and the values are lower away from 0. Using Property 4 in theorem 2.21, it is straightforward to construct filters with similar frequency responses, but centered around different frequencies: We simply need to multiply the filter coefficients with a complex exponential, in order to obtain a filter where the frequency response has been shifted to the left or right. In the MP3 standard, this observation is used to construct 32 filters, each having a frequency response which is a shifted copy of that of the prototype filter, so that all filters together cover the entire frequency range. 5 of these frequency responses are shown in Figure 3.8(b). To understand the effects of the different filters, let us apply them to our sample sound. If you apply all fil-

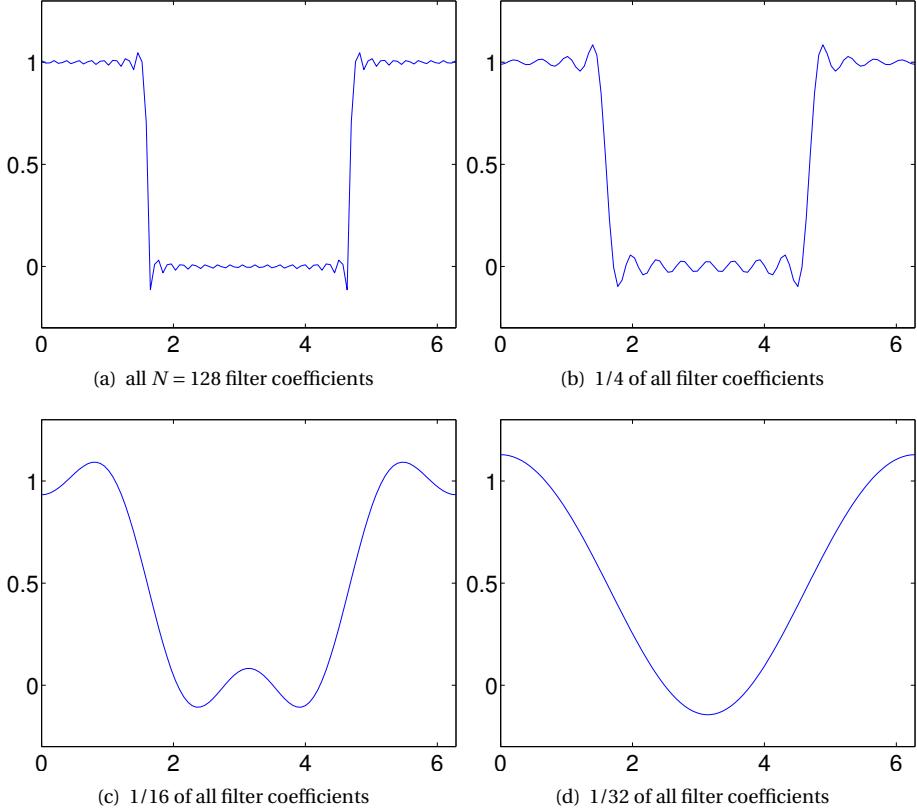


Figure 3.7: The frequency response which results by omitting the last filter coefficients for the ideal lowpass filter.

ters in the MP3 standard in successive order with the most lowpass filters first, the result will sound like this. You should interpret the result as low frequencies first, followed by the high frequencies. π corresponds to the frequency 22.05KHz (i.e. the highest representable frequency equals half the sampling rate on 44.1KHz . The different filters are concentrated on $1/32$ of these frequencies each, so that the angular frequencies you here are $[\pi/64, 3\pi/64]$, $[3\pi/64, 5\pi/64]$, $[5\pi/64, 7\pi/64]$, and so on, in that order.

In Section 3.3.1 we mentioned that the psychoacoustic model of the MP3 standard applied a window to the sound data, followed by an FFT to that data. This is actually performed in parallel on the same sound data. Applying two different operations in parallel to the sound data may seem strange. In the MP3 standard [17] (p. 109) this is explained by “the lack of spectral selectivity obtained at low frequencies” by the filters above. In other words, the FFT can give more precise frequency information than the filters can. This more precise information is then used to compute psychoacoustic information such as masking thresholds, and this information

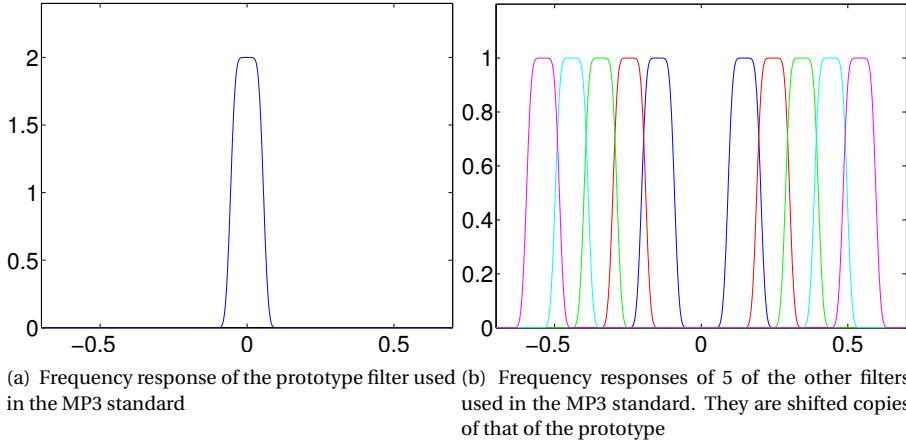


Figure 3.8: Frequency responses of some filters used in the MP3 standard.

is applied to the output of the filters. ♣

Example 3.37 (Reducing the treble II). When reducing the treble it is reasonable to let the middle sample x_i count more than the neighbours in the average, so an alternative is to compute the average by instead writing

$$z(t) = (x(t-1) + 2x(t) + x(t+1))/4;$$

The coefficients 1, 2, 1 here have been taken from row 2 in Pascal's triangle. It turns out that this is a good choice of coefficients. Also if we take averages of more numbers it will turn out that higher rows of Pascals triangle are good choices. Let us take a look at why this is the case. Let S be the moving average filter of two elements, i.e.

$$(S\mathbf{x})_n = \frac{1}{2}(x_{n-1} + x_n).$$

In Example 3.32 we had an odd number of filter coefficients. Here we have only two. We see that the frequency response in this case is

$$\lambda_S(\omega) = \frac{1}{2}(1 + e^{-i\omega}) = e^{-i\omega/2} \cos(\omega/2).$$

The frequency response is complex now, since the filter is not symmetric in this case. Let us now apply this filter k times, and denote by S_k the resulting filter. Theorem 3.20 gives us that the frequency response of S_k is

$$\lambda_{S^k}(\omega) = \frac{1}{2^k}(1 + e^{-i\omega})^k = e^{-ik\omega/2} \cos^k(\omega/2),$$

which is a polynomial in $e^{-i\omega}$ with the coefficients taken from Pascal's triangle (remember that the values in Pascals triangle are the coefficients of x in the expression

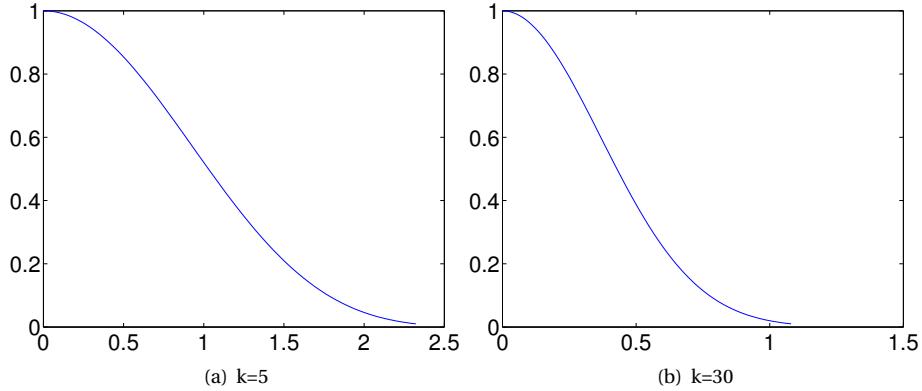


Figure 3.9: The frequency response of filters corresponding to a moving average filter convolved with itself k times.

$(1+x)^k$, i.e. the binomial coefficients $\binom{k}{r}$ for $0 \leq r \leq k$). At least, this partially explains how filters with coefficients taken from Pascal's triangle appear. The reason why these are more desirable than moving average filters, and are used much for smoothing abrupt changes in images and in sound, is the following: Since we take a k 'th power with k large, λ_{S^k} is more square-like near 0, i.e. it becomes more and more like a bandpass filter near 0. In Figure 3.9 we have plotted the magnitude of the frequency response when $k = 5$, and when $k = 30$. This behaviour near 0 is not so easy to see from the figure. Note that we have zoomed in on the frequency response to the area where it actually decreases to 0.

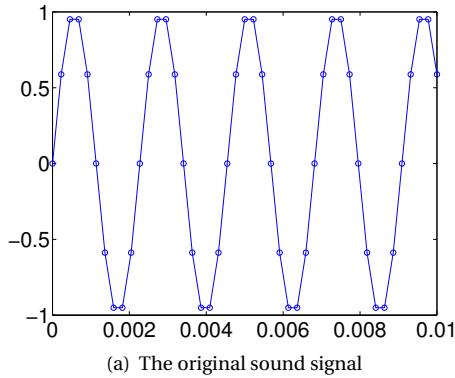
If we pick coefficients from row 4 of Pascals triangle instead, we would write

```
for t=3:(N-2)
    y(t)=(x(t-2)+4*x(t-1)+6*x(t)+4*x(t+1)+x(t+2))/16;
end
```

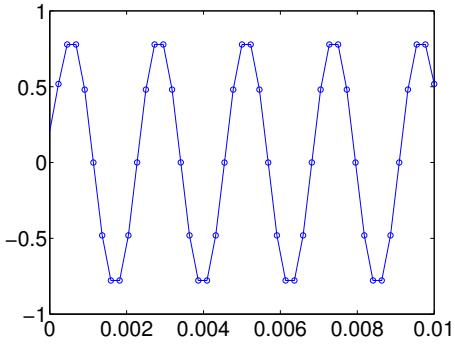
Here we have dropped the first and last part, which have special expressions due to the circulant structure of the matrix. Picking coefficients from a row in Pascal's triangle works better the longer the filter is:

Observation 3.38. Let x be the samples of a digital sound, and let $\{c_i\}_{i=1}^{2k+1}$ be the numbers in row $2k$ of Pascal's triangle. Then the sound with samples y given by

```
y=zeros(length(x));
for t=(k+1):(N-k)
    for j=1:(2*k+1)
        y(t)=y(t)+c(j)*x(t+k+1-j))/2^k;
    end
end
```



(a) The original sound signal

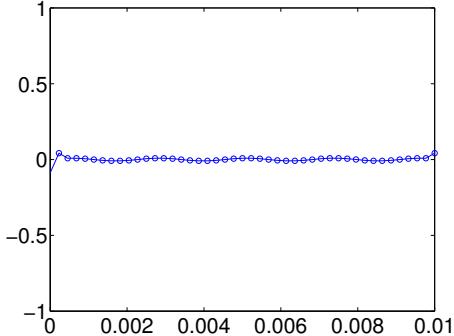


(b) The result of applying the filter from row 4 of
Pascal's triangle

Figure 3.10: Reducing the treble.

An example of the result of smoothing is shown in Figure 3.10. (a) shows the samples of the pure sound with frequency 440Hz (with sampling frequency $f_s = 4400\text{Hz}$). (b) shows the result of applying the averaging process by using row 4 of Pascals triangle. We see that the oscillations have been reduced. In Exercise 4 you will be asked to implement reducing the treble in our sample audio file. If you do this you should hear that the sound gets softer when you increase k : For $k = 32$ the sound will be like this, for $k = 256$ it will be like this. ♣

Another common option in an audio system is reducing the bass. This corresponds to reducing the low frequencies in the sound, or equivalently, the slow variations in the sample values. It turns out that this can be accomplished by simply changing the sign of the coefficients used for reducing the treble. Let us explain why this is the case. Let S be a filter with filter coefficients s_k , and let us consider the filter



(a) The result of applying the bass-reducing filter deduced from row 4 in Pascals triangle to the pure sound of figure 3.10(a).

Figure 3.11: Reducing the bass.

T with filter coefficient $(-1)^k s_k$. The frequency response of T is

$$\begin{aligned}\lambda_T(\omega) &= \sum_k (-1)^k s_k e^{-i\omega k} = \sum_k (e^{-i\pi})^k s_k e^{-i\omega k} \\ &= \sum_k e^{-i\pi k} s_k e^{-i\omega k} = \sum_k s_k e^{-i(\omega+\pi)k} = \lambda_S(\omega + \pi).\end{aligned}$$

where we have set $-1 = e^{-i\pi}$ (note that this is nothing but Property 4. in Theorem 2.21, with $d = N/2$). Now, for a lowpass filter S , $\lambda_S(\omega)$ has large values when ω is close to 0 (the low frequencies), and values near 0 when ω is close to π (the high frequencies). For a highpass filter T , $\lambda_T(\omega)$ has values near 0 when ω is close to 0 (the low frequencies), and large values when ω is close to π (the high frequencies). When T is obtained by adding an alternating sign to the filter coefficients of S , The relation $\lambda_T(\omega) = \lambda_S(\omega + \pi)$ thus says that T is a highpass filter when S is a lowpass filter, and vice versa:

Observation 3.39. Assume that T is obtained by adding an alternating sign to the filter coefficients of S . If S is a lowpass filter, then T is a highpass filter. If S is a highpass filter, then T is a lowpass filter.

The following example elaborates further on this.

Example 3.40 (Reducing the bass). Consider the bass-reducing filter deduced from the fourth row in Pascals triangle:

$$y(t) = (x(t-2) - 4*x(t-1) + 6*x(t) - 4*x(t+1) + x(t+2))/16;$$

The result of applying this filter to the sound in Figure 3.10 is shown in Figure 3.11. We observe that the samples oscillate much more than the samples of the original sound. If we play a sound after such a bass-reducing filter has been applied to it, the bass will be reduced.

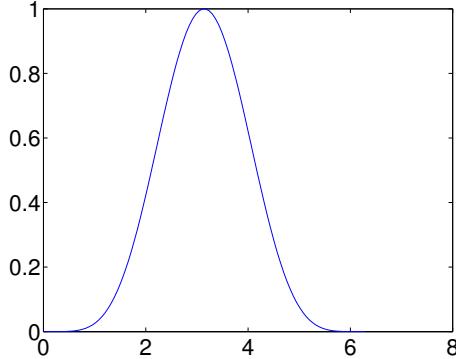


Figure 3.12: The frequency response of the bass reducing filter, which corresponds to row 5 of Pascal's triangle.

Observation 3.41. Let \mathbf{x} be the samples of a digital sound, and let $\{c_i\}_{i=1}^{2k+1}$ be the numbers in row $2k$ of Pascal's triangle. Then the sound with samples \mathbf{y} given by

```
y=zeros(length(x));
for t=(k+1):(N-k)
    for j=1:(2*k+1)
        y(t)=x(t)+(-1)^(k+1-j)*c(j)*x(t+j-k-1))/2^k;
    end
end
```

has reduced bass compared to the sound given by the samples \mathbf{y} .

In Exercise 4 you will be asked to implement reducing the bass in our sample audio file. The new sound will be difficult to hear for large k , and we will explain why later. For $k = 1$ the sound will be like this, for $k = 2$ it will be like this. Even if the sound is quite low, you can hear that more of the bass has disappeared for $k = 2$.

The frequency response we obtain from using row 5 of Pascal's triangle is shown in Figure 3.12. It is just the frequency response of the corresponding treble-reducing filter shifted with π . The alternating sign can also be achieved if we write the frequency response $\frac{1}{2^k}(1 + e^{-i\omega})^k$ from Example 3.37 as $\frac{1}{2^k}(1 - e^{-i\omega})^k$, which corresponds to applying the filter $S(\mathbf{x}) = \frac{1}{2}(-x_{n-1} + x_n)$ k times. ♣

What you should have learnt in this section

Simple examples of filters, such as time delay filters and filters which add echo. Low-pass and highpass filters and their frequency responses, and their interpretation as treble- and bass-reducing filters. Moving average filters, and filters arising from rows in Pascal's triangle, as examples of such filters. How to pass between the two by

adding an alternating sign to the filter coefficients.

Exercises for Section 3.5

1. Let E_{d_1} and E_{d_2} be two time delay filters. Show that $E_{d_1}E_{d_2} = E_{d_1+d_2}$ (i.e. that the composition of two time delays again is a time delay) in two different ways

- a. Give a direct argument which uses no computations.
 - b. By using Property 3 in Theorem 2.21, i.e. by using a property for the Discrete Fourier Transform.
2. In this exercise, we will experiment with adding echo to a signal.

- a. Write a function

```
function playwithecho(c,d)
```

which plays the sound samples of `castanets.wav` with echo added for damping constant c and delay d as described in Example 3.30.

- b. Generate the sound from Example 3.30, and verify that it is the same as the one you heard there.
 - c. Listen to the sound samples for different values of d and c . For which range of d is the echo distinguishable from the sound itself? How low can you choose c in order to still hear the echo?
3. Consider the two filters $S_1 = \{1, 0, \dots, 0, c\}$ and $S_2 = \{1, 0, \dots, 0, -c\}$. Both of these can be interpreted as filters which add an echo. Show that $\frac{1}{2}(S_1 + S_2) = I$. What is the interpretation of this relation in terms of echos?
4. In this exercise, we will experiment with increasing and reducing the treble and bass in a signal as in examples 3.37 and 3.40.

- a. Write functions

```
function reducetreble(k)
function reducebass(k)
```

which reduces bass and treble in the ways described above for the sound from the file `castanets.wav`, and plays the result, when row number $2k$ in Pascal's triangle is used to construct the filters. Look into the Matlab function `conv` to help you to find the values in Pascal's triangle.

- b. Generate the sounds you heard in examples 3.37 and 3.40, and verify that they are the same.
 - c. In your code, it will not be necessary to scale the values after reducing the treble, i.e. the values are already between -1 and 1 . Explain why this is the case.

- d.** How high must k be in order for you to hear difference from the actual sound? How high can you choose k and still recognize the sound at all?
- 5.** Consider again Example 3.34. Find an expression for a filter so that only frequencies so that $|\omega - \pi| < \omega_c$ are kept, i.e. the filter should only keep angular frequencies close to π (i.e. here we construct a highpass filter).
- 6.** In this exercise we will investigate how we can combine lowpass and highpass filters to produce other filters
- Assume that S_1 and S_2 are lowpass filters. What kind of filter is $S_1 S_2$? What if both S_1 and S_2 are highpass filters?
 - Assume that one of S_1, S_2 is a highpass filter, and that the other is a low-pass filter. What kind of filter $S_1 S_2$ in this case?
- 7.** A filter S_1 has the frequency response $\frac{1}{2}(1 + \cos \omega)$, and another filter has the frequency response $\frac{1}{2}(1 + \cos(2\omega))$.
- Is $S_1 S_2$ a lowpass filter, or a highpass filter?
 - What does the filter $S_1 S_2$ do with angular frequencies close to $\omega = \pi/2$.
 - Find the filter coefficients of $S_1 S_2$.
- Hint: Use Theorem 3.20 to compute the frequency response of $S_1 S_2$ first.
- Write down the matrix of the filter $S_1 S_2$ for $N = 8$.
- 8.** An operation describing some transfer of data in a system is defined as the composition of the following three filters:

- First a time delay filter with delay $d_1 = 2$, due to internal transfer of data in the system,
- then the treble-reducing filter $T = \{1/4, \underline{1/2}, 1/4\}$,
- finally a time delay filter with delay $d_2 = 4$ due to internal transfer of the filtered data.

We denote by $T_2 = E_{d_2} T E_{d_1} = E_4 T E_2$ the operation which applies these filters in succession.

- Explain why T_2 also is a digital filter. What is (the magnitude of) the frequency response of E_{d_1} ? What is the connection between (the magnitude of) the frequency response of T and T_2 ?

- Show that $T_2 = \{0, 0, 0, 0, 0, 1/4, 1/2, 1/4\}$.

Hint: Use the expressions $(E_{d_1} \mathbf{x})_n = x_{n-d_1}$, $(T \mathbf{x})_n = \frac{1}{4}x_{n+1} + \frac{1}{2}x_n + \frac{1}{4}x_{n-1}$, $(E_{d_2} \mathbf{x})_n = x_{n-d_2}$, and compute first $(E_{d_1} \mathbf{x})_n$, then $(T E_{d_1} \mathbf{x})_n$, and finally $(T_2 \mathbf{x})_n = (E_{d_2} T E_{d_1} \mathbf{x})_n$. From the last expression you should be able to read out the filter coefficients.

- c.** Assume that $N = 8$. Write down the 8×8 -circulant Toeplitz matrix for the filter T_2 .
- 9.** In Example 3.36, we mentioned that the filters used in the MP3-standard were constructed from a lowpass prototype filter by multiplying the filter coefficients with a complex exponential. Clearly this means that the new frequency response is a shift of the old one. The disadvantage is, however, that the new filter coefficients are complex. It is possible to address this problem as follows. Assume that t_k are the filter coefficients of a filter S_1 , and that S_2 is the filter with filter coefficients $\cos(2\pi kn/N) t_k$, where $n \in \mathbb{N}$. Show that

$$\lambda_{S_2}(\omega) = \frac{1}{2} (\lambda_{S_1}(\omega - 2\pi n/N) + \lambda_{S_1}(\omega + 2\pi n/N)).$$

In other words, when we multiply (modulate) the filter coefficients with a cosine, the new frequency response can be obtained by shifting the old frequency response with $2\pi n/N$ in both directions, and taking the average of the two.

- 10. a.** Explain what the code below does, line by line.

```
[S,fs]=wavread('castanets.wav');
S=S(:,1);
N=length(S);
newS=zeros(N,1);
for k=2:(N-1)
    newS(k)=2*S(k+1) + 4*S(k) + 2*S(k-1);
end
newS(1)=2*S(2) + 4*S(1) + 2*S(N);
newS(N)=2*S(1) + 4*S(N) + 2*S(N-1);
newS=newS/max(abs(newS));
playerobj=audioplayer(newS,fs);
playblocking(playerobj)
```

Comment in particular on what happens in the three lines directly after the `for`-loop, and why we do this. What kind of changes in the sound do you expect to hear?

- b.** Write down the compact filter notation for the filter which is used in the code, and write down a 5×5 circulant Toeplitz matrix which corresponds to this filter. Plot the (continuous) frequency response. Is the filter a lowpass- or a highpass filter?

- c.** Another filter is given by the circulant Toeplitz matrix

$$\begin{pmatrix} 4 & -2 & 0 & 0 & -2 \\ -2 & 4 & -2 & 0 & 0 \\ 0 & -2 & 4 & -2 & 0 \\ 0 & 0 & -2 & 4 & -2 \\ -2 & 0 & 0 & -2 & 4 \end{pmatrix}.$$

Express a connection between the frequency responses of this filter and the filter from b. Is the new filter a lowpass- or a highpass filter?

3.6 Time-invariance of filters

As we have seen, filters can be characterized by their eigenvectors, and also as circulant Toeplitz matrices. To make the picture more complete, we will in this section state a third, important characterization of filters. This characterization is stated in terms of the following concept:

Definition 3.42 (Time-invariance). A linear transformation from \mathbb{R}^N to \mathbb{R}^N is said to be time-invariant if, for any d , the output of the *delayed* input vector \mathbf{z} defined by $z_n = x_{n-d}$ is the delayed output vector \mathbf{w} defined by $w_n = y_{n-d}$.

We have the following result:

Theorem 3.43. A linear transformation S is a digital filter if and only if it is time-invariant.

Proof: Let $\mathbf{y} = S\mathbf{x}$, and \mathbf{z}, \mathbf{w} as defined above. We have that

$$\begin{aligned} w_n &= (S\mathbf{x})_{n-d} = \sum_{k=0}^{N-1} S_{n-d,k} x_k \\ &= \sum_{k=0}^{N-1} S_{n,k+d} x_k = \sum_{k=0}^{N-1} S_{n,k} x_{k-d} \\ &= \sum_{k=0}^{N-1} S_{n,k} z_k = (S\mathbf{z})_n \end{aligned}$$

This proves that $S\mathbf{z} = \mathbf{w}$, so that S is time-invariant. ■

By Example 3.29, delaying a vector with d elements corresponds to applying the filter E_d . Clearly, that S is time-invariant is the same as $SE_d = E_dS$ for any d . That all filters are time-invariant follows thus also immediately from the fact that all filters commute.

Due to Theorem 3.43, digital filters are also called *LTI filters* (LTI stands for Linear, Time-Invariant). By combining the definition of a digital filter with theorems 3.8 and 3.43 we get the following theorem which summarizes our findings.

Theorem 3.44 (Characterizations of digital filters). The following are equivalent characterizations of a digital filter:

1. $S = (F_N)^H D F_N$ for a diagonal matrix D , i.e. the Fourier basis is a basis of eigenvectors for S .
2. S is a circulant Toeplitz matrix.
3. S is linear and time-invariant.

Exercises for Section 3.6

1. In Example 2.7 we looked at time reversal as an operation on digital sound. In \mathbb{R}^N this can be defined as the linear mapping which sends the vector e_k to e_{N-1-k} for all $0 \leq k \leq N-1$.

- a. Write down the matrix for the time reversal linear mapping, and explain from this why time reversal is not a digital filter.
- b. Prove directly that time reversal is not a time-invariant operation.

3.7 More general filters

The starting point for defining filters at the beginning of this chapter was equations on the form

$$z_n = \sum_k t_k x_{n-k}.$$

The important point here was that we had a limited number of nonzero t_k , and this enabled us to compute the filter on a computer. In practice, however, there exist many filtering operations with an infinite number of filter coefficients. The ideal lowpass filter from Example 3.34 is one example. It turns out that many such cases can be made computable if we change our procedure slightly. The old procedure for computing a filter is to compute $\mathbf{z} = S\mathbf{x}$. Consider the following alternative:

Idea 3.45 (More general filters (1)). Let \mathbf{x} be the input to a filter, and let T be a filter. By solving the system $T\mathbf{z} = \mathbf{x}$ for \mathbf{z} we get another filter, which we denote by S .

Of course T must then be the inverse of S (which also is a filter), but the point is that the inverse of a filter may have a finite number of filter coefficients, even if the filter itself does not. In such cases this new procedure is more attractive than the old one, since the equation system can be solved with few arithmetic operations when T has few filter coefficients.

It turns out that there also are highly computable filters where neither the filter nor its inverse have a finite number of filter coefficients. Consider the following idea:

Idea 3.46 (More general filters (2)). Let \mathbf{x} be the input to a filter, and let U and V be filters. By solving the system $U\mathbf{z} = V\mathbf{x}$ for \mathbf{z} we get another filter, which we denote by S . The filter S can be implemented in two steps: first we compute the right hand side $\mathbf{y} = V\mathbf{x}$, and then we solve the equation $U\mathbf{z} = \mathbf{y}$.

If both U and V are invertible we have that the filter is $S = U^{-1}V$, and this is invertible with inverse $S^{-1} = V^{-1}U$. The point is that, when U and V have a finite number of filter coefficients, both S and its inverse will typically have an infinite

number of filter coefficients. The filters from this idea are thus more general than the ones from the previous idea, and the new idea makes a wider class of filters implementable using row reduction of sparse matrices. Computing a filter by solving $Uz = Vx$ may also give meaning when the matrices U and V are singular: The matrix system can have a solution even if U is singular. Therefore we should be careful in using the form $T = U^{-1}V$.

We have the following result concerning the frequency responses:

Theorem 3.47. Assume that S is the filter defined from the equation $Uz = Vx$. Then we have that $\lambda_S(\omega) = \frac{\lambda_V(\omega)}{\lambda_U(\omega)}$ whenever $\lambda_U(\omega) \neq 0$.

Proof: Set $x = \phi_n$. We have that $Uz = \lambda_{U,n}\lambda_{S,n}\phi_n$, and $Vx = \lambda_{V,n}\phi_n$. If the expressions are equal we must have that $\lambda_{U,n}\lambda_{S,n} = \lambda_{V,n}$, so that $\lambda_{S,n} = \frac{\lambda_{V,n}}{\lambda_{U,n}}$ for all n . By the definition of the continuous frequency response this means that $\lambda_S(\omega) = \frac{\lambda_V(\omega)}{\lambda_U(\omega)}$ whenever $\lambda_U(\omega) \neq 0$. \blacksquare

The following example clarifies the points made above, and how one may construct U and V from S . The example also shows that, in addition to making some filters with infinitely many filter coefficients computable, the procedure $Uz = Vx$ for computing a filter can also reduce the complexity in some filters where we already have a finite number of filter coefficients.

Example 3.48. Consider again the moving average filter S from Example 3.32:

$$z_n = \frac{1}{2L+1}(x_{n+L} + \dots + x_n + \dots + x_{n-L}).$$

If we implemented this directly, $2L$ additions would be needed for each n , so that we would need a total of $2NL$ additions. However, we can also write

$$\begin{aligned} z_{n+1} &= \frac{1}{2L+1}(x_{n+1+L} + \dots + x_{n+1} + \dots + x_{n+1-L}) \\ &= \frac{1}{2L+1}(x_{n+L} + \dots + x_n + \dots + x_{n-L}) + \frac{1}{2L+1}(x_{n+1+L} - x_{n-L}) \\ &= z_n + \frac{1}{2L+1}(x_{n+1+L} - x_{n-L}). \end{aligned}$$

This means that we can also compute the output from the formula

$$z_{n+1} - z_n = \frac{1}{2L+1}(x_{n+1+L} - x_{n-L}),$$

which can be written on the form $Uz = Vx$ with $U = \{1, -1\}$ and $V = \frac{1}{2L+1}\{1, 0, \dots, 0, -1\}$ where the 1 is placed at index $-L-1$ and the -1 is placed at index L . We now perform only $2N$ additions in computing the right hand side, and solving the equation system requires only $2(N-1)$ additions. The total number of additions is thus $2N + 2(N-1) = 4N-2$, which is much less than the previous $2LN$ when L is large.

A perhaps easier way to find U and V is to consider the frequency response of the moving average filter, which is

$$\begin{aligned}\frac{1}{2L+1}(e^{-Li\omega} + \dots + e^{Li\omega}) &= \frac{1}{2L+1}e^{-Li\omega} \frac{1 - e^{(2L+1)i\omega}}{1 - e^{i\omega}} \\ &= \frac{\frac{1}{2L+1}(-e^{(L+1)i\omega} + e^{-Li\omega})}{1 - e^{i\omega}},\end{aligned}$$

where we have used the formula for the sum of a geometric series. From here we easily see the frequency responses of U and V from the numerator and the denominator. ♣

Filters with an infinite number of filter coefficients are also called *IIR filters* (IIR stands for *Infinite Impulse Response*). Thus, we have seen that some IIR filters may still have efficient implementations.

Exercises for Section 3.7

1. A filter is defined by demanding that $z_{n+2} - z_{n+1} + z_n = x_{n+1} - x_n$.
 - a. Compute and plot the frequency response of the filter.
 - b. Use Matlab to compute the output of the vector $\mathbf{x} = (1, \dots, 10)$. In order to do this you should write down two 10×10 -circulant Toeplitz matrices, and use these in Matlab.

3.8 Implementation of filters

As we saw in Example 3.48, a filter with many filter coefficients could be factored into the application of two simpler filters, and this could be used as a basis for an efficient implementation. There are also several other possible efficient implementations of filters. In this section we will consider two such techniques. The first technique considers how we can use the DFT to speed up the computation of filters. The second technique considers how we can factorize a filter into a product of simpler filters.

3.8.1 Implementation of filters using the DFT

If there are k filter coefficients, a direct implementation of a filter would require kN multiplications. Since filters are diagonalized by the DFT, one can also compute the filter as the product $S = F_N^H D F_N$. This would instead require $O(N \log_2 N)$ complex multiplications when we use the FFT algorithm, which may be a higher number of multiplications. We will however see that, by slightly changing our algorithm, we may end up with a DFT-based implementation of the filter which requires fewer multiplications.

The idea is to split the computation of the filter into smaller parts. Assume that we compute M elements of the filter at a time. If the nonzero filter coefficients of S

are $t_{-k_0}, \dots, t_{k-k_0-1}$, we have that

$$(S\mathbf{x})_t = \sum_r t_r x_{s-r} = t_{-k_0} x_{t+k_0} + \dots + t_{k-k_0-1} x_{t-(k-k_0-1)}.$$

From this it is clear that $(S\mathbf{x})_t$ only depends on $x_{t-(k-k_0-1)}, \dots, x_{t+k_0}$. This means that, if we restrict the computation of S to $x_{t-(k-k_0-1)}, \dots, x_{t+M-1+k_0}$, the outputs x_t, \dots, x_{t+M-1} will be the same as without this restriction. This means that we can compute the output M elements at a time, at each step multiplying with a circulant Toeplitz matrix of size $(M+k-1) \times (M+k-1)$. If we choose M so that $M+k-1 = 2^r$, we can use the FFT and IFFT algorithms to compute $S = F_N^H D F_N$, and we require $O(r2^r)$ multiplications for every block of length M . The total number of multiplications is $\frac{Nr2^r}{M} = \frac{Nr2^r}{2^r - k + 1}$. If $k = 128$, you can check on your calculator that the smallest value is for $r = 10$ with value $11.4158 \times N$. Since the direct implementation gives kN multiplications, this clearly gives a benefit for the new approach, it gives a 90% decrease in the number of multiplications.

3.8.2 Factoring a filter into several filters

In practice, filters are often applied in hardware, and applied in real-time scenarios where performance is a major issue. The most CPU-intensive tasks in such applications often have few memory locations available. These tasks are thus not compatible with filters with many filter coefficients, since for each output sample we then need access to many input samples and filter coefficients. A strategy which addresses this is to factorize the filter into the product of several smaller filters, and then applying each filter in turn. Since the frequency response of the product of filters equals the product of the frequency responses, we get the following idea:

Idea 3.49. Let S be a filter with real coefficients. Assume that

$$\lambda_S(\omega) = K e^{ik\omega} (e^{i\omega} - a_1) \dots (e^{i\omega} - a_m) (e^{2i\omega} + b_1 e^{i\omega} + c_1) \dots (e^{2i\omega} + b_n e^{i\omega} + c_n). \quad (3.19)$$

Then we can write $S = K E_k A_1 \dots A_m B_1 \dots B_n$, where $A_i = \{1, \underline{-a_i}\}$ and $B_i = \{1, b_i, \underline{c_i}\}$.

Note that in Equation 3.19 a_i correspond to the real roots of the frequency response, while b_i, c_i are obtained by pairing the complex conjugate roots. Clearly the frequency responses of A_i, B_i equal the factors in the frequency response of S , which in any case can be factored into the product of filters with 2 and 3 filter coefficients, followed by a time-delay.

Note that, even though this procedure factorizes a filter into smaller parts (which is attractive for hardware implementations since smaller filters require fewer locations in memory), the number of arithmetic operations is usually not reduced. However, consider Example 3.37, where we factorized the treble-reducing filters into a product of moving average filters of length 2 (all roots in the previous idea are real, and equal). Each application of a moving average filter of length 2 does not really

require any multiplications, since multiplication with $\frac{1}{2}$ corresponds to a bitshift. Therefore, the factorization of Example 3.37 removes the need for doing any multiplications at all, while keeping the number of additions the same. There are computational savings in this case, due to the special filter structure here.

Exercises for Section 3.8

1. Write a function

```
function filterdftimpl(t,k0)
```

which takes the filter coefficients t and the value k_0 from Section 3.8.1 as input, computes the optimal M and implements the filter as described in that section.

2. Factor the filter $S = \{1, 5, 10, 6\}$ into a product of two filters, one with two filter coefficients, and one with three filter coefficients.

3.9 Summary

We defined digital filters, which do the same job for digital sound as analog filters do for (continuous) sound. Digital filters turned out to be linear transformations diagonalized by the DFT. We proved several other equivalent characterizations of digital filters as well, such as being time-invariant, and having a matrix which is circulant and Toeplitz. Just as for continuous sound, digital filters are characterized by their frequency response, which explains how the filter treats the different frequencies. We also went through several important examples of filters, some of which corresponded to meaningful operations on sound, such as adjustment of bass and treble, and adding echo. We also explained that there exist filters with useful implementations which have an infinite number of filter coefficients, and we considered techniques for implementing filters efficiently. Most of the topics covered on that can also be found in [28]. We also took a look at the role of filters in the MP3 standard for compression of sound.

In signal processing literature, the assumption that vectors are periodic is often not present, and filters are thus not defined as finite-dimensional operations. With matrix notation they would then be viewed as infinite matrices which have the Toeplitz structure (i.e. constant values on the diagonals), but with no circulation. The circulation in the matrices, as well as the restriction to finite vectors, come from the assumption of a periodic vector. There are, however, also some books which view filters as circulant Toeplitz matrices as we have done, such as [13].

Symmetric filters and the DCT

In Chapter 1 we approximated a signal of finite duration with trigonometric functions. Since these are all periodic, there are some undesirable effects near the boundaries of the signal (at least when the values at the boundaries are different), and this resulted in a slowly converging Fourier series. This was addressed by instead considering the symmetric extension of the function, for which we obtained a more precise Fourier representation, as fewer Fourier basis vectors were needed in order to get a precise approximation.

This chapter is dedicated to addressing these thoughts for vectors. We will start by defining symmetric extensions of vectors, similarly to how we defined these for functions. Just as the Fourier series of a symmetric function was a cosine series, we will see that the symmetric extension can be viewed as a cosine vector. This gives rise to a different change of coordinates than the DFT, which we will call the DCT, which enables us to express a symmetric vector as a sum of cosine-vectors (instead of the non-symmetric complex exponentials). Since a cosine also can be associated with a given frequency, the DCT is otherwise similar to the DFT, in that it extracts the frequency information in the vector. The advantage is that the DCT can give more precise frequency information than the DFT, since it avoids the discontinuity problem of the Fourier series. This makes the DCT very practical for applications, and we will explain some of these applications. We will also show that the DCT has a very efficient implementation, comparable with the FFT.

In this chapter we will also see that the DCT has a very similar role as the DFT when it comes to filters: just as the DFT diagonalized filters, we will see that symmetric filters can be diagonalized by the DCT, when we apply the filter to the symmetric extension of the input. We will actually show that the filters which preserve our symmetric extensions are exactly the symmetric filters.

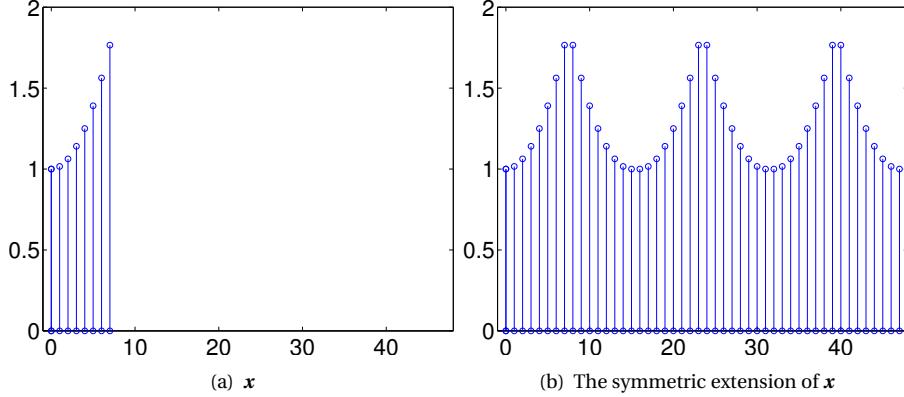


Figure 4.1: A vector and its symmetric extension.

4.1 Symmetric vectors and the DCT

As in Chapter 1, vectors can also be extended in a symmetric manner, besides the simple periodic extension procedure from Figure 2.1. In Figure 4.1 we have shown such an extension of a vector \mathbf{x} . It has \mathbf{x} as its first half, and a copy of \mathbf{x} in reverse order as its second half. We will call this the symmetric extension of \mathbf{x} :

Definition 4.1 (Symmetric extension of a vector). By the *symmetric extension* of $\mathbf{x} \in \mathbb{R}^N$, we mean the symmetric vector $\check{\mathbf{x}} \in \mathbb{R}^{2N}$ defined by

$$\check{\mathbf{x}}_k = \begin{cases} x_k & 0 \leq k < N \\ x_{2N-1-k} & N \leq k < 2N-1 \end{cases} \quad (4.1)$$

Clearly, the symmetric extension is symmetric around $N - 1/2$. This is not the only way to construct a symmetric extension, as we will return to later. As shown in Figure 4.1, we also repeat the vector in \mathbb{R}^{2N} in order to obtain a periodic vector. This is not included in Definition 4.1. Creating a symmetric extension is thus really a two-step process:

1. First, “mirror” the vector to obtain a vector in \mathbb{R}^{2N} ,
2. repeat this periodically to obtain a periodic vector.

The result from the first step lies in an N -dimensional subspace of all vectors in \mathbb{R}^{2N} , which we will call *the space of symmetric vectors*. To account for the fact that we a periodic vector also can have a different symmetry point than $N - 1/2$, let us make the following general definition:

Definition 4.2 (Symmetric vector). We say that a periodic vector \mathbf{x} is *symmetric* if there exists a number d so that $x_{d+k} = x_{d-k}$ for all k so that $d + k$ and $d - k$ are integers. d is called the *symmetry point* of \mathbf{x}

Due to the inherent periodicity of \mathbf{x} , it is clear that N must be an even number for symmetric vectors to exist at all. d can take any value, and it may not be an integer: It can also be an odd multiple of $1/2$, because then both $d + k$ and $d - k$ are integers when k also is an odd multiple of $1/2$. The symmetry point in symmetric extensions as defined in Definition 4.1 was $d = N - 1/2$. This is very common in the literature, and this is why we concentrate on this in this chapter. Later we will also consider symmetry around $N - 1$, as this also is much used.

We would like to find a basis for the N -dimensional space of symmetric vectors, and we would like this basis to be similar to the Fourier basis. Since the Fourier basis corresponds to the standard basis in the frequency domain, we are lead to studying the DFT of a symmetric vector. If the symmetry point is an integer, it is straightforward to prove the following:

Theorem 4.3 (Symmetric vectors with integer symmetry points). Let d be an integer. The following are equivalent

1. \mathbf{x} is real and symmetric with d as symmetry point.
2. $(\hat{\mathbf{x}})_n = z_n e^{-2\pi i dn/N}$ where z_n are real numbers so that $z_n = z_{N-n}$.

Proof: Assume first that $d = 0$. It follows in this case from property 2(a) of Theorem 2.21 that $(\hat{\mathbf{x}})_n$ is a real vector. Combining this with property 1 of Theorem 2.21 we see that $\hat{\mathbf{x}}$, just as \mathbf{x} , also must be a real vector symmetric about 0. Since the DFT is one-to-one, it follows that \mathbf{x} is real and symmetric about 0 if and only if $\hat{\mathbf{x}}$ is. From property 3 of Theorem 2.21 it follows that, when d is an integer, \mathbf{x} is real and symmetric about d if and only if $(\hat{\mathbf{x}})_n = z_n e^{-2\pi i dn/N}$, where z_n is real and symmetric about 0. This completes the proof. ■

Symmetric extensions were here defined by having the non-integer symmetry point $N - 1/2$, however. For these we prove the following, which is slightly more difficult.

Theorem 4.4 (Symmetric vectors with non-integer symmetry points). Let d be an odd multiple of $1/2$. The following are equivalent

1. \mathbf{x} is real and symmetric with d as symmetry point.
2. $(\hat{\mathbf{x}})_n = z_n e^{-2\pi i dn/N}$ where z_n are real numbers so that $z_{N-n} = -z_n$.

Proof: When x is as stated we can write

$$\begin{aligned}
(\hat{x})_n &= \frac{1}{\sqrt{N}} \sum_{k=0}^{N-1} x_k e^{-2\pi i k n / N} \\
&= \frac{1}{\sqrt{N}} \left(\sum_{s \geq 0} x_{d+s} e^{-2\pi i (d+s) n / N} + \sum_{s \geq 0} x_{d-s} e^{-2\pi i (d-s) n / N} \right) \\
&= \frac{1}{\sqrt{N}} \sum_{s \geq 0} x_{d+s} \left(e^{-2\pi i (d+s) n / N} + e^{-2\pi i (d-s) n / N} \right) \\
&= \frac{1}{\sqrt{N}} e^{-2\pi i d n / N} \sum_{s \geq 0} x_{d+s} \left(e^{-2\pi i s n / N} + e^{2\pi i s n / N} \right) \\
&= \frac{1}{\sqrt{N}} e^{-2\pi i d n / N} \sum_{s \geq 0} 2x_{d+s} \cos(2\pi s n / N).
\end{aligned}$$

Here s runs through odd multiples of $1/2$. Since $z_n = \frac{1}{\sqrt{N}} \sum_{s \geq 0} 2x_{d+s} \cos(2\pi s n / N)$ is a real number, we can write the result as $z_n e^{-2\pi i d n / N}$. Substituting $N - n$ for n , we get

$$\begin{aligned}
(\hat{x})_{N-n} &= \frac{1}{\sqrt{N}} e^{-2\pi i d (N-n) / N} \sum_{s \geq 0} 2x_{d+s} \cos(2\pi s (N-n) / N) \\
&= \frac{1}{\sqrt{N}} e^{-2\pi i d (N-n) / N} \sum_{s \geq 0} 2x_{d+s} \cos(-2\pi s n / N + 2\pi s) \\
&= -\frac{1}{\sqrt{N}} e^{-2\pi i d (N-n) / N} \sum_{s \geq 0} 2x_{d+s} \cos(2\pi s n / N) = -z_n e^{-2\pi i d (N-n) / N}.
\end{aligned}$$

This shows that $z_{N-n} = -z_n$, and this completes one way of the proof. The other way, we can write

$$x_k = \frac{1}{\sqrt{N}} \sum_{n=0}^{N-1} (\hat{x})_n e^{2\pi i k n / N}$$

if $(\hat{x})_n = z_n e^{-2\pi i d n / N}$ and $(\hat{x})_{N-n} = -z_n e^{-2\pi i d (N-n) / N}$, the sum of the n 'th term and the $N - n$ 'th term in the sum is

$$\begin{aligned}
&z_n e^{-2\pi i d n / N} e^{2\pi i k n / N} - z_n e^{2\pi i d (N-n) / N} e^{2\pi i k (N-n) / N} \\
&= z_n (e^{2\pi i (k-d)n / N} - e^{-2\pi i d + 2\pi i d n / N - 2\pi i k n / N}) \\
&= z_n (e^{2\pi i (k-d)n / N} + e^{2\pi i (d-k)n / N}) = 2z_n \cos(2\pi (k-d)n / N).
\end{aligned}$$

This is real, so that all x_k are real. If we set $k = d + s$, $k = d - s$ here we get

$$\begin{aligned}
2z_n \cos(2\pi((d+s)-d)n / N) &= 2z_n \cos(2\pi s n / N) \\
2z_n \cos(2\pi((d-s)-d)n / N) &= 2z_n \cos(-2\pi s n / N) = 2z_n \cos(2\pi s n / N).
\end{aligned}$$

By adding terms together and comparing we must have that $x_{d+s} = x_{d-s}$, and the proof is done. \blacksquare

Now, let us specialize to symmetric extensions as defined in Definition 4.1, i.e. where $d = N - 1/2$. The following result gives us an orthonormal basis for the symmetric extensions, which are very simple in the frequency domain:

Theorem 4.5. The set of all \mathbf{x} symmetric around $N - 1/2$ is a vector space of dimension N , and we have that

$$\left\{ \mathbf{e}_0, \left\{ \frac{1}{\sqrt{2}} \left(e^{\pi i n/(2N)} \mathbf{e}_n + e^{-\pi i n/(2N)} \mathbf{e}_{2N-n} \right) \right\}_{n=1}^{N-1} \right\}$$

is an orthonormal basis for $\hat{\mathbf{x}}$ where \mathbf{x} is symmetric around $N - 1/2$.

Proof: For a vector \mathbf{x} symmetric about $d = N - 1/2$ we know that

$$(\hat{\mathbf{x}})_n = z_n e^{-2\pi i(N-1/2)n/(2N)},$$

and the only requirement on the vector \mathbf{z} is the antisymmetry condition $z_{2N-n} = -z_n$. The vectors $\mathbf{z}_i = \frac{1}{\sqrt{2}}(\mathbf{e}_i - \mathbf{e}_{2N-i})$, $1 \leq i \leq N - 1$, together with the vector $\mathbf{z}_0 = \mathbf{e}_0$, are clearly orthonormal and satisfies the antisymmetry condition. From these we obtain that

$$\left\{ \mathbf{e}_0, \left\{ \frac{1}{\sqrt{2}} \left(e^{-2\pi i(N-1/2)n/(2N)} \mathbf{e}_n - e^{-2\pi i(N-1/2)(2N-n)/(2N)} \mathbf{e}_{2N-n} \right) \right\}_{n=1}^{N-1} \right\}$$

is an orthonormal basis for the $\hat{\mathbf{x}}$ with \mathbf{x} symmetric. We can write

$$\begin{aligned} & \frac{1}{\sqrt{2}} \left(e^{-2\pi i(N-1/2)n/(2N)} \mathbf{e}_n - e^{-2\pi i(N-1/2)(2N-n)/(2N)} \mathbf{e}_{2N-n} \right) \\ &= \frac{1}{\sqrt{2}} \left(e^{-\pi i n} e^{\pi i n/(2N)} \mathbf{e}_n + e^{\pi i n} e^{-\pi i n/(2N)} \mathbf{e}_{2N-n} \right) \\ &= \frac{1}{\sqrt{2}} e^{\pi i n} \left(e^{\pi i n/(2N)} \mathbf{e}_n + e^{-\pi i n/(2N)} \mathbf{e}_{2N-n} \right). \end{aligned}$$

This also means that

$$\left\{ \mathbf{e}_0, \left\{ \frac{1}{\sqrt{2}} \left(e^{\pi i n/(2N)} \mathbf{e}_n + e^{-\pi i n/(2N)} \mathbf{e}_{2N-n} \right) \right\}_{n=1}^{N-1} \right\}$$

is an orthonormal basis. ■

We immediately get the following result:

Theorem 4.6.

$$\left\{ \frac{1}{\sqrt{2N}} \cos \left(2\pi \frac{0}{2N} \left(k + \frac{1}{2} \right) \right), \left\{ \frac{1}{\sqrt{N}} \cos \left(2\pi \frac{n}{2N} \left(k + \frac{1}{2} \right) \right) \right\}_{n=1}^{N-1} \right\} \quad (4.2)$$

is an orthonormal basis for the set of vectors symmetric around $N - 1/2$ in \mathbb{R}^{2N} . Moreover, the n 'th vector in this basis has frequency contribution only from the indices n and $2N - n$.

Proof: Since the IDFT is unitary, the IDFT applied to the vectors above gives an orthonormal basis for the set of symmetric extensions. We get that

$$(F_{2N})^H(\mathbf{e}_0) = \left(\frac{1}{\sqrt{2N}}, \frac{1}{\sqrt{2N}}, \dots, \frac{1}{\sqrt{2N}} \right) = \frac{1}{\sqrt{2N}} \cos\left(2\pi \frac{0}{2N} \left(k + \frac{1}{2}\right)\right).$$

We also get that

$$\begin{aligned} & (F_{2N})^H\left(\frac{1}{\sqrt{2}}\left(e^{\pi i n/(2N)} \mathbf{e}_n + e^{-\pi i n/(2N)} \mathbf{e}_{2N-n}\right)\right) \\ &= \frac{1}{\sqrt{2}}\left(e^{\pi i n/(2N)} \frac{1}{\sqrt{2N}} e^{2\pi i nk/(2N)} + e^{-\pi i n/(2N)} \frac{1}{\sqrt{2N}} e^{2\pi i (2N-n)k/(2N)}\right) \\ &= \frac{1}{\sqrt{2}}\left(e^{\pi i n/(2N)} \frac{1}{\sqrt{2N}} e^{2\pi i nk/(2N)} + e^{-\pi i n/(2N)} \frac{1}{\sqrt{2N}} e^{-2\pi i nk/(2N)}\right) \\ &= \frac{1}{2\sqrt{N}}\left(e^{2\pi i (n/(2N))(k+1/2)} + e^{-2\pi i (n/(2N))(k+1/2)}\right) = \frac{1}{\sqrt{N}} \cos\left(2\pi \frac{n}{2N} \left(k + \frac{1}{2}\right)\right). \end{aligned}$$

Since F_{2N} is unitary, and thus preserves the scalar product, the given vectors are orthonormal. \blacksquare

We need to address one final thing before we can define the DCT: The vector \mathbf{x} we start with is in \mathbb{R}^N , but the vectors above are in \mathbb{R}^{2N} . We would like to have orthonormal vectors in \mathbb{R}^N , so that we can use them to decompose \mathbf{x} . It is possible to show with a direct argument that, when we restrict the vectors above to the first N elements, they are still orthogonal. We will, however, apply a more instructive argument to show this, which gives us some intuition into the connection with symmetric filters. We start with the following result, which shows that a filter preserves symmetric vectors if and only if the filter is symmetric.

Theorem 4.7. Let S be a filter. The following are equivalent

1. S preserves symmetric vectors (i.e. $S\mathbf{x}$ is a symmetric vector whenever \mathbf{x} is).
2. The set of filter coefficients of S is a symmetric vector.

Also, when S preserves symmetric vectors, the following hold:

1. The vector of filter coefficients has an integer symmetry point if and only if the input and output have the same type (integer or non-integer) of symmetry point.
2. The input and output have the same symmetry point if and only if the filter is symmetric.

Proof: Assume that the filter S maps a symmetric vector with symmetry at d_1 to another symmetric vector. Let \mathbf{x} be the symmetric vector so that $(\hat{\mathbf{x}})_n = e^{-2\pi i d_1 n/N}$ for $n < N/2$. Since the output is a symmetric vector, we must have that

$$\lambda_{S,n} e^{-2\pi i d_1 n/N} = z_n e^{-2\pi i d_2 n/N}$$

for some d_2 , z_n and for $n < N/2$. But this means that $\lambda_{S,n} = y_n e^{-2\pi i(d_2 - d_1)n/N}$. Similar reasoning applies for $n > N/2$, so that $\lambda_{S,n}$ clearly equals \hat{s} for some symmetric vector s from Theorems 4.3 and 4.4. This vector equals (up to multiplication with \sqrt{N}) the filter coefficients of S , which therefore is a symmetric. Moreover, it is clear that the filter coefficients have an integer symmetry point if and only if the input and output vector either both have an integer symmetry point, or both a non-integer symmetry point. ■

Since the filter coefficients of a filter which preserves symmetric vectors also is a symmetric vector, this means that its frequency response takes the form $\lambda_{S,n} = z_n e^{-2\pi i d n / N}$, where z is a real vector. This means that the phase (argument) of the frequency response is $-2\pi d n / N$ or $\pi - 2\pi d n / N$, depending on the sign of z_n . In other words, the phase is linear in n . Filters which preserve symmetric vectors are therefore also called *linear phase filters*.

Note also that the case $d = 0$ or $d = N - 1/2$ corresponds to symmetric filters. An example of linear phase filters which are not symmetric are smoothing filters where the coefficients are taken from odd rows in Pascal's triangle. We continue with the following result:

Theorem 4.8. (Characterization of filters which preserve vectors symmetric about $N - 1/2$) Assume that S is a symmetric filter, and let $y = Sx$, where $x, y \in \mathbb{R}^{2N}$. The mapping $S_r : \mathbb{R}^N \rightarrow \mathbb{R}^N$ which sends (x_0, \dots, x_{N-1}) to (y_0, \dots, y_{N-1}) is well-defined and symmetric. The restriction of S to vectors symmetric about $N - 1/2$ is uniquely characterized by S_r . Moreover, if we write $S = \begin{pmatrix} S_1 & S_2 \\ S_3 & S_4 \end{pmatrix}$, we have that $S_r = S_1 + (S_2)^f$, where $(S_2)^f$ is the matrix S_2 with the columns reversed.

Proof: With S as in the text of the theorem, we compute

$$\begin{pmatrix} y_0 \\ \vdots \\ y_{N-1} \end{pmatrix} = \begin{pmatrix} S_1 & S_2 \end{pmatrix} \begin{pmatrix} x_0 \\ \vdots \\ x_{N-1} \\ \hline x_N \\ \vdots \\ x_{2N-1} \end{pmatrix} = S_1 \begin{pmatrix} x_0 \\ \vdots \\ x_{N-1} \end{pmatrix} + S_2 \begin{pmatrix} x_N \\ \vdots \\ x_{2N-1} \end{pmatrix}.$$

When x is a symmetric vector we can rewrite this as

$$\begin{aligned} & S_1 \begin{pmatrix} x_0 \\ \vdots \\ x_{N-1} \end{pmatrix} + S_2 \begin{pmatrix} x_{N-1} \\ \vdots \\ x_0 \end{pmatrix} \\ &= S_1 \begin{pmatrix} x_0 \\ \vdots \\ x_{N-1} \end{pmatrix} + (S_2)^f \begin{pmatrix} x_0 \\ \vdots \\ x_{N-1} \end{pmatrix} = \left(S_1 + (S_2)^f \right) \begin{pmatrix} x_0 \\ \vdots \\ x_{N-1} \end{pmatrix}. \end{aligned}$$

This shows that the mapping $S_r : (x_0, \dots, x_{N-1}) \rightarrow (y_0, \dots, y_{N-1})$ is well-defined, and that $S_r = S_1 + (S_2)^f$. Note that S_2 is not symmetric, but it is constant on all diagonals since S is a filter. But then $(S_2)^f$ is constant on all anti-diagonals, and is in particular symmetric. But then $S_r = S_1 + (S_2)^f$ is also symmetric. This completes the proof. ■

Note that S_r is not a digital filter, since its matrix is not circulant. In particular, its eigenvectors are not pure tones. In the block matrix factorization of S , S_2 contains the circulant part of the matrix, and forming $(S_2)^f$ means that the circulant parts switch corners. With the help of Theorem 4.8 we can finally establish the orthogonality of the cosine-vectors in \mathbb{R}^N .

Corollary 4.9 (Basis of eigenvectors for S_r). Let S be a symmetric filter, and let S_r be the mapping defined in Theorem 4.8. Define

$$d_{n,N} = \begin{cases} \sqrt{\frac{1}{N}}, & n=0 \\ \sqrt{\frac{2}{N}}, & 1 \leq n < N \end{cases}$$

and $\mathbf{d}_n = d_{n,N} \cos(2\pi \frac{n}{2N} (k + \frac{1}{2}))$ for $0 \leq n \leq N-1$, then $\{\mathbf{d}_0, \mathbf{d}_1, \dots, \mathbf{d}_{N-1}\}$ is an orthonormal basis of eigenvectors for S_r .

Proof: Let S be a symmetric filter of length $2N$. We know then that $\lambda_{S,n} = \lambda_{S,2N-n}$, so that

$$\begin{aligned} & S \left(\cos \left(2\pi \frac{n}{2N} \left(k + \frac{1}{2} \right) \right) \right) \\ &= S \left(\frac{1}{2} \left(e^{2\pi i(n/(2N))(k+1/2)} + e^{-2\pi i(n/(2N))(k+1/2)} \right) \right) \\ &= \frac{1}{2} \left(e^{\pi i n / (2N)} S \left(e^{2\pi i n k / (2N)} \right) + e^{-\pi i n / (2N)} S \left(e^{-2\pi i n k / (2N)} \right) \right) \\ &= \frac{1}{2} \left(e^{\pi i n / (2N)} \lambda_{S,n} e^{2\pi i n k / (2N)} + e^{-\pi i n / (2N)} \lambda_{S,2N-n} e^{-2\pi i n k / (2N)} \right) \\ &= \frac{1}{2} \left(\lambda_{S,n} e^{2\pi i(n/(2N))(k+1/2)} + \lambda_{S,2N-n} e^{-2\pi i(n/(2N))(k+1/2)} \right) \\ &= \lambda_{S,n} \frac{1}{2} \left(e^{2\pi i(n/(2N))(k+1/2)} + e^{-2\pi i(n/(2N))(k+1/2)} \right) \\ &= \lambda_{S,n} \cos \left(2\pi \frac{n}{2N} \left(k + \frac{1}{2} \right) \right), \end{aligned}$$

where we have used that $e^{2\pi i n k / (2N)}$ is an eigenvector of S with eigenvalue $\lambda_{S,n}$, and $e^{-2\pi i n k / (2N)} = e^{2\pi i(2N-n)k/(2N)}$ is an eigenvector of S with eigenvalue $\lambda_{S,2N-n}$. This shows that the vectors are eigenvectors for symmetric filters of length $2N$. It is also clear that the first half of the vectors must be eigenvectors for S_r with the same eigenvalue, since when $\mathbf{y} = S\mathbf{x} = \lambda_{S,n}\mathbf{x}$, we also have that

$$(y_0, y_1, \dots, y_{N-1}) = S_r(x_0, x_1, \dots, x_{N-1}) = \lambda_{S,n}(x_0, x_1, \dots, x_{N-1}).$$

To see why these vectors are orthogonal, choose at the outset a symmetric filter where $\{\lambda_{S,n}\}_{n=0}^{N-1}$ are distinct. Then the cosine-vectors of length N are also eigenvec-

tors with distinct eigenvalues, and they must be orthogonal since S_r is symmetric. Moreover, since

$$\begin{aligned}
& \sum_{k=0}^{2N-1} \cos^2 \left(2\pi \frac{n}{2N} \left(k + \frac{1}{2} \right) \right) \\
&= \sum_{k=0}^{N-1} \cos^2 \left(2\pi \frac{n}{2N} \left(k + \frac{1}{2} \right) \right) + \sum_{k=N}^{2N-1} \cos^2 \left(2\pi \frac{n}{2N} \left(k + \frac{1}{2} \right) \right) \\
&= \sum_{k=0}^{N-1} \cos^2 \left(2\pi \frac{n}{2N} \left(k + \frac{1}{2} \right) \right) + \sum_{k=0}^{N-1} \cos^2 \left(2\pi \frac{n}{2N} \left(k + N + \frac{1}{2} \right) \right) \\
&= \sum_{k=0}^{N-1} \cos^2 \left(2\pi \frac{n}{2N} \left(k + \frac{1}{2} \right) \right) + (-1)^{2n} \sum_{k=0}^{N-1} \cos^2 \left(2\pi \frac{n}{2N} \left(k + \frac{1}{2} \right) \right) \\
&= 2 \sum_{k=0}^{N-1} \cos^2 \left(2\pi \frac{n}{2N} \left(k + \frac{1}{2} \right) \right),
\end{aligned}$$

where we used that $\cos(x + n\pi) = (-1)^n \cos x$. This means that

$$\left\| \left\{ \cos \left(2\pi \frac{n}{2N} \left(k + \frac{1}{2} \right) \right) \right\}_{k=0}^{2N-1} \right\|^{2N-1} = \sqrt{2} \left\| \left\{ \cos \left(2\pi \frac{n}{2N} \left(k + \frac{1}{2} \right) \right) \right\}_{k=0}^{N-1} \right\|^{N-1}.$$

Thus, in order to make the vectors orthonormal when we consider the first N elements instead of all $2N$ elements, we need to multiply with $\sqrt{2}$. This gives us the vectors \mathbf{d}_n as defined in the text of the theorem. This completes the proof. ■

We now clearly see the analogy between symmetric functions and vectors: while the first can be written as a sum of cosine-functions, the second can be written as a sum of cosine-vectors. The orthogonal basis we have found is given its own name:

Definition 4.10 (DCT basis). We denote by \mathcal{D}_N the orthogonal basis $\{\mathbf{d}_0, \mathbf{d}_1, \dots, \mathbf{d}_{N-1}\}$. We also call \mathcal{D}_N the N -point DCT basis.

Using the DCT basis instead of the Fourier basis we can make the following definitions, which parallel those for the DFT:

Definition 4.11 (Discrete Cosine Transform). The change of coordinates from the standard basis of \mathbb{R}^N to the DCT basis \mathcal{D}_N is called the *discrete cosine transform* (or DCT). The $N \times N$ matrix D_N that represents this change of basis is called the (N -point) DCT matrix. If \mathbf{x} is a vector in \mathbb{R}^N , its coordinates $\mathbf{y} = (y_0, y_1, \dots, y_{N-1})$ relative to the DCT basis are called the DCT coefficients of \mathbf{x} (in other words, $\mathbf{y} = D_N \mathbf{x}$).

Note that we can also write

$$D_N = \sqrt{\frac{2}{N}} \begin{pmatrix} 1/\sqrt{2} & 0 & \cdots & 0 \\ 0 & 1 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 1 \end{pmatrix} (\cos(2\pi \frac{n}{2N}(k+1/2))). \quad (4.3)$$

Since this matrix is orthogonal, it is immediate that

$$\left(\cos\left(2\pi\frac{n}{2N}(k+1/2)\right)\right)^{-1} = \frac{2}{N} \left(\cos\left(2\pi\frac{n+1/2}{2N}k\right)\right) \begin{pmatrix} 1/2 & 0 & \cdots & 0 \\ 0 & 1 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 1 \end{pmatrix} \quad (4.4)$$

$$\left(\cos\left(2\pi\frac{n+1/2}{2N}k\right)\right)^{-1} = \frac{2}{N} \begin{pmatrix} 1/2 & 0 & \cdots & 0 \\ 0 & 1 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 1 \end{pmatrix} \left(\cos\left(2\pi\frac{n}{2N}(k+1/2)\right)\right). \quad (4.5)$$

In other words, not only can D_N be directly expressed in terms of a cosine-matrix, but our developments helped us to express the inverse of a cosine matrix in terms of other cosine-matrices. In the literature different types of cosine-matrices have been useful:

- I** Cosine-matrices with entries $\cos(2\pi nk/(2(N-1)))$.
- II** Cosine-matrices with entries $\cos(2\pi n(k+1/2)/(2N))$.
- III** Cosine-matrices with entries $\cos(2\pi(n+1/2)k/(2N))$.
- IV** Cosine-matrices with entries $\cos(2\pi(n+1/2)(k+1/2)/(2N))$.

We will call these type-I, type-II, type-III, and type-IV cosine-matrices, respectively. What we did above handles the case of type-II cosine-matrices. It will turn out that not all of these cosine-matrices are orthogonal, but that we in all cases, as we did above for type-II cosine matrices, can express the inverse of a cosine-matrix of one type in terms of a cosine-matrix of another type, and that any cosine-matrix is easily expressed in terms of an orthogonal matrix. These orthogonal matrices will be called $D_N^{(I)}$, $D_N^{(II)}$, $D_N^{(III)}$, and $D_N^{(IV)}$, respectively, and they are all called DCT-matrices. The D_N we constructed above is thus $D_N^{(II)}$. The type-II DCT matrix is the most commonly used, and the type is therefore often dropped when referring to these. We will consider the other cases of cosine-matrices in Section ?? (type-I, in connection with a different extension strategy used for wavelets). Type-IV cosine-matrices are considered in an exercise of this section.

As with the Fourier basis vectors, the DCT basis vectors are called synthesis vectors, since we can write

$$\mathbf{x} = y_0 \mathbf{d}_0 + y_1 \mathbf{d}_1 + \cdots + y_{N-1} \mathbf{d}_{N-1} \quad (4.6)$$

in the same way as for the DFT. Following the same reasoning as for the DFT, D_N^{-1} is the matrix where the \mathbf{d}_n are columns. But since these vectors are real and orthonormal, D_N must be the matrix where the \mathbf{d}_n are rows. Moreover, since Theorem 4.8 also states that the same vectors are eigenvectors for filters which preserve symmetric extensions, we can state the following:

Theorem 4.12. D_N is the orthogonal matrix where the rows are \mathbf{d}_n . Moreover, for any digital filter S which preserves symmetric extensions, $(D_N)^T$ diagonalizes S_r , i.e. $S_r = D_N^T D D_N$ where D is a diagonal matrix.

Let us also make the following definition:

Definition 4.13 (IDCT). We will call $\mathbf{x} = (D_N)^T \mathbf{y}$ the inverse DCT or (IDCT) of \mathbf{x} .

Example 4.14. As with Example 2.19, exact expressions for the DCT can be written down just for a few specific cases. It turns out that the case $N = 4$ as considered in Example 2.19 does not give the same type of nice, exact values, so let us instead consider the case $N = 2$. We have that

$$D_4 = \begin{pmatrix} \frac{1}{\sqrt{2}} \cos(0) & \frac{1}{\sqrt{2}} \cos(0) \\ \cos\left(\frac{\pi}{2}\left(0 + \frac{1}{2}\right)\right) & \cos\left(\frac{\pi}{2}\left(1 + \frac{1}{2}\right)\right) \end{pmatrix} = \begin{pmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} \end{pmatrix}$$

The DCT of the same vector as in Example 2.19 can now be computed as:

$$D_2 \begin{pmatrix} 1 \\ 2 \end{pmatrix} = \begin{pmatrix} \frac{3}{\sqrt{2}} \\ -\frac{1}{\sqrt{2}} \end{pmatrix}.$$



Example 4.15. A direct implementation of the DCT could be made as follows:

```
function y=DCTImpl(x)
N=length(x);
DN=zeros(N);
DN(1,:)=ones(1,N)/sqrt(N);
for n=1:N
    DN(n,:)=cos(2*pi*((n-1)/(2*N))*((0:N-1)+1/2))*sqrt(2/N);
end
y=DN*x;
```

Matlab's functions for computing the DCT and IDCT are called `dct`, and `idct`, respectively. These are defined in Matlab exactly as they are here, contrary to the case for the FFT (where a different normalizing factor was used). In the next chapter we will see that one also can make a much more efficient implementation of the DCT than this.

With these functions we can repeat examples 2.29- 2.31, by simply replacing the calls to `DFTImpl/IDFTImpl` with calls to the DCT counterparts. You may not here much improvements in these simple experiments, but in theory the DCT should be able to approximate sound better. ♣

Similarly to Theorem 2.25 for the DFT, one can think of the DCT as a least squares approximation and the unique representation of a function having the same sample values, but this time in terms of sinusoids instead of complex exponentials:

Theorem 4.16 (Interpolation with the DCT basis). Let f be a function defined on the interval $[0, T]$, and let \mathbf{x} be the sampled vector given by

$$x_k = f((2k+1)T/(2N)) \quad \text{for } k = 0, 1, \dots, N-1.$$

There is exactly one linear combination $g(t)$ on the form

$$\sum_{n=0}^{N-1} y_n d_{n,N} \cos(2\pi(n/2)t/T)$$

which satisfies the conditions

$$g((2k+1)T/(2N)) = f((2k+1)T/(2N)), \quad k = 0, 1, \dots, N-1,$$

and its coefficients are determined by $\mathbf{y} = D_N \mathbf{x}$.

Proof: This follows by inserting $t = (2k+1)T/(2N)$ in the equation

$$g(t) = \sum_{n=0}^{N-1} y_n d_{n,N} \cos(2\pi(n/2)t/T)$$

to arrive at the equations

$$f(kT/N) = \sum_{n=0}^{N-1} y_n d_{n,N} \cos\left(2\pi \frac{n}{2N} \left(k + \frac{1}{2}\right)\right) \quad 0 \leq k \leq N-1.$$

This gives us an equation system for finding the y_n with the invertible DCT matrix as coefficient matrix, and the result follows. \blacksquare

Due to this there is a slight difference to how we applied the DFT, due to the subtle change in the sample points, from kT/N for the DFT, to $(2k+1)T/(2N)$ for the DCT. The sample points for the DCT are thus the midpoints on the intervals in a uniform partition of $[0, T]$ into N intervals, while they for the DFT are the start points on the intervals. Also, the frequencies are divided by 2. In Figure 4.2 we have plotted the sinusoids of Theorem 4.16 for $T = 1$, as well as the sample points used in that theorem. The sample points in (a) correspond to the first column in the DCT matrix, the sample points in (b) to the second column of the DCT matrix, and so on (up to normalization with $d_{n,N}$). As n increases, the functions oscillate more and more. As an example, y_5 says how much content of maximum oscillation there is. In other words, the DCT of an audio signal shows the proportion of the different frequencies in the signal, and the two formulas $\mathbf{y} = D_N \mathbf{x}$ and $\mathbf{x} = (D_N)^T \mathbf{y}$ allow us to switch back and forth between the time domain representation and the frequency domain representation of the sound. In other words, once we have computed $\mathbf{y} = D_N \mathbf{x}$, we can analyse the frequency content of \mathbf{x} . If we want to reduce the bass we can decrease the y -values with small indices and if we want to increase the treble we can increase the y -values with large indices.

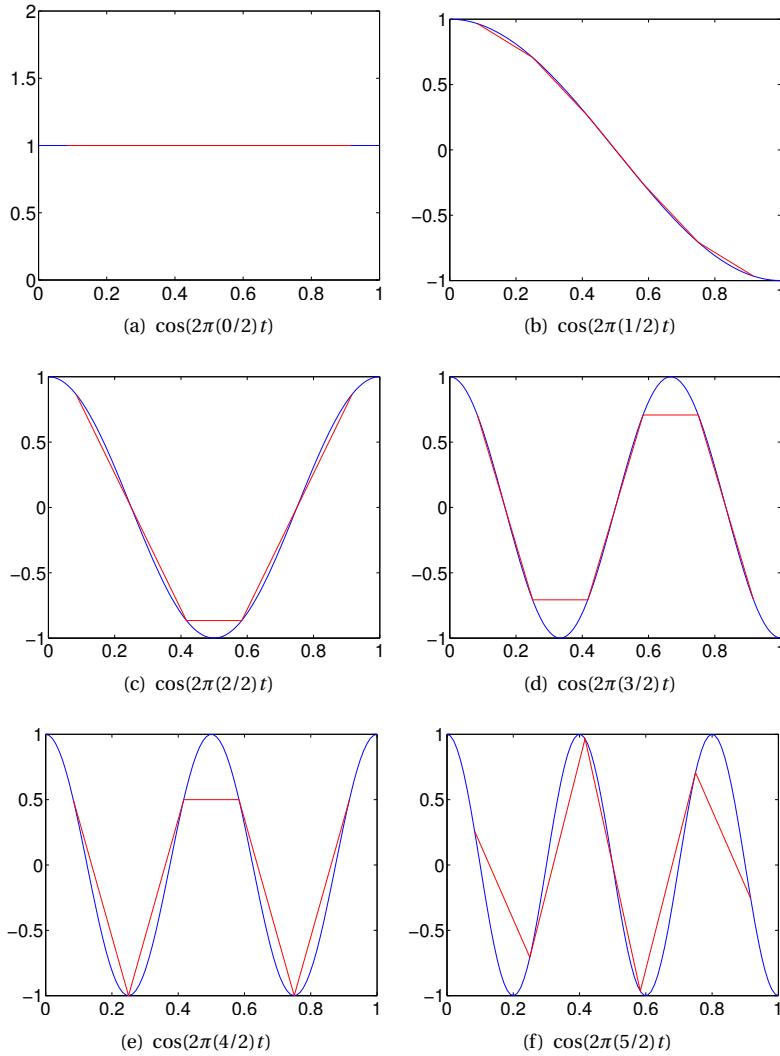


Figure 4.2: The 6 different sinusoids used in DCT for $N = 6$, i.e. $\cos(2\pi(n/2)t)$, $0 \leq n < 6$. The plots also show piecewise linear functions (in red) between the sample points $\frac{2k+1}{2N}$, $0 \leq k < 6$, since only the values at these points are used in Theorem 4.16.

Exercises for Section 4.1

- Consider the matrix

$$S = \frac{1}{3} \begin{pmatrix} 2 & 1 & 0 & 0 & 0 & 0 \\ 1 & 1 & 1 & 0 & 0 & 0 \\ 0 & 1 & 1 & 1 & 0 & 0 \\ 0 & 0 & 1 & 1 & 1 & 0 \\ 0 & 0 & 0 & 1 & 1 & 1 \\ 0 & 0 & 0 & 0 & 1 & 2 \end{pmatrix}$$

- a.** Compute the eigenvalues and eigenvectors of S using the results of this section. You should only need to perform one DFT or one DCT in order to achieve this.
- b.** Use Matlab to compute the eigenvectors and eigenvalues of S also. What are the differences from what you found in a.?
- c.** Find a filter T so that $S = T_r$. What kind of filter is T ?
- 2.** Consider the averaging filter $S = \{\frac{1}{4}, \frac{1}{2}, \frac{1}{4}\}$. Write down the matrix S_r for the case when $N = 4$.
- 3.** As in Example 4.14, state the exact cartesian form of the DCT matrix for the case $N = 3$.
- 4.** Show that the vectors $\left\{ \cos\left(2\pi \frac{n+\frac{1}{2}}{2N} (k + \frac{1}{2})\right) \right\}_{n=0}^{N-1}$ in \mathbb{R}^N are orthogonal, with lengths $\sqrt{N/2}$. This means that the matrix with entries $\sqrt{\frac{2}{N}} \cos\left(2\pi \frac{n+\frac{1}{2}}{2N} (k + \frac{1}{2})\right)$ is orthogonal. Since this matrix also is symmetric, it is its own inverse. This is the DCT-IV, which we denote by $D_N^{(\text{IV})}$. Although we will not consider this, the DCT-IV also has an efficient implementation. Hint: Compare with the orthogonal vectors \mathbf{d}_n , used in the DCT.
- 5.** The MDCT is defined as the $N \times (2N)$ -matrix M with elements $M_{n,k} = \cos(2\pi(n+1/2)(k+1/2+N/2)/(2N))$. This exercise will take you through the details of the transformation which corresponds to multiplication with this matrix. The MDCT is very useful, and is also used in the MP3 standard and in more recent standards.

a. Show that

$$M = \sqrt{\frac{N}{2}} D_N^{(\text{IV})} \begin{pmatrix} \mathbf{0} & A \\ B & \mathbf{0} \end{pmatrix}$$

where A and B are the $(N/2) \times N$ -matrices

$$A = \begin{pmatrix} \dots & \dots & 0 & -1 & -1 & 0 & \dots & \dots \\ \vdots & \vdots \\ 0 & -1 & \dots & \dots & \dots & \dots & -1 & 0 \\ -1 & 0 & \dots & \dots & \dots & \dots & 0 & -1 \end{pmatrix} = \begin{pmatrix} -I_{N/2}^f & -I_{N/2} \\ I_{N/2} & -I_{N/2}^f \end{pmatrix}$$

$$B = \begin{pmatrix} 1 & 0 & \dots & \dots & \dots & \dots & 0 & -1 \\ 0 & 1 & \dots & \dots & \dots & \dots & -1 & 0 \\ \vdots & \vdots \\ \dots & \dots & 0 & 1 & -1 & 0 & \dots & \dots \end{pmatrix} = \begin{pmatrix} I_{N/2} & -I_{N/2}^f \\ I_{N/2}^f & -I_{N/2} \end{pmatrix}.$$

Due to this expression, any algorithm for the DCT-IV can be used to compute the MDCT.

b. The MDCT is not invertible, since it is not a square matrix. We will show here that it still can be used in connection with invertible transformations. We first define the IMDCT as the matrix M^T/N . Transposing the matrix expression we obtained in a. gives

$$\frac{1}{\sqrt{2N}} \begin{pmatrix} \mathbf{0} & B^T \\ A^T & \mathbf{0} \end{pmatrix} D_N^{(IV)}$$

for the IMDCT, which thus also has an efficient implementation. Show that if

$$\mathbf{x}_0 = (x_0, \dots, x_{N-1}) \quad \mathbf{x}_1 = (x_N, \dots, x_{2N-1}) \quad \mathbf{x}_2 = (x_{2N}, \dots, x_{3N-1})$$

and

$$\mathbf{y}_{0,1} = M \begin{pmatrix} \mathbf{x}_0 \\ \mathbf{x}_1 \end{pmatrix} \quad \mathbf{y}_{1,2} = M \begin{pmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{pmatrix}$$

(i.e. we compute two MDCT's where half of the data overlap), then

$$\mathbf{x}_1 = \{\text{IMDCT}(\mathbf{y}_{0,1})\}_{k=N}^{2N-1} + \{\text{IMDCT}(\mathbf{y}_{1,2})\}_{k=0}^{N-1}.$$

Even though the MDCT itself is not invertible, the input can still be recovered from overlapping MDCT's.

4.2 Improvements from using the DCT to interpolate functions and approximate analog filters

Recall that, in Section 3.2.1, we explained how to approximate an analog filter from the samples. It turns out that, when an analog filter is symmetric, we can use symmetric extensions to create a better approximation from the samples.

Assume that s is an analog filter, and that we apply it to a general function f . Denote as before the symmetric extension of f by \check{f} . We start with the following observation, which follows from the continuity of s .

Observation 4.17. Since $(\check{f})_N$ is a better approximation to \check{f} , compared to what f_N is to f , $s((\check{f})_N)$ is a better approximation to $s(\check{f})$, compared to what $s(f_N)$ is to $s(f)$.

Since $s(\check{f})$ agrees with $s(f)$ except near the boundaries, we can thus conclude that $s((\check{f})_N)$ is a better approximation to $s(f)$ than what $s(f_N)$ is.

We have seen that the restriction of s to $V_{M,T}$ is equivalent to an $N \times N$ digital filter S , where $N = 2M + 1$. Let \mathbf{x} be the samples of f , $\check{\mathbf{x}}$ the samples of \check{f} . Turning around the fact that $(\check{f})_N$ is a better approximation to \check{f} , compared to what f_N is to f , the following is clear.

Observation 4.18. The samples $\check{\mathbf{x}}$ are a better approximation to the samples of $(\check{f})_N$, than the samples \mathbf{x} are to the samples of f_N .

Now, let $\mathbf{z} = S\mathbf{x}$, and $\check{\mathbf{z}} = S\check{\mathbf{x}}$. The following is also clear from the preceding observation, due to continuity of the digital filter S .

Observation 4.19. $\check{\mathbf{z}}$ is a better approximation to $S(\text{samples of } (\check{f})_N) = \text{samples of } s((\check{f})_N)$, than \mathbf{z} is to $S(\text{samples of } f_N) = \text{samples of } s(f_N)$.

Since by Observation 4.17 $s((\check{f})_N)$ is a better approximation to the output $s(f)$, we conclude that $\check{\mathbf{z}}$ is a better approximation than \mathbf{z} to the samples of the output of the filter.

Observation 4.20. $S\check{\mathbf{x}}$ is a better approximation to the samples of $s(f)$ than $S\mathbf{x}$ is (\mathbf{x} are the samples of f).

Now, let us also bring in the assumption that s is symmetric. Then the corresponding digital filter S is also symmetric, and we know then that we can view its restriction to symmetric extensions in \mathbb{R}^{2N} in terms of the mapping $S_r : \mathbb{R}^N \rightarrow \mathbb{R}^N$. We can thus specialize Figure 3.1 to symmetric filters by adding the step of creating the symmetric extension, and replacing S with S_r . We have summarized these remarks in Figure 4.3. The DCT here appears, since we have used Theorem 4.16 to interpolate with the DCT basis, instead of the Fourier basis. Note that this also requires that the sampling is performed as required in that theorem, i.e. the samples are the midpoints on all intervals. This new sampling procedure is not indicated in Figure 4.3. Figure 4.3 can be further simplified to that shown in Figure 4.4.

Note that the assumption that s is symmetric only helped us to implement the approximation $s(\check{f})$ in a more efficient way, since S_r has N points and S has $2N$ points. $s(\check{f})$ can in any way be used as an approximation, even if s is not symmetric. But this approximation is actually even better when s is symmetric: Since s is symmetric, $s(\check{f})$ is a symmetric function (since \check{f} is a symmetric function). The N 'th order Fourier series of $s(\check{f})$ is $s((\check{f})_N) = (s(\check{f}))_N$, and this is a better approximation to $s(\check{f})$ since $s(\check{f})$ is a symmetric function. Since the procedure above obtained an approximation to (the samples of) $(s(\check{f}))_N$, it follows that the approximations are better when s is symmetric.

As mentioned in Section 3.2, interpolation of a function from its samples can be seen as a special case. This can thus be illustrated as in Figure 4.5. Note that the approximation lies in $V_{2M,2T}$ (i.e. it is in a higher order Fourier space), but the point is that the same number of samples is used.

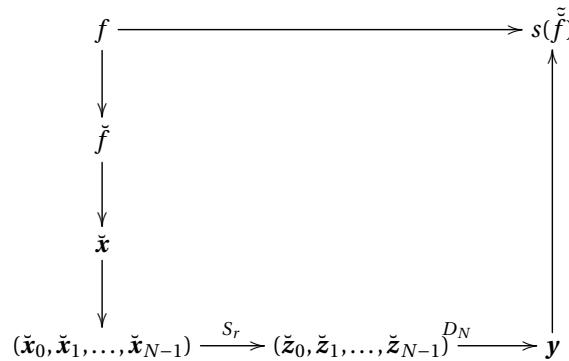


Figure 4.3: The connections between the new mapping S_r , sampling, and interpolation. The right vertical arrow represents interpolation with the DCT, i.e. that we compute $\sum_{n=0}^{N-1} y_n d_{n,N} \cos(2\pi(n/2)t/T)$ for values of t .

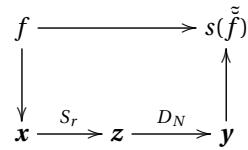


Figure 4.4: Simplification of Figure 4.3. The left vertical arrow represents sampling as dictated by the DCT.

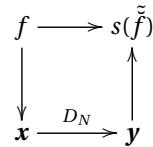


Figure 4.5: How we can approximate a function from its samples with the DCT.

4.2.1 Implementations of symmetric filters

Symmetric filters are also important for applications since they can be implemented efficiently. To see this, we can write

$$\begin{aligned}
 (\mathbf{S}\mathbf{x})_n &= \sum_{k=0}^{N-1} s_k x_{(n-k) \bmod N} \\
 &= s_0 x_n + \sum_{k=1}^{(N-1)/2} s_k x_{(n-k) \bmod N} + \sum_{k=(N+1)/2}^{N-1} s_k x_{(n-k) \bmod N} \\
 &= s_0 x_n + \sum_{k=1}^{(N-1)/2} s_k x_{(n-k) \bmod N} + \sum_{k=1}^{(N-1)/2} s_k x_{(N-(n-k)) \bmod N} \\
 &= s_0 x_n + \sum_{k=1}^{(N-1)/2} s_k (x_{(n-k) \bmod N} + x_{(k-n) \bmod N}). \tag{4.7}
 \end{aligned}$$

If we compare the first and last expressions here, we need the same number of summations, but the number of multiplications needed in the latter expression has been halved.

Observation 4.21. Assume that a symmetric filter has $2s + 1$ filter coefficients. The filter applied to a vector of length N can then be implemented using $(s + 1)N$ multiplications and $2sN$ additions. This gives a reduced number of arithmetic operations when compared to a filter with the same number of coefficients which is not symmetric, where a direct implementation requires $(2s + 1)N$ multiplications and $2sN$ additions.

Similarly to as in Section 3.8.2, a symmetric filter can be factored into a product of symmetric filters. To see how, note first that a real polynomial is symmetric if and only if $1/a$ is a root whenever a is. If we pair together the factors for the roots $a, 1/a$ when a is real we get a component in the frequency response of degree 2. If we pair the factors for the roots $a, 1/a, \bar{a}, 1/\bar{a}$ when a is complex, we get a component in the frequency response of degree 4. We thus get the following idea:

Idea 4.22. Let S be a symmetric filter with real coefficients. There exist constants $K, a_1, \dots, a_m, b_1, c_1, \dots, b_n, c_n$ so that

$$\begin{aligned}
 \lambda_S(\omega) &= K(a_1 e^{i\omega} + 1 + a_1 e^{-i\omega}) \dots (a_m e^{i\omega} + 1 + a_m e^{-i\omega}) \\
 &\quad \times (b_1 e^{2i\omega} + c_1 e^{i\omega} + 1 + c_1 e^{-i\omega} + b_1 e^{-2i\omega}) \dots (b_n e^{2i\omega} + c_n e^{i\omega} + 1 + c_n e^{-i\omega} + b_n e^{-2i\omega}).
 \end{aligned}$$

We can write $S = KA_1 \dots A_m B_1 \dots B_n$, where $A_i = \{a_i, 1, a_i\}$ and $B_i = \{b_i, c_i, 1, c_i, b_i\}$.

In any case we see that the component filters have 3 and 5 filter coefficients.

Exercises for Section 4.2

1. Recall that in Exercise 3.1.3 we wrote down formulas for the output of a filter. Using the results of this section these formulas can be written in a way which reduces the number of arithmetic operations. Assume that $S = t_{-E}, \dots, \underline{t_0}, \dots, t_E$ is a symmetric filter. Use Equation (4.7) to show that $z_n = (Sx)_n$ in this case can be split into the following different formulas, depending on n :

a. $0 \leq n < E$:

$$z_n = t_0 x_n + \sum_{k=1}^n t_k (x_{n+k} + x_{n-k}) + \sum_{k=n+1}^E t_k (x_{n+k} + x_{n-k+N}). \quad (4.8)$$

b. $E \leq n < N - E$:

$$z_n = t_0 x_n + \sum_{k=1}^E t_k (x_{n+k} + x_{n-k}). \quad (4.9)$$

c. $N - E \leq n < N$:

$$z_n = t_0 x_n + \sum_{k=1}^{N-1-n} t_k (x_{n+k} + x_{n-k}) + \sum_{k=N-1-n+1}^E t_k (x_{n+k-N} + x_{n-k}). \quad (4.10)$$

2. Assume that $S = t_{-E}, \dots, \underline{t_0}, \dots, t_E$ is a symmetric filter, and let $s = (t_0, \dots, t_E)$. Write a function

```
function z=filterS(s,x)
```

which takes the vector s as input, and returns $z = Sx$ using the formulas from Exercise 1.

3. Repeat Exercise 2.2.4 by reimplementing the functions `reducetreble` and `reducebass` using the function `filterS` from the previous exercise. The resulting sound files should sound the same, since the only difference is that we have modified the way we handle the beginning and end portion of the sound samples.

4.3 Efficient implementations of the DCT

When we defined the DCT in the preceding section, we converted to the frequency domain, so that the DFT is involved. This enables us to use efficient implementations of the DFT, such as the one we considered in Section 2.8. However, the way we defined the DCT, there is a penalty in that we need to compute a DFT of twice the length (the length doubles when we instead consider the symmetric extension). We are also forced to use complex arithmetic (note that any complex multiplication corresponds to 4 real multiplications, and that any complex addition corresponds to 2 real additions). Is there a way to get around these penalties, so that we can get an implementation of the DCT which is more efficient, and uses less additions and multiplications than the one you made in Exercise ???. The following theorem states an expression of the DCT which achieves this. This expression is, together with a

similar result for the DFT in the next section, much used in practical implementations:

Theorem 4.23 (DCT algorithm). Let $\mathbf{y} = D_N \mathbf{x}$ be the N -point DCT of the vector \mathbf{x} . Then we have that

$$y_n = c_{n,N} \left(\cos\left(\pi \frac{n}{2N}\right) \Re((F_N \mathbf{x}^{(1)}))_n + \sin\left(\pi \frac{n}{2N}\right) \Im((F_N \mathbf{x}^{(1)}))_n \right), \quad (4.11)$$

where $c_{0,N} = 1$ and $c_{n,N} = \sqrt{2}$ for $n \geq 1$, and where $\mathbf{x}^{(1)} \in \mathbb{R}^N$ is defined by

$$\begin{aligned} (\mathbf{x}^{(1)})_k &= x_{2k} \text{ for } 0 \leq k \leq N/2 - 1 \\ (\mathbf{x}^{(1)})_{N-k-1} &= x_{2k+1} \text{ for } 0 \leq k \leq N/2 - 1, \end{aligned}$$

Proof: The N -point DCT of \mathbf{x} is

$$y_n = d_{n,N} \sum_{k=0}^{N-1} x_k \cos\left(2\pi \frac{n}{2N} \left(k + \frac{1}{2}\right)\right).$$

Splitting this sum into two sums, where the indices are even and odd, we get

$$\begin{aligned} y_n &= d_{n,N} \sum_{k=0}^{N/2-1} x_{2k} \cos\left(2\pi \frac{n}{2N} \left(2k + \frac{1}{2}\right)\right) \\ &\quad + d_{n,N} \sum_{k=0}^{N/2-1} x_{2k+1} \cos\left(2\pi \frac{n}{2N} \left(2k + 1 + \frac{1}{2}\right)\right). \end{aligned}$$

If we reverse the indices in the second sum, this sum becomes

$$d_{n,N} \sum_{k=0}^{N/2-1} x_{N-2k-1} \cos\left(2\pi \frac{n}{2N} \left(N - 2k - 1 + \frac{1}{2}\right)\right).$$

If we then also shift the indices with $N/2$ in this sum, we get

$$\begin{aligned} &d_{n,N} \sum_{k=N/2}^{N-1} x_{2N-2k-1} \cos\left(2\pi \frac{n}{2N} \left(2N - 2k - 1 + \frac{1}{2}\right)\right) \\ &= d_{n,N} \sum_{k=N/2}^{N-1} x_{2N-2k-1} \cos\left(2\pi \frac{n}{2N} \left(2k + \frac{1}{2}\right)\right), \end{aligned}$$

where we used that \cos is symmetric and periodic with period 2π . We see that we now have the same cos-terms in the two sums. If we thus define the vector $\mathbf{x}^{(1)}$ as in

the text of the theorem, we see that we can write

$$\begin{aligned}
y_n &= d_{n,N} \sum_{k=0}^{N-1} (\mathbf{x}^{(1)})_k \cos\left(2\pi \frac{n}{2N} \left(2k + \frac{1}{2}\right)\right) \\
&= d_{n,N} \Re\left(\sum_{k=0}^{N-1} (\mathbf{x}^{(1)})_k e^{-2\pi i n(2k+\frac{1}{2})/(2N)}\right) \\
&= \sqrt{N} d_{n,N} \Re\left(e^{-\pi i n/(2N)} \frac{1}{\sqrt{N}} \sum_{k=0}^{N-1} (\mathbf{x}^{(1)})_k e^{-2\pi i nk/N}\right) \\
&= c_{n,N} \Re\left(e^{-\pi i n/(2N)} (F_N \mathbf{x}^{(1)})_n\right) \\
&= c_{n,N} \left(\cos\left(\pi \frac{n}{2N}\right) \Re((F_N \mathbf{x}^{(1)})_n) + \sin\left(\pi \frac{n}{2N}\right) \Im((F_N \mathbf{x}^{(1)})_n) \right),
\end{aligned}$$

where we have recognized the N -point DFT, and where $c_{n,N} = \sqrt{N} d_{n,N}$. Inserting the values for $d_{n,N}$, we see that $c_{0,N} = 1$ and $c_{n,N} = \sqrt{2}$ for $n \geq 1$, which agrees with the definition of $c_{n,N}$ in the theorem. This completes the proof. ■

With the result above we have avoided computing a DFT of double size. If we in the proof above define the $N \times N$ -diagonal matrix Q_N by $Q_{n,n} = c_{n,N} e^{-\pi i n/(2N)}$, the result can also be written on the more compact form

$$\mathbf{y} = D_N \mathbf{x} = \Re(Q_N F_N \mathbf{x}^{(1)}).$$

We will, however, not use this form, since there is complex arithmetic involved, contrary to Equation(4.11). Let us see how we can use (4.11) to implement the DCT, once we already have implemented the DFT in terms of the function FFTImpl as in Section 2.8:

```

function y = DCTImpl(x)
N = length(x);
if N == 1
    y = x;
else
    x1 = [x(1:2:(N-1)); x(N:(-2):2)];
    y = FFTImpl(x1);
    rp = real(y);
    ip = imag(y);
    y = cos(pi*((0:(N-1))'')/(2*N)).*rp + sin(pi*((0:(N-1))'')/(2*N)).*ip;
    y(2:N) = sqrt(2)*y(2:N);
end

```

In the code, the vector $\mathbf{x}^{(1)}$ is created first by rearranging the components, and it is sent as input to FFTImpl. After this we take real parts and imaginary parts, and multiply with the cos- and sin-terms in Equation (4.11).

4.3.1 Efficient implementations of the IDCT

As with the FFT, it is straightforward to modify the DCT implementation so that it returns the IDCT. To see how we can do this, write from Theorem 4.23, for $n \geq 1$

$$\begin{aligned} y_n &= c_{n,N} \left(\cos\left(\pi \frac{n}{2N}\right) \Re((F_N \mathbf{x}^{(1)})_n) + \sin\left(\pi \frac{n}{2N}\right) \Im((F_N \mathbf{x}^{(1)})_n) \right) \\ y_{N-n} &= c_{N-n,N} \left(\cos\left(\pi \frac{N-n}{2N}\right) \Re((F_N \mathbf{x}^{(1)})_{N-n}) + \sin\left(\pi \frac{N-n}{2N}\right) \Im((F_N \mathbf{x}^{(1)})_{N-n}) \right) \\ &= c_{n,N} \left(\sin\left(\pi \frac{n}{2N}\right) \Re((F_N \mathbf{x}^{(1)})_n) - \cos\left(\pi \frac{n}{2N}\right) \Im((F_N \mathbf{x}^{(1)})_n) \right), \end{aligned} \quad (4.12)$$

where we have used the symmetry of F_N for real signals. These two equations enable us to determine $\Re((F_N \mathbf{x}^{(1)})_n)$ and $\Im((F_N \mathbf{x}^{(1)})_n)$ from y_n and y_{N-n} . We get

$$\begin{aligned} \cos\left(\pi \frac{n}{2N}\right) y_n + \sin\left(\pi \frac{n}{2N}\right) y_{N-n} &= c_{n,N} \Re((F_N \mathbf{x}^{(1)})_n) \\ \sin\left(\pi \frac{n}{2N}\right) y_n - \cos\left(\pi \frac{n}{2N}\right) y_{N-n} &= c_{n,N} \Im((F_N \mathbf{x}^{(1)})_n). \end{aligned}$$

Adding we get

$$\begin{aligned} c_{n,N} (F_N \mathbf{x}^{(1)})_n &= \cos\left(\pi \frac{n}{2N}\right) y_n + \sin\left(\pi \frac{n}{2N}\right) y_{N-n} + i(\sin\left(\pi \frac{n}{2N}\right) y_n - \cos\left(\pi \frac{n}{2N}\right) y_{N-n}) \\ &= (\cos\left(\pi \frac{n}{2N}\right) + i \sin\left(\pi \frac{n}{2N}\right))(y_n - i y_{N-n}) = e^{\pi i n / (2N)} (y_n - i y_{N-n}). \end{aligned}$$

This means that $(F_N \mathbf{x}^{(1)})_n = \frac{1}{c_{n,N}} e^{\pi i n / (2N)} (y_n + i y_{N-n})$ for $n \geq 1$. Since $\Im((F_N \mathbf{x}^{(1)})_0) = 0$ we have that $(F_N \mathbf{x}^{(1)})_0 = \frac{1}{c_{0,N}} y_0$. This means that $\mathbf{x}^{(1)}$ can be recovered by taking the IDFT of the vector with component 0 being $\frac{1}{c_{0,N}} y_0$, and the remaining components being $\frac{1}{c_{n,N}} e^{\pi i n / (2N)} (y_n - i y_{N-n})$:

Theorem 4.24 (IDCT algorithm). Let $\mathbf{x} = (D_N)^T \mathbf{y}$ be the IDCT of \mathbf{y} , and let \mathbf{z} be the vector with component 0 being $\frac{1}{c_{0,N}} y_0$, and the remaining components being $\frac{1}{c_{n,N}} e^{\pi i n / (2N)} (y_n - i y_{N-n})$. Then we have that

$$\mathbf{x}^{(1)} = (F_N)^H \mathbf{z},$$

where $\mathbf{x}^{(1)}$ is defined as in Theorem 4.23.

The implementation of IDCT can thus go as follows:

```
function x = IDCTImpl(y)
N = length(y);
if N == 1
    x = y(1);
else
    Q=exp(pi*i*((0:(N-1))')/(2*N));
```

```

Q(2:N)=Q(2:N)/sqrt(2);
yrev=y(N:(-1):2);
toapply=[ y(1); Q(2:N).*((y(2:N)-1i*yrev) ] ;
x1=IFFTImpl(toapply);
x=zeros(N,1);
x(1:2:(N-1))=x1(1:(N/2));
x(2:2:N)=x1(N:(-1):(N/2+1));
end

```

4.3.2 Reduction in the number of multiplications with the DCT

Let us also state a result which confirms that the DCT and IDCT implementations we have described give the same type of reductions in the number of multiplications as the FFT and IFFT:

Theorem 4.25. (Number of multiplications required by the DCT and IDCT algorithms) Assume that the FFT algorithm from Section 2.8 is used. Both the N -point DCT and IDCT factorizations given by Theorem 4.23 and Theorem 4.24 then require $O(N \log_2 N + 3N/2)$ real multiplications. In comparison, the number of real multiplications required by a direct implementation of the N -point DCT and IDCT is N^2 .

Proof: By Theorem 2.40, the number of multiplications required by the FFT algorithm from Section 2.8 is $O(N \log_2 N)$. By Theorem 4.23, two additional multiplications are needed for each index. But in Exercise 1 we show that the joint computation of y_n and y_{N-n} can be done with 3 multiplications rather than 4, so that we can get away with $3N/2$ additional multiplications instead. Thus, a total number of $O(N \log_2 N + 3N/2)$ multiplications are required by the DCT. Since the number of multiplications for the IFFT is the same as for the FFT, we only need to count the additional multiplications needed in forming the vector $z = \frac{1}{c_{n,N}} e^{\pi i n/(2N)} (y_n - iy_{N-n})$. The same exercise shows that this can be done by $3N/2$ additional multiplications, so that the number of multiplications needed by the IDCT is $O(N \log_2 N + 3N/2)$ also. It is clear that the direct implementation of the DCT and IDCT needs N^2 multiplications, since only real arithmetic is involved. ■

Since the DCT and IDCT can be implemented using the FFT and IFFT, it has the same advantages as the FFT when it comes to parallel computing. Much literature is devoted to reducing the number of multiplications in the DFT and the DCT even further than what we have done (see [18] for one of the most recent developments). Some more notes on computational complexity are in order. For instance, we have not counted the operations sin and cos in the DCT. The reason is that these values can be precomputed, since we take the sine and cosine of a specific set of values for each DCT or DFT of a given size. This is contrary to multiplication and addition, since these include the input values, which are only known at runtime. We have, however, not written down that we use precomputed arrays for sine and cosine in our algorithms: This is an issue to include in more optimized algorithms. Another

point has to do with the multiplication of $\frac{1}{\sqrt{N}}$. As long as $N = 2^{2r}$, multiplication with N need not be considered as a multiplication, since it can be implemented using a bitshift.

Exercises for Section 4.3

1. In this exercise we will take a look at a small trick which reduces the number of multiplications needed by the DCT algorithm from Theorem 4.23. This is very similar to Exercise 2.8.9, where a simple observation lead us to halve the number of multiplications from $O(2N \log_2 N)$ to $O(N \log_2 N)$ in the FFT and IFFT algorithms.

a. Assume that \mathbf{x} is a real signal. Equation (4.12), which said that

$$y_n = c_{n,N} \left(\cos\left(\pi \frac{n}{2N}\right) \Re((F_N \mathbf{x}^{(1)}))_n + \sin\left(\pi \frac{n}{2N}\right) \Im((F_N \mathbf{x}^{(1)}))_n \right)$$

$$y_{N-n} = c_{n,N} \left(\sin\left(\pi \frac{n}{2N}\right) \Re((F_N \mathbf{x}^{(1)}))_n - \cos\left(\pi \frac{n}{2N}\right) \Im((F_N \mathbf{x}^{(1)}))_n \right)$$

for the n 'th and $N - n$ 'th coefficient of the DCT. This can also be rewritten as

$$y_n = c_{n,N} (\Re((F_N \mathbf{x}^{(1)}))_n + \Im((F_N \mathbf{x}^{(1)}))_n) \cos\left(\pi \frac{n}{2N}\right) - \Im((F_N \mathbf{x}^{(1)}))_n (\cos\left(\pi \frac{n}{2N}\right) - \sin\left(\pi \frac{n}{2N}\right))$$

$$y_{N-n} = c_{n,N} (-(\Re((F_N \mathbf{x}^{(1)}))_n + \Im((F_N \mathbf{x}^{(1)}))_n) \cos\left(\pi \frac{n}{2N}\right) + \Re((F_N \mathbf{x}^{(1)}))_n (\sin\left(\pi \frac{n}{2N}\right) + \cos\left(\pi \frac{n}{2N}\right))).$$

Explain that the first two equations require 4 multiplications to compute y_n and y_{N-n} , and that the last two equations require 3 multiplications to compute y_n and y_{N-n} (Do not count the multiplications with $c_{n,N}$).

b. Explain why the trick in a. reduces the number of multiplications in a DCT, so that it needs only $M_N = 3N/2$ multiplications plus the number of multiplications needed by an N -point DFT.

c. Explain why the trick in a. can be used to reduce the number of multiplications in an IDCT also to $M_N = 3N/2$ plus the number of multiplications needed by an N -point IDFT.

Hint: match the expression $e^{\pi i n / (2N)} (y_n - i y_{N-n})$ you encountered in the IDCT with the rewriting you did in b.

d. Show that the penalty of the trick we here have used to reduce the number of multiplications, is an increase in the number of additions. Why can this trick still be useful? (here we thus have similarities with

2. (An efficient joint implementation of the DCT and the FFT). In this exercise we will explain a joint implementation of the DFT and the DCT, which further reduces the number of arithmetic operations needed. It also has the additional benefit that

it avoids complex arithmetic altogether, and has a very structured implementation (this is not always the case for the quickest FFT implementations. Not surprisingly, one often sacrifices clarity of code when one pursues higher computational speed). For further details, the reader is referred to [35]

- a.** Let $\mathbf{y} = F_N \mathbf{x}$ be the N -point DFT of the real vector \mathbf{x} . Show that

$$\Re(y_n) = \begin{cases} \frac{1}{\sqrt{2}} \Re((F_{N/2} \mathbf{x}^{(e)})_n) + (ED_{N/4} \mathbf{z})_n & 0 \leq n \leq N/4 - 1 \\ \frac{1}{\sqrt{2}} \Re((F_{N/2} \mathbf{x}^{(e)})_n) & n = N/4 \\ \frac{1}{\sqrt{2}} \Re((F_{N/2} \mathbf{x}^{(e)})_n) - (ED_{N/4} \mathbf{z})_{N/2-n} & N/4 + 1 \leq n \leq N/2 - 1 \end{cases} \quad (4.13)$$

$$\Im(y_n) = \begin{cases} \frac{1}{\sqrt{2}} \Im((F_{N/2} \mathbf{x}^{(e)})_n) & n = 0 \\ \frac{1}{\sqrt{2}} \Im((F_{N/2} \mathbf{x}^{(e)})_n) + (ED_{N/4} \mathbf{w})_{N/4-n} & 1 \leq n \leq N/4 - 1 \\ \frac{1}{\sqrt{2}} \Im((F_{N/2} \mathbf{x}^{(e)})_n) + (ED_{N/4} \mathbf{w})_{n-N/4} & N/4 \leq n \leq N/2 - 1 \end{cases} \quad (4.14)$$

where $\mathbf{x}^{(e)}$ is as defined in Theorem 2.34, where $\mathbf{z}, \mathbf{w} \in \mathbb{R}^{N/4}$ defined by

$$\begin{aligned} z_k &= x_{2k+1} + x_{N-2k-1} & 0 \leq k \leq N/4 - 1, \\ w_k &= (-1)^k (x_{N-2k-1} - x_{2k+1}) & 0 \leq k \leq N/4 - 1, \end{aligned}$$

and where E is a diagonal matrix with diagonal entries $E_{0,0} = \frac{1}{2}$ and $E_{n,n} = \frac{1}{2\sqrt{2}}$ for $n \geq 1$.

- a.** says nothing about the coefficients y_n for $n > \frac{N}{2}$. These are obtained in the same way as before through symmetry. **a.** also says nothing about $y_{N/2}$. This can be obtained with the same formula as in Theorem 2.34.

- b.** Write down a matrix interpretation of what you found in a.

- c.** Write down an algorithm which implements the result in a..

While our previous FFT algorithm says that an FFT of length N can be computed as two FFT's of length $N/2$ (on $\mathbf{x}^{(e)}$ and $\mathbf{x}^{(o)}$), our revised algorithm says that an FFT of length N can be computed as one FFT of length $N/2$ (on $\mathbf{x}^{(e)}$), together with two DCT's of length $N/4$ (on \mathbf{z} and \mathbf{w}). Let us now compute the number of arithmetic operations our revised algorithm needs. Denote by the number of real multiplications needed by the revised N -point FFT algorithm

- d.** Explain from the algorithm in a. that

$$M_N = 2(M_{N/4} + 3N/8) + M_{N/2} = 3N/4 + M_{N/2} + 2M_{N/4} \quad (4.15)$$

when you employ the trick from the previous exercise.

- e.** Explain why $x_r = M_{2^r}$ is the solution to the difference equation

$$x_{r+2} - x_{r+1} - 2x_r = 3 \times 2^r,$$

and show that the general solution to this

$$x_r = \frac{1}{2}r2^r + C2^r + D(-1)^r.$$

f. Explain why $M_N = O(\frac{1}{2}N \log_2 N)$ both for the revised FFT and the revised DCT (you do not need to write down the initial conditions for the difference equation in order to find the particular solution to it).

g. Explain that, if you had not employed the trick from the previous exercise, the difference equation for obtaining the number of multiplications would instead have been

$$M_N = 2(M_{N/4} + N/2) + M_{N/2} = N + M_{N/2} + 2M_{N/4}.$$

This leads to the equation $x_{r+2} - x_{r+1} - 2x_r = 4 \times 2^r$, from which one can conclude that the number of multiplications is $O(\frac{2}{3}N \log_2 N)$. In conclusion, the trick we employed for our first FFT implementation reduced the number of multiplications needed from $O(2N \log_2 N)$ to $O(N \log_2 N)$, while the trick we employed for the revised FFT implementation reduced the number of multiplications needed from $O(\frac{2}{3}N \log_2 N)$ to $O(\frac{1}{2}N \log_2 N)$.

h. Show that both the revised FFT and DCT algorithms require $M_N = O(\frac{N}{2} \log_2 N)$ multiplications.

This has given us a big reduction from the $O(2N \log_2 N)$ required by the FFT algorithm from Theorem 2.34, and also from the $O(N \log_2 N)$ required from the FFT algorithm after it had been modified with the trick from Exercise 2.8.9.

3. We did not write down corresponding algorithms for the revised IFFT and IDCT algorithms. We will consider this in this exercise.

a. Using equations (4.13)-(4.14), show that

$$\begin{aligned}\Re(y_n) - \Re(y_{N/2-n}) &= 2(ED_{N/4}\mathbf{z})_n \\ \Im(y_n) + \Im(y_{N/2-n}) &= 2(ED_{N/4}\mathbf{w})_{N/4-n}\end{aligned}$$

for $1 \leq n \leq N/4 - 1$. Explain how one can compute \mathbf{z} and \mathbf{w} from this using two IDCT's of length $N/4$.

b. Using equations (4.13)-(4.14), show that

$$\begin{aligned}\Re(y_n) + \Re(y_{N/2-n}) &= \sqrt{2}\Re((F_{N/2}\mathbf{x}^{(e)})_n) \\ \Im(y_n) - \Im(y_{N/2-n}) &= \sqrt{2}\Im((F_{N/2}\mathbf{x}^{(e)})_n),\end{aligned}$$

and explain how one can compute $\mathbf{x}^{(e)}$ from this using an IFFT of length $N/2$.

c. Write Matlab functions which implement the revised IFFT and IDCT algorithms using this procedure, and explain why they can be written so that both require $O(\frac{N}{2} \log_2 N)$ multiplications.

Hint: Use what you found in Exercise 1.c.

4.4 Summary

We started this chapter by extending a previous result which had to do with that the Fourier series of a symmetric function converged quicker. To build on this we first needed to define symmetric extensions of vectors and symmetric vectors, before we classified symmetric extensions in the frequency domain. From this we could find a nice, orthonormal basis for the symmetric extensions, which lead us to the definition of the DCT. We also saw a connection with symmetric filters: These are exactly the filters which preserve symmetric extensions, and we could characterize symmetric filters restricted to symmetric extension as an N -dimensional mapping. We also showed that it is smart to replace the DFT with the DCT when we work with filters which are known to be symmetric. Among other things, this lead to a better way of approximating analog filters, and better interpolation of functions.

We also showed how to obtain an efficient implementation of the DCT, which could reuse the FFT implementation. The DCT has an important role in the MP3 standard. As we have explained, the MP3 standard applies several filters to the sound, in order to split it into bands concentrating on different frequency ranges. In Section 8.5 we will look closer at how these filters can be implemented and constructed. The implementation can use transforms similar to the MDCT, as explained in Exercise 4.1.5. The MDCT is also used in the more advanced version of the MP3 standard (layer III). Here it is applied to the filtered data to obtain a higher spectral resolution of the sound. The MDCT is applied to groups of 576 (in special circumstances 192) samples. The MP3 standard document [17] does not dig into the theory for this, only representing what is needed in order to make an implementation. It is somewhat difficult to read this document, since it is written in quite a different language, familiar mainly to those working with international standards.

The different type of cosine-matrices can all be associated with some extension strategy for the signal. [25] contains a review of these.

The DCT is particularly popular for processing sound data before they are compressed with lossless techniques such as Huffman coding or arithmetic coding. The reason is, as mentioned, that the DCT provides a better approximation from a low-dimensional space than the DFT does, and that it has a very efficient implementation. Libraries exist which goes into lengths to provide efficient implementation of the FFT and the DCT. FFTW, short for *Fastest Fourier Transform in the West* [15], is perhaps the best known of these.

Signal processing literature often does not motivate digital filters in explaining where they come from, and where the input to the filters come from. Using analog filters to motivate this, and to argue for improvements in using the DCT and symmetric extensions, is not that common. Much literature simply says that the property of linear phase is good, without elaborating on this further.

Part II

Wavelets and applications to image processing

Chapter 5

Motivation for wavelets and some simple examples

In Part I our focus was to approximate functions or vectors with trigonometric functions. We saw that the Discrete Fourier transform could be used to obtain a representation of a vector in terms of such functions, and that computations could be done efficiently with the FFT algorithm. This was useful for analyzing, filtering, and compressing sound and other discrete data. The approach with trigonometric functions has some limitations, however. One of these is that, in a representation with trigonometric functions, the frequency content is fixed over time. This is in contrast with most sound data, where the characteristics are completely different in different parts. We have also seen that, even if a sound has a simple representation in terms of trigonometric functions on two different parts, the representation of the entire sound may not be simple. In particular, if the function is nonzero only on a very small interval, a representation of it in terms of trigonometric functions is not so simple.

In this chapter we are going to introduce the basic properties of an alternative to Fourier analysis for representing functions. This alternative is called wavelets. Similar to Fourier analysis, wavelets are also based on the idea of expressing a function in some basis. But in contrast to Fourier analysis, where the basis is fixed, wavelets provide a general framework with many different types of bases. In this chapter we first give a motivation for wavelets, before we continue by introducing some very simple wavelets. The first wavelet we look at can be interpreted as an approximation scheme based on piecewise constant functions. The next wavelet we look at is similar, but with piecewise linear functions used instead. Following these examples we will establish a more general framework, based on experiences from the simple wavelets. In the following chapters we will interpret this framework in terms of filters, and use this connection to construct even more interesting wavelets.

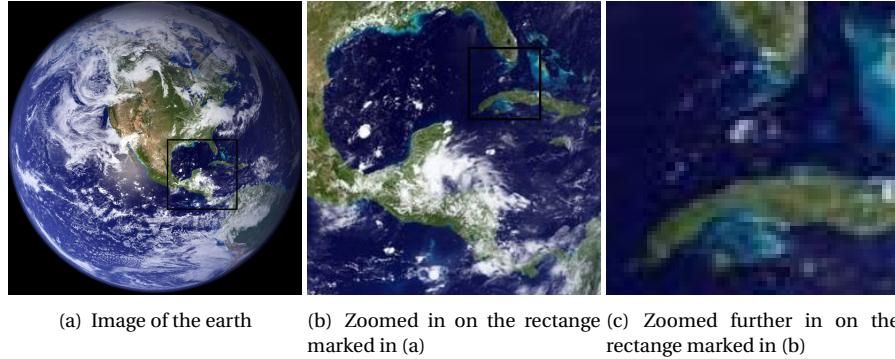


Figure 5.1: A view of Earth from space, together with versions of the image where we have zoomed in.

5.1 Why wavelets?

Figure 5.1(a) shows a view of the entire Earth. The startup image in Google Earth™, a program for viewing satellite images, maps and other geographic information, is very similar to this. In (b) we have zoomed in on the Mexican Gulff, as marked with a rectangle in (a). Similarly, in (c) we have further zoomed in on Cuba and a small portion of Florida, as marked with a rectangle in (b). There is clearly an amazing amount of information available behind a program like Google Earth™, since we there can zoom further in, and obtain enough detail to differentiate between buildings and even trees or cars all over the Earth. So, when the Earth is spinning in the opening screen of Google Earth™, all the Earth's buildings appear to be spinning with it! If this was the case the Earth would not be spinning on the screen, since there would just be so much information to process that a laptop would not be able to display a rotating Earth.

There is a simple reason that the globe can be shown spinning in spite of the huge amounts of information that need to be handled. We are going to see later that a digital image is just a rectangular array of numbers that represent the colour at a dense set of points. As an example, the images in Figure 5.1 are made up of a grid of 1064×1064 points, which gives a total of 1 132 096 points. The colour at a point is represented by three eight-bit integers, which means that the image files contain a total of 3 396 288 bytes each. So regardless of how close to the surface of the Earth our viewpoint is, the resulting image always contains the same number of points. This means that when we are far away from the Earth we can use a very coarse model of the geographic information that is being displayed, but as we zoom in, we need to display more details and therefore need a more accurate model.

Observation 5.1. When discrete information is displayed in an image, there is no need to use a mathematical model that contains more detail than what is visible in the image.

A consequence of Observation 5.1 is that for applications like Google Earth™ we should use a mathematical model that makes it easy to switch between different levels of detail, or different resolutions. Such models are called *multiresolution models*, and wavelets are prominent examples of this kind of models. We will see that multiresolution models also provide us with means of approximating functions, just as Taylor series and Fourier series. Our new approximation scheme differs from these in one important respect, however: When we approximate with Taylor series and Fourier series, the error must be computed at the same data points as well, so that the error contains just as much information as the approximating function, and the function to be approximated. Multiresolution models on the other hand will be defined in such a way that the error and the “approximating function” each contain half of the information from the function we approximate, i.e. their amount of data is reduced. This property makes multiresolution models attractive for the problems at hand, when compared to approaches such as Taylor series and Fourier series.

When we zoom in with Google Earth™, it seems that this is done continuously. The truth is probably that the program only has representations at some given resolutions (since each representation requires memory), and that one interpolates between these to give the impression of a continuous zoom. In the coming chapters we will first look at how we can represent the information at different resolutions, so that only new information at each level is included.

We will now turn to how wavelets are defined more formally, and construct the simplest wavelet we have. Its construction goes in the following steps: First we introduce what we call resolution spaces, and the corresponding scaling function. Then we introduce the detail spaces, and the corresponding mother wavelet. These two functions will give rise to certain bases for these spaces, and we will define the Discrete Wavelet Transform as a change of coordinates between these bases.

5.2 A wavelet based on piecewise constant functions

Our starting point will be the space of piecewise constant functions on an interval $[0, N]$. This will be called a resolution space.

Definition 5.2 (The resolution space V_0). Let N be a natural number. The resolution space V_0 is defined as the space of functions defined on the interval $[0, N]$ that are constant on each subinterval $[n, n + 1]$ for $n = 0, \dots, N - 1$.

Note that this also corresponds to piecewise constant functions which are periodic with period N . We will, just as we did in Fourier analysis, identify a function defined on $[0, N]$ with its (period N) periodic extension. An example of a function in V_0 for $N = 10$ is shown in Figure 5.2. It is easy to check that V_0 is a linear space, and for computations it is useful to know the dimension of the space and have a basis.

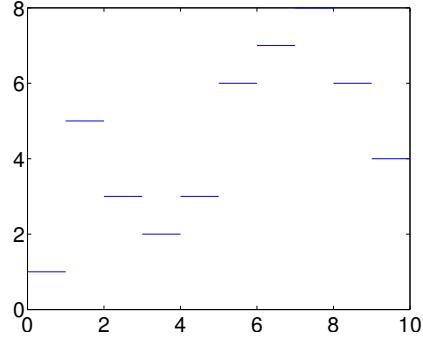


Figure 5.2: A piecewise constant function.

Lemma 5.3. Define the function $\phi(t)$ by

$$\phi(t) = \begin{cases} 1, & \text{if } 0 \leq t < 1; \\ 0, & \text{otherwise;} \end{cases} \quad (5.1)$$

and set $\phi_n(t) = \phi(t - n)$ for any integer n . The space V_0 has dimension N , and the N functions $\{\phi_n\}_{n=0}^{N-1}$ form an orthonormal basis for V_0 with respect to the standard inner product

$$\langle f, g \rangle = \int_0^N f(t)g(t) dt. \quad (5.2)$$

In particular, any $f \in V_0$ can be represented as

$$f(t) = \sum_{n=0}^{N-1} c_n \phi_n(t) \quad (5.3)$$

for suitable coefficients $(c_n)_{n=0}^{N-1}$. The function ϕ_n is referred to as the *characteristic* function of the interval $[n, n + 1]$

Note the small difference between the inner product we define here from the inner product we used for functions previously: Here there is no scaling $1/T$ involved. Also, for wavelets we will only consider real functions, and the inner product will therefore not be defined for complex functions. Two examples of the basis functions defined in Lemma 5.5 are shown in Figure 5.3.

Proof: Two functions ϕ_{n_1} and ϕ_{n_2} with $n_1 \neq n_2$ clearly satisfy $\int \phi_{n_1}(t)\phi_{n_2}(t) dt = 0$ since $\phi_{n_1}(t)\phi_{n_2}(t) = 0$ for all values of t . It is also easy to check that $\|\phi_n\| = 1$ for all n . Finally, any function in V_0 can be written as a linear combination the functions $\phi_0, \phi_1, \dots, \phi_{N-1}$, so the conclusion of the lemma follows. ■

In our discussion of Fourier analysis, the starting point was the function $\sin(2\pi t)$ that has frequency 1. We can think of the space V_0 as being analogous to this function: The function $\sum_{n=0}^{N-1} (-1)^n \phi_n(t)$ is (part of the) square wave that we discussed in

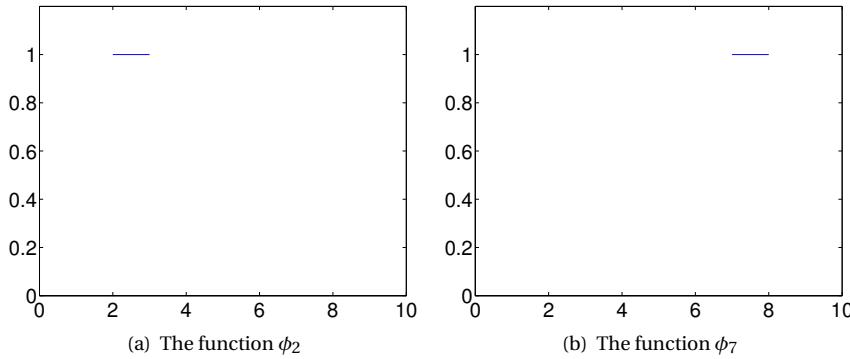


Figure 5.3: Two examples of the basis functions in V_0 .

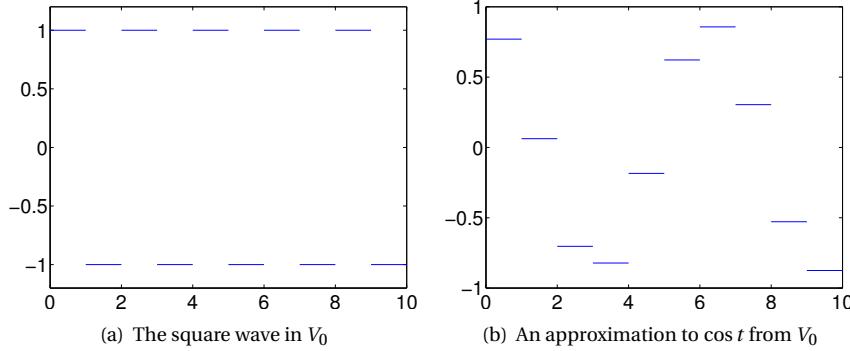


Figure 5.4: Examples of functions from V_0

Chapter 1, and which also oscillates regularly like the sine function, see Figure 5.4 (a). The difference is that we have more flexibility since we have a whole space at our disposal instead of just one function — Figure 5.4 (b) shows another function in V_0 .

In Fourier analysis we obtained a linear space of possible approximations by including sines of frequency 1, 2, 3, ..., up to some maximum. We use a similar approach for constructing wavelets, but we double the frequency each time and label the spaces as V_0, V_1, V_2, \dots

Definition 5.4 (Refined resolution spaces). The space V_m for the interval $[0, N]$ is the space of piecewise linear functions defined on $[0, N]$ that are constant on each subinterval $[n/2^m, (n+1)/2^m]$ for $n = 0, 1, \dots, 2^m N - 1$.

Some examples of functions in the spaces V_1, V_2 and V_3 for the interval $[0, 10]$ are shown in Figure 5.5. As m increases, we can represent smaller details. In particular, the function in (d) is a piecewise constant function that oscillates like $\sin(2\pi 2^2 t)$ on

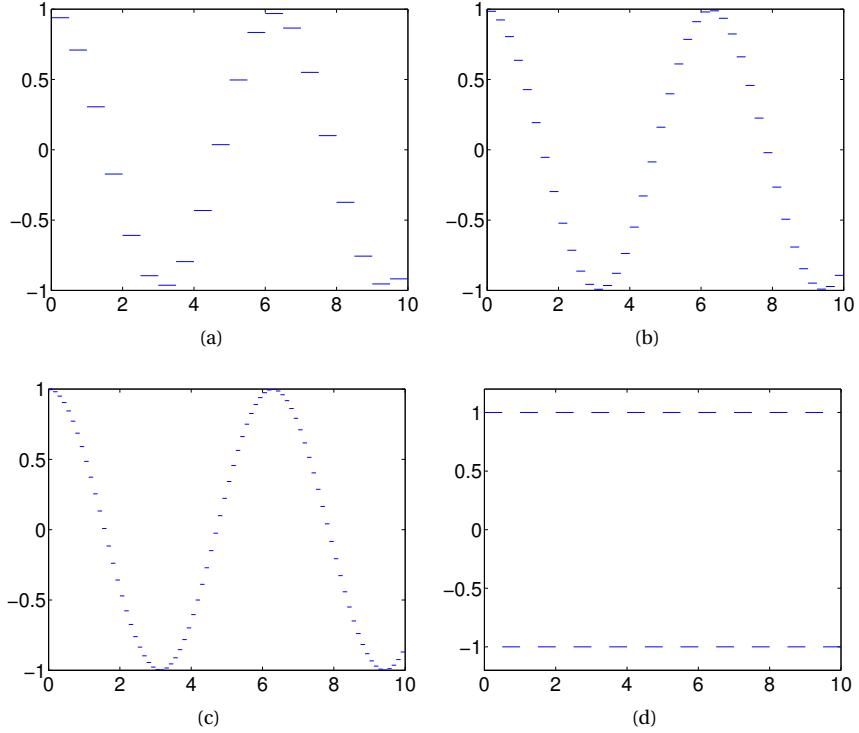


Figure 5.5: Piecewise constant approximations to $\cos t$ on the interval $[0, 10]$ in the spaces V_1 (a), V_2 (b), and V_3 (c). The plot in (d) shows the square wave in V_2 .

the interval $[0, 10]$.

It is easy to find a basis for V_m , we just use the characteristic functions of each subinterval.

Lemma 5.5. Let $[0, N)$ be a given interval with N some positive integer, and let V_m denote the resolution space of piecewise constant functions for some integer $m \geq 0$. Then the dimension of V_m is $2^m N$. Define the functions

$$\phi_{m,n}(t) = 2^{m/2} \phi(2^m t - n), \quad \text{for } n = 0, 1, \dots, 2^m N - 1, \quad (5.4)$$

where ϕ is the characteristic function of the interval $[0, 1]$. The functions $\{\phi_{m,n}\}_{n=0}^{2^m N - 1}$ form an orthonormal basis for V_m , and any function $f \in V_m$ can be represented as

$$f(t) = \sum_{n=0}^{2^m N - 1} c_n \phi_{m,n}(t)$$

for suitable coefficients $(c_n)_{n=0}^{2^m N - 1}$.

Proof: The functions given in (5.18) are exactly the characteristic functions of the subintervals $[n/2^m, (n+1)/2^m]$ which we referred to in Definition 5.4, so the proof is very similar to the proof of Lemma 5.5. The one mysterious thing may be the normalisation factor $2^{-m/2}$. This comes from the fact that

$$\int_0^N \phi(2^m t - n)^2 dt = \int_{n/2^m}^{(n+1)/2^m} \phi(2^m t - n)^2 dt = 2^{-m} \int_0^1 \phi(u)^2 du = 2^{-m}.$$

The normalisation therefore ensures that $\|\phi_{m,n}\| = 1$ for all m . ■

We will denote by ϕ_m the basis $\{\phi_{m,n}\}_{n=0}^{2^m N-1}$. Note that our definition restricts the dimensions of the spaces we study to be on the form $N2^m$. In Chapter 6 we will explain how this restriction can be dropped, but until then the dimensions will be assumed to be on this form. In the theory of wavelets, the function ϕ is also called a *scaling function*. The origin behind this name is that the scaled (and translated) functions $\phi_{m,n}$ of ϕ are used as basis functions for the refined resolution spaces. Later on we will see that other scaling functions ϕ can be chosen, where the scaled versions $\phi_{m,n}$ will be used to define similar resolution spaces, with slightly different properties.

5.2.1 Function approximation property

Each time m is increased by 1, the dimension of V_m doubles, and the subinterval on which the functions in V_m are constant are halved in size. It therefore seems reasonable that, for most functions, we can find good approximations in V_m provided m is big enough.

Theorem 5.6. Let f be a given function that is continuous on the interval $[0, N]$. Given $\epsilon > 0$, there exists an integer $m \geq 0$ and a function $g \in V_m$ such that

$$|f(t) - g(t)| \leq \epsilon$$

for all t in $[0, N]$.

Proof: Since f is (uniformly) continuous on $[0, N]$, we can find an integer m so that $|f(t_1) - f(t_2)| \leq \epsilon$ for any two numbers t_1 and t_2 in $[0, N]$ with $|t_1 - t_2| \leq 2^{-m}$. Define the approximation g by

$$g(t) = \sum_{n=0}^{2^m N-1} f(t_{m,n+1/2}) \phi_{m,n}(t),$$

where $t_{m,n+1/2}$ is the midpoint of the subinterval $[n2^{-m}, (n+1)2^{-m}]$,

$$t_{m,n+1/2} = (n + 1/2)2^{-m}.$$

For t in this subinterval we then obviously have $|f(t) - g(t)| \leq \epsilon$, and since these intervals cover $[0, N]$, the conclusion holds for all $t \in [0, N]$. ■

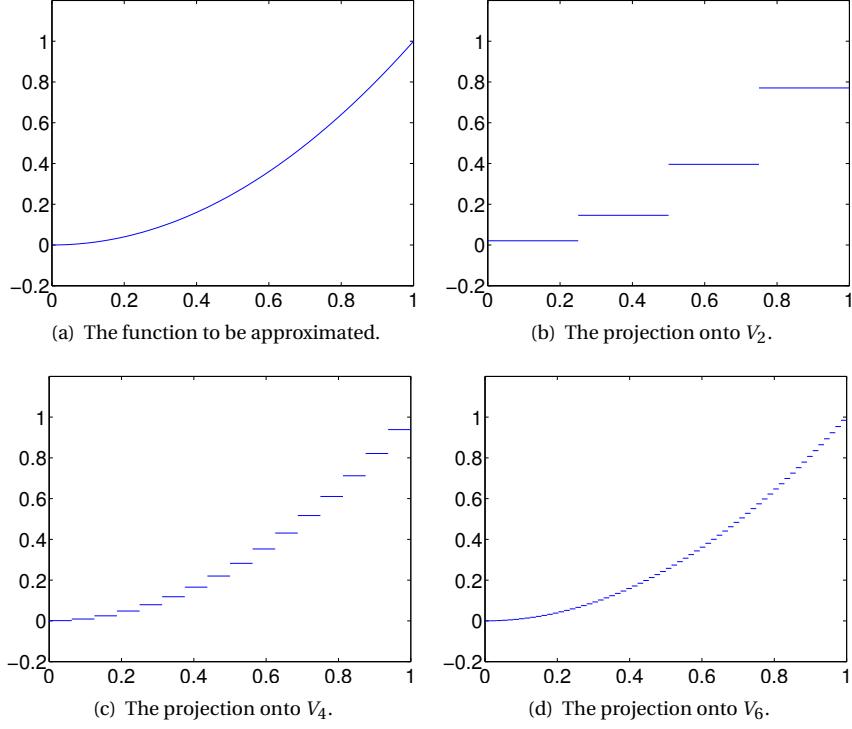


Figure 5.6: Comparison of the function defined by $f(t) = t^2$ on $[0, 1]$ with the projection onto different spaces V_m .

Theorem 5.6 does not tell us how to find the approximation g although the proof makes use of an approximation that interpolates f at the midpoint of each subinterval. Note that if we measure the error in the L^2 -norm, we have

$$\|f - g\|^2 = \int_0^N |f(t) - g(t)|^2 dt \leq N\epsilon^2,$$

so $\|f - g\| \leq \epsilon\sqrt{N}$. We therefore have the following corollary.

Corollary 5.7. Let f be a given continuous function on the interval $[0, N]$ and let $\text{proj}_{V_m}(f)$ denote the best approximation to f from V_m . Then

$$\lim_{m \rightarrow \infty} \|f - \text{proj}_{V_m}(f)\| = 0.$$

Figure 5.6 illustrates how some of the approximations of the function $f(x) = x^2$ from the resolution spaces for the interval $[0, 1]$ improve with increasing m .

5.2.2 Detail spaces and wavelets

So far we have described a family of function spaces that allow us to determine arbitrarily good approximations to a continuous function. The next step is to introduce the so-called detail spaces and the wavelet functions. For this we focus on the two spaces V_0 and V_1 .

We start by observing that since

$$[n, n+1) = [2n/2, (2n+1)/2) \cup [(2n+1)/2, (2n+2)/2)$$

we have

$$\phi_{0,n} = \frac{1}{\sqrt{2}}\phi_{1,2n} + \frac{1}{\sqrt{2}}\phi_{1,2n+1}. \quad (5.5)$$

This provides a formal proof of the intuitive observation that $V_0 \subset V_1$. For if $g \in V_0$, then we can write

$$g(t) = \sum_{n=0}^{N-1} c_n \phi_{0,n}(t) = \sum_{n=0}^{N-1} c_n (\phi_{1,2n} + \phi_{1,2n+1})/\sqrt{2}.$$

The right-hand side clearly lies in V_1 . A similar argument shows that $V_k \subset V_{k+1}$ for any integer $k \geq 0$.

Lemma 5.8. The spaces $V_0, V_1, \dots, V_m, \dots$ are nested,

$$V_0 \subset V_1 \subset V_2 \subset \dots \subset V_m \dots$$

The next step is to investigate what happens if we start with a function g_1 in V_1 and project this to an approximation g_0 in V_0 .

Lemma 5.9. Let proj_{V_0} denote the orthogonal projection onto the subspace V_0 . Then the projection of a basis function $\phi_{1,n}$ is given by

$$\text{proj}_{V_0}(\phi_{1,n}) = \begin{cases} \phi_{0,n/2}/\sqrt{2}, & \text{if } n \text{ is even;} \\ \phi_{0,(n-1)/2}/\sqrt{2}, & \text{if } n \text{ is odd.} \end{cases} \quad (5.6)$$

If $g_1 \in V_1$ is given by

$$g_1 = \sum_{n=0}^{2N-1} c_{1,n} \phi_{1,n}, \quad (5.7)$$

then

$$\text{proj}_{V_0}(g_1) = g_0 = \sum_{n=0}^{N-1} c_{0,n} \phi_{0,n}$$

where $c_{0,n}$ is given by

$$c_{0,n} = \frac{c_{1,2n} + c_{1,2n+1}}{\sqrt{2}}. \quad (5.8)$$

Proof: We first observe that $\phi_{1,n}(t) \neq 0$ if and only if $n/2 \leq t < (n+1)/2$. Suppose that n is even. Then the intersection

$$\left[\frac{n}{2}, \frac{n+1}{2} \right] \cap [n_1, n_1 + 1] \quad (5.9)$$

is nonempty only if $n_1 = \frac{n}{2}$. Using the orthogonal decomposition formula we get

$$\begin{aligned} \text{proj}_{V_0}(\phi_{1,n}) &= \sum_{k=0}^{N-1} \langle \phi_{1,n}, \phi_{0,k} \rangle \phi_{0,k} = \langle \phi_{1,n}, \phi_{0,n_1} \rangle \phi_{0,n_1} \\ &= \int_{n/2}^{(n+1)/2} \sqrt{2} dt \phi_{0,n/2} = \frac{1}{\sqrt{2}} \phi_{0,n/2}. \end{aligned}$$

When n is odd, the intersection (5.9) is nonempty only if $n_1 = (n-1)/2$, which gives the second formula in (5.6) in the same way.

We project the function g_1 in V_1 using the formulas in (5.6). We split the sum in (5.7) into even and odd values of n ,

$$g_1 = \sum_{n=0}^{2N-1} c_{1,n} \phi_{1,n} = \sum_{n=0}^{N-1} c_{1,2n} \phi_{1,2n} + \sum_{n=0}^{N-1} c_{1,2n+1} \phi_{1,2n+1}. \quad (5.10)$$

We can now apply the two formulas in (5.6),

$$\begin{aligned} \text{proj}_{V_0}(g_1) &= \text{proj}_{V_0} \left(\sum_{n=0}^{N-1} c_{1,2n} \phi_{1,2n} + \sum_{n=0}^{N-1} c_{1,2n+1} \phi_{1,2n+1} \right) \\ &= \sum_{n=0}^{N-1} c_{1,2n} \text{proj}_{V_0}(\phi_{1,2n}) + \sum_{n=0}^{N-1} c_{1,2n+1} \text{proj}_{V_0}(\phi_{1,2n+1}) \\ &= \sum_{n=0}^{N-1} c_{1,2n} \phi_{0,n}/\sqrt{2} + \sum_{n=0}^{N-1} c_{1,2n+1} \phi_{0,n}/\sqrt{2} \\ &= \sum_{n=0}^{N-1} \frac{c_{1,2n} + c_{1,2n+1}}{\sqrt{2}} \phi_{0,n} \end{aligned}$$

which proves (5.8) ■

When $g_1 \in V_1$ is projected onto V_0 , the result $g_0 = \text{proj}_{V_0} g_1$ is in general different from g_0 . We can write $g_1 = g_0 + e_0$, where $e_0 = g_1 - g_0$ represents the error we have committed in making this projection. e_0 lies in the orthogonal complement of V_0 in V_1 (in particular, $e_0 \in V_1$).

Definition 5.10. We will denote by W_0 the orthogonal complement of V_0 in V_1 . We also call W_0 a *detail space*

The name detail space is used since $e_0 \in W_0$ can be considered as the detail which is left out when considering g_0 instead of g_1 (due to the expression $g_1 = g_0 + e_0$). Since V_0 and W_0 are mutually orthogonal spaces they are also linearly independent spaces, so that we can use the following definition to rewrite the space V_0 :

Definition 5.11 (Direct sum of vector spaces). Assume that $U, V \subset W$ are vector spaces, and that U and V are mutually linearly independent. By $U \oplus V$ we mean the vector space consisting of all vectors of the form $\mathbf{u} + \mathbf{v}$, where $\mathbf{u} \in U$, $\mathbf{v} \in V$. We will also call $U \oplus V$ the *direct sum* of U and V .

This definition also makes sense if we have several vector spaces, since the direct sum clearly obeys the associate law $U \oplus (V \oplus W) = (U \oplus V) \oplus W$, i.e. we can define $U \oplus V \oplus W = U \oplus (V \oplus W)$. We will have use for this use of direct sum of several vector space in the next section.

We can now write $V_1 = V_0 \oplus W_0$, i.e. the resolution space V_1 is the direct sum of the lower order resolution space V_0 , and the detail space W_0 . The expression $g_1 = g_0 + e_0$ is thus a decomposition into a low-resolution approximation, and the details which are left out in this approximation. In the context of our Google Earth™ example, in Figure 5.1 you should interpret g_0 as the image in (a), g_1 as the image in (b), and e_0 as the additional details which are needed to reproduce (b) from (a). While Lemma 5.12 explained how we can compute the low level approximation g_0 from g_1 , the next result states how we can compute the detail/error e_0 from g_1 .

Lemma 5.12. With W_0 the orthogonal complement of V_0 in V_1 , set

$$\hat{\psi}_{0,n} = \frac{\phi_{1,2n} - \phi_{1,2n+1}}{2}$$

for $n = 0, 1, \dots, N-1$. Then $\hat{\psi}_{0,n} \in W_0$ and

$$\text{proj}_{W_0}(\phi_{1,n}) = \begin{cases} \hat{\psi}_{0,n/2}, & \text{if } n \text{ is even;} \\ -\hat{\psi}_{0,(n-1)/2}, & \text{if } n \text{ is odd.} \end{cases} \quad (5.11)$$

If $g_1 \in V_1$ is given by $g_1 = \sum_{n=0}^{2N-1} c_{1,n} \phi_{1,n}$, then

$$\text{proj}_{W_0}(g_1) = e_0 = \sum_{n=0}^{N-1} (c_{1,2n} - c_{1,2n+1}) \hat{\psi}_{0,n}.$$

Proof: We start by determining the error when $\phi_{1,n}$, for n even, is projected onto V_0 . The error is then

$$\begin{aligned} \text{proj}_{W_0}(\phi_{1,n}) &= \phi_{1,n} - \frac{\phi_{0,n/2}}{\sqrt{2}} = \phi_{1,n} - \frac{1}{\sqrt{2}} \left(\frac{1}{\sqrt{2}} \phi_{1,n} + \frac{1}{\sqrt{2}} \phi_{1,n+1} \right) \\ &= \frac{1}{2} \phi_{1,n} - \frac{1}{2} \phi_{1,n+1} = \hat{\psi}_{0,n/2}. \end{aligned}$$

Here we used the relation (5.6) in the second equation. When n is odd we have

$$\begin{aligned} \text{proj}_{W_0}(\phi_{1,n}) &= \phi_{1,n} - \frac{\phi_{0,(n-1)/2}}{\sqrt{2}} = \phi_{1,n} - \frac{1}{\sqrt{2}} \left(\frac{1}{\sqrt{2}} \phi_{1,n-1} + \frac{1}{\sqrt{2}} \phi_{1,n} \right) \\ &= \frac{1}{2} \phi_{1,n} - \frac{1}{2} \phi_{1,n-1} = -\hat{\psi}_{0,(n-1)/2}. \end{aligned}$$

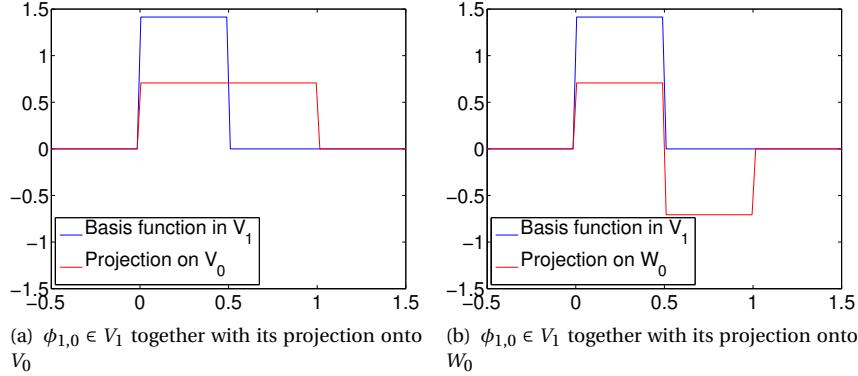


Figure 5.7: The projection of a basis function in V_1 onto V_0 and W_0 .

For a general function g_1 we first split the sum into even and odd terms as in (5.10) and then project each part onto W_0 ,

$$\begin{aligned}\text{proj}_{W_0}(g_1) &= \text{proj}_{W_0} \left(\sum_{n=0}^{N-1} c_{1,2n} \phi_{1,2n} + \sum_{n=0}^{N-1} c_{1,2n+1} \phi_{1,2n+1} \right) \\ &= \sum_{n=0}^{N-1} c_{1,2n} \text{proj}_{W_0}(\phi_{1,2n}) + \sum_{n=0}^{N-1} c_{1,2n+1} \text{proj}_{W_0}(\phi_{1,2n+1}) \\ &= \sum_{n=0}^{N-1} c_{1,2n} \hat{\psi}_{0,n} - \sum_{n=0}^{N-1} c_{1,2n+1} \hat{\psi}_{0,n} = \sum_{n=0}^{N-1} (c_{1,2n} - c_{1,2n+1}) \hat{\psi}_{0,n}.\end{aligned}$$

In Figure 5.7 we have used lemmas 5.9 and 5.12 to plot the projections of $\phi_{1,0} \in V_1$ onto V_0 and W_0 . It is an interesting exercise to see from the plots why exactly these functions should be least-squares approximations of $\phi_{1,n}$. It is also an interesting exercise to prove the following from lemmas 5.9 and 5.12:

Proposition 5.13. Let $f(t) \in V_1$, and let $f_{n,1}$ be the value f attains on $[n, n + 1/2]$, and $f_{n,2}$ the value f attains on $[n + 1/2, n + 1]$. Then $\text{proj}_{V_0}(f)$ is the function in V_0 which equals $(f_{n,1} + f_{n,2})/2$ on the interval $[n, n + 1]$. Moreover, $\text{proj}_{W_0}(f)$ is the function in W_0 which is $(f_{n,1} - f_{n,2})/2$ on $[n, n + 1/2]$, and $-(f_{n,1} - f_{n,2})/2$ on $[n + 1/2, n + 1]$.

In other words, the projection on V_0 is constructed by averaging on two subintervals, while the projection on W_0 is constructed by taking the difference from the mean. This sounds like a reasonable candidate for the least-squares approximations. In the exercise we generalize these observations.

Consider the functions $\hat{\psi}_{0,n} = (\phi_{1,2n} - \phi_{1,2n+1})/2$ from Lemma 5.12. They are clearly orthogonal since their nonzero parts do not overlap. We also note that $\|\hat{\psi}_{0,n}\| =$

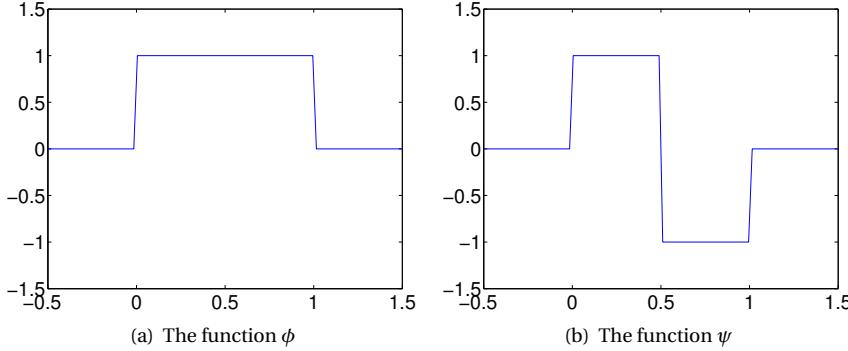


Figure 5.8: The functions we used to analyse the space of piecewise constant functions

$\sqrt{2}/2$, since it has absolute value $\sqrt{2}/2$ on two intervals of length $1/2$. The functions defined by $\psi_{0,n}(t) = \sqrt{2} \hat{\psi}_{0,n}(t)$ will therefore form an orthonormal set.

Lemma 5.14. Define the function ψ by

$$\psi(t) = (\phi_{1,0}(t) - \phi_{1,1}(t)) / \sqrt{2} = \phi(2t) - \phi(2t-1) \quad (5.12)$$

and set

$$\psi_{0,n}(t) = \psi(t-n) = (\phi_{1,2n}(t) - \phi_{1,2n+1}(t)) / \sqrt{2} \quad \text{for } n = 0, 1, \dots, N-1. \quad (5.13)$$

Then the set $\{\psi_{0,n}\}_{n=0}^{N-1}$ is an orthonormal basis for W_0 , the orthogonal complement of V_0 in V_1 .

Later we will encounter other functions, which also will be denoted by ψ , and have similar properties as stated in Lemma 5.14. In the theory of wavelets, such ψ are called *mother wavelets*. In Figure 5.8 we have plotted the functions ϕ and ψ . There is one important property of ψ , which we will return to:

Observation 5.15. We have that $\int_0^N \psi(t) dt = 0$.

This can be seen directly from the plot in Figure 5.8, since the parts of the graph above and below the x -axis cancel. In general we say that ψ has k vanishing moments if the integrals $\int t^l \psi(t) dt = 0$ for all $0 \leq l \leq k-1$. Due to Observation 5.15, ψ has one vanishing moment. In Chapter 7 we will show that mother wavelets with many vanishing moments are very desirable when it comes to approximation of functions.

We now have all the tools needed to define the Discrete Wavelet Transform. If $\{\mathcal{B}_i\}_{i=1}^n$ are mutually independent bases, we will in the following write $(\mathcal{B}_1, \mathcal{B}_2, \dots, \mathcal{B}_n)$ for the basis where the basis vectors from \mathcal{B}_i are included before \mathcal{B}_j when $i < j$.

Theorem 5.16 (Discrete Wavelet Transform). The space V_1 can be decomposed as the orthogonal sum $V_1 = V_0 \oplus W_0$ where W_0 is the orthogonal complement of V_0 in V_1 , and V_1 therefore has the two bases

$$\boldsymbol{\phi}_1 = (\phi_{1,n})_{n=0}^{2N-1} \quad \text{and} \quad (\boldsymbol{\phi}_0, \boldsymbol{\psi}_0) = ((\phi_{0,n})_{n=0}^{N-1}, (\psi_{0,n})_{n=0}^{N-1}).$$

The Discrete Wavelet Transform (DWT) is the change of coordinates from $\boldsymbol{\phi}_1$ to $(\boldsymbol{\phi}_0, \boldsymbol{\psi}_0)$. If $g_1 = g_0 + e_0$ with

$$g_1 = \sum_{n=0}^{2N-1} c_{1,n} \phi_{1,n} \in V_1$$

$$g_0 = \sum_{n=0}^{N-1} c_{0,n} \phi_{0,n} \in V_0 \quad e_0 = \sum_{n=0}^{N-1} w_{0,n} \psi_{0,n} \in W_0,$$

then the DWT is given by

$$\begin{pmatrix} c_{0,n} \\ w_{0,n} \end{pmatrix} = \begin{pmatrix} 1/\sqrt{2} & 1/\sqrt{2} \\ 1/\sqrt{2} & -1/\sqrt{2} \end{pmatrix} \begin{pmatrix} c_{1,2n} \\ c_{1,2n+1} \end{pmatrix} \quad (5.14)$$

Conversely, the Inverse Discrete Wavelet Transform (IDWT) is the change of coordinates from $(\boldsymbol{\phi}_0, \boldsymbol{\psi}_0)$ to $\boldsymbol{\phi}_1$, and is given by

$$\begin{pmatrix} c_{1,2n} \\ c_{1,2n+1} \end{pmatrix} = \begin{pmatrix} 1/\sqrt{2} & 1/\sqrt{2} \\ 1/\sqrt{2} & -1/\sqrt{2} \end{pmatrix}^{-1} \begin{pmatrix} c_{0,n} \\ w_{0,n} \end{pmatrix} \quad (5.15)$$

Proof: We have that

$$\begin{aligned} \phi_{1,2n} &= \text{proj}_{V_0}(\phi_{1,2n}) + \text{proj}_{W_0}(\phi_{1,2n}) = \phi_{0,n}/\sqrt{2} + \hat{\psi}_{0,n} \\ &= \phi_{0,n}/\sqrt{2} + \psi_{0,n}/\sqrt{2} \\ \phi_{1,2n+1} &= \text{proj}_{V_0}(\phi_{1,2n+1}) + \text{proj}_{W_0}(\phi_{1,2n+1}) = \phi_{0,n}/\sqrt{2} - \hat{\psi}_{0,n} \\ &= \phi_{0,n}/\sqrt{2} - \psi_{0,n}/\sqrt{2}, \end{aligned}$$

where we have used Equations (5.6) and (5.11). From this it follows that $(\phi_{1,2n}, \phi_{1,2n+1})$ and $(\phi_{0,n}, \psi_{0,n})$ span the same space, and that the corresponding change of coordinate matrix is given by $\begin{pmatrix} 1/\sqrt{2} & 1/\sqrt{2} \\ 1/\sqrt{2} & -1/\sqrt{2} \end{pmatrix}$. This proves Equation (5.14). Equation (5.15) follows immediately since this matrix equals its inverse. ■

What you should have learnt in this section

Definition of resolution spaces and scaling function for the wavelet based on piecewise constant functions. The nesting of resultion spaces, and how one can project

from one resolution space onto another. Definition of details spaces and mother wavelet, and how one can project onto the detail spaces. The definition of the Discrete Wavelet Transform. How to compute the coordinate vector of a piecewise constant function in the basis ϕ_m .

Exercises for Section 5.2

1. Show that the coordinate vector for $f \in V_0$ in the basis $\{\phi_{0,0}, \phi_{0,1}, \dots, \phi_{0,N-1}\}$ is $(f(0), f(1), \dots, f(N-1))$.
2. Prove Proposition 5.13.
3. In this exercise we will consider the two projections from V_1 onto V_0 and W_0 .
 - a. Consider the projection proj_{V_0} of V_1 onto V_0 . Use lemma 5.9 to write down the matrix for proj_{V_0} relative to the bases ϕ_1 and ϕ_0 .
 - b. Consider the projection proj_{W_0} of V_1 onto V_0 . Use lemma 5.12 to write down the matrix for proj_{W_0} relative to the bases ϕ_1 and ψ_0 .
4. Consider again the projection proj_{V_0} of V_1 onto V_0 .
 - a. Explain why $\text{proj}_{V_0}(\phi) = \phi$ and $\text{proj}_{V_0}(\psi) = 0$.
 - b. Show that the matrix of proj_{V_0} relative to (ϕ_0, ψ_0) is given by the diagonal matrix where the first half of the entries on the diagonal are 1, the second half 0.
 - c. Show in a similar way that the projection of V_1 onto W_0 has a matrix relative to (ϕ_0, ψ_0) given by the diagonal matrix where the first half of the entries on the diagonal are 0, the second half 1.

5. Show that

$$\text{proj}_{V_0}(f) = \sum_{n=0}^{N-1} \left(\int_n^{n+1} f(t) dt \right) \phi_{0,n}(t) \quad (5.16)$$

for any f . Show also that the first part of Proposition 5.13 follows from this.

6. Show that

$$\left\| \sum_n \left(\int_n^{n+1} f(t) dt \right) \phi_{0,n}(t) - f \right\|^2 = \langle f, f \rangle - \sum_n \left(\int_n^{n+1} f(t) dt \right)^2.$$

This, together with the previous exercise, gives us an expression for the least-squares error for f from V_0 (at least after taking square roots).

7. Show that

$$\text{proj}_{W_0}(f) = \sum_{n=0}^{N-1} \left(\int_n^{n+1/2} f(t) dt - \int_{n+1/2}^{n+1} f(t) dt \right) \psi_{0,n}(t) \quad (5.17)$$

for any f . Show also that the second part of Proposition 5.13 follows from this.

5.3 m -level Discrete Wavelet Transform

In Section 5.2 we introduced the important decomposition $V_1 = V_0 \oplus W_0$ which lets us rewrite a function in V_1 as an approximation in V_0 and the corresponding error in W_0 , orthogonal to the approximation. The resolution spaces V_m were in fact defined for all integers $m \geq 0$. It turns out that all these higher order resolution spaces can be decomposed in the same way as V_1 .

Definition 5.17 (Higher order detail spaces). The orthogonal complement of V_{m-1} in V_m is denoted W_{m-1} . All the spaces $\{W_k\}_k$ are also called detail spaces.

The first question we will try to answer is how we can, for $f \in V_m$, extract the corresponding detail in W_{m-1} . We first need to define $\psi_{m,n}$ in terms of ψ , similarly to how we defined $\phi_{m,n}$ in terms of ϕ ,

$$\psi_{m,n}(t) = 2^{m/2} \psi(2^m t - n), \quad \text{for } n = 0, 1, \dots, 2^m N - 1. \quad (5.18)$$

As in Lemma 5.14, it is straightforward to prove that $\psi_m = \{\psi_{m,n}\}_{n=0}^{2^m N - 1}$ is an orthonormal basis for W_m . Moreover, we have the following result, which is completely analogous to Theorem 5.16.

Theorem 5.18. The space V_m can be decomposed as the orthogonal sum $V_m = V_{m-1} \oplus W_{m-1}$ where W_{m-1} is the orthogonal complement of V_{m-1} in V_m , and V_m has the two bases

$$\phi_m = (\phi_{m,n})_{n=0}^{2^m N - 1}, \text{ and } (\phi_{m-1}, \psi_{m-1}) = ((\phi_{m-1,n})_{n=0}^{2^{m-1} N - 1}, (\psi_{m-1,n})_{n=0}^{2^{m-1} N - 1}).$$

If $g_m = g_{m-1} + e_{m-1}$ with

$$g_m = \sum_{n=0}^{2^m N - 1} c_{m,n} \phi_{m,n} \in V_m,$$

$$g_{m-1} = \sum_{n=0}^{2^{m-1} N - 1} c_{m-1,n} \phi_{m-1,n} \in V_{m-1} \quad e_{m-1} = \sum_{n=0}^{2^{m-1} N - 1} w_{m-1,n} \psi_{m-1,n} \in W_{m-1},$$

then the change of coordinates from ϕ_m to (ϕ_{m-1}, ψ_{m-1}) is given by

$$\begin{pmatrix} c_{m-1,n} \\ w_{m-1,n} \end{pmatrix} = \begin{pmatrix} 1/\sqrt{2} & 1/\sqrt{2} \\ 1/\sqrt{2} & -1/\sqrt{2} \end{pmatrix} \begin{pmatrix} c_{m,2n} \\ c_{m,2n+1} \end{pmatrix} \quad (5.19)$$

Conversely, the change of coordinates from (ϕ_{m-1}, ψ_{m-1}) to ϕ_m is given by

$$\begin{pmatrix} c_{m,2n} \\ c_{m,2n+1} \end{pmatrix} = \begin{pmatrix} 1/\sqrt{2} & 1/\sqrt{2} \\ 1/\sqrt{2} & -1/\sqrt{2} \end{pmatrix} \begin{pmatrix} c_{m-1,n} \\ w_{m-1,n} \end{pmatrix} \quad (5.20)$$

We will omit the proof of Theorem 5.18, and only remark that it can be proved by making the substitution $t \rightarrow 2^m u$ in Lemma 5.9 and Lemma 5.12, and then following the proof of Theorem 5.16.

Let us return to our interpretation of the Discrete Wavelet Transform as writing a function $g_1 \in V_1$ as a sum of a function $g_0 \in V_0$ at low resolution, and a detail function $e_0 \in W_0$. Theorem 5.18 states similarly how we can write $g_m \in V_m$ as a sum of a function $g_{m-1} \in V_{m-1}$ at lower resolution, and a detail function $e_{m-1} \in W_{m-1}$. The same decomposition can of course be applied to g_{m-1} in V_{m-1} , then to the resulting approximation g_{m-2} in V_{m-2} , and so on,

$$\begin{aligned} V_m &= V_{m-1} \oplus W_{m-1} \\ &= V_{m-2} \oplus W_{m-2} \oplus W_{m-1} \\ &\vdots \\ &= V_0 \oplus W_0 \oplus W_1 \oplus \cdots \oplus W_{m-2} \oplus W_{m-1}. \end{aligned} \tag{5.21}$$

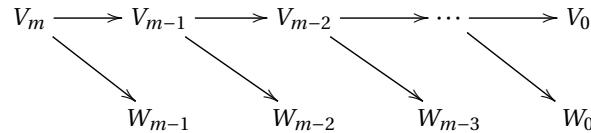
This corresponding change of coordinates corresponds to replacing as many ϕ -functions as we can with ψ -functions, i.e. replacing the original function with a sum of as much detail at different resolutions as possible. It is also called a Discrete Wavelet Transform.

Definition 5.19 (m -level Discrete Wavelet Transform). The change of coordinates from ϕ_m to the base

$$(\phi_0, \psi_0, \psi_1, \dots, \psi_{m-2} \oplus \psi_{m-1})$$

for $V_0 \oplus W_0 \oplus W_1 \oplus \cdots \oplus W_{m-2} \oplus W_{m-1}$ is called a (m -level) *Discrete Wavelet Transform*, or a DWT. After this change of coordinates, the resulting coordinates are called *wavelet coefficients*. The change of coordinates the opposite way is called an (m -level) *Inverse Discrete Wavelet Transform*, or IDWT.

Clearly, this generalizes the Discrete Wavelet Transform defined in Section 5.2. At each level in a DWT, V_k is split into one low-resolution part from V_{k-1} , and one detail part from W_{k-1} . We can visualize this with the following figure, where the arrows represent changes of coordinates:



The part from W_{k-1} is not subject to further transformation. This is seen in the figure since W_{m-1} is a leaf node, i.e. there are no arrows going out from W_{m-1} . In a similar illustration for the IDWT, the arrows would go the opposite way. The Discrete Wavelet Transform is the analogue in a wavelet setting to the Discrete Fourier transform. When applying the DFT to a vector of length N , one starts by viewing

this vector as coordinates relative to the standard basis. When applying the DWT to a vector of length N , one instead views the vector as coordinates relative to the basis ϕ_m . This makes sense in light of Exercise 5.2.1.

The DWT is what is used in practice when transforming a signal using wavelets, and it is straightforward to implement: One simply needs to iterate Equation (5.19) for $m, m-1, \dots, 1$, also at each step, the coordinates in ϕ_{m-1} should be placed before the ones in ψ_{m-1} . At each step, only the first coordinates are further transformed. The following function, called `DWTHaarImpl`, follows this procedure. It takes as input the number of levels m , as well as the input vector x , runs the m -level DWT on x , and returns the result:

```
function xnew=DWTHaarImpl(x,m)
    xnew=x;
    for mres=m:(-1):1
        len=length(xnew)/2^(m-mres);
        c=(xnew(1:2:(len-1))+xnew(2:2:len))/sqrt(2);
        w=(xnew(1:2:(len-1))-xnew(2:2:len))/sqrt(2);
        xnew(1:len)=[c w];
    end
```

For simplicity this code assumes that the length of the vector can be written as an appropriate power of 2. Note that this implementation is not recursive, contrary to the FFT. The for-loop here runs through the different resolutions. Inside the loop we perform the change of coordinates from ϕ_k to (ϕ_{k-1}, ψ_{k-1}) by applying Equation (5.19). This works on the first coordinates, since the coordinates from ϕ_k are stored first in

$$V_k \oplus W_k \oplus W_{k+1} \oplus \cdots \oplus W_{m-2} \oplus W_{m-1}.$$

Finally, the c -coordinates are stored before the w -coordinates. In this implementation, note that the first levels require the most multiplications, since the latter levels leave an increasing part of the coordinates unchanged. Note also that the change of coordinates matrix is a very sparse matrix: At each level a coordinate can be computed from only two of the other coordinates, so that this matrix has only two nonzero elements in each row/column. The algorithm clearly shows that there is no need to perform a full matrix multiplication to perform the change of coordinates.

The corresponding function for the IDWT, called `IDWTHaarImpl`, goes as follows:

```
function x=IDWTHaarImpl(xnew,m)
    x=xnew;
    for mres=1:m
        len=length(x)/2^(m-mres);
        ev=(x(1:(len/2))+x((len/2+1):len))/sqrt(2);
        od=(x(1:(len/2))-x((len/2+1):len))/sqrt(2);
        x(1:2:(len-1))=ev;
        x(2:2:len)=od;
    end
```

Here the steps are simply performed in the reverse order, and by iterating Equation (5.20).

You may be puzzled by the names `DWTHaarImpl` and `IDWTHaarImpl`. In the next sections we will consider other cases, where the underlying function ϕ may be a different function, not necessarily piecewise constant. It will turn out that much of the analysis we have done makes sense for other functions ϕ as well, giving rise to other structures which we also will refer to as wavelets. The wavelet resulting from piecewise constant functions is thus simply one example out of many, and it is commonly referred to as the *Haar wavelet*.

Example 5.20. It is important to keep in mind that the DWT always is a change of coordinates. To visualize this, let us take the simple vector \mathbf{x} of length $2^{10} = 1024$ defined by

$$x_n = \begin{cases} 1 & \text{for } n < 512 \\ 0 & \text{for } n \geq 512, \end{cases}$$

and let us compute the 10-level DWT of this vector by first visualizing the function with these coordinates. Since $m = 10$ here, we should view \mathbf{x} as coordinates in the basis ϕ_{10} of a function $f(t) \in V_{10}$. This is $f(t) = \sum_{n=0}^{511} \phi_{10,n}$, and since $\phi_{10,n}$ is supported on $[2^{-10}n, 2^{-10}(n+1))$, the support of f has width $512 \times 2^{-10} = 1/2$ (512 translates, each with width 2^{-10}). Moreover, since $\phi_{10,n}$ is $2^{10/2} = 2^5 = 32$ on $[2^{-10}n, 2^{-10}(n+1))$ and 0 elsewhere, it is clear that

$$f(t) = \begin{cases} 32 & \text{for } 0 \leq t < 1/2 \\ 0 & \text{for } 0 < t \geq 1/2. \end{cases}$$

This is by definition a function in V_1 : f must in fact be a multiplum of $\phi_{1,0}$, since this also is supported on $[0, 1/2)$. We can thus write $f(t) = c\phi_{1,0}(t)$ for some c . We can find c by setting $t = 0$. This gives that $32 = 2^{1/2}c$ (since $f(0) = 32$, $\phi_{1,0}(0) = 2^{1/2}$), so that $c = 32/\sqrt{2}$. This means that $f(t) = \frac{32}{\sqrt{2}}\phi_{1,0}(t)$.

When we run a 10-level DWT we make a change of coordinates from ϕ_{10} to $(\phi_0, \psi_0, \dots, \psi_9)$. We have seen that f actually is in V_1 , so that what remains is to perform an additional change of coordinates from ϕ_1 to (ϕ_0, ψ_0) . This can be achieved by substituting $\phi_{1,0} = \frac{1}{\sqrt{2}}(\phi_{0,0} + \psi_{0,0})$, so that we get

$$f(t) = \frac{32}{\sqrt{2}}\phi_{1,0}(t) = \frac{32}{\sqrt{2}}\frac{1}{\sqrt{2}}(\phi_{0,0} + \psi_{0,0}) = 16\phi_{0,0} + 16\psi_{0,0}.$$

From this we see that the coordinate vector of f in $(\phi_0, \psi_0, \dots, \psi_9)$ is $(16, 16, 0, 0, \dots, 0)$, which thus also is the 10-level DWT of \mathbf{x} . Note that here V_0 and W_0 are both 1-dimensional, since V_{10} was assumed to be of dimension 2^{10} (in particular, $N = 1$).

It is straightforward to verify what we found using the algorithm above:

```
DWTHaarImpl([ones(512,1);zeros(512,1)],10);
```



Example 5.21. When you run a DWT you may be led to believe that coefficients from the lower order resolution spaces may correspond to lower frequencies. This

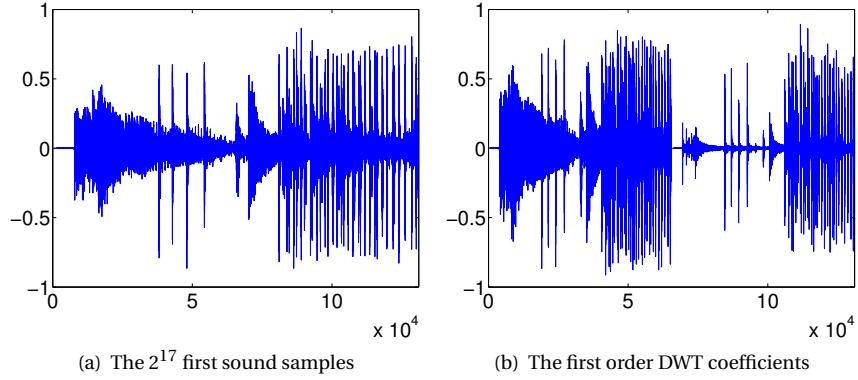


Figure 5.9: The sound samples and the DWT coefficients of the sound `castanets.wav`.

sounds reasonable, since the functions $\phi(2^m t - n) \in V_m$ change more quickly than $\phi(t - n) \in V_0$. However, the functions $\phi_{m,n}$ do not correspond to pure tones in the setting of wavelets. But we can still listen to sound from the different resolution spaces. In Exercise 10 you will be asked to implement a function which runs an m -level DWT on the first samples of the sound file `castanets.wav`, extracts the coefficients from the lower order resolution spaces, transforms the values back to sound samples with the IDWT, and plays the result. When you listen to the result the sound is clearly recognizable for lower values of m , but is degraded for higher values of m . The explanation is that too much of the detail is omitted when you use a higher m . To be more precise, when listening to the sound by throwing away wverything from the detail spaces W_0, W_1, \dots, W_{m-1} , we are left with a 2^{-m} share of the data. Note that this procedure is mathematically not the same as setting some DFT coefficients to zero, since the DWT does not operate on pure tones.

It is of interest to plot the samples of our test audio file `castanets.wav`, and compare it with the first order DWT coefficients of the same samples. This is shown in Figure 5.9. The first half part of the plot represents the low-resolution approximation of the sound, the second half part represents the detail/error. We see that the detail is quite significant in this case. This means that the first order wavelet approximation does not give a very good approximation to the sound. In the exercises we will experiment more on this.

It is also interesting to plot only the detail/error in the sound, for different resolutions. For this, we must perform a DWT so that we get a representation in the basis $V_0 \oplus W_0 \oplus W_1 \oplus \dots \oplus W_{m-1}$, set the coefficients from V_0 to zero, and transform back with the IDWT. In figure 5.10 the error is shown for the test audio file `castanets.wav` for $m = 1, m = 2$. This clearly shows that the error is larger when two levels of the DWT are performed, as one would suspect. It is also seen that the error is larger in the part of the file where there are bigger variations. This also sounds reasonable. ♣

The previous example illustrates that wavelets as well may be used to perform

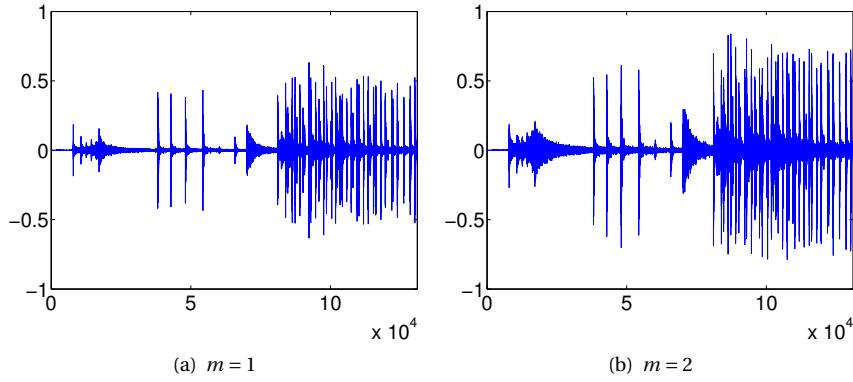


Figure 5.10: The error (i.e. the contribution from $W_0 \oplus W_1 \oplus \cdots \oplus W_{m-1}$) in the sound file castanets.wav, for different values of m .

operations on sound. As we will see later, however, our main application for wavelets will be images, where they have found a more important role than for sound. Images typically display variations which are less abrupt than the ones found in sound. Just as the functions above had smaller errors in the corresponding resolution spaces than the sound had, images are thus more suited for use with wavelets. The main idea behind why wavelets are so useful comes from the fact that the detail, i.e., wavelet coefficients corresponding to the spaces W_k , are often very small. After a DWT one is therefore often left with a couple of significant coefficients, while most of the coefficients are small. The approximation from V_0 can be viewed as a good approximation, even though it contains much less information. This gives another reason why wavelets are popular for images: Detailed images can be very large, but when they are downloaded to a web browser, the browser can very early show a low-resolution of the image, while waiting for the rest of the details in the image to be downloaded. When we later look at how wavelets are applied to images, we will need to handle one final hurdle, namely that images are two-dimensional.

Example 5.22. Above we plotted the DWT coefficients of a sound, as well as the detail/error. We can also experiment with samples generated from a mathematical function. Figure 5.11 plots the error for different functions, with $N = 1024$. In these cases, we see that we require large m before the detail/error becomes significant. We see also that there is no error for the square wave. The reason is that the square wave is a piecewise constant function, so that it can be represented exactly by the ϕ -functions. For the other functions, however, this is not the case, so we here get an error. ♣

Above we used the functions `DWTHaarImpl`, `IDWTHaarImpl` to plot the error. For the functions we plotted in the previous example it is also possible to compute the wavelet coefficients, which we previously have denoted by $w_{m,n}$, exactly. You will be asked to do this in exercises 13 and 14. The following example shows the general procedure which can be used for this:

Example 5.23. Let us compute the wavelet coefficients $w_{m,n}$ for the function $f(t) =$

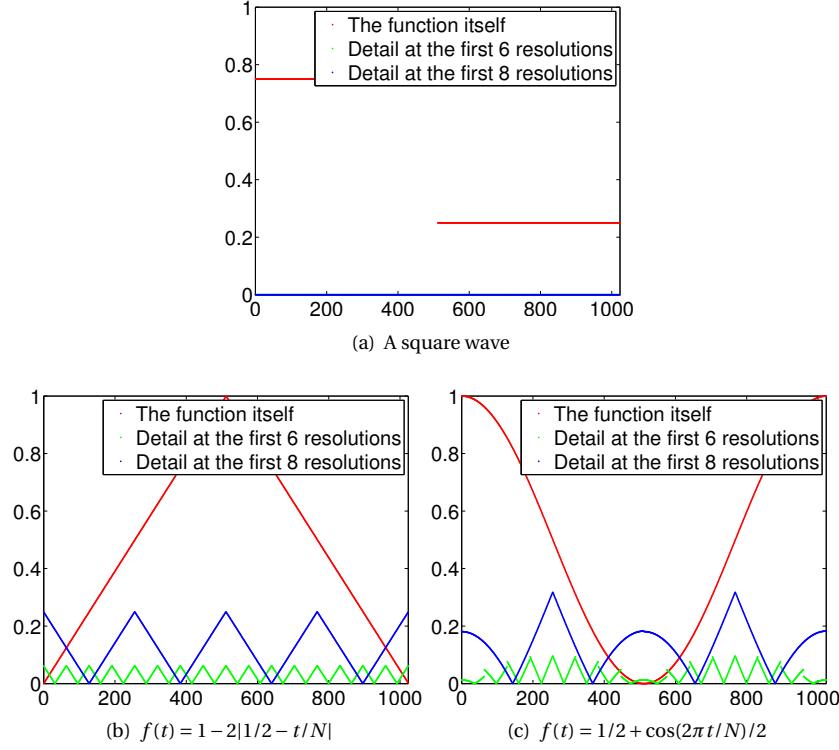


Figure 5.11: The error (i.e. the contribution from $W_0 \oplus W_1 \oplus \dots \oplus W_{m-1}$) for $N = 1024$ for different functions $f(t)$, for different values of m .

$1 - t/N$. This function decreases linearly from 1 to 0 on $[0, N]$. Since the $w_{m,n}$ are coefficients in the basis $\{\psi_{m,n}\}$, it follows by the orthogonal decomposition formula that $w_{m,n} = \langle f, \psi_{m,n} \rangle = \int_0^N f(t) \psi_{m,n}(t) dt$. Using the definition of $\psi_{m,n}$ we get that

$$w_{m,n} = \int_0^N (1 - t/N) \psi_{m,n}(t) dt = 2^{m/2} \int_0^N (1 - t/N) \psi(2^m t - n) dt.$$

Moreover $\psi_{m,n}$ is nonzero only on $[2^{-m}n, 2^{-m}(n+1))$, and is 1 on $[2^{-m}n, 2^{-m}(n+1))$.

$1/2)$), and -1 on $[2^{-m}(n+1/2), 2^{-m}(n+1))$. We can therefore write

$$\begin{aligned}
w_{m,n} &= 2^{m/2} \int_{2^{-m}n}^{2^{-m}(n+1/2)} (1-t/N) dt - 2^{m/2} \int_{2^{-m}(n+1/2)}^{2^{-m}(n+1)} (1-t/N) dt \\
&= 2^{m/2} \left[t - \frac{t^2}{2N} \right]_{2^{-m}n}^{2^{-m}(n+1/2)} - 2^{m/2} \left[t - \frac{t^2}{2N} \right]_{2^{-m}(n+1/2)}^{2^{-m}(n+1)} \\
&= 2^{m/2} \left(2^{-m}(n+1/2) - \frac{2^{-2m}(n+1/2)^2}{2N} - 2^{-m}n + \frac{2^{-2m}n^2}{2N} \right) \\
&\quad - 2^{m/2} \left(2^{-m}(n+1) - \frac{2^{-2m}(n+1)^2}{2N} - 2^{-m}(n+1/2) + \frac{2^{-2m}(n+1/2)^2}{2N} \right) \\
&= 2^{m/2} \left(\frac{2^{-2m}n^2}{2N} - \frac{2^{-2m}(n+1/2)^2}{N} + \frac{2^{-2m}(n+1)^2}{2N} \right) \\
&= \frac{2^{-3m/2}}{2N} (n^2 - 2(n+1/2)^2 + (n+1)^2) = \frac{1}{N2^{2+3m/2}}.
\end{aligned}$$

We see in particular that $w_{m,n} \rightarrow 0$ when $m \rightarrow \infty$. We see also that there were a lot of computations even in this very simple example. For most functions we therefore usually do not compute $w_{m,n}$ exactly. Instead we use implementations like DWTHaarImpl, IDWTHaarImpl, and run them on a computer. ♣

What you should have learnt in this section

Definition of the m -level Discrete Wavelet Transform. Matlab implementation of the Haar wavelet transform and its inverse. Experimentation with wavelets on sound. Direct computation of the wavelet coefficients.

Exercises for Section 5.3

1. Generalize exercise 4 to the projections from V_{m+1} onto V_m and W_m .
2. Show that $f(t) \in V_m$ if and only if $g(t) = f(2t) \in V_{m+1}$.
3. Let C_1, C_2, \dots, C_n be independent vector spaces, and let $T_i : C_i \rightarrow C_i$ be linear transformations. The direct sum of T_1, T_2, \dots, T_n , written as $T_1 \oplus T_2 \oplus \dots \oplus T_n$, denotes the linear transformation from $C_1 \oplus C_2 \oplus \dots \oplus C_n$ to itself defined by

$$T_1 \oplus T_2 \oplus \dots \oplus T_n(\mathbf{c}_1 + \mathbf{c}_2 + \dots + \mathbf{c}_n) = T_1(\mathbf{c}_1) + T_2(\mathbf{c}_2) + \dots + T_n(\mathbf{c}_n)$$

when $\mathbf{c}_1 \in C_1, \mathbf{c}_2 \in C_2, \dots, \mathbf{c}_n \in C_n$. Similarly, when A_1, A_2, \dots, A_n are square matrices, $A_1 \oplus A_2 \oplus \dots \oplus A_n$ is defined as the block matrix where the blocks along the diagonal are A_1, A_2, \dots, A_n , and where all other blocks are 0. Show that, if \mathcal{B}_i is a basis for C_i then

$$[T_1 \oplus T_2 \oplus \dots \oplus T_n]_{(\mathcal{B}_1, \mathcal{B}_2, \dots, \mathcal{B}_n)} = [T_1]_{\mathcal{B}_1} \oplus [T_2]_{\mathcal{B}_2} \oplus \dots \oplus [T_n]_{\mathcal{B}_n},$$

Here two new concepts are used: a direct sum of matrices, and a direct sum of linear transformations.

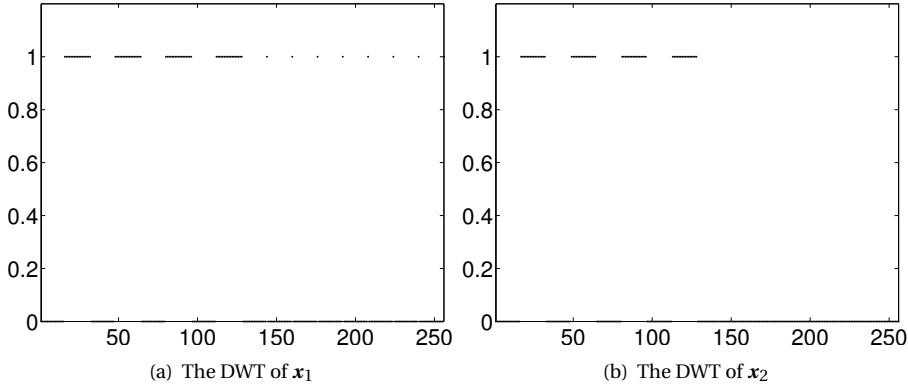


Figure 5.12: 2 vectors which seem equal, but where the DWT's are very different

4. Assume that T_1 and T_2 are matrices, and that the eigenvalues of T_1 are equal to those of T_2 . What are the eigenvalues of $T_1 \oplus T_2$? Can you express the eigenvectors of $T_1 \oplus T_2$ in terms of those of T_1 and T_2 ?
5. Assume that A and B are square matrices which are invertible. Show that $A \oplus B$ is invertible, and that $(A \oplus B)^{-1} = A^{-1} \oplus B^{-1}$.
6. Let A, B, C, D be square matrices of the same dimensions. Show that $(A \oplus B)(C \oplus D) = (AC) \oplus (BD)$.
7. Assume that you run an m -level DWT on a vector of length r . What value of N does this correspond to? Note that an m -level DWT performs a change of coordinates from V_m to $V_0 \oplus W_0 \oplus W_1 \oplus \dots \oplus W_{m-2} \oplus W_{m-1}$.
8. In Figure 5.12 we have plotted the DWT's of two vectors x_1 and x_2 . In both vectors we have 16 ones followed by 16 zeros, and this pattern repeats cyclically so that the length of both vectors is 256. The only difference is that the second vector is obtained by delaying the first vector with one element. You see that the two DWT's are very different: For the first vector we see that there is much detail present (the second part of the plot), while for the second vector there is no detail present. Attempt to explain why this is the case. Based on your answer, also attempt to explain what can happen if you change the point of discontinuity for the piecewise constant function in Figure 5.22(a) to something else.
9. Run a 2-level DWT on the first 2^{17} sound samples of the audio file `castanets.wav`, and plot the values of the resulting DWT-coefficients. Compare the values of the coefficients from V_0 with those from W_0 and W_1 .
10. In this exercise we will experiment with applying an m -level DWT to a sound file.
 - a. Write a function

```
function playDWTlower(m)
```

which

1. reads the audio file `castanets.wav`,
 2. performs an m -level DWT to the first 2^{17} sound samples of \mathbf{x} using the function `DWTHaarImpl`,
 3. sets all wavelet coefficients representing detail to zero (i.e. keep only wavelet coefficients from V_0 in the decomposition $V_0 \oplus W_0 \oplus W_1 \oplus \dots \oplus W_{m-2} \oplus W_{m-1}$),
 4. performs an IDWT on the resulting coefficients using the function `IDWTHaarImpl`,
 5. plays the resulting sound.
- b.** Run the function `playDWTlower` for different values of m . For which m can you hear that the sound gets degraded? How does it get degraded? Compare with what you heard through the function `playDFTlower` in Example 2.29, where you performed a DFT on the sound sample instead, and set some of the DFT coefficients to zero.
- c.** Do the sound samples returned by `playDWTlower` lie in $[-1, 1]$?

- 11.** Attempt to construct a (nonzero) sound where the function `playDWTlower` from the previous exercise does not change the sound for $m = 1, 2$.
- 12.** Repeat Exercise 10, but this time instead keep only wavelet coefficients from the detail spaces W_0, W_1, \dots . Call the new function `playDWTlowerdifference`. What kind of sound do you hear? Can you recognize the original sound in what you hear?
- 13.** Compute the wavelet detail coefficients analytically for the functions in Example 5.22, i.e. compute the quantities $w_{m,n} = \int_0^N f(t)\psi_{m,n}(t)dt$ similarly to how this was done in Example 5.23.
- 14.** Compute the wavelet detail coefficients analytically for the functions $f(t) = (\frac{t}{N})^k$, i.e. compute the quantities $w_{m,n} = \int_0^N (\frac{t}{N})^k \psi_{m,n}(t)dt$ similarly to how this was done in Example 5.23. How do these compare with the coefficients from the Exercise 13?
- 15.** Suppose that we have the vector \mathbf{x} with length $2^{10} = 1024$, defined by $x_n = 1$ for n even, $x_n = -1$ for n odd. What will be the result if you run a 10-level DWT on \mathbf{x} ? Use the function `DWTHaarImpl` to verify what you have found.
Hint: We defined ψ by $\psi(t) = (\phi_{1,0}(t) - \phi_{1,1}(t))/\sqrt{2}$. From this connection it follows that $\psi_{9,n} = (\phi_{10,2n} - \phi_{10,2n+1})/\sqrt{2}$, and thus $\phi_{10,2n} - \phi_{10,2n+1} = \sqrt{2}\psi_{9,n}$. Try to couple this identity with the alternating sign you see in \mathbf{x} .

5.4 A wavelet based on piecewise linear functions

Unfortunately, piecewise constant functions are too simple to provide good approximations. In this section we are going to extend the construction of wavelets to piecewise linear functions. The advantage is that piecewise linear functions are better for approximating smooth functions and data than piecewise constants, which should

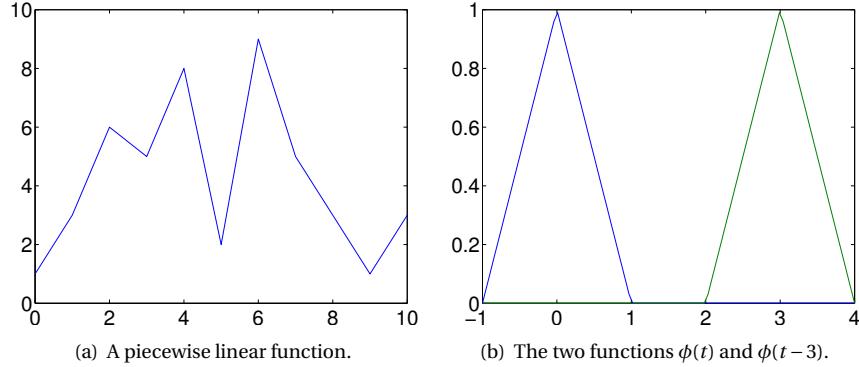


Figure 5.13: Some piecewise linear functions.

translate into smaller components (errors) in the detail spaces in many practical situations. As an example, this would be useful if we are interested in compression. In this new setting it turns out that we loose the orthonormality we had for the Haar wavelet. On the other hand, we will see that the new scaling functions and mother wavelets are symmetric functions. We will later see that this implies that the corresponding DWT and IDWT have simple implementations with higher precision. Our experience from deriving Haar wavelets will guide us in the construction of piecewise linear wavelets. The first task is to define the new resolution spaces.

Definition 5.24 (Resolution spaces of piecewise linear functions). The space V_m is the subspace of continuous functions on \mathbb{R} which are periodic with period N , and linear on each subinterval of the form $[n2^{-m}, (n+1)2^{-m}]$.

Any $f \in V_m$ is uniquely determined by its values on $[0, N]$. Figure 5.13(a) shows an example of a piecewise linear function in V_0 on the interval $[0, 10]$. We note that a piecewise linear function in V_0 is completely determined by its value at the integers, so the functions that are 1 at one integer and 0 at all others are particularly simple and therefore interesting, see Figure 5.13(b). These simple functions are all translates of each other and can therefore be built from one scaling function, as is required for a multiresolution analysis.

Recall that the *support* of a function f defined on a subset I of \mathbb{R} is given by the closure of the set of points where the function is nonzero,

$$\text{supp}(f) = \overline{\{t \in I \mid f(t) \neq 0\}}.$$

Lemma 5.25. Let the function ϕ be defined by

$$\phi(t) = \begin{cases} 1+t, & \text{if } -1 \leq t < 0; \\ 1-t, & \text{if } 0 \leq t < 1; \\ 0, & \text{otherwise;} \end{cases} \quad (5.22)$$

and for any $m \geq 0$ set

$$\phi_{m,n}(t) = 2^{m/2} \phi(2^m t - n) \quad \text{for } n = 0, 1, \dots, 2^m N - 1,$$

or in vector notation

$$\boldsymbol{\phi}_m = (\phi_{m,0}, \phi_{m,1}, \dots, \phi_{m,2^m N - 1}).$$

The functions $\{\phi_{m,n}\}_{n=0}^{2^m N - 1}$, restricted to the interval $[0, N]$, form a basis for the space V_m for this interval. In other words, the function ϕ is a scaling function for the spaces V_0, V_1, \dots . Moreover, the function $\phi_{0,n}(t)$ is the function in V_0 with smallest support that is nonzero at $t = n$.

Proof: The proof is similar for all the resolution spaces, so it is sufficient to consider the proof in the case of V_0 . The function ϕ is clearly linear between each pair of neighbouring integers, and it is also easy to check that it is continuous. Its restriction to $[0, N]$ therefore lies in V_0 . And as we noted above $\phi_{0,n}(t)$ is 0 at all the integers except at $t = n$ where its value is 1.

A general function f in V_0 is completely determined by its values at the integers in the interval $[0, N]$ since all straight line segments between neighbouring integers are then fixed. Note that we can also write f as

$$f(t) = \sum_{n=0}^{N-1} f(n) \phi_{0,n}(t) \quad (5.23)$$

since this function agrees with f at the integers in the interval $[0, N]$ and is linear on each subinterval between two neighbouring integers. This means that V_0 is spanned by the functions $\{\phi_{0,n}\}_{n=0}^{N-1}$. On the other hand, if f is identically 0, all the coefficients in (5.23) are also 0, so $\{\phi_{0,n}\}_{n=0}^{N-1}$ are linearly independent and therefore a basis for V_0 .

Suppose that the function $g \in V_0$ has smaller support than $\phi_{0,n}$, but is nonzero at $t = n$. Then g must be identically zero either on $[n-1, n]$ or on $[n, n+1]$, since a straight line segment cannot be zero on just a part of an interval between integers. But then g cannot be continuous, which contradicts the fact that it lies in V_0 . ■

The function ϕ and its translates and dilates are often referred to as hat functions for obvious reasons. Note that the new function ϕ is nonzero for small negative x -values, contrary to the ϕ we defined in Chapter 5. If we plotted the function on $[0, N]$, we would see the nonzero parts at the beginning and end of this interval, due to the period N , but we will mostly plot on an interval around zero, since such an interval captures the entire support of the function. A formula like (5.23) is also valid for functions in V_m :

Lemma 5.26. A function $f \in V_m$ may be written as

$$f(t) = \sum_{n=0}^{2^m N - 1} f(n/2^m) 2^{-m/2} \phi_{m,n}(t). \quad (5.24)$$

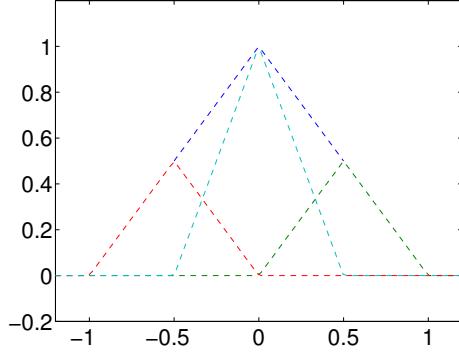


Figure 5.14: How $\phi(t)$ can be decomposed as a linear combination of $\phi_{1,-1}$, $\phi_{1,0}$, and $\phi_{1,1}$.

An essential property also here is that the spaces are nested.

Lemma 5.27. The piecewise linear resolution spaces are nested,

$$V_0 \subset V_1 \subset \dots \subset V_m \subset \dots$$

Proof: We only need to prove that $V_0 \subset V_1$ since the other inclusions are similar. But this is immediate since any function in V_0 is continuous, and linear on any subinterval in the form $[n/2, (n+1)/2]$. ■

In the piecewise constant case, we saw in Lemma 5.5 that the scaling functions were automatically orthogonal since their supports did not overlap. This is not the case in the linear case, but we could orthogonalise the basis ϕ_m with the Gram-Schmidt process from linear algebra. The disadvantage is that we lose the nice local behaviour of the scaling functions and end up with basis functions that are nonzero over all of $[0, N]$. And for most applications, orthogonality is not essential; we just need a basis. The next step in the derivation of wavelets is to find formulas that let us express a function given in the basis ϕ_0 for V_0 in terms of the basis ϕ_1 for V_1 .

Lemma 5.28. The function $\phi_{0,n}$ satisfies the relation

$$\phi = \frac{1}{\sqrt{2}} \left(\frac{1}{2} \phi_{1,-1} + \phi_{1,0} + \frac{1}{2} \phi_{1,1} \right). \quad (5.25)$$

Proof: Since $\phi_{0,n}$ is in V_0 it may be expressed in the basis ϕ_1 with formula (5.24),

$$\phi_{0,n}(t) = 2^{-1/2} \sum_{k=0}^{2N-1} \phi_{0,n}(k/2) \phi_{1,k}(t).$$

The relation (5.25) now follows since

$$\phi_{0,n}((2n-1)/2) = \phi_{0,n}((2n+1)/2) = 1/2, \quad \phi_{0,n}(2n/2) = 1,$$

and $\phi_{0,n}(k/2) = 0$ for all other values of k . ■

The relationship given by Equation 5.25 is shown in Figure 5.14.

5.4.1 Detail spaces and wavelets

The next step in our derivation of wavelets for piecewise linear functions is the definition of the detail spaces. We need to determine a space W_0 that is linearly independent from V_0 , and so that $V_1 = V_0 \oplus W_0$. In the case of piecewise constant functions we started with a function g_1 in V_1 , computed the least squares approximation g_0 in V_0 , and then defined the error function $e_0 = g_1 - g_0$, with $e_0 \in W_0$ and W_0 as the orthogonal complement of V_0 in V_1 .

It turns out that this strategy is less appealing in the case of piecewise linear functions. The reason is that the functions $\phi_{0,n}$ are not orthogonal anymore (see Exercise 1). Due to this we have no simple, orthogonal basis for the set of piecewise linear functions, so that the orthogonal decomposition theorem fails to give us the projection onto V_0 in a simple way. It is therefore no reason to use the orthogonal complement of V_0 in V_1 as our error space, since it is hard to write a piecewise linear function as a sum of two other piecewise linear functions which are orthogonal. Instead of using projections to find low-resolution approximations, and orthogonal complements to find error functions, we will attempt the following simple approximation method:

Definition 5.29. Let g_1 be a function in V_1 given by

$$g_1 = \sum_{n=0}^{2N-1} c_{1,n} \phi_{1,n}. \quad (5.26)$$

The approximation $g_0 = S(g_1)$ in V_0 is defined as the unique function in V_0 which has the same values as g_1 at the integers, i.e.

$$g_0(n) = g_1(n), \quad n = 0, 1, \dots, N-1. \quad (5.27)$$

It is easy to show that $S(g_1)$ actually is different from the projection of g_1 onto V_0 : If $g_1 = \phi_{1,1}$, then g_1 is zero at the integers, and then clearly $S(g_1) = 0$. But in Exercise 2 you will be asked to compute the projection onto V_0 using different means than the orthogonal decomposition theorem, and the result will be seen to be nonzero. It is also very easy to see that the coordinates of g_0 in ϕ_0 can be obtained by dropping every second coordinate of g_0 in ϕ_1 . To be more precise, the following holds:

Lemma 5.30. We have that

$$S(\phi_{1,n}) = \begin{cases} \sqrt{2}\phi_{0,n/2}, & \text{if } n \text{ is an even integer;} \\ 0, & \text{otherwise.} \end{cases}$$

Once the method of approximation is determined, it is straightforward to determine the detail space as the space of error functions. With the notation from Definition 5.29, the error is given by $e_0 = g_1 - g_0$. Since g_0 interpolates g_1 at the integers, the error is 0 there,

$$e_0(n) = 0, \quad \text{for } n = 0, 1, \dots, N-1.$$

Conversely, any function in V_1 which is 0 at the integers may be viewed as an error function in the above sense. This provides the basis for a precise description of the error functions.

Lemma 5.31. Suppose that $g_1 \in V_1$ and that $g_0 = S(g_1)$. Then the error $e_0 = g_1 - g_0$ lies in the space W_0 defined by

$$W_0 = \{f \in V_1 \mid f(n) = 0, \quad \text{for } n = 0, 1, \dots, N-1\}.$$

A basis for W_0 is given by $\{\psi_{0,n}\}_{n=0}^{N-1}$, where the mother wavelet ψ is defined by

$$\psi(t) = \frac{1}{\sqrt{2}}\phi_{1,1}(t). \quad (5.28)$$

Proof: Since $g_0(n) = g_1(n)$ for all integers n , $e_0(n) = 0$. This means that the error functions are precisely the piecewise linear functions at the half intervals which are zero at the integers. Clearly this is the same space as that spanned by the $\psi_{0,n} = \frac{1}{\sqrt{2}}\phi_{1,2n+1}$. \blacksquare

We now have all the ingredients to formulate an analog of Theorem 5.19 that describes how V_m can be expressed as a direct sum of V_{m-1} and W_{m-1} . The formulas for $m = 1$ generalise without change, except that the upper bound on the summation indices must be adjusted.

Theorem 5.32. The space V_m can be decomposed as the direct sum $V_m = V_{m-1} \oplus W_{m-1}$ where W_{m-1} is the space of all functions in V_m that are zero at the points $\{n/2^{m-1}\}_{n=0}^{N2^{m-1}-1}$. The space V_m has the two bases

$$\boldsymbol{\phi}_m = (\phi_{m,n})_{n=0}^{2^m N-1}, \text{ and } (\boldsymbol{\phi}_{m-1}, \boldsymbol{\psi}_{m-1}) = ((\phi_{m-1,n})_{n=0}^{2^{m-1} N-1}, (\psi_{m-1,n})_{n=0}^{2^{m-1} N-1}).$$

Equation (5.25) and (5.28) written together gives

$$\begin{aligned}\phi &= \frac{1}{2\sqrt{2}}\phi_{1,-1} + \frac{1}{\sqrt{2}}\phi_{1,0} + \frac{1}{2\sqrt{2}}\phi_{1,1} \\ \psi &= \frac{1}{\sqrt{2}}\phi_{1,1}.\end{aligned}\quad (5.29)$$

These two relations together give all columns in the change of coordinate matrix $P_{\phi_1 \leftarrow (\phi_0, \psi_0)}$ (i.e. the IDWT), when the spaces ϕ_m, ψ_m instead are defined in terms of the function ψ and ϕ . Similarly we can compute the change of coordinate matrix the opposite way by rewriting equations (5.29) as

$$\begin{aligned}\frac{1}{\sqrt{2}}\phi_{1,2n} &= \phi_{0,n} - \frac{1}{2\sqrt{2}}\phi_{1,2n-1} - \frac{1}{2\sqrt{2}}\phi_{1,2n+1} \\ \frac{1}{\sqrt{2}}\phi_{1,2n+1} &= \psi_{0,n},\end{aligned}$$

from which it follows that

$$\phi_{1,2n} = \sqrt{2}\phi_{0,n} - \frac{1}{2}\phi_{1,2n-1} - \frac{1}{2}\phi_{1,2n+1} = -\frac{\sqrt{2}}{2}\psi_{0,n-1} + \sqrt{2}\phi_{0,n} - \frac{\sqrt{2}}{2}\psi_{0,n} \quad (5.30)$$

$$\phi_{1,2n+1} = \sqrt{2}\psi_{0,n}. \quad (5.31)$$

This gives us the columns in the change of coordinate matrix $P_{(\phi_0, \psi_0) \rightarrow \phi_1}$ as well, i.e. the DWT.

Example 5.33. In Chapter 6 we will construct an algorithm which performs DWT/IDWT, for a general wavelet. In particular, this algorithm can be used for the wavelet we constructed in this section. Let us also for this wavelet plot the detail/error in the test audio file `castanets.wav` for different resolutions, as we did in Example 5.21. The result is shown in Figure 5.15. When comparing with Figure 5.10 we see much of the same, but it seems here that the error is bigger than before. In the next section we will try to explain why this is the case, and construct another wavelet based on piecewise linear functions which remedies this. ♣

Example 5.34. Let us also repeat Exercise 5.22, where we plotted the detail/error at different resolutions, for the samples of a mathematical function. Figure 5.16 shows the new plot. With the square wave we see now that there is an error. The reason is that a piecewise constant function can not be represented exactly by piecewise linear functions, due to discontinuity. For the second function we see that there is no error. The reason is that this function is piecewise linear, so there is no error when we represent the function from the space V_0 . With the third function, however, we see an error. ♣

What you should have learnt in this section

Definition of scaling function, mother wavelet, resolution spaces, and detail spaces for the wavelet of piecewise linear functions.

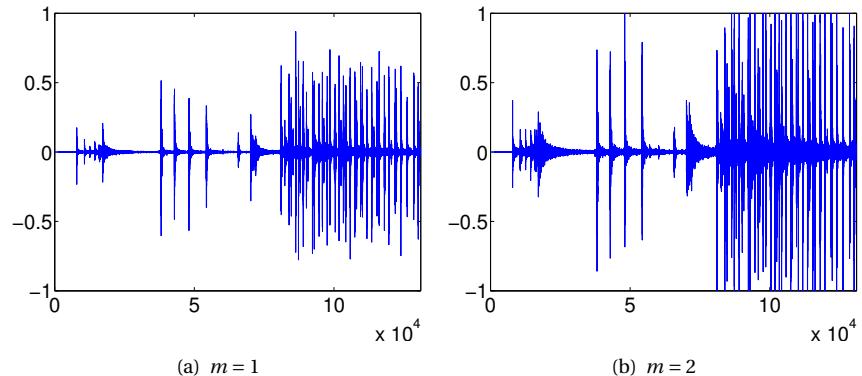


Figure 5.15: The error (i.e. the contribution from $W_0 \oplus W_1 \oplus \cdots \oplus W_{m-1}$) in the sound file castanets.wav, for different values of m .

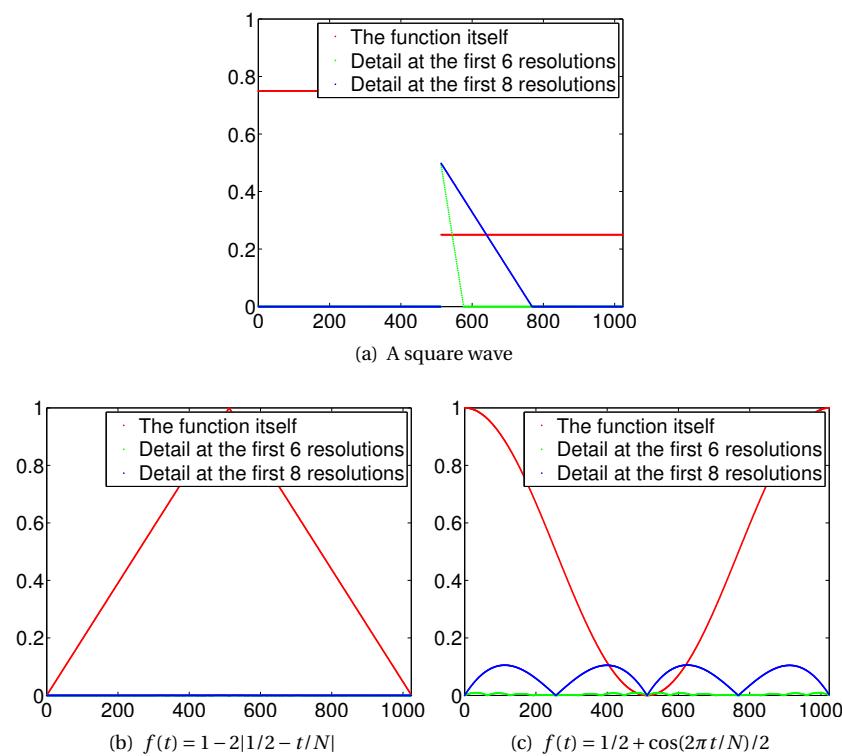


Figure 5.16: The error (i.e. the contribution from $W_0 \oplus W_1 \oplus \dots \oplus W_{m-1}$) for $N = 1024$ for different functions $f(t)$, for different values of m .

Exercises for Section 5.4

1. Show that, for $f \in V_0$ we have that $[f]_{\phi_0} = (f(0), f(1), \dots, f(N-1))$. This generalizes the result for piecewise constant functions.
2. Show that the projection of $\phi_{1,1}$ onto V_0 is ..
3. Show that

$$\langle \phi_{0,n}, \phi_{0,n} \rangle = \frac{2}{3} \quad \langle \phi_{0,n}, \phi_{0,n \pm 1} \rangle = \frac{1}{6} \quad \langle \phi_{0,n}, \phi_{0,n \pm k} \rangle = 0 \text{ for } k > 1.$$

As a consequence, the $\{\phi_{0,n}\}_n$ are neither orthogonal, nor have norm 1.

4. The convolution of two functions defined on $(-\infty, \infty)$ is defined by

$$(f * g)(x) = \int_{-\infty}^{\infty} f(t)g(x-t)dt.$$

Show that we can obtain the piecewise linear ϕ we have defined as $\phi = \chi_{[-1/2,1/2]} * \chi_{[-1/2,1/2]}$ (recall that $\chi_{[-1/2,1/2]}$ is the function which is 1 on $[-1/2, 1/2]$ and 0 elsewhere). This gives us a nice connection between the piecewise constant scaling function (which is similar to $\chi_{[-1/2,1/2]}$) and the piecewise linear scaling function in terms of convolution.

5.5 Alternative wavelet based on piecewise linear functions

For the scaling function used for piecewise linear functions, $\{\phi(t-n)\}_{0 \leq n < N}$ were not orthogonal anymore, contrary to the case for piecewise constant functions. We were still able to construct what we could call resolution spaces and detail spaces. We also mentioned that having many vanishing moments is desirable for a mother wavelet, and that the scaling function used for piecewise constant functions had one vanishing moment. It is easily checked, however, that the mother wavelet we now introduced for piecewise linear functions (i.e. $\psi(t) = \frac{1}{\sqrt{2}}\phi_{1,1}(t)$) has no vanishing moments. Therefore, this is not a very good choice of mother wavelet. We will attempt the following adjustment strategy to construct an alternative mother wavelet $\hat{\psi}$ which has two vanishing moments, i.e. one more than the Haar wavelet.

Idea 5.35. Adjust the wavelet construction in Theorem 5.32 so that $\{\hat{\psi}_{m,n}\}_{n=0}^{N2^m-1}$ is a basis for W_m , where the new mother wavelet $\hat{\psi}$ satisfies

$$\int_0^N \hat{\psi}(t) dt = \int_0^N t\hat{\psi}(t) dt = 0. \quad (5.32)$$

From this we see that we need to enforce two conditions for each wavelet function. If we adjust the wavelets in Theorem 5.32 by adding multiples of the two neighbouring hat functions, we have two free parameters,

$$\hat{\psi} = \psi - \alpha\phi_{0,0} - \beta\phi_{0,1} \quad (5.33)$$

that we may determine so that the two conditions in (5.32) are enforced. In Exercise 1 you are taken through the details of this, which gives the following result:

Lemma 5.36. The function

$$\hat{\psi}(t) = \psi(t) - \frac{1}{4}(\phi_{0,0}(t) + \phi_{0,1}(t)) \quad (5.34)$$

satisfies the conditions

$$\int_0^N \hat{\psi}(t) dt = \int_0^N t\hat{\psi}(t) dt = 0.$$

Using Equation (5.25), which stated that

$$\phi(t) = \frac{1}{\sqrt{2}}\left(\frac{1}{2}\phi_{1,-1} + \phi_{1,0} + \frac{1}{2}\phi_{1,1}\right) \quad (5.35)$$

we get

$$\begin{aligned} \hat{\psi}(t) &= \psi(t) - \frac{1}{4}(\phi_{0,0}(t) + \phi_{0,1}(t)) \\ &= \frac{1}{\sqrt{2}}\phi_{1,1}(t) - \frac{1}{4}\frac{1}{\sqrt{2}}\left(\frac{1}{2}\phi_{1,-1} + \phi_{1,0} + \frac{1}{2}\phi_{1,1} + \frac{1}{2}\phi_{1,1} + \phi_{1,2} + \frac{1}{2}\phi_{1,3}\right) \\ &= -\frac{1}{8\sqrt{2}}\phi_{1,-1} - \frac{1}{4\sqrt{2}}\phi_{1,0} + \frac{3}{4\sqrt{2}}\phi_{1,1} - \frac{1}{4\sqrt{2}}\phi_{1,2} - \frac{1}{8\sqrt{2}}\phi_{1,3} \end{aligned} \quad (5.36)$$

Note that what we did here is equivalent to finding the coordinates of $\hat{\psi}$ in the basis ϕ_1 : Equation (5.34) says that

$$[\hat{\psi}]_{(\phi_0, \psi_0)} = (-1/4, -1/4, 0, \dots, 0) \oplus (1, 0, \dots, 0). \quad (5.37)$$

Since the IDWT is the change of coordinates from (ϕ_0, ψ_0) to ϕ_1 , we could also have computed $[\hat{\psi}]_{\phi_1}$ by taking the IDWT of $(-1/4, -1/4, 0, \dots, 0) \oplus (1, 0, \dots, 0)$. In the next section we will consider more general implementations of the DWT and the IDWT, which we thus can use instead of performing the computation above.

In summary we have

$$\begin{aligned} \phi(t) &= \frac{1}{2\sqrt{2}}\phi_{1,-1} + \frac{1}{\sqrt{2}}\phi_{1,0} + \frac{1}{2\sqrt{2}}\phi_{1,1} \\ \hat{\psi}(t) &= -\frac{1}{8\sqrt{2}}\phi_{1,-1} - \frac{1}{4\sqrt{2}}\phi_{1,0} + \frac{3}{4\sqrt{2}}\phi_{1,1} - \frac{1}{4\sqrt{2}}\phi_{1,2} - \frac{1}{8\sqrt{2}}\phi_{1,3}, \end{aligned} \quad (5.38)$$

which gives us the change of coordinate matrix $P_{\phi_1 \leftarrow (\phi_0, \psi_0)}$. The new function $\hat{\psi}$ is plotted in Figure 5.17. We see that $\hat{\psi}$ has support $(-1, 2)$, and consists of four linear segments glued together. This is in contrast with the old ψ , which was simpler in that it had the shorter support $(0, 1)$, and consisted of only two linear segments glued together. It may therefore seem surprising that $\hat{\psi}$ is better suited for approximating functions than ψ . This is indeed a more complex fact, which may not be deduced by simply looking at plots of the functions.

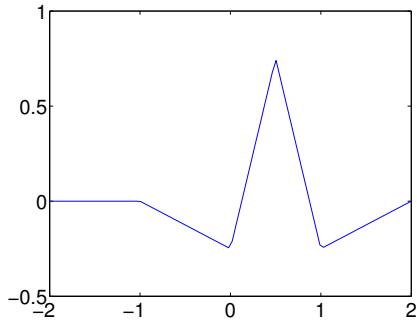


Figure 5.17: The function $\hat{\psi}$ we constructed as an alternative wavelet for piecewise linear functions.

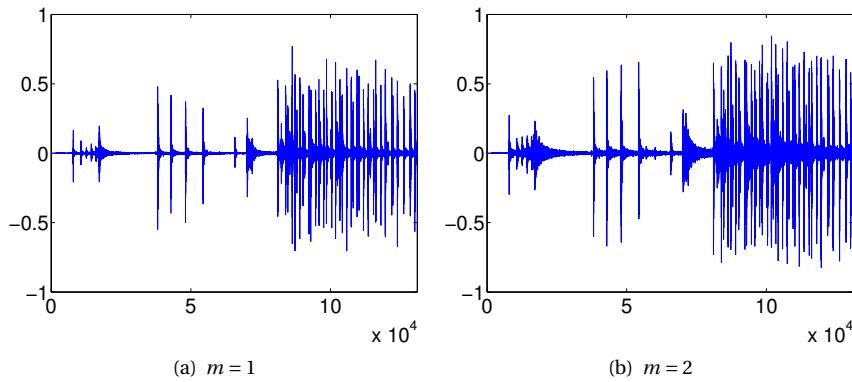
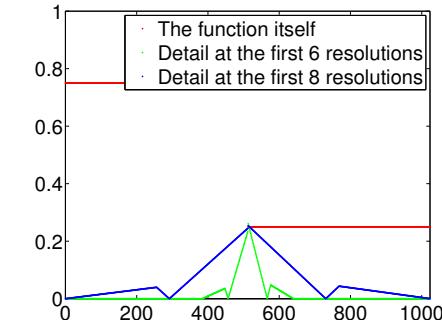


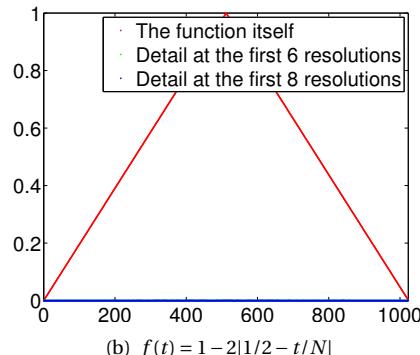
Figure 5.18: The error (i.e. the contribution from $W_0 \oplus W_1 \oplus \dots \oplus W_{m-1}$) in the sound file `castanets.wav`, for different values of m .

Example 5.37. Let us also plot the detail/error in the test audio file `castanets.wav` for different resolutions for our alternative wavelet, as we did in Example 5.21. The result is shown in Figure 5.18. Again, when comparing with Figure 5.10 we see much of the same. It is difficult to see an improvement from this figure. However, this figure also clearly shows a smaller error than the wavelet of the preceding section. A partial explanation is that the wavelet we now have constructed has two vanishing moments, while the previous one had not. ♣

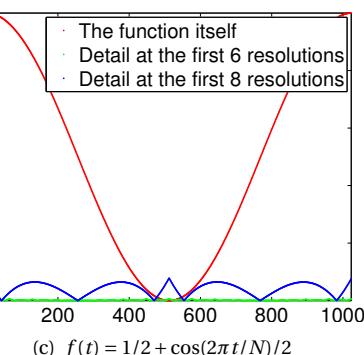
Example 5.38. Let us also repeat Exercise 5.22 for our alternative wavelet, where we plotted the detail/error at different resolutions, for the samples of a mathematical function. Figure 5.19 shows the new plot. Again for the square wave there is an error, which seems to be slightly lower than for the previous wavelet. For the second function we see that there is no error, as before. The reason is the same as before, since the function is piecewise linear. With the third function there is an error. The error seems to be slightly lower than for the previous wavelet, which fits well with



(a) A square wave



(b) $f(t) = 1 - 2|1/2 - t/N|$



(c) $f(t) = 1/2 + \cos(2\pi t/N)/2$

Figure 5.19: The error (i.e. the contribution from $W_0 \oplus W_1 \oplus \dots \oplus W_{m-1}$) for $N = 1024$ for different functions $f(t)$, for different values of m .

the fact that this new wavelet has a bigger number of vanishing moments. ♣

What you should have learnt in this section

How one alters the mother wavelet for piecewise linear functions, in order to add a vanishing moment.

Exercises for Section 5.5

1. In this exercise we will show that there is a unique function on the form (5.33) which has two vanishing moments.

- a.** Show that, when $\hat{\psi}$ is defined by (5.33), we have that

$$\hat{\psi}(t) = \begin{cases} -\alpha t - \alpha & \text{for } -1 \leq t < 0 \\ (2 + \alpha - \beta)t - \alpha & \text{for } 0 \leq t < 1/2 \\ (\alpha - \beta - 2)t - \alpha + 2 & \text{for } 1/2 \leq t < 1 \\ \beta t - 2\beta & \text{for } 1 \leq t < 2 \\ 0 & \text{for all other } t \end{cases}$$

- b.** Show that

$$\int_0^N \hat{\psi}(t) dt = \frac{1}{2} - \alpha - \beta, \quad \int_0^N t \hat{\psi}(t) dt = \frac{1}{4} - \beta.$$

- c.** Explain why there is a unique function on the form (5.33) which has two vanishing moments, and that this function is given by Equation (5.34).

- 2.** In the previous exercise we ended up with a lot of calculations to find α, β in Equation (5.33). Let us try to make a program which does this for us, and which also makes us able to generalize the result.

- a.** Define

$$a_k = \int_{-1}^1 t^k (1 - |t|) dt, \quad b_k = \int_0^2 t^k (1 - |t - 1|) dt, \quad e_k = \int_0^1 t^k (1 - 2|t - 1/2|) dt,$$

for $k \geq 0$. Explain why finding α, β so that we have two vanishing moments in Equation 5.33 is equivalent to solving the following equation:

$$\begin{pmatrix} a_0 & b_0 \\ a_1 & b_1 \end{pmatrix} \begin{pmatrix} \alpha \\ \beta \end{pmatrix} = \begin{pmatrix} e_0 \\ e_1 \end{pmatrix}$$

Write a program which sets up and solves this system of equations, and use this program to verify the values for α, β we previously have found.

Hint: recall that you can integrate functions in Matlab with the function quad. As an example, the function $\phi(t)$, which is nonzero only on $[-1, 1]$, can be integrated as follows:

```
quad(@(t)t.^k.*(1-abs(t)), -1, 1)
```

- b.** The procedure where we set up a matrix equation in a. allows for generalization to more vanishing moments. Define

$$\hat{\psi} = \psi_{0,0} - \alpha \phi_{0,0} - \beta \phi_{0,1} - \gamma \phi_{0,-1} - \delta \phi_{0,2}. \quad (5.39)$$

We would like to choose $\alpha, \beta, \gamma, \delta$ so that we have 4 vanishing moments. Define also

$$g_k = \int_{-2}^0 t^k (1 - |t + 1|) dt, \quad d_k = \int_1^3 t^k (1 - |t - 2|) dt$$

for $k \geq 0$. Show that $\alpha, \beta, \gamma, \delta$ must solve the equation

$$\begin{pmatrix} a_0 & b_0 & g_0 & d_0 \\ a_1 & b_1 & g_1 & d_1 \\ a_2 & b_2 & g_2 & d_2 \\ a_3 & b_3 & g_3 & d_3 \end{pmatrix} \begin{pmatrix} \alpha \\ \beta \\ \gamma \\ \delta \end{pmatrix} = \begin{pmatrix} e_0 \\ e_1 \\ e_2 \\ e_3 \end{pmatrix},$$

and solve this with Matlab.

- c.** Plot the function defined by (5.39), which you found in b.

Hint: If t is the vector of t -values, and you write

$(t \geq 0) .*(t \leq 1) .*(1 - 2 * \text{abs}(t - 0.5))$

you get the points $\phi_{1,1}(t)$.

- d.** Explain why the coordinate vector of $\hat{\psi}$ in the basis (ϕ_0, ψ_0) is

$$[\hat{\psi}]_{(\phi_0, \psi_0)} = (-\alpha, -\beta, -\delta, 0, \dots, 0 - \gamma) \oplus (1, 0, \dots, 0).$$

Hint: You can also compare with Equation (5.37) here. The placement of $-\gamma$ may seem a bit strange here, and has to do with that $\phi_{0,-1}$ is not one of the basis functions $\{\phi_{0,n}\}_{n=0}^{N-1}$. However, we have that $\phi_{0,-1} = \phi_{0,N-1}$, i.e. $\phi(t+1) = \phi(t-N+1)$, since we always assume that the functions we work with have period N .

- e.** Sketch a more general procedure than the one you found in b., which can be used to find wavelet bases where we have even more vanishing moments.

- 3.** Let $\phi(t)$ be the function we used when we defined the Haar-wavelet.

- a.** Compute $\text{proj}_{V_0}(f(t))$, where $f(t) = t^2$, and where f is defined on $[0, N]$.

- b.** Find constants α, β so that $\hat{\psi}(t) = \psi(t) - \alpha\phi_{0,0}(t) - \beta\phi_{0,1}(t)$ has two vanishing moments, i.e. so that $\langle \hat{\psi}, 1 \rangle = 0$, and $\langle \hat{\psi}, t \rangle = 0$. Plot also the function $\hat{\psi}$.

Hint: Start with computing the integrals $\int \psi(t) dt$, $\int t\psi(t) dt$, $\int \phi_{0,0}(t) dt$, $\int \phi_{0,1}(t) dt$, and $\int t\phi_{0,0}(t) dt$, $\int t\phi_{0,1}(t) dt$.

- c.** Express ϕ and $\hat{\psi}$ with the help of functions from ϕ_1 , and use this to write down the change of coordinate matrix from (ϕ_0, ψ_0) to ϕ_1 .

- 4.** It is also possible to add more vanishing moments to the Haar wavelet. Define

$$\hat{\psi} = \psi_{0,0} - a_0\phi_{0,0} - \cdots - a_{k-1}\phi_{0,k-1}.$$

Define also $c_{r,l} = \int_l^{l+1} t^r dt$, and $e_r = \int_0^1 t^r \psi(t) dt$.

- a.** Show that $\hat{\psi}$ has k vanishing moments if and only if a_0, \dots, a_{k-1} solves the equation

$$\begin{pmatrix} c_{0,0} & c_{0,1} & \cdots & c_{0,k-1} \\ c_{1,0} & c_{1,1} & \cdots & c_{1,k-1} \\ \vdots & \vdots & \vdots & \vdots \\ c_{k-1,0} & c_{k-1,1} & \cdots & c_{k-1,k-1} \end{pmatrix} \begin{pmatrix} a_0 \\ \hat{E}a_1 \\ \vdots \\ a_{k-1} \end{pmatrix} = \begin{pmatrix} e_0 \\ e_1 \\ \vdots \\ e_{k-1} \end{pmatrix} \quad (5.40)$$

- b.** Write a function

```
function a=vanishingmomshaar(k)
```

which solves Equation 5.40, and returns a_0, a_1, \dots, a_{k-1} in the vector \mathbf{a} .

5.6 Multiresolution analysis: A generalization

Let us summarize the properties of the spaces V_m . In both our examples we showed that they were nested, i.e.

$$V_0 \subset V_1 \subset V_2 \subset \cdots \subset V_m \subset \cdots$$

We also showed that continuous functions could be approximated arbitrarily well from V_m , as long as m was chosen large enough. Moreover, the space V_0 is closed under all translates, at least if we view the functions in V_0 as periodic with period N . In the following we will always identify a function with this periodic extension, just as we did in Fourier analysis. When performing this identification, we also saw that $f(t) \in V_m$ if and only if $g(t) = f(2t) \in V_{m+1}$. We have therefore shown that the scaling functions we have considered fit into the following general framework.

Definition 5.39 (Multiresolution analysis). A Multiresolution analysis, or MRA, is a nested sequence of function spaces

$$V_0 \subset V_1 \subset V_2 \subset \cdots \subset V_m \subset \cdots, \quad (5.41)$$

called resolution spaces, so that

1. Any function can be approximated arbitrarily well from V_m , as long as m is large enough,
2. $f(t) \in V_0$ if and only if $f(2^m t) \in V_m$,
3. $f(t) \in V_0$ if and only if $f(t - n) \in V_0$ for all n .
4. There is a function ϕ , called a scaling function, so that $\phi = \{\phi(t - n)\}_{0 \leq n < N}$ is a basis for V_0 .

When ϕ is an orthonormal basis we say that the MRA is *orthonormal*.

The wavelet of piecewise constant functions was an orthonormal MRA, while the wavelets for piecewise linear functions were not. Although the definition above states that any function can be approximated with MRA's, in practice one needs to restrict to certain functions: Certain pathological functions may be difficult to approximate. In the literature one typically requires that the function is in $L^2(\mathbb{R})$, and also that the scaling function and the spaces V_m are in $L^2(\mathbb{R})$. MRA's are much used, and one can find a wide variety of functions ϕ , not only piecewise constant functions, which give rise to MRA's.

In the examples we have considered we also chose a mother wavelet. The term wavelet is used in very general terms. However, the term mother wavelet is quite concrete, and is what gives rise to the theory of wavelets. This was necessary in order to efficiently decompose the $g_m \in V_m$ into a low resolution approximation $g_{m-1} \in V_{m-1}$, and a detail/error e_{m-1} in a detail space we called W_{m-1} . We have freedom in how we define these detail spaces, as well as how we define a mother wavelet whose translates span the detail space (in general we choose a mother wavelet which simplifies the computation of the decomposition $g_m = g_{m-1} + e_{m-1}$, but we will see later that it also is desirable to choose a ψ with other properties). Once we agree on the detail spaces and the mother wavelet, we can perform a change of coordinates to find detail and low resolution approximations. We thus have the following general recipe.

Idea 5.40 (Recipe for constructing wavelets). In order to construct MRA's which are useful for practical purposes, we need to do the following:

1. Find a function ϕ which can serve as the scaling function for an MRA,
2. Find a function ψ so that $\boldsymbol{\psi} = \{\psi(t - n)\}_{0 \leq n < N}$ and $\boldsymbol{\phi} = \{\phi(t - n)\}_{0 \leq n < N}$ together form an orthonormal basis for V_1 . The function ψ is also called a mother wavelet.

With V_0 the space spanned by $\boldsymbol{\phi} = \{\phi(t - n)\}_{0 \leq n < N}$, and W_0 the space spanned by $\boldsymbol{\psi} = \{\psi(t - n)\}_{0 \leq n < N}$, ϕ and ψ should be chosen so that we easily can compute the decomposition of $g_1 \in V_1$ into $g_0 + e_0$, where $g_0 \in V_0$ and $e_0 \in W_0$. If we can achieve this, the Discrete Wavelet Transform is defined as the change of coordinates from $\boldsymbol{\phi}_1$ to $(\boldsymbol{\phi}_0, \boldsymbol{\psi}_0)$.

More generally, if

$$f(t) = \sum_n c_{m,n} \phi_{m,n} = \sum_n c_{0,n} \phi_{0,n} + \sum_{m' < m, n} w_{m',n} \psi_{m',n},$$

then the m -level DWT is defined by $\text{DWT}(\mathbf{c}_m) = (\mathbf{c}_0, \mathbf{w}_0, \dots, \mathbf{w}_{m-1})$. It is useful to interpret m as frequency, n as time, and $w_{m,n}$ as the contribution at frequency m and time n . In this sense, wavelets provide a *time-frequency representation* of signals. This is what can make them more useful than Fourier analysis, which only provides frequency representations.

While there are in general many possible choices of detail spaces, in the case of an orthonormal wavelet we saw that it was natural to choose the detail space W_{m-1}

as the orthogonal complement of V_{m-1} in V_m , and obtain the mother wavelet by projecting the scaling function onto the detail space. Thus, for orthonormal MRA's, the low-resolution approximation and the detail can be obtained by computing projections, and the least squares approximation of f from V_m can be computed as

$$\text{proj}_{V_m}(f) = \sum_n \langle f, \phi_{m,n} \rangle \phi_{m,n}(t).$$

5.6.1 Working with the samples of f instead of f

In the MRA-setting it helps to think about the continuous-time function $f(t)$ as the model for an image, which is the object under study. f itself may not be in any V_m , however (this corresponds to that detail is present in the image for infinitely many m), and increasing m corresponds to that we also include the detail we see when we zoom in on the image. But how can we obtain useful approximations to f from V_m ? In case of an orthonormal MRA we can compute the least squares approximation as above, but we then need to compute the integrals $\langle f, \phi_{m,n} \rangle$, so that all function values are needed. However, as before we have only access to some samples $f(2^{-m}n)$, $0 \leq n < 2^m N$. These are called pixel values in the context of images, so that we can only hope to obtain a good approximation to $f^{(m)}$ (and thus f) from the pixel values. The following result explains how we can obtain this.

Theorem 5.41. If f is continuous, and ϕ has compact support, we have that, for all t ,

$$f(t) = \lim_{m \rightarrow \infty} \sum_{n=0}^{2^m N - 1} \frac{2^{-m}}{\int_0^N \phi_{m,0}(t) dt} f(n/2^m) \phi_{m,n}(t).$$

Proof: We have that

$$2^{-m} \sum_{n=0}^{2^m N - 1} \phi_{m,n} = \sum_{n=0}^{2^m N - 1} 2^{-m} \phi_{m,0}(t - 2^{-m}n).$$

We recognize this as a Riemann sum for the integral $\int_0^N \phi_{m,0}(t) dt$. Therefore, $\lim_{m \rightarrow \infty} \sum_{n=0}^{2^m N - 1} 2^{-m} \phi_{m,n} = \int_0^N \phi_{m,0}(t) dt$. Also, finitely many n contribute in this sum since ϕ has compact support. We now get that

$$\begin{aligned} \sum_{n=0}^{2^m N - 1} 2^{-m} f(n/2^m) \phi_{m,n}(t) &= \sum_{\substack{n \text{ so that } 2^{-m}n \approx t}} 2^{-m} f(n/2^m) \phi_{m,n}(t) \\ &\approx \sum_{\substack{n \text{ so that } 2^{-m}n \approx t}} 2^{-m} f(t) \phi_{m,n}(t) \\ &= f(t) \sum_{\substack{n \text{ so that } 2^{-m}n \approx t}} 2^{-m} \phi_{m,n}(t) \approx f(t) \int_0^N \phi_{m,0}(t) dt. \end{aligned}$$

where we have used the continuity of f and that $\lim_{m \rightarrow \infty} \sum_{n=0}^{2^m N - 1} 2^{-m} \phi_{m,n} = \int_0^N \phi_{m,0}(t) dt$. The result follows. Note that here we have not used the requirement that $\{\phi(t-n)\}_n$ are orthogonal. ■

The coordinate vector $\mathbf{x} = \left(\frac{2^{-m}}{\int_0^N \phi_{m,0}(t) dt} f(n/2^m) \right)_{n=0}^{2^m N - 1}$ in ϕ_m is therefore a candidate to an approximation of both f and $f^{(m)}$ from V_m , using only the pixel values. Normally one drops the leading constant $\frac{2^{-m}}{\int_0^N \phi_{m,0}(t) dt}$, so that one simply considers the sample values $f(n/2^m)$ as a coordinate vector in ϕ_m . This is used as the input to the DWT.

5.6.2 Increasing the precision of the DWT

Even though the samples of f give a good approximation to f as above, the approximation and f are still different, so that we obtain different output from the DWT. In Section 7.1 we will argue that the output from the DWT is equivalent to sampling the output from certain analog filters. We would like the difference in the output from these analog filters to be as small as possible. If the functions ϕ, ψ are symmetric around 0, we will also see that the analog filters are symmetric (a filter is symmetric if and only if the convolution kernel is symmetric around 0), in which case we know that such a high precision implementation is possible using the simple technique of symmetric extension. Let us summarize this as the following idea.

Idea 5.42. If the functions ϕ, ψ in a wavelet are symmetric around 0, then we can obtain an implementation of the DWT with higher precision when we consider symmetric extensions of the input.

Unfortunately, the piecewise constant scaling function we encountered was not symmetric. However, the piecewise linear scaling function was, and so are also many other interesting scaling functions we will encounter later. For a symmetric function, denote as before the symmetric extension of the input f with \check{f} . If the input \mathbf{x} to the DWT are the samples $(f(n/2^m))_{n=0}^{2^m N - 1}$, we create a vector $\check{\mathbf{x}}$ representing the samples of \check{f} . It is clear that this vector should be

$$\check{\mathbf{x}} = \left((f(n/2^m))_{n=0}^{2^m N - 1}, \lim_{t \rightarrow N_-} f(t), (f((2^m N - n)/2^m))_{n=1}^{2^m N - 1} \right).$$

In this vector there is symmetry around entry $2^m N$, so that the vector is determined from the $2^m N + 1$ first elements. Also the boundary is not duplicated, contrary to the previous symmetric extension given by Definition 4.1. We are thus lead to define a symmetric extension in the following way instead:

Definition 5.43 (Symmetric extension of a vector). By the *symmetric extension* of $\mathbf{x} \in \mathbb{R}^N$, we mean $\check{\mathbf{x}} \in \mathbb{R}^{2N-2}$ defined by

$$\check{\mathbf{x}}_k = \begin{cases} x_k & 0 \leq k < N \\ x_{2N-2-k} & N \leq k < 2N-3 \end{cases} \quad (5.42)$$

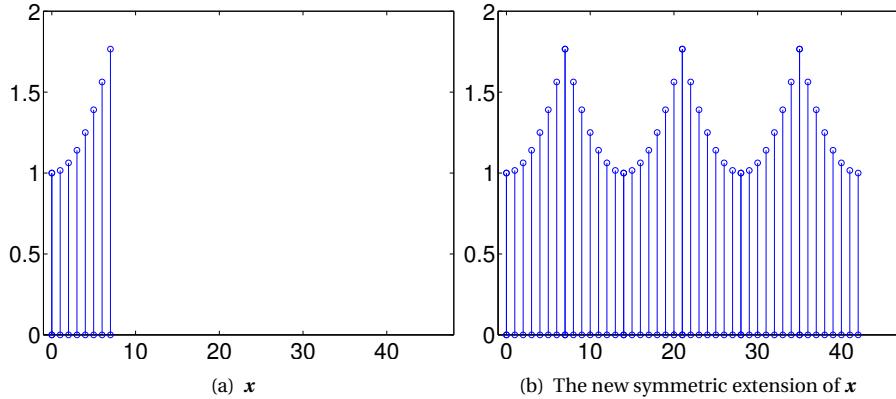


Figure 5.20: A vector and its symmetric extension. Note that the period of the vector is now $2N - 2$, while it was $2N$ for the vector shown in Figure 4.1.

With this notation, $N - 1$ is the symmetry point in all symmetric extensions. This is illustrated in Figure 5.20 From Chapter 4 it follows that symmetric filters preserve the symmetry around $N - 1$ when applied to such vectors. Applying a filter to this kind of symmetric extension in \mathbb{R}^{2N-2} can therefore be viewed as a mapping from \mathbb{R}^N to \mathbb{R}^N . If S is a symmetric filter, we will as before write S_r for this operation. We have the following analog to Theorem 4.8 for these new mappings, which is proved in very much the same way.

Theorem 5.44. With $S = \begin{pmatrix} S_1 & S_2 \\ S_3 & S_4 \end{pmatrix}$ a symmetric filter, we have that $S_r = S_1 + (0 \quad (S_2)^f \quad 0)$.

The proof is left as an exercise. With the Haar wavelet we succeeded in finding a function ψ which could be used in the recipe above. Note, however, that there may be many other ways to define a function ψ which can be used in the recipe. In the next chapter we will follow the recipe in order to construct other wavelets, and we will try to express which pairs of functions ϕ, ψ are most interesting, and which resolution spaces are most interesting.

Exercises for Section 5.6

1. Prove Theorem 5.44. Use the proof of Theorem 4.8 as a guide.

5.7 Summary

We started this chapter by motivating the theory of wavelets as a different function approximation scheme, which solved some of the shortcomings of Fourier series.

While one approximates functions with trigonometric functions in Fourier theory, with wavelets one instead approximates a function in several stages, where one at each stage attempts to capture information at a given resolution, using a function prototype. This prototype is localized in time, contrary to the Fourier basis functions, and this makes the theory of wavelets suitable for time-frequency representations of signals. We used an example based on Google Earth™ to illustrate that the wavelet-based scheme can represent an image at different resolutions in a scalable way, so that passing from one resolution to another simply amounts to adding some detail information to the lower resolution version of the image. This also made wavelets useful for compression, since the images at different resolutions can serve as compressed versions of the image.

We defined the simplest wavelet, the Haar wavelet, which is a function approximation scheme based on piecewise constant functions, and deduced its properties. We defined the Discrete Wavelet Transform (DWT) as a change of coordinates corresponding to the function spaces we defined. This transform is the crucial object to study when it comes to more general wavelets also, since it is the object which makes wavelets useful for computation. In the following chapters, we will see that reordering of the source and target bases of the DWT will aid in expressing connections between wavelets and filters, and in constructing optimized implementations of the DWT.

We then defined another wavelet, which corresponded to a function approximation scheme based on piecewise linear functions, instead of piecewise constant functions. There were several differences with the new wavelet when compared to the previous one. First of all, the basis functions were not orthonormal, and we did not attempt to make them orthonormal. The resolution spaces we now defined were not defined in terms of orthogonal bases, and we had some freedom on how we defined the detail spaces, since they are not defined as orthogonal complements anymore. Similarly, we had some freedom on how we define the mother wavelet, and we mentioned that we could define it so that it is more suitable for approximation of functions, by adding what we called vanishing moments.

From these examples of wavelets and their properties we made a generalization to what we called a multiresolution analysis (MRA). In an MRA we construct successively refined spaces of functions that may be used to approximate functions arbitrarily well. We will continue in the next chapter to construct even more general wavelets, within the MRA framework.

The book [21] goes through developments for wavelets in detail. While wavelets have been recognized for quite some time, it was with the important work of Daubechies [8, 9] that they found new arenas in the 80's. Since then they found important applications. The main application we will focus on in later chapters is image processing.

The filter representation of wavelets

Previously we have seen that analog filters restricted to the Fourier spaces gave rise to digital filters. These digital filters sent the samples of the input function to the samples of the output function, and they are easily implementable in contrast to the analog filters. We have also seen that wavelets give rise to analog filters. This leads us to believe that the DWT also can be implemented in terms of digital filters. In this chapter we will prove that this is in fact the case. The starting point for this is a reordering of the wavelet bases. There are some differences, however.

First of all, the DWT is not constructed by looking at the samples of a function, but rather by looking at coordinates in a given basis. Secondly, the function spaces we work in (i.e. V_m) are different from the Fourier spaces. Thirdly, we saw that the wavelet transform gave rise to two different types of analog filters: The filter defined by Equation (7.11) for obtaining $c_{m,n}$, and the filter defined by Equation (7.12) for obtaining $w_{m,n}$. We want both to correspond to digital filters. Due to these differences, the way we realize wavelet transformations in terms of digital filters is a bit different from before.

Despite the differences, this chapter will make it clear that the output of a DWT can be interpreted as the combined output of two different filters, and each filter will have an interpretation in terms of frequency representations. We will also see that the IDWT has a similar interpretation in terms of filters.

We will also define transforms which are useful generalizations of the filter interpretation of wavelets. In these many different filters are applied to the input, and the output is simply assembly of the results of these. It is fruitful to think about each filter as concentrating on a particular frequency range, and that these transforms thus simply splits the input into different frequency bands. Such transforms have important applications to the processing and compression of sound, and we will show that the much used MP3 standard for compression of sound takes use of such transforms.

6.1 The filters of a wavelet transformation

We will make the connection with digital filters by looking again at the different examples of wavelet bases we have seen: The one for piecewise constant functions, and the one for piecewise linear functions. We start by reordering the basis vectors in (ϕ_0, ψ_0) as

$$\mathcal{C}_1 = \{\phi_{0,0}, \psi_{0,0}, \phi_{0,1}, \psi_{0,1}, \dots, \phi_{0,N-1}, \psi_{0,N-1}\}. \quad (6.1)$$

The subscript 1 is used since \mathcal{C}_1 is a basis for V_1 . It turns out to be useful to reorder of the basis functions since it makes it easier to write down the change of coordinates matrices. To be more precise, let us first consider the Haar wavelet. From formula (5.15) it is apparent that $P_{\phi_1 \leftarrow \mathcal{C}_1}$ is the matrix where

$$\begin{pmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} \end{pmatrix}$$

is repeated along the main diagonal N times. Also, from formula (5.14) it is apparent that $P_{\mathcal{C}_1 \leftarrow \phi_1}$ is the same matrix. Such matrices are called *block diagonal matrices*. This particular block diagonal matrix is clearly orthogonal, since it transforms one orthonormal base to another.

For higher m we also reorder the basis vectors for (ϕ_{m-1}, ψ_{m-1}) as in Equation (6.1), i.e. we define

$$\mathcal{C}_m = \{\phi_{m-1,0}, \psi_{m-1,0}, \phi_{m-1,1}, \psi_{m-1,1}, \dots, \phi_{m-1,2^{m-1}N-1}, \psi_{m-1,2^{m-1}N-1}\}. \quad (6.2)$$

The bases ϕ_m and \mathcal{C}_m are both referred to as *wavelet bases*. Again, both change of coordinates matrices $P_{\phi_m \leftarrow \mathcal{C}_m}$, $P_{\mathcal{C}_m \leftarrow \phi_m}$ can be obtained by repeating the matrix from Equation (6.1) along the diagonal, but this time it is repeated $2^{m-1}N$ times.

For the piecewise linear wavelet, if we define the bases \mathcal{C}_m again by Equation (6.1) (i.e. with ϕ and ψ replaced). Equation (5.29) gives that the first two columns in $P_{\phi_1 \leftarrow \mathcal{C}_1}$ take the form

$$\frac{1}{\sqrt{2}} \begin{pmatrix} 1 & 0 \\ 1/2 & 1 \\ 0 & 0 \\ \vdots & \vdots \\ 0 & 0 \\ 1/2 & 0 \end{pmatrix}. \quad (6.3)$$

The remaining columns are obtained by shifting this, as in a circulant Toeplitz matrix. Similarly, Equation (5.31) gives that the first two columns in $P_{\mathcal{C}_1 \leftarrow \phi_1}$ take the form

$$\sqrt{2} \begin{pmatrix} 1 & 0 \\ -1/2 & 1 \\ 0 & 0 \\ \vdots & \vdots \\ 0 & 0 \\ -1/2 & 0 \end{pmatrix}. \quad (6.4)$$

Also here, the remaining columns are obtained by shifting this, as in a circulant Toeplitz matrix.

For the alternative piecewise linear wavelet, Equation (5.38) give all columns in the change of coordinate matrix $P_{\phi_1 \leftarrow \mathcal{C}_1}$ also, when the spaces ϕ_m, \mathcal{C}_m instead are defined in terms of the function $\hat{\psi}$, and ϕ . In particular, the first two columns in this matrix are

$$\frac{1}{\sqrt{2}} \begin{pmatrix} 1 & -1/4 \\ 1/2 & 3/4 \\ 0 & -1/4 \\ 0 & -1/8 \\ 0 & 0 \\ \vdots & \vdots \\ 0 & 0 \\ 1/2 & -1/8 \end{pmatrix}. \quad (6.5)$$

The first column is the same as before, since there was no change in the definition of ϕ . The remaining columns are obtained by shifting this, as in a circulant Toeplitz matrix. Similarly we could compute the change of coordinate matrix the opposite way, $P_{\mathcal{C}_1 \leftarrow \phi_1}$. We will explain shortly how this can be done.

In each case above it turned out that the change of coordinate matrices $P_{\phi_m \leftarrow \mathcal{C}_m}$, $P_{\mathcal{C}_m \leftarrow \phi_m}$ had a special structure: They were obtained by repeating the first two columns in a circulant way, similarly to how we did in a circulant Toeplitz matrix. The matrices were not exactly circulant Toeplitz matrices, however, since there are two different columns repeating. The change of coordinate matrices occurring in the stages in a DWT are thus not digital filters, but they seem to be related. Let us start by giving these new matrices names:

Definition 6.1 (MRA-matrices). An $N \times N$ -matrix T , with N even, is called an MRA-matrix if the columns are translates of the first two columns in alternating order, in the same way as the columns of a circulant Toeplitz matrix.

From our previous calculations it is clear that, once ϕ and ψ are given through an MRA, the corresponding change of coordinate matrices will always be MRA-matrices. The MRA-matrices is our connection between filters and wavelets. Let us make the following definition:

Definition 6.2. We denote by H_0 the (unique) filter with the same first row as $P_{\mathcal{C}_m \leftarrow \phi_m}$, and by H_1 the (unique) filter with the same second row as $P_{\mathcal{C}_m \leftarrow \phi_m}$. H_0 and H_1 are also called the *DWT filter components*.

Using this definition it is clear that

$$\begin{aligned} (P_{\mathcal{C}_m \leftarrow \phi_m} \mathbf{c}_m)_k &= (H_0 \mathbf{c}_m)_k && \text{when } k \text{ is even} \\ (P_{\mathcal{C}_m \leftarrow \phi_m} \mathbf{c}_m)_k &= (H_1 \mathbf{c}_m)_k && \text{when } k \text{ is odd} \end{aligned}$$

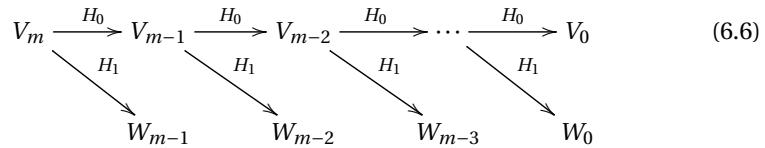
since the left hand side depends only on row k in the matrix $P_{\mathcal{C}_m \leftarrow \phi_m}$, and this is equal to row k in H_0 (when k is even) or row k in H_1 (when k is odd). This means that $P_{\mathcal{C}_m \leftarrow \phi_m} \mathbf{c}_m$ can be computed with the help of H_0 and H_1 as follows:

Theorem 6.3 (DWT expressed in terms of filters). Let \mathbf{c}_m be the coordinates in ϕ_m , and let H_0, H_1 be defined as above. Any stage in a DWT can be implemented in terms of filters as follows:

1. Compute $H_0 \mathbf{c}_m$. The even-indexed entries in the result are the coordinates \mathbf{c}_{m-1} in ϕ_{m-1} .
2. Compute $H_1 \mathbf{c}_m$. The odd-indexed entries in the result are the coordinates \mathbf{w}_{m-1} in ψ_{m-1} .

This gives an important connection between wavelets and filters: The DWT corresponds to applying two filters, H_0 and H_1 , and the result from the DWT is produced by assembling half of the coordinates from each. Keeping only every second coordinate is called *downsampling* (with a factor of two). Had we not performed downsampling, we would have ended up with twice as many coordinates as we started with. Downsampling with a factor of two means that we end up with the same number of samples as we started with. We also say that the output of the two filters is *critically sampled*. Due to the critical sampling, it is inefficient to compute the full application of the filters. We will return to the issue of making efficient implementations of critically sampled filter banks later.

We can now complement Figure 5.3 by giving names to the arrows as follows:



Let us make a similar analysis for the IDWT, and let us first make the following definition:

Definition 6.4. We denote by G_0 the (unique) filter with the same first column as $P_{\phi_m \leftarrow \mathcal{C}_m}$, and by G_1 the (unique) filter with the same second column as $P_{\phi_m \leftarrow \mathcal{C}_m}$. G_0 and G_1 are also called the *IDWT filter components*.

These filters are uniquely determined, since any filter is uniquely determined

from one of its columns. We can now write

$$\begin{aligned}
P_{\phi_m \leftarrow \mathcal{C}_m} & \begin{pmatrix} c_{m-1,0} \\ w_{m-1,0} \\ c_{m-1,1} \\ w_{m-1,1} \\ \dots \\ c_{m-1,2^{m-1}N-1} \\ w_{m-1,2^{m-1}N-1} \end{pmatrix} = P_{\phi_m \leftarrow \mathcal{C}_m} \begin{pmatrix} c_{m-1,0} \\ 0 \\ c_{m-1,1} \\ 0 \\ \dots \\ c_{m-1,2^{m-1}N-1} \\ 0 \end{pmatrix} + \begin{pmatrix} 0 \\ w_{m-1,0} \\ 0 \\ w_{m-1,1} \\ \dots \\ 0 \\ w_{m-1,2^{m-1}N-1} \end{pmatrix} \\
& = P_{\phi_m \leftarrow \mathcal{C}_m} \begin{pmatrix} c_{m-1,0} \\ 0 \\ c_{m-1,1} \\ 0 \\ \dots \\ c_{m-1,2^{m-1}N-1} \\ 0 \end{pmatrix} + P_{\phi_m \leftarrow \mathcal{C}_m} \begin{pmatrix} 0 \\ w_{m-1,0} \\ 0 \\ w_{m-1,1} \\ \dots \\ 0 \\ w_{m-1,2^{m-1}N-1} \end{pmatrix} \\
& = G_0 \begin{pmatrix} c_{m-1,0} \\ 0 \\ c_{m-1,1} \\ 0 \\ \dots \\ c_{m-1,2^{m-1}N-1} \\ 0 \end{pmatrix} + G_1 \begin{pmatrix} 0 \\ w_{m-1,0} \\ 0 \\ w_{m-1,1} \\ \dots \\ 0 \\ w_{m-1,2^{m-1}N-1} \end{pmatrix}.
\end{aligned}$$

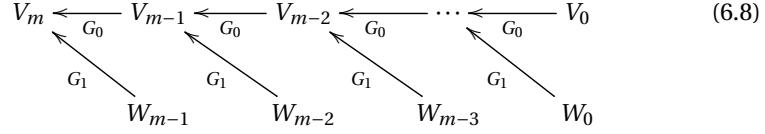
Here we have split a vector into its even-indexed and odd-indexed elements, which correspond to the coefficients from ϕ_{m-1} and ψ_{m-1} , respectively. In the last equation, we replaced with G_0, G_1 , since the multiplications with $P_{\phi_m \leftarrow \mathcal{C}_m}$ depend only on the even and odd columns in that matrix (due to the zeros inserted), and these columns are equal in G_0, G_1 . We can now state the following characterization of the inverse Discrete Wavelet transform:

Theorem 6.5 (IDWT expressed in terms of filters). Let G_0, G_1 be defined as above. Any stage in an IDWT can be implemented in terms of filters as follows:

$$\mathbf{c}_m = G_0 \begin{pmatrix} c_{m-1,0} \\ 0 \\ c_{m-1,1} \\ 0 \\ \dots \\ c_{m-1,2^{m-1}N-1} \\ 0 \end{pmatrix} + G_1 \begin{pmatrix} 0 \\ w_{m-1,0} \\ 0 \\ w_{m-1,1} \\ \dots \\ 0 \\ w_{m-1,2^{m-1}N-1} \end{pmatrix}. \quad (6.7)$$

Making a new vector where zeroes have been inserted in this way is also called *upsampling* (with a factor of two). We can now also complement Figure 5.3 for the

IDWT with named arrows as follows:



Note that the filters G_0, G_1 were defined in terms of the columns of $P_{\phi_m \leftarrow \mathcal{C}_m}$, while the filters H_0, H_1 were defined in terms of the rows of $P_{\mathcal{C}_m \leftarrow \phi_m}$. This difference is seen from the computations above to come from that the change of coordinates one way splits the coordinates into two parts, while the inverse change of coordinates performs the opposite. Let us summarize what we have found as follows.

Fact 6.6. The DWT can be computed with the help of two filters H_0, H_1 , as described by Theorem 6.3. The IDWT can be computed with the help of the two filters G_0, G_1 as described by Theorem 6.5. The filter coefficients in these four filters can be found from the relations between the bases ϕ_1 and (ϕ_0, ψ_0) .

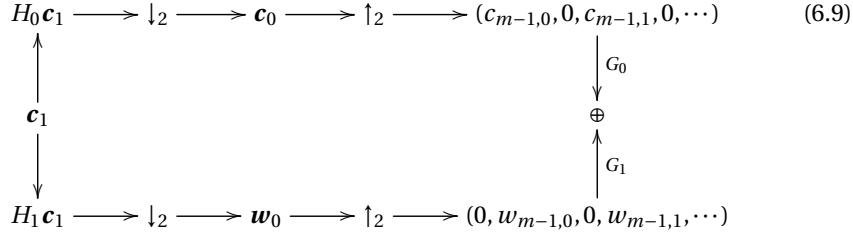
More generally, any transformation where the output is produced from the input using two filters H_0, H_1 as in a DWT, is called a *forward filter bank transform*. Any transformation which is defined from filters G_0, G_1 as in an IDWT is called a *reverse filter bank transform*. In particular, the filters H_0, H_1, G_0, G_1 may not stem from a wavelet setting, but from some signal processing setting, and in this setting it may be meaningful to allow for the reverse transform not to invert the forward transform exactly. It may also be meaningful that the reverse filter bank transform does not originate from change of coordinate matrices between certain function bases. In the setting of wavelets, however, the reverse filter bank transform is always associated with a change of coordinates between certain function bases, and the forward and reverse filter banks transforms are always inverses. This can be summed up as $GH = I$, and if we transpose this expression we obtain $H^T G^T = I$. Clearly H^T is a reverse filter bank transform with filters $(H_0)^T, (H_1)^T$, and G^T is a forward filter bank transform with filters $(G_0)^T, (G_1)^T$. These are called the dual filter bank transforms:

Definition 6.7 (Dual filter bank transforms). Assume that H_0, H_1 are the filters of a forward filter bank transform, and that G_0, G_1 are the filters of a reverse filter bank transform. By the *dual transforms* we mean the forward filter bank transform with filters $(G_0)^T, (G_1)^T$, and the reverse filter bank transform with filters $(H_0)^T, (H_1)^T$.

Note that, even though the reverse filter bank transform G can be associated with certain function bases, it is not clear if the reverse filter bank transform H^T also can be associated with such bases. We will see in the next chapter that such bases can in many cases be found. We will also denote these bases with the term dual.

Note that Figure 6.6 and 6.8 do not indicate the additional downsampling and upsampling steps described in Theorem 6.3 and 6.5. If we indicate downsampling

with \downarrow_2 , and upsampling with \uparrow_2 , the algorithms given in Theorem 6.3 and 6.5 can also be summarized as follows:



Her \oplus represents summing the elements which point inwards to the plus sign. In this figure, the left side represents the DWT, the right side the IDWT. In the literature, wavelet transforms are more often illustrated in this way using filters, since it makes all steps involved in the process more clear. This type of figure also opens for generalization. We will shortly look into this.

There are several reasons why it is smart to express a wavelet transformation in terms of filters. First of all, it enables us to reuse theoretical results from the world of filters in the world of wavelets, and to give useful interpretations of the wavelet transform in terms of frequencies. Secondly, and perhaps most important, it enables us to reuse efficient implementations of filters in order to compute wavelet transformations. A lot of work has been done in order to establish efficient implementations of filters, due to their importance.

In Example 5.21 we argued that the elements in V_{m-1} correspond to frequencies at lower frequencies than those in V_m , since $V_0 = \text{Span}(\phi_{0,n})$ should be interpreted as content of lower frequency than the $\phi_{1,n}$, with $W_0 = \text{Span}(\psi_{0,n})$ the remaining high frequency detail. To elaborate more on this, we have that

$$\phi(t) = \sum_{n=0}^{2N-1} (G_0)_{n,0} \phi_{1,n}(t) \quad (6.10)$$

$$\psi(t) = \sum_{n=0}^{2N-1} (G_1)_{n,1} \phi_{1,n}(t), \quad (6.11)$$

where $(G_k)_{i,j}$ are the entries in the matrix G_k . Similar equations are true for $\phi(t-k), \psi(t-k)$. Due to Equation (6.10), the filter G_0 should have lowpass characteristics, since it extracts the information at lower frequencies. Similarly, G_1 should have highpass characteristics due to Equation (6.11). Let us verify this for the different wavelets we have considered up to now by computing the frequency responses for the filters of these wavelets. We start with the Haar wavelet.

Example 6.8. For the Haar wavelet we saw that, in $P_{\phi_m \leftarrow \mathcal{C}_m}$, the matrix

$$\begin{pmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} \end{pmatrix} \quad (6.12)$$

repeated along the diagonal. The filters G_0 and G_1 can be found directly from these

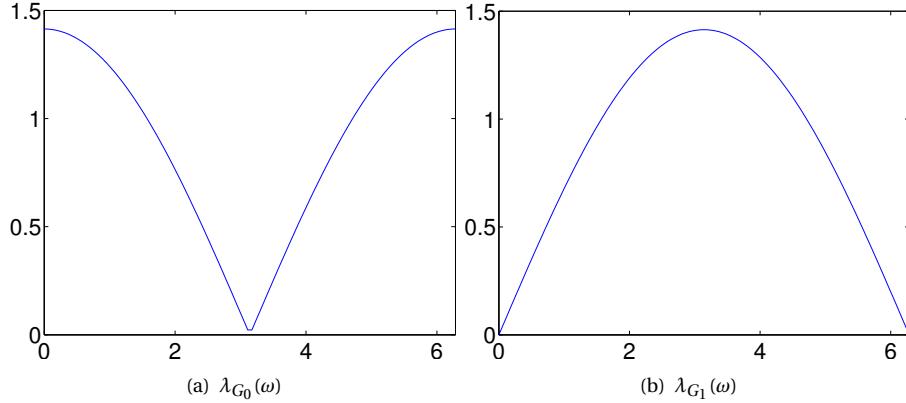


Figure 6.1: The frequency responses for the MRA of piecewise constant functions.

columns:

$$\begin{aligned} G_0 &= \{1/\sqrt{2}, 1/\sqrt{2}\} \\ G_1 &= \{1/\sqrt{2}, -1/\sqrt{2}\}. \end{aligned}$$

We have seen these filters previously: G_0 is a moving average filter of two elements (up to multiplication with a constant). This is a lowpass filter. G_1 is a bass-reducing filter, which is a highpass filter. Up to a constant, this is also an approximation to the derivative. Since G_1 is constructed from G_0 by adding an alternating sign to the filter coefficients, we know from before that G_1 is the highpass filter corresponding to the lowpass filter G_0 , so that the frequency response of the second is given by a shift of frequency with π in the first. The frequency responses are

$$\begin{aligned} \lambda_{G_0}(\omega) &= \frac{1}{\sqrt{2}} + \frac{1}{\sqrt{2}}e^{-i\omega} = \sqrt{2}e^{-i\omega/2} \cos(\omega/2) \\ \lambda_{G_1}(\omega) &= \frac{1}{\sqrt{2}}e^{i\omega} - \frac{1}{\sqrt{2}} = \sqrt{2}ie^{i\omega/2} \sin(\omega/2). \end{aligned}$$

The magnitude of these are plotted in Figure 6.1, where the lowpass/highpass characteristics are clearly seen. By considering the filters where the rows in Equation (6.12), it is clear that

$$\begin{aligned} H_0 &= \{1/\sqrt{2}, 1/\sqrt{2}\} \\ H_1 &= \{-1/\sqrt{2}, 1/\sqrt{2}\}, \end{aligned}$$

so that the frequency responses for the DWT have the same lowpass/highpass characteristics. ♣

It turns out that this connection between G_0 and G_1 as lowpass and highpass filters corresponding to each other can be found in all orthonormal wavelets. We will prove this in the next chapter.

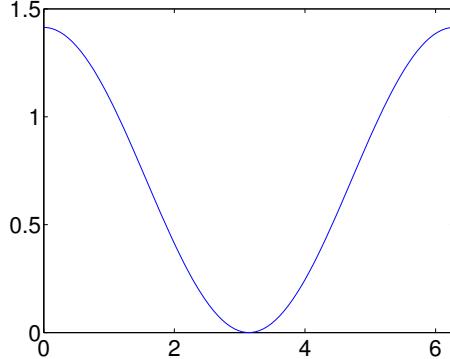


Figure 6.2: The frequency response $\lambda_{G_0}(\omega)$ for the first choice of wavelet for piecewise linear functions

Example 6.9. For the first wavelet for piecewise linear functions we looked at in the previous section, Equation (6.3) gives that

$$\begin{aligned} G_0 &= \frac{1}{\sqrt{2}}\{1/2, \underline{1}, 1/2\} \\ G_1 &= \frac{1}{\sqrt{2}}\{\underline{1}\}. \end{aligned} \quad (6.13)$$

G_0 is again a filter we have seen before: Up to multiplication with a constant, it is the treble-reducing filter with values from row 2 of Pascal's triangle. We see something different here when compared to the Haar wavelet, in that the filter G_1 is not the highpass filter corresponding to G_0 . The frequency responses are now

$$\begin{aligned} \lambda_{G_0}(\omega) &= \frac{1}{2\sqrt{2}}e^{i\omega} + \frac{1}{\sqrt{2}} + \frac{1}{2\sqrt{2}}e^{-i\omega} = \frac{1}{\sqrt{2}}(\cos\omega + 1) \\ \lambda_{G_1}(\omega) &= \frac{1}{\sqrt{2}}. \end{aligned}$$

$\lambda_{G_1}(\omega)$ thus has magnitude $\frac{1}{\sqrt{2}}$ at all points. The magnitude of $\lambda_{G_0}(\omega)$ is plotted in Figure 6.2. Comparing with Figure 6.1 we see that here also the frequency response has a zero at π . The frequency response seems also to be flatter around π . For the DWT, Equation (6.4) gives us

$$\begin{aligned} H_0 &= \sqrt{2}\{\underline{1}\} \\ H_1 &= \sqrt{2}\{-1/2, \underline{1}, -1/2\}. \end{aligned} \quad (6.14)$$

Even though G_1 was not the highpass filter corresponding to G_0 , we see that, up to a constant, H_1 is (it is a bass-reducing filter with values taken from row 2 of Pascals triangle). ♣

Note that the role of H_1 as the highpass filter corresponding to G_0 is the case in both previous examples. We will prove in the next chapter that this is a much more general result which holds for all wavelets, not only for the orthonormal ones.

For the alternative wavelet for piecewise linear functions, we are only able to find expressions for the filters G_0, G_1 at this stage (these can be extracted from Equation (6.5)). In the next chapter we will learn a general technique of computing the transformations the opposite way from these, so this will be handled in the next chapter.

The scaling functions and mother wavelets we encounter will turn out to always be functions with compact support. An interesting consequence of equations (6.10) and (6.11) is that we can find the size of these supports from the number of filter coefficients in G_0 and G_1 :

Theorem 6.10. Assume that the filters G_0, G_1 have N_0, N_1 nonzero filter coefficients, respectively. Then the support size of ϕ is $N_0 - 1$, and the support size of ψ is $(N_0 + N_1)/2 - 1$.

Proof: Let q be the support size of ϕ . Then the functions $\phi_{1,n}$ all have support size $q/2$. On the right hand side of Equation (6.10) we thus add N_0 functions, all with support size $q/2$. These functions are translated with $1/2$ with respect to one another, so that the sum has support size $q/2 + (N_0 - 1)/2$. Comparing with the support of the left hand side we get the equation $q = q/2 + (N_0 - 1)/2$, so that $q = N_0 - 1$. Similarly, with Equation (6.11), the function on the right hand side has support size $q/2 + (N_1 - 1)/2 = (N_0 + N_1)/2 - 1$, which thus is the support of ψ . ■

Since both filters for the Haar wavelet have 2 filter coefficients, it follows from this that both ϕ and ψ have support size 1. Similarly, since G_0 has 3 filter coefficients in the piecewise linear wavelet, it follows that the scaling function for this wavelet has support size 2. This verifies what we already know. It is not too difficult to show that, when the filters in a wavelet are symmetric, the support of ϕ is centered around 0, and the support of ψ is centered around $1/2$.

What you should have learnt in this section

How one can find the filters of a wavelet transformation by considering its matrix and its inverse. Forward and reverse filter bank transforms. Plot of the frequency responses for the filters of the wavelets we have considered.

Exercises for Section 6.1

1. Write down the corresponding filters G_0 og G_1 for Exercise 5.5.3. Plot their frequency responses, and characterize the filters as lowpass- or highpass filters.
2. Find two symmetric filters, so that the corresponding MRA-matrix, constructed with alternating rows from these two filters, is not a symmetric matrix.
3. Assume that an MRA-matrix is symmetric. Are the corresponding filters H_0, H_1, G_0, G_1 also symmetric? If not, find a counterexample.

6.2 Properties of the filter bank transforms of a wavelet

We have now described the DWT/IDWT as linear transformations G, H so that $GH = I$, and where two filters G_0, G_1 characterize G , two filters H_0, H_1 characterize H . G and H are not Toeplitz matrices, however, so they are not filters. Since filters produce the same output frequency from an input frequency, we must have that G and H produce other (undesired) frequencies in the output than those that are present in the input. We will call this phenomenon *aliasing*. In order for $GH = I$, the undesired frequencies must cancel each other, so that we end up with what we started with. Thus, GH must have what we will refer to as *alias cancellation*. This is the same as saying that GH is a filter. In order for $GH = I$, alias cancellation is not enough: We also need that the amount at the given frequency is unchanged, i.e. that $GH\phi_n = \phi_n$ for any Fourier basis vector ϕ_n . We then say that we have *perfect reconstruction*. Perfect reconstruction is always the case for wavelets by construction, but in signal processing many interesting examples (G_0, G_1, H_0, H_1) exist, for which we do not have perfect reconstruction. Historically, forward and reverse filter bank transforms have been around long before they appeared in a wavelet context. Operations where $GH\phi_n = c\phi_n$ for all n may also be useful, in particular when c_n is close to 1 for all n . If c_n is real for all n , we say that we have *no phase distortion*. If we have no phase distortion, the output from GH has the same phase, even if we do not have perfect reconstruction. Such “near-perfect reconstruction systems” have also been around long before many perfect reconstruction wavelet systems were designed. In signal processing, these transforms also exist in more general variants, and we will define these later. Let us summarize as follows.

Definition 6.11. (Alias cancellation, phase distortion, and perfect reconstruction). We say that we have *alias cancellation* if, for any n ,

$$GH\phi_n = c_n\phi_n,$$

for some constant c_n (i.e. GH is a filter). If all c_n are real, we say that we *no phase distortion*. If $GH = I$ (i.e. $c_n = 1$ for all n) we say that we have *perfect reconstruction*. If all c_n are close to 1, we say that we have *near-perfect reconstruction*

In signal processing, one also says that we have perfect- or near-perfect reconstruction when GH equals E_d , or is close to E_d (i.e. the overall result is a delay). The reason why a delay occurs has to do with that the transforms are used in real-time processing, for which we may not be able to compute the output at a given time instance before we know some of the following samples. Clearly the delay is unproblematic, since one can still reconstruct the input from the output. We will encounter a useful example of near-perfect reconstruction soon in the MP3 standard.

Let us now find a criterium for alias cancellation: When do we have that $GHe^{2\pi i rk/N}$

is a multiplum of $e^{2\pi i r k / N}$, for any r ? We first remark that

$$H(e^{2\pi i r k / N}) = \begin{cases} \lambda_{H_0,r} e^{2\pi i r k / N} & k \text{ even} \\ \lambda_{H_1,r} e^{2\pi i r k / N} & k \text{ odd.} \end{cases}$$

The frequency response of $H(e^{2\pi i r k / N})$ is

$$\begin{aligned} & \sum_{k=0}^{N/2-1} \lambda_{H_0,r} e^{2\pi i r (2k) / N} e^{-2\pi i (2k) n / N} + \sum_{k=0}^{N/2-1} \lambda_{H_1,r} e^{2\pi i r (2k+1) / N} e^{-2\pi i (2k+1) n / N} \\ &= \sum_{k=0}^{N/2-1} \lambda_{H_0,r} e^{2\pi i r (r-n)(2k) / N} + \sum_{k=0}^{N/2-1} \lambda_{H_1,r} e^{2\pi i r (r-n)(2k+1) / N} \\ &= (\lambda_{H_0,r} + \lambda_{H_1,r} e^{2\pi i r (r-n) / N}) \sum_{k=0}^{N/2-1} e^{2\pi i r (r-n) k / (N/2)}. \end{aligned}$$

Clearly, $\sum_{k=0}^{N/2-1} e^{2\pi i r (r-n) k / (N/2)} = N/2$ if $n = r$ or $n = r + N/2$, and 0 else. The frequency response is thus the vector

$$\frac{N}{2}(\lambda_{H_0,r} + \lambda_{H_1,r}) \mathbf{e}_r + \frac{N}{2}(\lambda_{H_0,r} - \lambda_{H_1,r}) \mathbf{e}_{r+N/2},$$

so that

$$H(e^{2\pi i r k / N}) = \frac{1}{2}(\lambda_{H_0,r} + \lambda_{H_1,r}) e^{2\pi i r k / N} + \frac{1}{2}(\lambda_{H_0,r} - \lambda_{H_1,r}) e^{2\pi i (r+N/2) k / N}. \quad (6.15)$$

Let us now turn to the reverse filter bank transform. We can write

$$\begin{aligned} (e^{2\pi i r \cdot 0 / N}, 0, e^{2\pi i r \cdot 2 / N}, 0, \dots, e^{2\pi i r (N-2) / N}, 0) &= \frac{1}{2}(e^{2\pi i r k / N} + e^{2\pi i (r+N/2) k / N}) \\ (0, e^{2\pi i r \cdot 1 / N}, 0, e^{2\pi i r \cdot 3 / N}, \dots, 0, e^{2\pi i r (N-1) / N}) &= \frac{1}{2}(e^{2\pi i r k / N} - e^{2\pi i (r+N/2) k / N}). \end{aligned}$$

This means that

$$\begin{aligned} G(e^{2\pi i r k / N}) &= G_0 \left(\frac{1}{2} (e^{2\pi i r k / N} + e^{2\pi i (r+N/2) k / N}) \right) + G_1 \left(\frac{1}{2} (e^{2\pi i r k / N} - e^{2\pi i (r+N/2) k / N}) \right) \\ &= \frac{1}{2}(\lambda_{G_0,r} e^{2\pi i r k / N} + \lambda_{G_0,r+N/2} e^{2\pi i (r+N/2) k / N}) + \frac{1}{2}(\lambda_{G_1,r} e^{2\pi i r k / N} - \lambda_{G_1,r+N/2} e^{2\pi i (r+N/2) k / N}) \\ &= \frac{1}{2}(\lambda_{G_0,r} + \lambda_{G_1,r}) e^{2\pi i r k / N} + \frac{1}{2}(\lambda_{G_0,r+N/2} - \lambda_{G_1,r+N/2}) e^{2\pi i (r+N/2) k / N}. \end{aligned} \quad (6.16)$$

Now, if we combine equations (6.15) and (6.16), we get

$$\begin{aligned}
& GH(e^{2\pi i r k/N}) \\
&= \frac{1}{2}(\lambda_{H_0,r} + \lambda_{H_1,r})G(e^{2\pi i r k/N}) + \frac{1}{2}(\lambda_{H_0,r} - \lambda_{H_1,r})G(e^{2\pi i(r+N/2)k/N}) \\
&= \frac{1}{2}(\lambda_{H_0,r} + \lambda_{H_1,r})\left(\frac{1}{2}(\lambda_{G_0,r} + \lambda_{G_1,r})e^{2\pi i r k/N} + \frac{1}{2}(\lambda_{G_0,r+N/2} - \lambda_{G_1,r+N/2})e^{2\pi i(r+N/2)k/N}\right) \\
&\quad + \frac{1}{2}(\lambda_{H_0,r} - \lambda_{H_1,r})\left(\frac{1}{2}(\lambda_{G_0,r+N/2} + \lambda_{G_1,r+N/2})e^{2\pi i(r+N/2)k/N} + \frac{1}{2}(\lambda_{G_0,r} - \lambda_{G_1,r})e^{2\pi i r k/N}\right) \\
&= \frac{1}{4}\left((\lambda_{H_0,r} + \lambda_{H_1,r})(\lambda_{G_0,r} + \lambda_{G_1,r}) + (\lambda_{H_0,r} - \lambda_{H_1,r})(\lambda_{G_0,r} - \lambda_{G_1,r})\right)e^{2\pi i r k/N} \\
&\quad + \frac{1}{4}\left((\lambda_{H_0,r} + \lambda_{H_1,r})(\lambda_{G_0,r+N/2} - \lambda_{G_1,r+N/2}) + (\lambda_{H_0,r} - \lambda_{H_1,r})(\lambda_{G_0,r+N/2} + \lambda_{G_1,r+N/2})\right)e^{2\pi i(r+N/2)k/N} \\
&= \frac{1}{2}(\lambda_{H_0,r}\lambda_{G_0,r} + \lambda_{H_1,r}\lambda_{G_1,r})e^{2\pi i r k/N} + \frac{1}{2}(\lambda_{H_0,r}\lambda_{G_0,r+N/2} - \lambda_{H_1,r}\lambda_{G_1,r+N/2})e^{2\pi i(r+N/2)k/N}.
\end{aligned}$$

If we also replace with the continuous frequency response, we obtain the following:

Theorem 6.12. We have that

$$\begin{aligned}
GH(e^{2\pi i r k/N}) &= \frac{1}{2}(\lambda_{H_0,r}\lambda_{G_0,r} + \lambda_{H_1,r}\lambda_{G_1,r})e^{2\pi i r k/N} \\
&\quad + \frac{1}{2}(\lambda_{H_0,r}\lambda_{G_0,r+N/2} - \lambda_{H_1,r}\lambda_{G_1,r+N/2})e^{2\pi i(r+N/2)k/N}. \quad (6.17)
\end{aligned}$$

In particular, we have alias cancellation if and only if

$$\lambda_{H_0}(\omega)\lambda_{G_0}(\omega + \pi) = \lambda_{H_1}(\omega)\lambda_{G_1}(\omega + \pi). \quad (6.18)$$

We will refer to this as the *alias cancellation condition*. If in addition

$$\lambda_{H_0}(\omega)\lambda_{G_0}(\omega) + \lambda_{H_1}(\omega)\lambda_{G_1}(\omega) = 2, \quad (6.19)$$

we also have perfect reconstruction. We will refer to as the *condition for perfect reconstruction*.

No phase distortion means that we have alias cancellation, and that

$$\lambda_{H_0}(\omega)\lambda_{G_0}(\omega) + \lambda_{H_1}(\omega)\lambda_{G_1}(\omega) \text{ is real.}$$

Now let us turn to how we can construct wavelets/perfect reconstruction systems from FIR-filters. We will have use for some theorems which allow us to construct wavelets from prototype filters. In particular we show that, when G_0 and H_0 are given lowpass filters which satisfy a certain common property, we can define unique (up to a constant) highpass filters H_1 and G_1 so that the collection of these four filters can be used to implement a wavelet. We first state the following general theorem.

Theorem 6.13. The following statements are equivalent for FIR filters H_0, H_1, G_0, G_1 :

1. H_0, H_1, G_0, G_1 give perfect reconstruction,
2. there exist $\alpha \in \mathbb{R}$ and $d \in \mathbb{Z}$ so that

$$(H_1)_n = (-1)^n \alpha^{-1} (G_0)_{n-2d} \quad (6.20)$$

$$(G_1)_n = (-1)^n \alpha (H_0)_{n+2d} \quad (6.21)$$

$$2 = \lambda_{H_0,n} \lambda_{G_0,n} + \lambda_{H_0,n+N/2} \lambda_{G_0,n+N/2} \quad (6.22)$$

Let us translate this to continuous frequency responses. We first have that

$$\begin{aligned} \lambda_{H_1}(\omega) &= \sum_k (H_1)_k e^{-ik\omega} = \sum_k (-1)^k \alpha^{-1} (G_0)_{k-2d} e^{-ik\omega} \\ &= \alpha^{-1} \sum_k (-1)^k (G_0)_k e^{-i(k+2d)\omega} = \alpha^{-1} e^{-2id\omega} \sum_k (G_0)_k e^{-ik(\omega+\pi)} \\ &= \alpha^{-1} e^{-2id\omega} \lambda_{G_0}(\omega + \pi). \end{aligned}$$

We have a similar computation for $\lambda_{G_1}(\omega)$. We can thus state the following:

Theorem 6.14. The following statements are equivalent for FIR filters H_0, H_1, G_0, G_1 :

1. H_0, H_1, G_0, G_1 give perfect reconstruction,
2. there exist $\alpha \in \mathbb{R}$ and $d \in \mathbb{Z}$ so that

$$\lambda_{H_1}(\omega) = \alpha^{-1} e^{-2id\omega} \lambda_{G_0}(\omega + \pi) \quad (6.23)$$

$$\lambda_{G_1}(\omega) = \alpha e^{2id\omega} \lambda_{H_0}(\omega + \pi) \quad (6.24)$$

$$2 = \lambda_{H_0}(\omega) \lambda_{G_0}(\omega) + \lambda_{H_0}(\omega + \pi) \lambda_{G_0}(\omega + \pi) \quad (6.25)$$

Proof: Let us prove first that equations (6.23)-(6.25) for a FIR filter implies that we have perfect reconstruction. Equations (6.23)-(6.24) mean that the alias cancellation condition (6.18) is satisfied, since

$$\begin{aligned} \lambda_{H_1}(\omega) \lambda_{G_1}(\omega + \pi) &= \alpha^{-1} e^{-2id\omega} \lambda_{G_0}(\omega + \pi) (\alpha) e^{2id(\omega+\pi)} \lambda_{H_0}(\omega) \\ &= \lambda_{H_0}(\omega) \lambda_{G_0}(\omega + \pi). \end{aligned}$$

Inserting this in the perfect reconstruction condition (6.25), we get

$$\begin{aligned} 2 &= \lambda_{H_0}(\omega) \lambda_{G_0}(\omega) + \lambda_{G_0}(\omega + \pi) \lambda_{H_0}(\omega + \pi) \\ &= \lambda_{H_0}(\omega) \lambda_{G_0}(\omega) + \alpha^{-1} e^{-2id\omega} \lambda_{G_0}(\omega + \pi) \alpha e^{2id\omega} \lambda_{H_0}(\omega + \pi) \\ &= \lambda_{H_0}(\omega) \lambda_{G_0}(\omega) + \lambda_{H_1}(\omega) \lambda_{G_1}(\omega), \end{aligned}$$

which is Equation (6.19), so that equations (6.23)- (6.25) imply perfect reconstruction. We therefore only need to prove that any set of FIR filters which give perfect reconstruction, also satisfy these equations. Due to the calculation above, it is enough to prove that equations (6.23)-(6.24) are satisfied. The proof of this will wait till section 8.4, since it uses some techniques we have not introduced yet. ■

Note that, even though conditions (6.23) and (6.24) together ensure that the alias cancellation condition is satisfied, alias cancellation can occur also if these conditions are not satisfied. Conditions (6.23) and (6.24) thus give a stronger requirement than alias cancellation. We will be particularly concerned with wavelets where the filters are symmetric, for which we can state the following corollary.

Corollary 6.15. The following statements are equivalent:

1. H_0, H_1, G_0, G_1 are the filters of a symmetric wavelet,
2. $\lambda_{H_0}(\omega), \lambda_{H_1}(\omega), \lambda_{G_0}(\omega), \lambda_{G_1}(\omega)$ are real functions, and

$$\lambda_{H_1}(\omega) = \alpha^{-1} \lambda_{G_0}(\omega + \pi) \quad (6.26)$$

$$\lambda_{G_1}(\omega) = \alpha \lambda_{H_0}(\omega + \pi) \quad (6.27)$$

$$2 = \lambda_{H_0}(\omega) \lambda_{G_0}(\omega) + \lambda_{H_0}(\omega + \pi) \lambda_{G_0}(\omega + \pi). \quad (6.28)$$

The delay d is thus 0 for symmetric wavelets.

Proof: Since H_0 is symmetric, $(H_0)_n = (H_0)_{-n}$, and from equations (6.20) and (6.21) it follows that

$$\begin{aligned} (G_1)_{n-2d} &= (-1)^{n-2d} \alpha (H_0)_n = (-1)^n \alpha^{-1} (H_0)_{-n} \\ &= (-1)^{(-n-2d)} \alpha^{-1} (H_0)_{(-n-2d)+2d} = (G_1)_{-n-2d} \end{aligned}$$

This shows that G_1 is symmetric about both $-2d$, in addition to being symmetric about 0 (by assumption). We must thus have that $d = 0$, so that $(H_1)_n = (-1)^n \alpha (G_0)_n$ and $(G_1)_n = (-1)^n \alpha^{-1} (H_0)_n$. We now get that

$$\begin{aligned} \lambda_{H_1}(\omega) &= \sum_k (H_1)_k e^{-ik\omega} = \alpha^{-1} \sum_k (-1)^k (G_0)_k e^{-ik\omega} \\ &= \alpha^{-1} \sum_k e^{-ik\pi} (G_0)_k e^{-ik\omega} = \alpha^{-1} \sum_k (G_0)_k e^{-ik(\omega+\pi)} \\ &= \alpha^{-1} \lambda_{G_0}(\omega + \pi), \end{aligned}$$

which proves Equation (6.26). Equation (6.26) follows similarly. ■

When constructing a wavelet it may be that we know one of the two pairs (G_0, G_1) , (H_0, H_1) , and that we would like to construct the other two. This can be achieved if we can find the constants d and α from above. If the filters are symmetric we just saw that $d = 0$. If G_0, G_1 are known, it follows from equations (6.20) and (6.21) that

$$1 = \sum_n (G_1)_n (H_1)_n = \sum_n (G_1)_n \alpha^{-1} (-1)^n (G_0)_n = \alpha^{-1} \sum_n (-1)^n (G_0)_n (G_1)_n,$$

so that $\alpha = \sum_n (-1)^n (G_0)_n (G_1)_n$. On the other hand, if H_0, H_1 are known instead, we must have that

$$1 = \sum_n (G_1)_n (H_1)_n = \sum_n \alpha (-1)^n (H_0)_n (H_1)_n = \alpha \sum_n (-1)^n (H_0)_n (H_1)_n,$$

so that $\alpha = 1 / (\sum_n (-1)^n (H_0)_n (H_1)_n)$. Let us use these observations to state the filters for the alternative wavelet of piecewise linear functions, which is the only wavelet we have gone through we have not computed the filters and the frequency response for.

Example 6.16. In Equation (6.5) we wrote down the first two columns in $P_{\phi_m \leftarrow \psi_m}$ for the alternative piecewise linear wavelet. This gives us that the filters G_0 and G_1 are

$$\begin{aligned} G_0 &= \frac{1}{\sqrt{2}} \{1/2, \underline{1}, 1/2\} \\ G_1 &= \frac{1}{\sqrt{2}} \{-1/8, -1/4, \underline{3/4}, -1/4, -1/8\}. \end{aligned} \quad (6.29)$$

Here G_0 was as for the wavelet of piecewise linear functions since we use the same scaling function. G_1 was changed, however. From Theorem 6.10, ψ should have support size 3, which can be seen to be right from Figure 5.17. Let us use Theorem 6.14 and the remark above to compute the two remaining filters H_0 and H_1 . These filters are also symmetric, since G_0, G_1 were. From the simple computation above we get that

$$\alpha = \sum_n (-1)^n (G_0)_n (G_1)_n = \frac{1}{2} \left(-\frac{1}{2} \left(-\frac{1}{4} \right) + 1 \cdot \frac{3}{4} - \frac{1}{2} \left(-\frac{1}{4} \right) \right) = \frac{1}{2}.$$

Theorem 6.14 now gives

$$\begin{aligned} (H_0)_n &= \alpha^{-1} (-1)^n (G_1)_n = 2 (-1)^n (G_1)_n \\ (H_1)_n &= \alpha^{-1} (-1)^n (G_0)_n = 2 (-1)^n (G_0)_n, \end{aligned} \quad (6.30)$$

so that

$$\begin{aligned} H_0 &= \sqrt{2} \{-1/8, 1/4, \underline{3/4}, 1/4, -1/8\} \\ H_1 &= \sqrt{2} \{-1/2, \underline{1}, -1/2\}. \end{aligned} \quad (6.31)$$

We now have that

$$\begin{aligned} \lambda_{G_1}(\omega) &= -1/(8\sqrt{2})e^{2i\omega} - 1/(4\sqrt{2})e^{i\omega} + 3/(4\sqrt{2}) - 1/(4\sqrt{2})e^{-i\omega} - 1/(8\sqrt{2})e^{-2i\omega} \\ &= -\frac{1}{4\sqrt{2}} \cos(2\omega) - \frac{1}{2\sqrt{2}} \cos\omega + \frac{3}{4\sqrt{2}}. \end{aligned}$$

The magnitude of $\lambda_{G_1}(\omega)$ is plotted in Figure 6.3. Clearly, G_1 now has highpass characteristics, while the lowpass characteristic of G_0 has been preserved. The filters G_0, G_1, H_0, H_1 are particularly important in applications: Apart from the scaling factors $1/\sqrt{2}, \sqrt{2}$ in front, we see that the filter coefficients are all dyadic fractions, i.e.

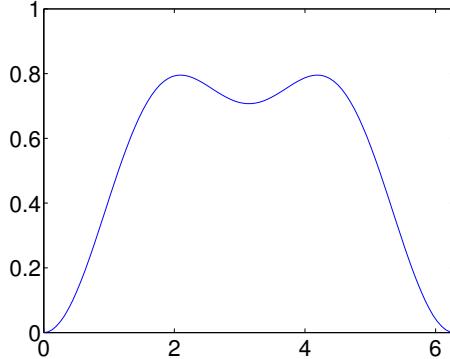


Figure 6.3: The frequency response $\lambda_{G_1}(\omega)$ for the alternative wavelet for piecewise linear functions.

they are on the form $\beta/2^j$. Arithmetic operations with dyadic fractions can be carried out exactly on a computer, due to representations as binary numbers in computers. These filters are thus important in applications, since they can be used as transformations for lossless coding. The same argument can be made for the Haar wavelet, but this wavelet had one less vanishing moment. ♣

In the literature, two particular cases of filter banks have been important. They are both referred to as *Quadrature Mirror Filter banks*, or QMF filter banks, and some confusion exists between the two. Let us therefore make precise definitions of the two.

Definition 6.17 (Classical QMF filter banks). In the classical definition of a QMF filter banks it is required that $G_0 = H_0$ and $G_1 = H_1$ (i.e. the filters in the forward and reverse transforms are equal), and that

$$\lambda_{H_1}(\omega) = \lambda_{H_0}(\omega + \pi). \quad (6.32)$$

It is straightforward to check that, for a classical QMF filter bank, the forward and reverse transforms are equal (i.e. $G = H$). It is easily checked that conditions (6.23) and (6.24) are satisfied with $\alpha = 1, d = 0$ for a classical QMF filter bank. In particular, the alias cancellation condition is satisfied. The perfect reconstruction condition can be written as

$$2 = \lambda_{H_0}(\omega)\lambda_{G_0}(\omega) + \lambda_{H_1}(\omega)\lambda_{G_1}(\omega) = \lambda_{H_0}(\omega)^2 + \lambda_{H_0}(\omega + \pi)^2. \quad (6.33)$$

Unfortunately, it is impossible to find non-trivial FIR-filters which satisfy this quadrature formula (Exercise 1). Therefore, classical QMF filter banks which give perfect reconstruction do not exist. Nevertheless, one can construct such filter banks which give close to perfect reconstruction [19], and this together with the fulfillment of the

alias cancellation condition still make them useful. In fact, we will see in Section 8.5 that the MP3 standard take use of such filters, and this explains our previous observation that the MP3 standard does not give perfect reconstruction. Note, however, that if the filters in a classical QMF filter bank are symmetric (so that $\lambda_{H_0}(\omega)$ is real), we have no phase distortion.

The second type of QMF filter bank is defined as follows.

Definition 6.18 (Alternative QMF filter banks). In the alternative definition of a QMF filter bank it is required that $G_0 = (H_0)^T$ and $G_1 = (H_1)^T$ (i.e. the filter coefficients in the forward and reverse transforms are reverse of one another), and that

$$\lambda_{H_1}(\omega) = \overline{\lambda_{H_0}(\omega + \pi)}. \quad (6.34)$$

The perfect reconstruction condition for an alternative QMF filter bank can be written as

$$\begin{aligned} 2 &= \lambda_{H_0}(\omega)\lambda_{G_0}(\omega) + \lambda_{H_1}(\omega)\lambda_{G_1}(\omega) = \lambda_{H_0}(\omega)\overline{\lambda_{H_0}(\omega)} + \overline{\lambda_{H_0}(\omega + \pi)}\lambda_{H_0}(\omega + \pi) \\ &= |\lambda_{H_0}(\omega)|^2 + |\lambda_{H_0}(\omega + \pi)|^2. \end{aligned}$$

We see that the perfect reconstruction property of the two definitions of QMF filter banks only differ in that the latter take absolute values. It turns out that the latter also has many interesting solutions, as we will see in Chapter 7. If we in condition (6.23) substitute $G_0 = (H_0)^T$ we get

$$\lambda_{H_1}(\omega) = \alpha^{-1} e^{-2id\omega} \lambda_{G_0}(\omega + \pi) = \alpha^{-1} e^{-2id\omega} \overline{\lambda_{H_0}(\omega + \pi)}.$$

If we set $\alpha = 1, d = 0$, we get equality here. A similar computation follows for Condition (6.24). In other words, also alternative QMF filter banks satisfy the alias cancellation condition. In the literature, a wavelet is called *orthonormal* if $G_0 = (H_0)^T, G_1 = (H_1)^T$. From our little computation it follows that alternative QMF filter banks with perfect reconstruction are examples of orthonormal wavelets, and correspond to orthonormal wavelets which satisfy $\alpha = 1, d = 0$.

For the Haar wavelet it is easily checked that $G_0 = (H_0)^T, G_1 = (H_1)^T$, but it does not satisfy the relation $\lambda_{H_1}(\omega) = \overline{\lambda_{H_0}(\omega + \pi)}$. Instead it satisfies the relation $\lambda_{H_1}(\omega) = -\overline{\lambda_{H_0}(\omega + \pi)}$. In other words, the Haar wavelet is not an alternative QMF filter bank—the way we have defined them. The difference lies only in a sign, however. This is the reason why the Haar wavelet is still listed as an alternative QMF filter bank in the literature. The additional sign leads to orthonormal wavelets which satisfy $\alpha = -1, d = 0$ instead.

The following is clear for orthonormal wavelets.

Theorem 6.19. A DWT is the transpose of an IDWT if and only if the filters satisfy $G_0 = (H_0)^T, G_1 = (H_1)^T$, i.e. if and only if the MRA equals the dual MRA.

This can be proved simply by observing that, if we transpose the DWT-matrix, Theorem 6.24 says that we get an IDWT matrix with filters $(H_0)^T, (H_1)^T$, and this is

equal to the IDWT if and only if $G_0 = (H_0)^T$, $G_1 = (H_1)^T$. It follows that QMF filter banks with perfect reconstruction give rise to orthonormal wavelets.

In practice we want to apply the wavelet transform to a symmetric extension, since then symmetric filters can give a better approximation to the underlying analog filters. In order to achieve this, the following result says that we only need to replace the filters H_0 , H_1 , G_0 , and G_1 in the wavelet transform with $(H_0)_r$, $(H_1)_r$, $(G_0)_r$, and $(G_1)_r$.

Theorem 6.20. If the filters H_0 , H_1 , G_0 , and G_1 in a wavelet transform are symmetric, then the DWT/IDWT preserve symmetric extensions (as defined in Definition 5.43). Also, applying the filters H_0 , H_1 , G_0 , and G_1 to $\tilde{\mathbf{x}} \in \mathbb{R}^{2N-2}$ in the DWT/IDWT is equivalent to applying $(H_0)_r$, $(H_1)_r$, $(G_0)_r$, and $(G_1)_r$ to $\mathbf{x} \in \mathbb{R}^N$ in the same way.

Proof: Since H_0 and H_1 are symmetric, their output from $\tilde{\mathbf{x}}$ is also a symmetric vector, and by assembling their outputs as the even- and odd-indexed entries, we see that the output $(c_0, w_0, c_1, w_1, \dots)$ of the MRA-matrix H also is a symmetric vector. The same then applies for the matrix G , since it inverts the first. This proves the first part.

Now, assume that $\mathbf{x} \in \mathbb{R}^N$. By definition of $(H_i)_r$, $(H_i \tilde{\mathbf{x}})_n = ((H_i)_r \mathbf{x})_n$ for $0 \leq n \leq N-1$. This means that we get the same first N output elements in a wavelet transform if we replace H_0, H_1 with $(H_0)_r, (H_1)_r$. Since the vectors $(c_0, 0, c_1, 0, \dots)$ and $(0, w_0, 0, w_1, \dots)$ also are symmetric vectors when $(c_0, w_0, c_1, w_1, \dots)$ is, it follows that $(G_0)_r, (G_1)_r$ will reproduce the same first N elements as G_0, G_1 also. In conclusion, for symmetric vectors, the wavelet transform restricted to the first N elements produces the same result when we replace H_0, H_1, G_0 , and G_1 with $(H_0)_r, (H_1)_r, (G_0)_r$, and $(G_1)_r$. This proves the result. ■

As in Chapter 4, it follows that when the filters of a wavelet are symmetric, applying $(H_0)_r, (H_1)_r, (G_0)_r$, and $(G_1)_r$ to the input better approximates an underlying analog filter.

Exercises for Section 6.2

1. Show that it is impossible to find a non-trivial FIR-filter which satisfies Equation (6.33).
2. Show that the Haar wavelet satisfies $\lambda_{H_1}(\omega) = -\overline{\lambda_{H_0}(\omega + \pi)}$, and $G_0 = (H_0)^T$, $G_1 = (H_1)^T$. The Haar wavelet can thus be considered as an alternative QMF filter bank.
3. In this exercise we will establish an orthonormal basis for the symmetric extensions, as defined by Definition 5.43. This parallels Theorem 4.6.
 - a. Explain why, if $\mathbf{x} \in \mathbb{R}^{2N-2}$ is a symmetric extension (according to definition 4.1), then $(\tilde{\mathbf{x}})_n = z_n e^{-\pi i n}$, where \mathbf{z} is a real vectors which satisfies $z_n = z_{2N-2-n}$

b. Show that

$$\left\{ \mathbf{e}_0, \left\{ \frac{1}{\sqrt{2}}(\mathbf{e}_i + \mathbf{e}_{2N-2-i}) \right\}_{n=1}^{N-2}, \mathbf{e}_{N-1} \right\} \quad (6.35)$$

is an orthonormal basis for the vectors on the form $\hat{\mathbf{x}}$ with $\mathbf{x} \in \mathbb{R}^{2N-2}$ a symmetric extension.

c. Show that

$$\begin{aligned} & \frac{1}{\sqrt{2N-2}} \cos\left(2\pi \frac{0}{2N-2} k\right) \\ & \left\{ \frac{1}{\sqrt{N-1}} \cos\left(2\pi \frac{n}{2N-2} k\right) \right\}_{n=1}^{N-2} \\ & \frac{1}{\sqrt{2N-2}} \cos\left(2\pi \frac{N-1}{2N-2} k\right) \end{aligned} \quad (6.36)$$

is an orthonormal basis for the symmetric extensions in \mathbb{R}^{2N-2} .

d. Assume that S is symmetric. Show that the vectors listed in (6.36) are eigenvectors for S_r , when the vectors are viewed as vectors in \mathbb{R}^N , and that they are linearly independent. This shows that S_r is diagonalizable.

4. Let us explain how the matrix S_r can be diagonalized, similarly to how we previously diagonalized using the DCT. In Exercise 3 we showed that the vectors

$$\left\{ \cos\left(2\pi \frac{n}{2N-2} k\right) \right\}_{n=0}^{N-1} \quad (6.37)$$

in \mathbb{R}^N is a basis of eigenvectors for S_r when S is symmetric. S_r itself is not symmetric, however, so that this basis can not possibly be orthogonal (S is symmetric if and only if it is orthogonally diagonalizable). However, when the vectors are viewed in \mathbb{R}^{2N-2} we showed in Exercise 3.c an orthogonality statement which can be written as

$$\sum_{k=0}^{2N-3} \cos\left(2\pi \frac{n_1}{2N-2} k\right) \cos\left(2\pi \frac{n_2}{2N-2} k\right) = (N-1) \times \begin{cases} 2 & \text{if } n_1 = n_2 \in \{0, N-1\} \\ 1 & \text{if } n_1 = n_2 \notin \{0, N-1\} \\ 0 & \text{if } n_1 \neq n_2 \end{cases}. \quad (6.38)$$

a. Show that

$$\begin{aligned} & (N-1) \times \begin{cases} 1 & \text{if } n_1 = n_2 \in \{0, N-1\} \\ \frac{1}{2} & \text{if } n_1 = n_2 \notin \{0, N-1\} \\ 0 & \text{if } n_1 \neq n_2 \end{cases} \\ & = \frac{1}{\sqrt{2}} \cos\left(2\pi \frac{n_1}{2N-2} \cdot 0\right) \frac{1}{\sqrt{2}} \cos\left(2\pi \frac{n_2}{2N-2} \cdot 0\right) \\ & + \sum_{k=1}^{N-2} \cos\left(2\pi \frac{n_1}{2N-2} k\right) \cos\left(2\pi \frac{n_2}{2N-2} k\right) \\ & + \frac{1}{\sqrt{2}} \cos\left(2\pi \frac{n_1}{2N-2} (N-1)\right) \frac{1}{\sqrt{2}} \cos\left(2\pi \frac{n_2}{2N-2} (N-1)\right). \end{aligned}$$

Hint: Use that $\cos x = \cos(2\pi - x)$ to pair the summands k and $2N-2-k$.

Now, define the vector $\mathbf{d}_n^{(I)}$ as

$$d_{n,N} \left(\frac{1}{\sqrt{2}} \cos \left(2\pi \frac{n}{2N-2} \cdot 0 \right), \left\{ \cos \left(2\pi \frac{n}{2N-2} k \right) \right\}_{k=1}^{N-2}, \frac{1}{\sqrt{2}} \cos \left(2\pi \frac{n}{2N-2} (N-1) \right) \right),$$

and define $d_{0,N}^{(I)} = d_{N-1,N}^{(I)} = 1/\sqrt{N-1}$, and $d_{n,N}^{(I)} = \sqrt{2/(N-1)}$ when $n > 1$. The orthogonal $N \times N$ matrix where the rows are $\mathbf{d}_n^{(I)}$ is called the DCT-I, and we will denote it by $D_N^{(I)}$. DCT-I is also much used, just as the DCT-II of Chapter 4. The main difference from the previous cosine vectors is that $2N$ has been replaced by $2N-2$.

- b.** Explain that the vectors $\mathbf{d}_n^{(I)}$ are orthonormal, and that the matrix

$$\sqrt{\frac{2}{N-1}} \begin{pmatrix} 1/\sqrt{2} & 0 & \cdots & 0 & 0 \\ 0 & 1 & \cdots & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & \cdots & 1 & 0 \\ 0 & 0 & \cdots & 0 & 1/\sqrt{2} \end{pmatrix} (\cos(2\pi \frac{n}{2N-2} k)) \begin{pmatrix} 1/\sqrt{2} & 0 & \cdots & 0 & 0 \\ 0 & 1 & \cdots & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & \cdots & 1 & 0 \\ 0 & 0 & \cdots & 0 & 1/\sqrt{2} \end{pmatrix}$$

is orthogonal.

- c.** Explain from b. that $(\cos(2\pi \frac{n}{2N-2} k))^{-1}$ can be written as

$$\frac{2}{N-1} \begin{pmatrix} 1/2 & 0 & \cdots & 0 & 0 \\ 0 & 1 & \cdots & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & \cdots & 1 & 0 \\ 0 & 0 & \cdots & 0 & 1/2 \end{pmatrix} (\cos(2\pi \frac{n}{2N-2} k)) \begin{pmatrix} 1/2 & 0 & \cdots & 0 & 0 \\ 0 & 1 & \cdots & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & \cdots & 1 & 0 \\ 0 & 0 & \cdots & 0 & 1/2 \end{pmatrix}$$

With the expression we found in c., S_r can now be diagonalized as

$$(\cos(2\pi \frac{n}{2N-2} k)) D (\cos(2\pi \frac{n}{2N-2} k))^{-1}.$$

- 5.** In Exercise 5.3.5 we computed the DWT of two very simple vectors \mathbf{x}_1 and \mathbf{x}_2 , using the Haar wavelet.

- a.** Compute $H_0 \mathbf{x}_1$, $H_1 \mathbf{x}_1$, $H_0 \mathbf{x}_2$, and $H_1 \mathbf{x}_2$, where H_0 and H_1 are the filters used by the Haar wavelet.
- b.** Compare the odd-indexed elements in $H_1 \mathbf{x}_1$ with the odd-indexed elements in $H_1 \mathbf{x}_2$. From this comparison, attempt to find an explanation to why the two vectors have very different detail components.
- 6.** Suppose that we run the following algorithm on the sound represented by the vector \mathbf{S} :

```

l=length(S);
c=(S(1:2:(l-1))+S(2:2:l))/sqrt(2);
w=(S(1:2:(l-1))-S(2:2:l))/sqrt(2);

newS=[c w];
newS=newS/max(abs(newS));
playerobj=audioplayer(newS,44100);
playblocking(playerobj)

```

- a.** Comment the code and explain what happens. Which wavelet is used? What do the vectors c and w represent? Describe the sound you believe you will hear.
- b.** Assume that we add lines in the code above which sets the elements in the vector w to 0 before we compute the inverse operation. What will you hear if you play the new sound you then get?

6.3 Filter-based algorithm for the DWT and the IDWT

When we apply wavelets in practice, they are often defined in terms of the filters H_0 , H_1 , G_0 , G_1 . We therefore need an implementation of the DWT and the IDWT which takes as input the filter coefficients of these filters. All wavelets we look at, except for the Haar wavelet (which we have already implemented), will have symmetric filters, so we will restrict our implementations to symmetric filters. You will be spared writing these implementations: you can assume that the function

```
xnew=DWTImpl(h0,h1,x,m)
```

takes as input symmetric filters h_0 and h_1 , a vector x , and m , and returns the result of the m -level DWT for us, i.e. it returns the coordinates of x in the basis $V_0 \oplus W_0 \oplus W_1 \oplus \dots \oplus W_{m-1}$. You can also assume that the function

```
x=IDWTImpl(g0,g1,xnew,m)
```

performs the m -level IDWT in a similar way, where g_0 and g_1 also are symmetric filters. You can assume that both these functions handle input of any size i.e. the length of the vector needs not be an appropriate power of 2 (as we saw that the code in the Haar wavelet transformation assumed).

Although you are spared writing these implementations, we will comment on some of the issues in writing them. Previously we mentioned that we could use Matlab's `conv` function to implement filters. Since the DWT and the IDWT now have been expressed in terms of filters, the `conv`-function can be used to implement the DWT and the IDWT as well. More precisely, they can be implemented by taking the convolution with the rows (since the entries in the first column equal the entries in the first row when the filter is symmetric) of the MRA-matrix. With the DWT we

have these rows, since the rows in $P_{\mathcal{C}_m \leftarrow \phi_m}$ are given by the rows of H_0, H_1 . With the IDWT we need a translation from the column representation of the MRA-matrix to the row representation of the matrix. This is a straightforward task, however a bit tedious. In the function `IDWTImpl` you will see some additional code at the beginning taking care of this.

Previously we have mentioned that we assume that vectors have lengths on the form $N2^m$, in order for us to apply m levels of the DWT to it. You can assume, however, that the functions `DWTImpl` and `IDWTImpl` also work for vectors of general lengths. Let us how we can generalize what we have done in order to address this. Instead of assuming that the functions in V_m have period N , let us assume that they have period $N2^{-m}$ (i.e. that the period is not an integer), and that the number N represents the dimension of V_m (instead of the period as before). This certainly affects functions in the spaces $V_{m'}$ for $m' < m$, since their supports may split in two near $t = T = N2^{-m}$. But the change of coordinates between the spaces with this new period are governed by the same formulas as before. This means that, even if the number of basis functions in ϕ_m is odd, the change of coordinates from $\phi_m = \{\phi_{m,0}, \phi_{m,1}, \dots, \phi_{m,N-1}\}$ to

$$\mathcal{C}_m = \{\phi_{m-1,0}, \psi_{m-1,0}, \dots, \phi_{m-1,(N-3)/2}, \psi_{m-1,(N-3)/2}, \phi_{m-1,(N-1)/2}\}$$

is given by the same repetition of the rows/columns in the MRA-matrix as before (the only difference is that the values in the first and last last columns are equal). We can therefore define the basis \mathcal{C}_m by including the desired number of functions (even if it is odd), and the MRA-matrix can be used for the change of coordinates as before.

Symmetric filters are very common in practice, since MRA-matrices with symmetric filters also preserve symmetric vectors. Due to this they share some of the desirable properties of symmetric filters (which lead us to the definition of the DCT). `DWTImpl` and `IDWTImpl` are written so that they are applied to the symmetric extension of the sound. Note also that, since these implementations call the `conv`-function, they compute all output elements of the filter. But Theorem 6.3 states that only half of the output entries for each filter need to be computed. The current implementation is therefore not optimal. In Chapter 8 we will describe a method which solves this issue, and which also provides further computational savings when compared to a direct implementation of the filters.

Finally, let us provide an example of how we write down the parameters h_0, h_1, g_0, g_1 to the implementations. They represent the nonzero filter coefficients of H_0, H_1, G_0, G_1 , respectively. If

$$H_0 = \{h_{0,-k_0}, \dots, h_{0,-1} \underline{h_{0,0}}, h_{0,1}, \dots, h_{0,k_0}\}$$

$$H_1 = \{h_{1,-k_1}, \dots, h_{1,-1} \underline{h_{1,0}}, h_{1,1}, \dots, h_{1,k_1}\},$$

then the vectors

$$h_0 = (h_{0,-k_0}, \dots, h_{0,-1}, h_{0,0}, h_{0,1}, \dots, h_{0,k_0})$$

$$h_1 = (h_{1,-k_1}, \dots, h_{1,-1}, h_{1,0}, h_{1,1}, \dots, h_{1,k_1})$$

should be input to the `DWTImpl` function.

Example 6.21. In Exercise 3 you will be asked to implement a function `playDWTfilterslower` which plays the low-resolution approximations of our audio test file, for any type of wavelet, using the functions we have described. With this function we can play the result for all the wavelets we have considered up to now, in succession, and at a given resolution, with the following code:

```
function playDWTall(m)
    disp('Haar wavelet');
    playDWTlower(m);
    disp('Wavelet for piecewise linear functions');
    playDWTfilterslower(m,sqrt(2)*[1],...
        sqrt(2)*[-1/2 1 -1/2],...
        [1/2 1 1/2]/sqrt(2),...
        [1]/sqrt(2));
    disp('Wavelet for piecewise linear functions, alternative version');
    playDWTfilterslower(m,...);
    sqrt(2)*[-1/8 1/4 3/4 1/4 -1/8],...
    sqrt(2)*[-1/2 1 -1/2],...
    [1/2 1 1/2]/sqrt(2),...
    [-1/8 -1/4 3/4 -1/4 -1/8]/sqrt(2));
```

The call to `playDWTlower` first plays the result, using the Haar wavelet. The code then moves on to the piecewise linear wavelet. From Equation (6.14) we first see that

$$h_0 = (h_{0,-k_0}, \dots, h_{0,-1}, h_{0,0}, h_{0,1}, \dots, h_{0,k_0}) = \sqrt{2}(1) \quad (6.39)$$

$$h_1 = (h_{1,-k_1}, \dots, h_{1,-1}, h_{1,0}, h_{1,1}, \dots, h_{1,k_1}) = \sqrt{2}(-1/2, 1, -1/2), \quad (6.40)$$

and from Equation (6.13) we see that

$$g_0 = (g_{0,0}, g_{0,1}, \dots, g_{0,l_0}) = \frac{1}{\sqrt{2}}(1/2, 1, 1/2) \quad (6.41)$$

$$g_1 = (g_{1,0}, g_{1,1}, \dots, g_{1,l_1}) = \frac{1}{\sqrt{2}}(1). \quad (6.42)$$

These explain the parameters to the call to `playDWTfilterslower` for the piecewise linear wavelet. The code then moves to the alternative piecewise linear wavelet. For this wavelet we have not computed all filter coefficients yet. This is delayed till Example 6.16, since we need some more techniques in order to compute these. From Equation (6.31) in that example we see that

$$h_0 = (h_{0,-k_0}, \dots, h_{0,-1}, h_{0,0}, h_{0,1}, \dots, h_{0,k_0}) = \sqrt{2}(-1/8, 1/4, 3/4, 1/4, -1/8)$$

$$h_1 = (h_{1,-k_1}, \dots, h_{1,-1}, h_{1,0}, h_{1,1}, \dots, h_{1,k_1}) = \sqrt{2}(-1/2, 1, -1/2),$$

and from Equation (6.29) we see that

$$g_0 = (g_{0,-l_0}, \dots, g_{0,-1}, g_{0,0}, g_{0,1}, \dots, g_{0,l_0}) = \frac{1}{\sqrt{2}}(1/2, 1, 1/2)$$

$$g_1 = (g_{1,-l_1}, \dots, g_{1,-1}, g_{1,0}, g_{1,1}, \dots, g_{1,l_1}) = \frac{1}{\sqrt{2}}(-1/8, -1/4, 3/4, -1/4, -1/8).$$

These explain the parameters to the call to `playDWTfilterslower` for the alternative piecewise linear wavelet. ♣

Exercises for Section 6.3

- 1.** In this exercise we will practice setting up the parameters `h0,h1,g0,g1` which are used in the calls to `DWTImpl` and `IDWTImpl`.

- a.** Assume that one stage in a DWT is given by the MRA-matrix

$$P_{\mathcal{C}_1 \leftarrow \phi_1} = \begin{pmatrix} 1/5 & 1/5 & 1/5 & 0 & 0 & 0 & \cdots & 0 & 1/5 & 1/5 \\ -1/3 & 1/3 & -1/3 & 0 & 0 & 0 & \cdots & 0 & 0 & 0 \\ 1/5 & 1/5 & 1/5 & 1/5 & 1/5 & 0 & \cdots & 0 & 0 & 0 \\ 0 & 0 & -1/3 & 1/3 & -1/3 & 0 & \cdots & 0 & 0 & 0 \\ \vdots & \vdots \end{pmatrix}$$

Write down the compact form for the corresponding filters H_0, H_1 , and compute and plot the frequency responses. Are the filters symmetric? If so, also write down the parameters `h0,h1` you would use for this matrix in a call to `DWTImpl`.

- b.** Assume that one stage in the IDWT is given by the MRA-matrix

$$P_{\phi_1 \leftarrow \mathcal{C}_1} = \begin{pmatrix} 1/2 & -1/4 & 0 & 0 & \cdots \\ 1/4 & 3/8 & 1/4 & 1/16 & \cdots \\ 0 & -1/4 & 1/2 & -1/4 & \cdots \\ 0 & 1/16 & 1/4 & 3/8 & \cdots \\ 0 & 0 & 0 & -1/4 & \cdots \\ 0 & 0 & 0 & 1/16 & \cdots \\ 0 & 0 & 0 & 0 & \cdots \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & 0 & 0 & \cdots \\ 1/4 & 1/16 & 0 & 0 & \cdots \end{pmatrix}$$

Write down the compact form for the filters G_0, G_1 , and compute and plot the frequency responses. Are the filters symmetric? If so, also write down the parameters `g0,g1` you would use for this matrix in a call to `IDWTImpl`.

- 2.** Let us also practice on writing down the change of coordinate matrices from the parameters `h0,h1,g0,g1`.

- a.** Assume that $\mathbf{h0}=[1/16 \ 1/4 \ 3/8 \ 1/4 \ 1/16]$ and $\mathbf{h1}=[-1/4 \ 1/2 \ -1/4]$. Write down the compact form for the filters H_0, H_1 . Plot the frequency responses and verify that H_0 is a lowpass filter, and that H_1 is a highpass filter. Also write down the change of coordinate matrix $P_{\mathcal{C}_1 \leftarrow \phi_1}$ for the wavelet corresponding to these filters.

- b.** Assume that $g0=[1/3 \ 1/3 \ 1/3]$ and $g1=[1/5 \ -1/5 \ 1/5 \ -1/5 \ 1/5]$. Write down the compact form for the filters G_0, G_1 . Plot the frequency responses and verify that G_0 is a lowpass filter, and that G_1 is a highpass filter. Also write down the change of coordinate matrix $P_{\phi_1 \leftarrow \mathcal{C}_1}$ for the wavelet corresponding to these filters.

3. Write a function

```
function playDWTfilterslower(m,h0,h1,g0,g1)
```

which reimplements the function `playDWTlower` from Exercise 5.3.10 so that it takes as input the coefficients of the four different filters as in Example 6.21. Listen to the result using the different wavelets we have encountered and for different m , using the code from Example 6.21. Can you hear any difference from the Haar wavelet? If so, which wavelet gives the best sound quality?

- 4.** In this exercise we will change the code in Example 6.21 so that it instead only plays the contribution from the detail spaces (i.e. $W_0 \oplus W_1 \oplus \dots \oplus W_{m-1}$).

- a.** Reimplement the function you made in Exercise 3 so that it instead plays the contribution from the detail spaces. Call the new function `playDWTfilterslowerdifference`.
- b.** In Exercise 5.3.12 we implemented a function `playDWTlowerdifference` for listening to the detail/error when the Haar wavelet is used. In the function `playDWTall` from Example 6.21, replace `playDWTlower` and `playDWTfilterslower` with `playDWTlowerdifference` and `playDWTfilterslowerdifference`. Describe the sounds you hear for different m . Try to explain why the sound seems to get louder when you increase m .

- 5.** Let us return to the piecewise linear wavelet from Exercise 5.5.2.

- a.** With $\hat{\psi}$ as defined as in Exercise 5.5.2 b., compute the coordinates of $\hat{\psi}$ in the basis ϕ_1 (i.e. $[\hat{\psi}]_{\phi_1}$) with $N = 8$, i.e. compute the IDWT of

$$[\hat{\psi}]_{(\phi_0, \psi_0)} = (-\alpha, -\beta, -\delta, 0, 0, 0, 0, -\gamma) \oplus (1, 0, 0, 0, 0, 0, 0, 0),$$

which is the coordinate vector you computed in Exercise 5.5.2 d.. For this, you should use the function `IDWTImpl`, with parameters being the filters G_0, G_1 , given as described by $g0, g1$ by equations (6.41)-(6.42) in Example 6.21.

- b.** If we redefine the basis \mathcal{C}_1 from $\{\phi_{0,0}, \psi_{0,0}, \phi_{0,1}, \psi_{0,1}, \dots\}$, to $\{\phi_{0,0}, \hat{\psi}_{0,0}, \phi_{0,1}, \hat{\psi}_{0,1}, \dots\}$, the vector you obtained in a. gives us an expression for the second column in $P_{\phi_1 \leftarrow \mathcal{C}_1}$. After redefining the basis like this, the corresponding filter G_1 has changed from that of the piecewise linear wavelet we started with. Use Matlab to so state the new filter G_1 with our compact filter notation. Also, plot its frequency response.

Hint: Here you are asked to find the unique filter with the same second column as $P_{\phi_1 \leftarrow \mathcal{C}_1}$, i.e. the vector from a..

c. Write code which uses Equation (6.30) to find H_0, H_1 from G_0, G_1 , and state these filters with our compact filter notation. Also, state the forms $h0, h1$, which should be used in calls to `DWTImpl` for our new wavelet. These replace the forms from equations (6.39)-(6.40) in Example 6.21, which we found for the first piecewise linear wavelet.

Hint: Note that the filter G_0 is unchanged from that of the first piecewise linear wavelet (since ϕ is unchanged when compared to the other wavelets for piecewise linear functions).

d. The filters you have found above should be symmetric, so that we can follow the procedure from Example 6.21 to listen to sound which has been wavelet-transformed by this wavelet. Write a program which plays our audio test file as in Example 6.21 for $m = 1, 2, 3, 4$ (i.e. plays the part in V_0), as well as the difference as in Exercise 4 (i.e. play the part from $W_0 \oplus W_1 \oplus \dots \oplus W_{m-1}$), where the new filters you have found are used. Listen to the sounds.

6. Repeat the previous exercise for the Haar wavelet as in exercise 4, and plot the corresponding frequency responses for $k = 2, 4, 6$.

6.4 A generalization of the filter representation, and its use in audio coding

It turns out that the filter representation, which we now have used for an alternative representation of a wavelet transformation, can be generalized in such a way that it also is useful for audio coding. In this section we will first define this generalization. We will then state how the MP3 standard encodes and decodes audio, and see how our generalization is connected to this. Much literature fails to elaborate on this connection. We will call our generalizations *filter bank transforms*, or simply *filter banks*. Just as for wavelets, filters are applied differently for the forward and reverse transforms. We start by defining the forward filter bank transform and its filters.

Definition 6.22 (Forward filter bank transforms). Let H_0, H_1, \dots, H_{M-1} be $N \times N$ -filters. A *forward filter bank transform* H produces output $\mathbf{z} \in \mathbb{R}^N$ from the input $\mathbf{x} \in \mathbb{R}^N$ in the following way:

1. $z_{iM} = (H_0 \mathbf{x})_{iM}$ for any i so that $0 \leq iM < N$.
2. $z_{iM+1} = (H_1 \mathbf{x})_{iM+1}$ for any i so that $0 \leq iM + 1 < N$.
3. ...
4. $z_{iM+(M-1)} = (H_{M-1} \mathbf{x})_{iM+(M-1)}$ for any i so that $0 \leq iM + (M - 1) < N$.

In other words, the output of a forward filter bank transform is computed by applying filters H_0, H_1, \dots, H_{M-1} to the input, and by downsampling and assembling these so that we obtain the same number of output samples as input samples

(also in this more general setting this is called critical sampling). H_0, H_1, \dots, H_{M-1} are also called *analysis filter components*, the output of filter H_i is called channel i *channel*, and M is called the number of channels. The output samples z_{iM+k} are also called the *subband samples* of channel k .

Clearly this definition generalizes the DWT and its analysis filters, since these can be obtained by setting $M = 2$. The DWT is thus a 2-channel forward filter bank transform. While the DWT produces the output $\begin{pmatrix} \mathbf{c}_{m-1} \\ \mathbf{w}_{m-1} \end{pmatrix}$ from the input \mathbf{c}_m , an M -channel forward filter bank transform splits the output into M components, instead of 2. Clearly, in the matrix of a forward filter bank transform the rows repeat cyclically with period M , similarly to MRA-matrices. In practice, the filters in a forward filter bank transform are chosen so that they concentrate on specific frequency ranges. This parallels what we saw for the filters of a wavelet, where one concentrated on high frequencies, one on low frequencies. Using a filter bank to split a signal into frequency components is also called *subband coding*. But the filters in a filter bank are usually not ideal bandpass filters. The introduction of frequencies outside the frequency band for a given filter is called *aliasing*. There exist a variety of different filter banks, for many different purposes [34, 30]. In Chapter 7 we will say more on how one can construct filter banks which can be used for subband coding.

Let us now turn to reverse filter bank transforms.

Definition 6.23 (Reverse filter bank transforms). Let G_0, G_1, \dots, G_{M-1} be $N \times N$ -filters. An *reverse filter bank transform* G produces $\mathbf{x} \in \mathbb{R}^N$ from $\mathbf{z} \in \mathbb{R}^N$ in the following way:

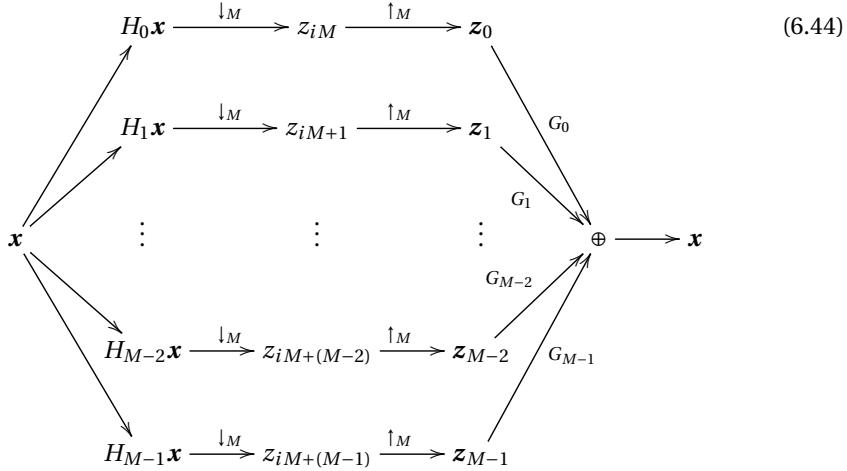
1. Define $\mathbf{z}_k \in \mathbb{R}^N$ as the vector where $(\mathbf{z}_k)_{iM+k} = z_{iM+k}$ for all i so that $0 \leq iM + k < N$, and $(\mathbf{z}_k)_s = 0$ for all other s .
- 2.

$$\mathbf{x} = G_0 \mathbf{z}_0 + G_1 \mathbf{z}_1 + \dots + G_{M-1} \mathbf{z}_{M-1}. \quad (6.43)$$

G_0, G_1, \dots, G_{M-1} are also called *synthesis filter components*.

Again, this generalizes the IDWT and its synthesis filters, and the IDWT can be seen as a 2-channel reverse filter bank transform. Also, in the matrix of a reverse filter bank transform, the columns repeat cyclically with period M , similarly to MRA-matrices. Also in this more general setting the filters G_i are in general different from the filters H_i . But we will see that, just as we saw for the Haar wavelet, there are important special cases where the analysis and synthesis filters are equal, and where their frequency responses are simply shifts of one another. It is clear that definitions 6.22 and 6.23 give the following diagram for computing forward and reverse

filter bank transforms:



Here \downarrow_M and \uparrow_M means that we extract every M 'th element in the vector, and add $M - 1$ zeros between the elements, respectively, similarly to how we previously defined \downarrow_2 and \uparrow_2 . Comparing Figure 6.9 with Figure 6.44 makes the similarities between wavelet transformations and the transformation used in the MP3 standard very visible: Although the filters used are different, they are subject to the same kind of processing, and can therefore be subject to the same implementations.

In general it may be that the synthesis filters do not invert exactly the analysis filters. If the synthesis system exactly inverts the analysis system, we say that we have a *perfect reconstruction filter bank*. Since the analysis system introduces undesired frequencies in the different channels, these have to cancel in the inverse transform, in order to reconstruct the input exactly. This is also called *alias cancellation*.

We will have use for the following simple connection between forward and reverse filter bank transforms, which follows immediately from the definitions.

Theorem 6.24 (Connection between forward and reverse filter bank transforms). Assume that H is a forward filter bank transform with filters H_0, \dots, H_{N-1} . Then H^T is a reverse filter bank transform with filters $G_0 = (H_0)^T, \dots, G_{N-1} = (H_{N-1})^T$.

6.4.1 Forward filter bank transform used for encoding in the MP3 standard

Now, let us turn to the MP3 standard. The MP3 standard document states that it applies a filter bank, and explains the following procedure for applying this filter bank, see p. 67 of the standard document (the procedure is slightly modified with mathematical terminology adapted to this book):

1. Input 32 audio samples at a time.

2. Build an input sample vector $X \in \mathbb{R}^{512}$, where the 32 new samples are placed first, all other samples are delayed with 32 elements. In particular the 32 last samples are taken out.
3. Multiply X componentwise with a vector C (this vector is defined through a table in the standard), to obtain a vector $Z \in \mathbb{R}^{512}$. The standard calls this *windowing*.
4. Compute the vector $Y \in \mathbb{R}^{64}$ where $Y_i = \sum_{j=0}^7 Z_{i+64j}$. The standard calls this a *partial calculation*.
5. Calculate $S = MY \in \mathbb{R}^{32}$, where M is the 32×64 - matrix where $M_{ik} = \cos((2i + 1)(k - 16)\pi/64)$. S is called the vector of output samples, or output subband samples. The standard calls this *matrixing*.

The standard does not motivate these steps, and does not put them into the filter bank transform framework which we have established. Also, the standard does not explain how the values in the vector C have been constructed.

Let us start by proving that the steps above really corresponds to applying a forward filter bank transform, and let us state the corresponding filters of this transform. The procedure computes 32 outputs in each iteration, and each of them is associated with a subband. Therefore, from the standard we would guess that we have $M = 32$ channels, and we would like to find the corresponding 32 filters H_0, H_1, \dots, H_{31} .

It may seem strange to use the name matrixing here, for something which obviously is matrix multiplication. The reason for this name must be that the at the origin of the procedure come from outside a linear algebra framework. The name windowing is a bit strange, too. This really does not correspond to applying a window to the sound samples as we explained in Section 3.3.1. We will see that it rather corresponds to applying a filter coefficient to a sound sample. A third and final thing which seems a bit strange is that the order of the input samples is reversed, since we are used to having the first sound samples in time with lowest index. This is perhaps more usual to do in an engineering context, and not so usual in a mathematical context. FIFO.

Clearly, the procedure above defines a linear transformation, and we need to show that this linear transformation coincides with the procedure we defined for a forward filter bank transform, for a set of 32 filters. The input to the transformation are the audio samples, which we will denote by a vector \mathbf{x} . At iteration s of the procedure above the input audio samples are $x_{32s-512}, x_{32s-510}, \dots, x_{32s-1}$, and $X_i = x_{32s-i-1}$ due to the reversal of the input samples. The output to the transformation at iteration s of the procedure are the S_0, \dots, S_{31} . We assemble these into a vector \mathbf{z} , so that the output at iteration s are $z_{32(s-1)} = S_0, z_{32(s-1)+1} = S_1, \dots, z_{32(s-1)+31} = S_{31}$.

We will have use for the following cosine-properties, which are easily verified:

$$\cos(2\pi(n + 1/2)(k + 2Nr)/(2N)) = (-1)^r \cos(2\pi(n + 1/2)k/(2N)) \quad (6.45)$$

$$\cos(2\pi(n + 1/2)(2N - k)/(2N)) = -\cos(2\pi(n + 1/2)k/(2N)). \quad (6.46)$$

With the terminology above and using Property (6.45) the transformation can be written as

$$\begin{aligned}
z_{32(s-1)+n} &= \sum_{k=0}^{63} \cos((2n+1)(k-16)\pi/64) Y_k = \sum_{k=0}^{63} \cos((2n+1)(k-16)\pi/64) \sum_{j=0}^7 Z_{k+64j} \\
&= \sum_{k=0}^{63} \sum_{j=0}^7 (-1)^j \cos((2n+1)(k+64j-16)\pi/64) Z_{k+64j} \\
&= \sum_{k=0}^{63} \sum_{j=0}^7 \cos((2n+1)(k+64j-16)\pi/64) (-1)^j C_{k+64j} X_{k+64j} \\
&= \sum_{k=0}^{63} \sum_{j=0}^7 \cos((2n+1)(k+64j-16)\pi/64) (-1)^j C_{k+64j} x_{32s-(k+64j)-1}.
\end{aligned}$$

Now, if we define $\{h_r\}_{r=0}^{511}$ by $h_{k+64j} = (-1)^j C_{k+64j}$, $0 \leq j < 8, 0 \leq k < 64$, and $h^{(n)}$ as the filter with coefficients $\{\cos((2n+1)(k-16)\pi/64) h_k\}_{k=0}^{511}$, the above can be simplified as

$$\begin{aligned}
z_{32(s-1)+n} &= \sum_{k=0}^{511} \cos((2n+1)(k-16)\pi/64) h_k x_{32s-k-1} = \sum_{k=0}^{511} (h^{(n)})_k x_{32s-k-1} \\
&= (h^{(n)} \mathbf{x})_{32s-1} = (E_{n-31} h^{(n)} \mathbf{x})_{32(s-1)+n}.
\end{aligned}$$

This means that the output of the procedure stated in the MP3 standard can be computed as a forward filter bank transform, and that we can choose the analysis filters as $H_n = E_{n-31} h^{(n)}$.

Theorem 6.25 (Forward filter bank transform for the MP3 standard). Define $\{h_r\}_{r=0}^{511}$ by $h_{k+64j} = (-1)^j C_{k+64j}$, $0 \leq j < 8, 0 \leq k < 64$, and $h^{(n)}$ as the filter with coefficients $\{\cos((2n+1)(k-16)\pi/64) h_k\}_{k=0}^{511}$. If we define $H_n = E_{n-31} h^{(n)}$, the procedure stated in the MP3 standard corresponds to applying the corresponding forward filter bank transform.

The filters H_n were shown in Example 3.36 as examples of filters which concentrate on specific frequency ranges. The h_k are the filter coefficients of what is called a prototype filter. This kind of filter bank is also called a *cosine-modulated filter bank*. The multiplication with $\cos(2\pi(n+1/2)(k-16)/(2N)) h_k$, modulated the filter coefficients so that the new filter has a frequency response which is simply shifted in frequency in a symmetric manner: In Exercise 3.5.9, we saw that, by multiplying with a cosine, we could construct new filters with real filter coefficients, which also corresponded to shifting a prototype filter in frequency. Of course, multiplication with a complex exponential would also shift the frequency response (such filter banks are called *DFT-modulated filter banks*), but the problem with this is that the new filter has complex coefficients: It will turn out that cosine-modulated filter banks can also be constructed so that they are invertible, and that one can find such filter banks where the inverse is easily found.

The effect of the delay in the definition of H_n is that, for each n , the multiplications with the vector \mathbf{x} are “aligned”, so that we can save a lot of multiplications by performing this multiplication first, and summing these. We actually save even more multiplications in the sum where j goes from 0 to 7, since we here multiply with the same cosines. The steps defined in the MP3 standard are clearly motivated by the desire to reduce the number of multiplications due to these facts. A simple arithmetic count illustrates these savings: For every 32 output samples, we have the following number of multiplications:

1. The first step computes 512 multiplications.
2. The second step computes 64 sums of 8 elements each, i.e. a total of $7 \times 64 = 448$ additions (note that $q = 512/64 = 8$).

The standard says nothing about how the matrix multiplication in the third step can be implemented. A direct multiplication would yield $32 \times 64 = 2048$ multiplications, leaving a total number of multiplications at 2560. In a direct implementation of the forward filter bank transform, the computation of 32 samples would need $32 \times 512 = 16384$ multiplications, so that the procedure sketched in the standard gives a big reduction.

The standard does not mention all possibilities for saving multiplications, however: We can reduce the number of multiplications even further, since clearly a DCT-type implementation can be used for the matrixing operation. We already have an efficient implementation for multiplication with a 32×32 type-III cosine matrix (this is simply the IDCT). We have seen that this implementation can be chosen to reduce the number of multiplications to $N \log_2 N/2 = 80$, so that the total number of multiplications is $512 + 80 = 592$. Clearly then, when we use the DCT, the first step is the computationally most intensive part.

6.4.2 Reverse filter bank transform used for decoding in the MP3 standard

Let us now turn to how decoding is specified in the MP3 standard, and see that we can associate this with a reverse filter bank transform. The MP3 standard also states the following procedure for decoding:

1. Input 32 new subband samples as the vector S .
2. Change vector $V \in \mathbb{R}^{512}$, so that all elements are delayed with 64 elements. In particular the 64 last elements are taken out.
3. Set the first 64 elements of V as $NS \in \mathbb{R}^{64}$, where N is the 64×32 -matrix where $N_{ik} = \cos((16+i)(2k+1)\pi/64)$. The standard also calls this matrixing.
4. Build the vector $U \in \mathbb{R}^{512}$ from V from the formulas $U_{64i+j} = V_{128i+j}$, $U_{64i+32+j} = V_{128i+96+j}$ for $0 \leq i \leq 7$ and $0 \leq j \leq 31$, i.e. U is the vector where V is first split into segments of length 132, and U is constructed by assembling the first and last 32 elements of each of these segments.

5. Multiply U componentwise with a vector D (this vector is defined in the standard), to obtain a vector $W \in \mathbb{R}^{512}$. The standard also calls this windowing.
6. Compute the 32 next sound samples as $\sum_{i=0}^{15} W_{32i+j}$.

To interpret this also in terms of filters, rewrite first steps 4 to 6 as

$$\begin{aligned}
x_{32(s-1)+j} &= \sum_{i=0}^{15} W_{32i+j} = \sum_{i=0}^{15} D_{32i+j} U_{32i+j} \\
&= \sum_{i=0}^7 D_{64i+j} U_{64i+j} + \sum_{i=0}^7 D_{64i+32+j} U_{64i+32+j} \\
&= \sum_{i=0}^7 D_{64i+j} V_{128i+j} + \sum_{i=0}^7 D_{64i+32+j} V_{128i+96+j}.
\end{aligned} \tag{6.47}$$

The elements in V are obtained by “matrixing” different segments of the vector z . More precisely, at iteration s we have that

$$\begin{pmatrix} V_{64r} \\ V_{64r+1} \\ \vdots \\ V_{64r+63} \end{pmatrix} = N \begin{pmatrix} z_{32(s-r-1)} \\ z_{32(s-r-1)+1} \\ \vdots \\ z_{32(s-r-1)+31} \end{pmatrix},$$

so that

$$V_{64r+j} = \sum_{k=0}^{31} \cos((16+j)(2k+1)\pi/64) z_{32(s-r-1)+k}$$

for $0 \leq j \leq 63$. Since also

$$V_{128i+j} = V_{64(2i)+j} \quad V_{128i+96+j} = V_{64(2i+1)+j+32},$$

we can rewrite Equation (6.47) as

$$\begin{aligned}
&\sum_{i=0}^7 D_{64i+j} \sum_{k=0}^{31} \cos((16+j)(2k+1)\pi/64) z_{32(s-2i-1)+k} \\
&+ \sum_{i=0}^7 D_{64i+32+j} \sum_{k=0}^{31} \cos((16+j+32)(2k+1)\pi/64) z_{32(s-2i-2)+k}.
\end{aligned}$$

Again using Relation (6.45), this can be written as

$$\begin{aligned}
&\sum_{k=0}^{31} \sum_{i=0}^7 (-1)^i D_{64i+j} \cos((16+64i+j)(2k+1)\pi/64) z_{32(s-2i-1)+k} \\
&+ \sum_{k=0}^{31} \sum_{i=0}^7 (-1)^i D_{64i+32+j} \cos((16+64i+j+32)(2k+1)\pi/64) z_{32(s-2i-2)+k}.
\end{aligned}$$

Now, if we define $\{g_r\}_{r=0}^{511}$ by $g_{64i+s} = (-1)^i C_{64i+s}$, $0 \leq i < 8$, $0 \leq s < 64$, and $g^{(k)}$ as the filter with coefficients $\{\cos((r+16)(2k+1)\pi/64)g_r\}_{r=0}^{511}$, the above can be simplified as

$$\begin{aligned} & \sum_{k=0}^{31} \sum_{i=0}^7 (g^{(k)})_{64i+j} z_{32(s-2i-1)+k} + \sum_{k=0}^{31} \sum_{i=0}^7 (g^{(k)})_{64i+j+32} z_{32(s-2i-2)+k} \\ &= \sum_{k=0}^{31} \left(\sum_{i=0}^7 (g^{(k)})_{32(2i)+j} z_{32(s-2i-1)+k} + \sum_{i=0}^7 (g^{(k)})_{32(2i+1)+j} z_{32(s-2i-2)+k} \right) \\ &= \sum_{k=0}^{31} \sum_{r=0}^{15} (g^{(k)})_{32r+j} z_{32(s-r-1)+k}, \end{aligned}$$

where we observed that $2i$ and $2i+1$ together run through the values from 0 to 15 when i runs from 0 to 7. Since z has the same values as z_k on the indices $32(s-r-1)+k$, this can be written as

$$\begin{aligned} &= \sum_{k=0}^{31} \sum_{r=0}^{15} (g^{(k)})_{32r+j} (z_k)_{32(s-r-1)+k} \\ &= \sum_{k=0}^{31} (g^{(k)} z_k)_{32(s-1)+j+k} = \sum_{k=0}^{31} ((E_{-k} g^{(k)}) z_k)_{32(s-1)+j}. \end{aligned}$$

By substituting a general s and j we see that $x = \sum_{k=0}^{31} (E_{-k} g^{(k)}) z_k$. We have thus proved the following.

Theorem 6.26 (Reverse filter bank transform for the MP3 standard). Define $\{g_r\}_{r=0}^{511}$ by $g_{64i+s} = (-1)^i C_{64i+s}$, $0 \leq i < 8$, $0 \leq s < 64$, and $g^{(k)}$ as the filter with coefficients $\{\cos((r+16)(2k+1)\pi/64)g_r\}_{r=0}^{511}$. If we define $G_k = E_{-k} g^{(k)}$, the procedure stated in the MP3 standard corresponds to applying the corresponding reverse filter bank transform.

In other words, both procedures for encoding and decoding stated in the MP3 standard both correspond to filter banks: A forward filter bank transform for the encoding, and a reverse filter bank transform for the decoding. Moreover, both filter banks can be constructed by cosine-modulating prototype filters, and the coefficients of these prototype filters are stated in the MP3 standard (up to multiplication with an alternating sign). Note, however, that the two prototype filters may be different. When we compare the two tables for these coefficients in the standard they do indeed seem to be different. At closer inspection, however, one sees a connection: If you multiply the values in the D -table with 32, and reverse them, you get the values in the C -table. This indicates that the analysis and synthesis prototype filters are the same, up to multiplication with a scalar. This connection will be explained in Section 8.5.

While the steps defined in the MP3 standard for decoding seem a bit more complex than the steps for encoding, they are clearly also motivated by the desire to reduce the number of multiplications. In both cases (encoding and decoding), the

window tables (C and D) are in direct connection with the filter coefficients of the prototype filter: one simply adds a sign which alternates for every 64 elements. The standard document does not mention this connection, and it is perhaps not so simple to find this connection in the literature (but see [26]).

The forward and reverse filter bank transforms are clearly very related. The following result clarifies this.

Theorem 6.27. (Connection between the forward and reverse filter bank transforms in the MP3 standard). Assume that a forward filter bank transform has filters on the form $H_i = E_{i-31}h^{(i)}$ for a prototype filter h . Then $G = E_{481}H^T$ is a reverse filter bank transform with filters on the form $G_k = E_{-k}g^{(k)}$, where g is a prototype filter where the elements equal the reverse of those in h . Vice versa, $H = E_{481}G^T$.

Proof: From Theorem 6.24 we know that H^T is a reverse filter bank transform with filters

$$(H_i)^T = (E_{i-31}h^{(i)})^T = E_{31-i}(h^{(i)})^T.$$

$(h^{(i)})^T$ has filter coefficients $\cos((2i+1)(-k-16)\pi/64))h_{-k}$. If we delay all $(H_i)^T$ with $481 = 512 - 31$ elements as in the theorem, we get a total delay of $512 - 31 + 31 - i = 512 - i$ elements, so that we get the filter

$$\begin{aligned} & E_{512-i}\{\cos((2i+1)(-k-16)\pi/64))h_{-k}\}_k \\ &= E_{-i}\{\cos((2i+1)(-(k-512)-16)\pi/64))h_{-(k-512)}\}_k \\ &= E_{-i}\{\cos((2i+1)(k+16)\pi/64))h_{-(k-512)}\}_k. \end{aligned}$$

Now, we define the prototype filter g with elements $g_k = h_{-(k-512)}$. This has, just as h , its support on $[1, 511]$, and consists of the elements from h in reverse order. If we define $g^{(i)}$ as the filter with coefficients $\cos((2i+1)(k+16)\pi/64))g_k$, we see that $E_{481}H^T$ is a reverse filter bank transform with filters $E_{-i}g^{(i)}$. Since $g^{(k)}$ now has been defined as for the MP3 standard, and its elements are the reverse of those in h , the result follows. ■

We will have use for this result in Section 8.5, when we find conditions on the prototype filter in order for the reverse transform to invert the forward transform. Preferably, the reverse filter bank transform inverts exactly the forward filter bank transform. In Exercise 2 we construct examples which show that this is not the case. In the same exercise we also find many examples where the reverse transform does what we would expect. These examples will also be explained in Section 8.5, where we also will see how one can get around this so that we obtain a system with perfect reconstruction. It may seem strange that the MP3 standard does not do this.

In the MP3 standard, the output from the forward filter bank transform is processed further, before the result is compressed using a lossless compression method.

Exercises for Section 6.4

1. The values C_q, D_q can be found by calling the functions `mp3ctable`, `mp3dtable` which can be found on the book's webpage.

- a. Use Matlab to verify the connection we stated between the tables C and D , i.e. that $D_i = 32C_i$ for all i .
 - b. Plot the frequency responses of the corresponding prototype filters, and verify that they both are lowpass filters. Use the connection from Theorem (6.25) to find the prototype filter coefficients from the C_q .
2. It is not too difficult to make implementations of the forward and reverse steps as explained in the MP3 standard. In this exercise we will experiment with this. In your code you can for simplicity assume that the input and output vectors to your methods all have lengths which are multiples of 32. Also, use the functions `mp3ctable`, `mp3dtable` mentioned in the previous exercise.

- a. Write a function

```
function z=mp3forwardfbt(x)
```

which implements the steps in the forward direction of the MP3 standard.

- b. Write also a function

```
function z=mp3reversefbt(x)
```

which implements the steps in the reverse direction.

6.5 Summary

We started this chapter by noting that, by reordering the target base of the DWT, the change of coordinate matrix took a particular form. From this form we understood that the DWT could be realized in terms of two filters H_0 and H_1 , and that the IDWT could be realized in a similar way in terms of two filters G_0 and G_1 . This gave rise to what we called the filter representation of wavelets. The filter representation gives an entirely different view on wavelets: instead of constructing function spaces with certain properties and deducing corresponding filters from these, we can instead construct filters with certain properties (such as alias cancellation and perfect reconstruction), and attempt to construct corresponding mother wavelets, scaling functions, and function spaces. This strategy, which replaces problems from function theory with discrete problems, will be the subject of the next chapter. In practice this is what is done.

We stated what is required for filter bank matrices to invert each other: The frequency responses of the lowpass filters needed to satisfy a certain equation, and once this is satisfied the highpass filters can easily be obtained in the same way we previously obtained highpass filters from lowpass filters. We will return to this equation in the next chapter.

A useful consequence of the filter representation was that we could reuse existing implementations of filters to implement the DWT and the IDWT, and reuse existing theory, such as symmetric extensions. For wavelets, symmetric extensions

are applied in a slightly different way, when compared to the developments which lead to the DCT. We looked at the frequency responses of the filters for the wavelets we have encountered upto now. From these we saw that G_0, H_0 were lowpass filters, and that G_1, H_1 were highpass filters, and we argued why this is typically the case for other wavelets as well. The filter representation was also easily generalized from 2 to $M > 2$ filters, and such transformations had a similar interpretation in terms of splitting the input into a uniform set of frequencies. Such transforms were generally called filter bank transforms, and we saw that the processing performed by the MP3 standard could be interpreted as a certain filter bank transform, called a cosine-modulated filter bank. This is just one of many possible filter banks. In fact, the filter bank of the MP3 standard is largely outdated, since it is too simple, and as we will see it does not even give perfect reconstruction (only alias cancellation and no phase distortion). It is merely chosen here since it is the simplest to present theoretically, and since it is perhaps the best known standard for compression of sound. Other filters banks with better properties have been constructed, and they are used in more recent standards. In many of these filter banks, the filters do not partition frequencies uniformly, and have been adapted to the way the human auditory system handles the different frequencies. Different construction methods are used to construct such filter banks. The motivation behind filter bank transforms is that their output is more suitable for further processing, such as compression, or playback in an audio system, and that they have efficient implementations.

We mentioned that the MP3 standard does not say how the prototype filters were chosen. We will have more to say on what dictates their choice in Section 8.5.

There are several differences between the use of wavelet transformations in wavelet theory, and the use of filter bank transforms in signal processing theory. One is that wavelet transforms are typically applied in stages, while filter bank transforms often are not. Nevertheless, such use of filter banks also has theoretical importance, and this gives rise to what is called *tree-structured filter banks* [34]. Another difference lies in the use of the term perfect reconstruction system. In wavelet theory this is a direct consequence of the wavelet construction, since the DWT and the IDWT correspond to change of coordinates to and from the same bases. The alternative QMF filter bank was used as an example of a filter bank which stems from signal processing, and which also shows up in wavelet transformation. In signal processing theory, one has a wider perspective, since one can design many useful systems with fast implementations when one replaces the perfect reconstruction requirement with a near perfect reconstruction requirement. One instead requires that the reverse transform gives alias cancellation. The classical QMF filter banks were an example of this. The original definition of classical QMF filter banks are from [7], and differ only in a sign from how they are defined here.

All filters we encounter in wavelets and filter banks in this book are FIR. This is just done to limit the exposition. Much useful theory has been developed using IIR-filters.

Constructing interesting wavelets

In the previous chapter we saw that, once we have an MRA with corresponding scaling function and mother wavelet, we could write down corresponding forward and reverse filter bank matrices. We also exemplified this in terms of the wavelets we have looked at. In practice, however, it is difficult to come up with good choices of such functions. Constructing such matrices which invert each other is a discrete problem, however, and may be more easily solved. It is therefore tempting to turn the problem around and ask the following question: Can one, from usable filters H_0, H_1, G_0, G_1 , where the associated forward and reverse filter bank transforms revert each other, also associate a scaling function ϕ and a mother wavelet ψ ? Can one also associate a scaling function and mother wavelet with the dual filter bank transforms? If so we will denote these by $\tilde{\phi}$ and $\tilde{\psi}$, respectively. And if we can do this, what are the appropriate conditions to put on the filters in order for them to be useful in approximation of functions?

To ensure that the mother wavelets have a given number of vanishing moments, or that the scaling functions are differentiable? Answers to these questions would certainly transfer even more of the theory of wavelets to the theory of filters.

In this chapter we will indeed see that we can find answers to these questions. We will first see how we can associate scaling functions and mother wavelets from filters, and for which filters this association can be done. There are two things we will be concerned with. First of all, which filters produce scaling functions and mother wavelets which are differentiable? This question is clearly important in cases where we know something about the regularity of the function we approximate. Secondly, which filters produce mother wavelets with a given number of vanishing moments? It will also be clear why it is desirable for the mother wavelet to have many vanishing moments, in terms of approximation of functions.

Finally, we will also solve for these properties to find the simplest filters with a given number of vanishing moments, or a given degree of differentiability. Several of these filters enjoy a widespread use in applications. We will look at two of these. These are used for lossless and lossy compression in JPEG2000, which is a

much used standard. These wavelets all have symmetric filters. We end the chapter by looking at a family of orthonormal wavelets with different number of vanishing moments.

7.1 From filters to scaling functions and mother wavelets

Let us turn to the first issue mentioned in the introduction: If we have filters H_0, H_1, G_0, G_1 which give perfect reconstruction as defined in the previous chapter, how can we from these associate scaling functions and mother wavelets for the transform and dual transform (see Definition 6.7)? Also, for which filters can this association be done? Denote by $\phi, \tilde{\phi}$ the scaling functions for the transform and the dual transform, respectively. Since $(H_0)^T$ is the lowpass filter in the dual reverse transform we require that (see Equation (6.10))

$$\phi(t) = \sum_{n=0}^{2N-1} (G_0)_{n,0} \phi_{1,n}(t) \quad (7.1)$$

$$\tilde{\phi}(t) = \sum_{n=0}^{2N-1} ((H_0)^T)_{n,0} \tilde{\phi}_{1,n}(t) = \sum_{n=0}^{2N-1} (H_0)_{0,n} \tilde{\phi}_{1,n}(t). \quad (7.2)$$

We define as before V_m as the space spanned by $\phi_m = \{\phi_{m,n}\}_n$, and we define also the dual space \tilde{V}_m as the space spanned by $\tilde{\phi}_m = \{\tilde{\phi}_{m,n}\}_n$. Hopefully $\phi, \tilde{\phi}$ are bases when $\phi, \tilde{\phi}$ are defined in this way. Similarly we denote by $\psi, \tilde{\psi}$ the mother wavelet for the transform and the dual transform, respectively. We require that (see Equation (6.11))

$$\psi(t) = \sum_{n=0}^{2N-1} (G_1)_{n,1} \phi_{1,n}(t) \quad (7.3)$$

$$\tilde{\psi}(t) = \sum_{n=0}^{2N-1} ((H_1)^T)_{n,1} \tilde{\phi}_{1,n}(t) = \sum_{n=0}^{2N-1} (H_1)_{1,n} \tilde{\phi}_{1,n}(t). \quad (7.4)$$

We define the detail space W_m as the space spanned by $\psi_m = \{\psi_{m,n}\}_n$, and the dual detail space \tilde{W}_m as the space spanned by $\tilde{\psi}_m = \{\tilde{\psi}_{m,n}\}_n$. As in the case of the wavelets of piecewise linear functions, these spaces are typically different from the orthogonal complement of V_{m-1} in V_m (\tilde{V}_{m-1} in \tilde{V}_m), but still constructed in such a way that $V_m = V_{m-1} \oplus W_m$ ($\tilde{V}_m = \tilde{V}_{m-1} \oplus \tilde{W}_m$). The DWT and the dual DWT can now be defined as before, assuming that these functions exist.

Example 7.1. Let us return to the alternative piecewise linear wavelet. In Example 6.16 we found the filters H_0, H_1 for this wavelet, and these determine the dual scaling function and the dual mother wavelet. We already know how the scaling function and the mother wavelet look, but how do these dual functions look? It turns out that there is usually no way to find analytical expressions for these dual functions (as is the case for the scaling function and the mother wavelet itself in most cases), but that there still is an algorithm we can apply in order to see how these functions look. This algorithm is called the *cascade algorithm*. The cascade algorithm works by approximating these functions from V_m or \tilde{V}_m . By increasing m we can clearly

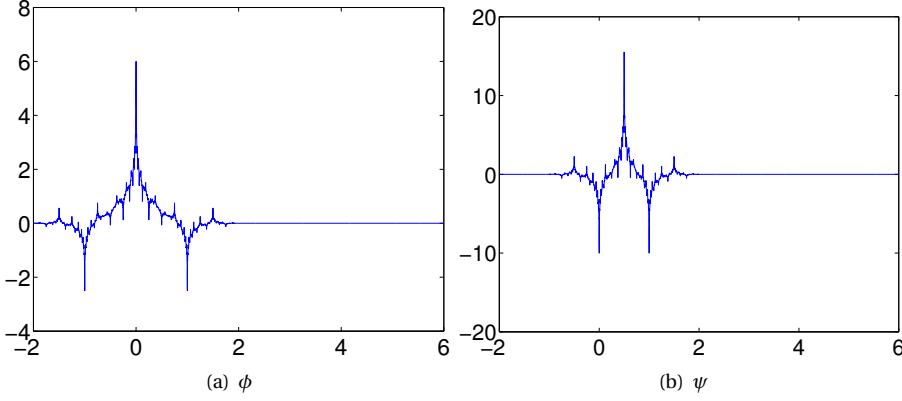


Figure 7.1: Dual scaling function and dual mother wavelet for the alternative piecewise linear wavelet.

obtain arbitrarily good such approximations, and we have previously argued that the coordinates in ϕ_m or $\tilde{\phi}_m$ of these approximations are good approximations to the samples of the function.

To be more specific, we start with the following observation:

1. the coordinates of $\tilde{\phi}$ in $(\tilde{\phi}_0, \tilde{\psi}_0, \tilde{\psi}_1 \dots)$ is the vector with 1 first, followed by only zeros,
2. the coordinates of $\tilde{\psi}$ in $(\tilde{\phi}_0, \tilde{\psi}_0, \tilde{\psi}_1 \dots)$ is the vector with N zeros first, then a 1, and then only zeros.

The length of these vectors is $N2^m$. The coordinates in $\tilde{\phi}_m$ for $\tilde{\phi}$ and $\tilde{\psi}$ can be obtained by applying the m -level dual IDWT (i.e. the filters $(H_0)^T$, $(H_1)^T$ are used) to these vectors. In Exercise 1 we will study code which uses this approach to approximate the scaling function and mother wavelet. In Figure 7.1 we have plotted the resulting coordinates in $\tilde{\phi}_{10}$, and thus a good approximation to $\tilde{\phi}$ and $\tilde{\psi}$. We see that these functions look very irregular. Also, they are very different from the original scaling function and mother wavelet. We will later argue that this is bad, it would be much better if $\phi \approx \tilde{\phi}$ and $\psi \approx \tilde{\psi}$. From Theorem 6.10 it follows that the support sizes of these dual functions are 3 and 2. This is verified in these figures. ♣

Let us formalize the cascade algorithm from the previous example as follows.

Definition 7.2 (The cascade algorithm). The *cascade algorithm* applies a change of coordinates for the functions $\tilde{\phi}, \tilde{\psi}$ from $(\tilde{\phi}_0, \tilde{\psi}_0, \tilde{\psi}_1 \dots)$ to $\tilde{\phi}_m$, and uses the new coordinates as an approximation to the function values of these functions.

While the previous example certainly shows that the dual functions can be visualized if they exist, we need to investigate when the functional equations defining

them are solvable. In this direction, we start by taking the CTFT of the left and right hand side of Equation (7.1) to obtain

$$\begin{aligned}
\hat{\phi}(\nu) &= \int_{-\infty}^{\infty} \phi(t) e^{-2\pi i \nu t} dt = \int_{-\infty}^{\infty} \left(\sum_n (G_0)_{n,0} \sqrt{2} \phi(2t-n) \right) e^{-2\pi i \nu t} dt \\
&= \sqrt{2} \sum_n \int_{-\infty}^{\infty} (G_0)_{n,0} \phi(2t-n) e^{-2\pi i \nu t} dt = \frac{1}{\sqrt{2}} \sum_n \int_{-\infty}^{\infty} (G_0)_{n,0} \phi(t) e^{-2\pi i \nu (t+n)/2} dt \\
&= \frac{1}{\sqrt{2}} \sum_n \int_{-\infty}^{\infty} (G_0)_{n,0} e^{-2\pi i \nu n/2} \phi(t) e^{-2\pi i (\nu/2)t} dt \\
&= \frac{1}{\sqrt{2}} \left(\sum_n (G_0)_{n,0} e^{-2\pi i \nu n/2} \right) \int_{-\infty}^{\infty} \phi(t) e^{-2\pi i (\nu/2)t} dt = \frac{\lambda_{G_0}(2\pi\nu/2)}{\sqrt{2}} \hat{\phi}(\nu/2).
\end{aligned}$$

We obtain similar expressions for the CTFT's of $\psi, \tilde{\phi}, \tilde{\psi}$, so that

$$\begin{aligned}
\hat{\phi}(\nu) &= \frac{\lambda_{G_0}(2\pi\nu/2)}{\sqrt{2}} \hat{\phi}(\nu/2) & \hat{\psi}(\nu) &= e^{-2\pi i \nu/2} \frac{\lambda_{G_1}(2\pi\nu/2)}{\sqrt{2}} \hat{\phi}(\nu/2) \\
\hat{\tilde{\phi}}(\nu) &= \frac{\lambda_{H_0}(-2\pi\nu/2)}{\sqrt{2}} \hat{\phi}(\nu/2) & \hat{\tilde{\psi}}(\nu) &= e^{-2\pi i \nu/2} \frac{\lambda_{H_1}(-2\pi\nu/2)}{\sqrt{2}} \hat{\phi}(\nu/2)
\end{aligned}$$

Clearly these expressions can be continued recursively, so that we obtain the following:

Theorem 7.3. Define

$$g_N(\nu) = \prod_{s=1}^N \frac{\lambda_{G_0}(2\pi\nu/2^s)}{\sqrt{2}} \chi_{[0,1]}(2^{-N}\nu) \quad h_N(\nu) = \prod_{s=1}^N \frac{\lambda_{H_0}(-2\pi\nu/2^s)}{\sqrt{2}} \chi_{[0,1]}(2^{-N}\nu). \quad (7.5)$$

Then on $[0, 2\pi 2^N]$ we have that

$$\begin{aligned}
\hat{\phi}(\nu) &= g_N(\nu) \hat{\phi}(\nu/2^N) & \hat{\psi}(\nu) &= e^{-2\pi i \nu/2} \frac{\lambda_{G_1}(2\pi\nu/2)}{\lambda_{G_0}(2\pi\nu/2)} \hat{\phi}(\nu) \\
\hat{\tilde{\phi}}(\nu) &= h_N(\nu) \hat{\tilde{\phi}}(\nu/2^N) & \hat{\tilde{\psi}}(\nu) &= e^{-2\pi i \nu/2} \frac{\lambda_{H_1}(-2\pi\nu/2)}{\lambda_{H_0}(-2\pi\nu/2)} \hat{\tilde{\phi}}(\nu)
\end{aligned}$$

Here we have used that

$$\begin{aligned}
\hat{\psi}(\nu) &= e^{-2\pi i \nu/2} \frac{\lambda_{G_1}(2\pi\nu/2)}{\lambda_{G_0}(2\pi\nu/2)} \frac{\lambda_{G_0}(2\pi\nu/2)}{\sqrt{2}} \hat{\phi}(\nu/2) = e^{-2\pi i \nu/2} \frac{\lambda_{G_1}(2\pi\nu/2)}{\lambda_{G_0}(2\pi\nu/2)} g_1(\nu) \hat{\phi}(\nu/2) \\
&= e^{-2\pi i \nu/2} \frac{\lambda_{G_1}(2\pi\nu/2)}{\lambda_{G_0}(2\pi\nu/2)} \hat{\phi}(\nu) \\
\hat{\tilde{\psi}}(\nu) &= e^{-2\pi i \nu/2} \frac{\lambda_{H_1}(-2\pi\nu/2)}{\lambda_{H_0}(-2\pi\nu/2)} \frac{\lambda_{H_0}(-2\pi\nu/2)}{\sqrt{2}} \hat{\tilde{\phi}}(\nu/2) = e^{-2\pi i \nu/2} \frac{\lambda_{H_1}(-2\pi\nu/2)}{\lambda_{H_0}(-2\pi\nu/2)} h_1(\nu) \hat{\tilde{\phi}}(\nu/2) \\
&= e^{-2\pi i \nu/2} \frac{\lambda_{H_1}(-2\pi\nu/2)}{\lambda_{H_0}(-2\pi\nu/2)} \hat{\tilde{\phi}}(\nu)
\end{aligned}$$

Note that $g_{N+1}(\nu) = \frac{\lambda_{G_0}(2\pi\nu/2)}{\sqrt{2}} g_N(\nu/2)$, $h_{N+1}(\nu) = \frac{\lambda_{H_0}(-2\pi\nu/2)}{\sqrt{2}} h_N(\nu/2)$. We will continue based on two assumptions:

1. There exist functions u_N, v_N defined on \mathbb{R} so that $\hat{u}_N = g_N, \hat{v}_N = h_N$.
2. The limits $\lim_{N \rightarrow \infty} u_N(t)$ and $\lim_{N \rightarrow \infty} v_N(t)$ exist for all t .

If these are fulfilled we will define the scaling functions by $\phi(t) = c \lim_{N \rightarrow \infty} u_N(t)$ and $\tilde{\phi} = \frac{1}{c} \lim_{N \rightarrow \infty} v_N(t)$, where c is a constant we can choose as we wish. Following the calculation above under these assumptions, we see that

$$u_{N+1}(t) = \sum_n (G_0)_{n,0} \sqrt{2} u_N(2t - n) \quad v_{N+1}(t) = \sum_n (H_0)_{0,n} \sqrt{2} v_N(2t - n).$$

Note first that $g_0(\nu) = h_0(\nu) = \chi_{[0,1]}(\nu)$, and that $= \chi_{[0,1]}(\nu)$ //her
Note first that, since $\langle u_0, v_0 \rangle = \langle g_0, h_0 \rangle$,

$$\int_{-\infty}^{-\infty} u_0(t) \overline{v_0(t-k)} dt = \int_{-\infty}^{-\infty} g_0(\nu) \overline{h_0(\nu) e^{2\pi i k \nu}} d\nu = \int_0^{2\pi} e^{-2\pi i k \nu} d\nu = \delta_{k,0}.$$

Now assume that we have proved that $\langle u_N(t), v_N(t-k) \rangle = \delta_{k,0}$. We then get that

$$\begin{aligned} \langle u_{N+1}(t), v_{N+1}(t-k) \rangle &= 2 \sum_{n_1, n_2} (G_0)_{n_1,0} (H_0)_{0,n_2} \langle u_N(2t-n_1), v_N(2(t-k)-n_2) \rangle \\ &= 2 \sum_{n_1, n_2} (G_0)_{n_1,0} (H_0)_{0,n_2} \langle u_N(t), v_N(t+n_1-n_2-2k) \rangle \\ &= \sum_{n_1, n_2 | n_1 - n_2 = 2k} (G_0)_{n_1,0} (H_0)_{0,n_2} = \sum_n (H_0)_{0,n-2k} (G_0)_{n,0} \\ &= \sum_n (H_0)_{2k,n} (G_0)_{n,0} = \sum_n H_{2k,n} G_{n,0} = (HG)_{2k,0} = I_{2k,0} = \delta_{k,0} \end{aligned}$$

where we did the change of variables $u = 2t - n_1$. Taking the limit as $N \rightarrow \infty$, we get that $\langle \phi, \tilde{\phi}_{0,k} \rangle = \delta_{k,0}$, regardless of c . From this it immediately follows that $\langle \phi_{m,k}, \tilde{\phi}_{m,l} \rangle = \delta_{k,l}$ for all k, l, m , and

$$\begin{aligned} \langle \psi_{0,k}, \tilde{\psi}_{0,l} \rangle &= \sum_{n_1, n_2} (G_1)_{n_1,1} (H_1)_{1,n_2} \langle \phi_{1,n_1+2k}(t) \tilde{\phi}_{1,n_2+2l}(t) \rangle \\ &= \sum_n (G_1)_{n,1} (H_1)_{1,n+2(k-l)} = \sum_n (H_1)_{1+2(l-k),n} (G_1)_{n,1} = \sum_n H_{1+2(l-k),n} G_{n,1} \\ &= (HG)_{1+2(l-k),1} = \delta_{k,0}. \end{aligned}$$

Similarly,

$$\begin{aligned} \langle \psi_{0,k} \tilde{\phi}_{0,l} \rangle &= \sum_{n_1, n_2} (G_1)_{n_1,1} (H_0)_{0,n_2} \langle \phi_{1,n_1+2k}(t) \tilde{\phi}_{1,n_2+2l}(t) \rangle = \sum_n (G_1)_{n,1} (H_0)_{0,n+2(k-l)} \\ &= \sum_n (H_0)_{2(l-k),n} (G_1)_{n,1} = \sum_n H_{2(l-k),n} G_{n,1} = (HG)_{2(l-k),1} = 0 \\ \langle \phi_{0,k} \tilde{\psi}_{0,l} \rangle &= \sum_{n_1, n_2} (G_0)_{n_1,0} (H_1)_{1,n_2} \langle \phi_{1,n_1+2k}(t) \tilde{\phi}_{1,n_2+2l}(t) \rangle = \sum_n (G_0)_{n,0} (H_1)_{1,n+2(k-l)} \\ &= \sum_n (H_1)_{1+2(l-k),n} (G_0)_{n,0} = \sum_n H_{1+2(l-k),n} G_{n,0} = (HG)_{1+2(l-k),0} = 0. \end{aligned}$$

From this we also get with a simple change of coordinates that

$$\langle \psi_{m,k}, \tilde{\psi}_{m,l} \rangle = \langle \psi_{m,k}, \tilde{\phi}_{m,l} \rangle = \langle \phi_{m,k}, \tilde{\psi}_{m,l} \rangle = 0.$$

Finally, if $m' < m$, $\phi_{m',k'}, \psi_{m',k}$ can be written as a linear combination of $\phi_{m,l}$, so that $\langle \phi_{m',k}, \tilde{\psi}_{m,l} \rangle = \langle \psi_{m',k}, \tilde{\psi}_{m,l} \rangle = 0$ due to what we showed above. Similarly, $\langle \tilde{\phi}_{m',k}, \psi_{m,l} \rangle = \langle \tilde{\psi}_{m',k}, \psi_{m,l} \rangle = 0$. Before we summarize, let us make the following definition.

Definition 7.4 (Biorthogonal bases). We say that two bases $\{f_n\}$, $\{g_m\}$ are *biorthogonal* if $\langle f_n, g_m \rangle = 0$ whenever $n \neq m$.

Using this concept the above deductions can be summarized as follows.

Theorem 7.5. Let H_0, H_1, G_0, G_1 be filters so that the corresponding forward and reverse filter bank matrices invert each other. Assume that the limits

$$\phi(t) = c \lim_{N \rightarrow \infty} u_N(t) \quad \tilde{\phi}(t) = \frac{1}{c} \lim_{N \rightarrow \infty} v_N(t) \quad (7.6)$$

exist, where $\hat{u}_N = g_N$, $\hat{v}_N = h_N$, and where g_N and h_N are defined by Equation (7.5). Define also $\psi, \tilde{\psi}$ by equations (7.3)-(7.4). Then Equation (7.1) is also fulfilled. Moreover,

1. the bases $\{\phi_{m,n}\}_n$ and $\{\tilde{\phi}_{m,n}\}_n$ are biorthogonal,
2. the bases $\{\psi_{m',n}, \phi_{0,n}\}_{n,m' < m}$ and $\{\tilde{\psi}_{m',n}, \tilde{\phi}_{0,n}\}_{n,m' < m}$ are biorthogonal.

Thus, the source bases of the DWT and the dual DWT are biorthogonal, and the target bases are also biorthogonal.

From this we can make the following definition as a generalization OF MRA.

Definition 7.6 (Dual multiresolution analysis). A Dual multiresolution analysis consists of two nested sequence of function spaces

$$V_0 \subset V_1 \subset V_2 \subset \cdots \subset V_m \subset \cdots \quad \tilde{V}_0 \subset \tilde{V}_1 \subset \tilde{V}_2 \subset \cdots \subset \tilde{V}_m \subset \cdots$$

so that

1. There exist functions $\phi, \tilde{\phi}$, called a scaling function and a dual scaling function, so that $\phi, \tilde{\phi}$ are bases for V_0 and \tilde{V}_0 , respectively, and so that these bases are biorthogonal.
2. The pairs $(\phi, V_m), (\tilde{\phi}, \tilde{V}_m)$ both satisfy requirements 1-3 from Definition 5.39.

From our construction it is also clear that a perfect reconstruction filter bank gives rise to a dual MRA, and how the associated mother wavelets are constructed.

Definition 7.6 has the original definition of an MRA (orthonormal MRA) as a special case, since $\tilde{\phi} = \phi$ then. Clearly, if $\phi = \tilde{\phi}$ and the basis $\{\phi(t-n)\}_{0 \leq n < N}$ is orthonormal, the approximations above coincide with the projection of f onto V_m . When $\phi \neq \tilde{\phi}$, however, there are no reasons to believe that these approximations equal the best approximations to f from V_m . For non-orthonormal MRA's, we have no procedure for computing best approximations. Anyway, we have mentioned that Theorem 5.41 is valid also when $\{\phi(t-n)\}_n$ are not orthogonal, i.e. $\frac{2^{-m}}{\int_0^N \phi_{m,0}(t) dt} f(n/2^m) \phi_{m,n}(t)$ has validity as an approximation to f also in this more general setting. The following is also clear from the theorem above.

Theorem 7.7. For all $f \in V_m$, $\tilde{f} \in \tilde{V}_m$, we can write

$$\begin{aligned} f(t) &= \sum_n \langle f(t), \tilde{\phi}_{m,n} \rangle \phi_{m,n} = \sum_n \langle f(t), \tilde{\phi}_{0,n} \rangle \phi_{0,n} + \sum_{m' < m, n} \langle f(t), \tilde{\psi}_{m',n} \rangle \psi_{m',n} \\ \tilde{f}(t) &= \sum_n \langle \tilde{f}(t), \phi_{m,n} \rangle \tilde{\phi}_{m,n} = \sum_n \langle \tilde{f}(t), \phi_{0,n} \rangle \tilde{\phi}_{0,n} + \sum_{m' < m, n} \langle \tilde{f}(t), \psi_{m',n} \rangle \tilde{\psi}_{m',n}. \end{aligned}$$

Thus,

1. The input to the DWT is $c_{m,n} = \langle f, \tilde{\phi}_{m,n} \rangle$, and the output of the DWT is $c_{0,n} = \langle f, \tilde{\phi}_{0,n} \rangle$ and $w_{m',n} = \langle f, \tilde{\psi}_{m',n} \rangle$
2. The input to the dual DWT is $\tilde{c}_{m,n} = \langle \tilde{f}, \phi_{m,n} \rangle$, and the output of the dual DWT is $\tilde{c}_{0,n} = \langle \tilde{f}, \phi_{0,n} \rangle$ and $\tilde{w}_{m',n} = \langle \tilde{f}, \psi_{m',n} \rangle$.

Since the DWT is the change of coordinates from ϕ_1 to (ϕ_0, ψ_0) , and similarly for the IDWT, inserting these basis functions for f we obtain the following corollary.

Corollary 7.8. All entries in both the DWT and the IDWT matrices (and their dual matrices, since these are obtained by transposing these matrices) can be written as inner products. More precisely,

1. in the DWT matrix, column k has entries $\langle \phi_{1,k}, \tilde{\phi}_{0,l} \rangle$, and $\langle \phi_{1,k}, \tilde{\psi}_{0,l} \rangle$,
2. in the IDWT matrix, column $2k$ has entries $\langle \phi_{0,k}, \tilde{\phi}_{1,l} \rangle$, and column $2k+1$ has entries $\langle \psi_{0,k}, \tilde{\phi}_{1,l} \rangle$.

Such expansions also make sense when the functions are not in V_m, \tilde{V}_m . In fact,

$$\sum_n \langle f(t), \tilde{\phi}_{m,n} \rangle \phi_{m,n}(t) \text{ and } \sum_n \langle f(t), \phi_{m,n} \rangle \tilde{\phi}_{m,n}(t)$$

serve as approximations to f from V_m and \tilde{V}_m , respectively (although they are typically not the best approximations). To see why, if $f(t) = \sum_n c_n \phi_{m,n}(t) + \epsilon(t)$ for $\epsilon(t)$ a

small function, then the first approximation is

$$\begin{aligned}\sum_n \left\langle \sum_{n'} c_{n'} \phi_{m,n'}(t) + \epsilon(t), \tilde{\phi}_{m,n} \right\rangle \phi_{m,n}(t) &= \sum_n c_n \phi_{m,n}(t) + \sum_n \langle \epsilon(t), \tilde{\phi}_{m,n} \rangle \phi_{m,n}(t) \\ &= f(t) + \sum_n \langle \epsilon(t), \tilde{\phi}_{m,n} \rangle \phi_{m,n}(t) - \epsilon(t).\end{aligned}$$

Clearly, the difference $\sum_n \langle \epsilon(t), \tilde{\phi}_{m,n} \rangle \phi_{m,n}(t) - \epsilon(t)$ from f is small, at least under suitable assumptions on the scaling function and the dual scaling function, and when ϵ is small.

Now we need to say something about when the products $g_N(v)$, $h_N(v)$ converge.

Theorem 7.9. Assume that H_0 , H_1 , G_0 , G_1 are FIR-filters, and that $\lambda_{G_0}(0) = \lambda_{H_0}(0) = \sqrt{2}$.

Then $\lim_{N \rightarrow \infty} g_N(v)$, $\lim_{N \rightarrow \infty} h_N(v)$ exist for all v , and the limits are infinitely differentiable functions. Also, the following hold.

1. The multiplicity of a zero at $v = 0$ for $\hat{\psi}(v)$ equals that of $\lambda_{G_1}(\omega)$ at $\omega = 0$.
2. The multiplicity of a zero at $v = 0$ for $\hat{\psi}(v)$ equals that of $\lambda_{H_1}(\omega)$ at $\omega = 0$.

Proof. Setting $v = 0$ in equations (7.5), we see that we must have that $\lambda_{G_0}(0) = \lambda_{H_0}(0) = \sqrt{2}$, otherwise the infinite products above will not converge. This explains the requirement $\lambda_{G_0}(0) = \lambda_{H_0}(0) = \sqrt{2}$. Since the filters are FIR, $\lambda_{H_0}(\omega)$ and $\lambda_{G_0}(\omega)$ are polynomials in $e^{i\omega}$, and they must thus have a bounded derivative on $[0, 2\pi]$. Since $\lambda_{G_0}(0) = \lambda_{H_0}(0) = \sqrt{2}$, there exists $c > 0$ so that

$$|\lambda_{G_0}(\omega)| \leq \sqrt{2} + c|\omega| \text{ and } |\lambda_{H_0}(\omega)| \leq \sqrt{2} + c|\omega|.$$

We have that

$$\begin{aligned}\left| \prod_{s=2}^{\infty} \frac{\lambda_{H_0}(-2\pi v/2^s)}{\sqrt{2}} \right| &\leq \prod_{s=2}^{\infty} \frac{\sqrt{2} + c(\pi v/2^{s-1})}{\sqrt{2}} = \prod_{s=2}^{\infty} 1 + c(\pi v/2^{s-3/2}) \\ &\leq \prod_{s=2}^{\infty} e^{c(\pi v/2^{s-3/2})} = e^{\sum_{s=2}^{\infty} c(\pi v/2^{s-3/2})} = e^{c\pi v 2/\sqrt{2}}.\end{aligned}$$

It follows that the product converges. Differentiability follows from the fact that the infinite product converges uniformly and absolutely on compact sets. The last statement follows from the expressions, and from that $\lambda_{G_0}(0) = \lambda_{H_0}(0) = \sqrt{2}$. ■

We also need to state why the g_N, h_N are CTFT's of functions u_N, v_N . Clearly, if $g(v)$ is integrable and decays faster than $|v|^{-1/2}$, then $g \in L^2(\mathbb{R})$, and it follows from standard Fourier analysis that there exists a unique $\phi \in L^2(\mathbb{R})$ so that $\hat{\phi}(v) = g(v)$. Moreover, ϕ can be recovered by

$$\phi(t) = \int_{-\infty}^{\infty} g(v) e^{2\pi i vt} dv.$$

If also g decays faster than v^{-1-z} for some integer $z \geq 0$, then ϕ is z times differentiable. To prove this, the integral above is bounded by $\int_{-\infty}^{\infty} v^{-1-z} e^{2\pi i vt} dv$, and if we

differentiate z times w.r.t. t we get $\int_{-\infty}^{\infty} v^{-1} e^{ivt} dv$, which converges. We therefore have the following result.

Theorem 7.10. It is possible to construct scaling functions $\phi, \tilde{\phi} \in L^2(\mathbb{R})$ from the filters G_0, H_0 when the infinite products $\lim_{N \rightarrow \infty} g_N(v), \lim_{N \rightarrow \infty} h_N(v)$ decay faster than $|v|^{-1/2}$. If they also decay faster than v^{-1-z} for some integer $z \geq 0$, then $\phi, \tilde{\phi}$ will also be z times differentiable.

We also prove the following result, which will turn out to be useful in the next section.

Theorem 7.11. Assume that

$$\lambda_{G_0}(v) = \left(\frac{1+e^{-iv}}{2} \right)^{N_1} r_1(v) \quad \lambda_{H_0}(v) = \left(\frac{1+e^{-iv}}{2} \right)^{N_2} r_2(v)$$

and that the $r_i(v)$ are bounded by $B_i < 2^{\gamma_i}$ on $[0, 2\pi]$. Then

$$g(v) \leq C_1(1+|v|)^{-(N_1-\gamma_1)-\epsilon_1} \quad h(v) \leq C_2(1+|v|)^{-(N_2-\gamma_2)-\epsilon_2}$$

where $\epsilon_i = \gamma_i - \log_2 B_i$. In particular, ϕ is $N_1 - \gamma_1$ times differentiable, and $\tilde{\phi}$ is $N_2 - \gamma_2$ times differentiable.

We will not prove this. We refer instead to [8] for a proof. In that paper it is also proved that there is a connection between the nonzero filter coefficients of H_0, G_0 , and the supports of $\phi, \tilde{\phi}$. Similarly, there exist $\psi, \tilde{\psi}$ so that $\lambda_2 = \lambda_{s-\psi(-t)}, \lambda_4 = \lambda_{s-\tilde{\psi}(-t)}$. Finally we prove the following.

7.2 Vanishing moments

The scaling functions and mother wavelets we constructed in Chapter 5 were very simple. They were however, enough to provide scaling functions which were differentiable. This may clearly be important for signal approximation, at least in cases where we know certain things about the regularity of the functions we approximate. However, there seemed to be nothing which dictated how the mother wavelet should be chosen in order to be useful. To see that this may pose a problem, consider the mother wavelet we chose for piecewise linear functions. Set $N = 1$ and consider the space V_{10} , which has dimension 2^{10} . When we apply a DWT, we start with a function $g_{10} \in V_{10}$. This may be a very good representation of the underlying data. However, when we compute g_{m-1} we just pick every other coefficient from g_m . By the time we get to g_0 we are just left with the first and last coefficient from g_{10} . In some situations this may be adequate, but usually not.

Idea 7.12. We would like a wavelet basis to be able to represent f efficiently. By this we mean that the approximation $f^{(m)} = \sum_n c_{0,n} \phi_{0,n} + \sum_{m' < m, n} w_{m',n} \psi_{m',n}$ to f from Observation 7.15 should converge quickly for the f we work with, as m increases. This means that, with relatively few $\psi_{m,n}$, we can create good approximations of f .

In this section we will address a property which the mother wavelet must fulfill in order to be useful in this respect. To motivate this property, let us first use Theorem 7.7 to write $f \in V_m$ as

$$f = \sum_{n=0}^{N-1} \langle f, \tilde{\phi}_{0,n} \rangle \phi_{0,n} + \sum_{r=0}^{m-1} \sum_{n=0}^{2^r N-1} \langle f, \tilde{\psi}_{r,n} \rangle \psi_{r,n}. \quad (7.7)$$

If f is s times differentiable, it can be represented as $f = P_s(x) + Q_s(x)$, where P_s is a polynomial of degree s , and Q_s is a function which is very small (P_s could for instance be a Taylor series expansion of f). If in addition $\langle t^k, \tilde{\psi} \rangle = 0$, for $k = 1, \dots, s$, we have also that $\langle t^k, \tilde{\psi}_{r,t} \rangle = 0$ for $r \leq s$, so that $\langle P_s, \tilde{\psi}_{r,t} \rangle = 0$ also. This means that Equation (7.7) can be written

$$\begin{aligned} f &= \sum_{n=0}^{N-1} \langle P_s + Q_s, \tilde{\phi}_{0,n} \rangle \phi_{0,n} + \sum_{r=0}^{m-1} \sum_{n=0}^{2^r N-1} \langle P_s + Q_s, \tilde{\psi}_{r,n} \rangle \psi_{r,n} \\ &= \sum_{n=0}^{N-1} \langle P_s + Q_s, \tilde{\phi}_{0,n} \rangle \phi_{0,n} + \sum_{r=0}^{m-1} \sum_{n=0}^{2^r N-1} \langle P_s, \tilde{\psi}_{r,n} \rangle \psi_{r,n} + \sum_{r=0}^{m-1} \sum_{n=0}^{2^r N-1} \langle Q_s, \tilde{\psi}_{r,n} \rangle \psi_{r,n} \\ &= \sum_{n=0}^{N-1} \langle f, \tilde{\phi}_{0,n} \rangle \phi_{0,n} + \sum_{r=0}^{m-1} \sum_{n=0}^{2^r N-1} \langle Q_s, \tilde{\psi}_{r,n} \rangle \psi_{r,n}. \end{aligned}$$

Here the first sum lies in V_0 . We see that the wavelet coefficients from W_r are $\langle Q_s, \tilde{\psi}_{r,n} \rangle$, which are very small since Q_s is small. This means that the detail in the different spaces W_r is very small, which is exactly what we aimed for. Let us summarize this as follows:

Theorem 7.13 (Vanishing moments). If a function $f \in V_m$ is r times differentiable, and $\tilde{\psi}$ has r vanishing moments, then f can be approximated well from V_0 . Moreover, the quality of this approximation improves when r increases.

Having many vanishing moments is thus very good for compression, since the corresponding wavelet basis is very efficient for compression. In particular, if f is a polynomial of degree less than or equal to $k-1$ and $\tilde{\psi}$ has k vanishing moments, then the detail coefficients $w_{m,n}$ are exactly 0. Since (ϕ, ψ) and $(\tilde{\phi}, \tilde{\psi})$ both are wavelet bases, it is equally important for both to have vanishing moments. We will in the following concentrate on the number of vanishing moments of ψ .

The Haar wavelet has one vanishing moment, since $\tilde{\psi} = \psi$ and $\int_0^N \psi(t) dt = 0$ as we noted in Observation 5.15. It is an exercise to see that the Haar wavelet has only one vanishing moment, i.e. $\int_0^N t \psi(t) dt \neq 0$.

Theorem 7.14. Assume that the filters are chosen so that the scaling functions exist. Then the following hold

1. The number of vanishing moments of $\tilde{\psi}$ equals the multiplicity of a zero at $\omega = \pi$ for $\lambda_{G_0}(\omega)$.
2. The number of vanishing moments of ψ equals the multiplicity of a zero at $\omega = \pi$ for $\lambda_{H_0}(\omega)$.

number of vanishing moments of ψ , $\tilde{\psi}$ equal the multiplicities of the zeros of the frequency responses $\lambda_{H_0}(\omega)$, $\lambda_{G_0}(\omega)$, respectively, at $\omega = \pi$.

In other words, the flatter the frequency responses $\lambda_{H_0}(\omega)$ and $\lambda_{G_0}(\omega)$ are near high frequencies ($\omega = \pi$), the better the wavelet functions are for approximation of functions. This is analogous to the smoothing filters we constructed previously, where the use of values from Pascals triangle resulted in filters which behaved like the constant function one at low frequencies. The frequency response for the Haar wavelet had just a simple zero at π , so that it cannot represent functions efficiently. The result also proves why we should consider G_0, H_0 as lowpass filters, G_1, H_1 as highpass filters.

Proof: We have that

$$\lambda_{s_{-\tilde{\psi}(-t)}}(v) = - \int_{-\infty}^{\infty} \tilde{\psi}(-t) e^{-2\pi i v t} dt. \quad (7.8)$$

By differentiating this expression k times w.r.t. v (differentiate under the integral sign) we get

$$(\lambda_{s_{-\tilde{\psi}(-t)}})^{(k)}(v) = - \int (-2\pi i t)^k \tilde{\psi}(t) e^{-2\pi i v t} dt. \quad (7.9)$$

Evaluating this at $v = 0$ gives

$$(\lambda_{s_{-\tilde{\psi}(-t)}})^{(k)}(0) = - \int (-2\pi i t)^k \tilde{\psi}(t) dt. \quad (7.10)$$

From this expression it is clear that the number of vanishing moments of $\tilde{\psi}$ equals the multiplicity of a zero at $v = 0$ for $\lambda_{s_{-\tilde{\psi}(-t)}}(v)$, which we have already shown equals the multiplicity of a zero at $\omega = 0$ for $\lambda_{H_1}(\omega)$. Similarly it follows that the number of vanishing moments of ψ equals the multiplicity of a zero at $\omega = 0$ for $\lambda_{G_1}(\omega)$. Since we know that $\lambda_{G_0}(\omega)$ has the same number of zeros at π as $\lambda_{H_1}(\omega)$ has at 0, and $\lambda_{H_0}(\omega)$ has the same number of zeros at π as $\lambda_{G_1}(\omega)$ has at 0, the result follows. ■

These results explain how we can construct $\phi, \psi, \tilde{\phi}, \tilde{\psi}$ from FIR-filters H_0, G_1, G_0, G_1 satisfying the perfect reconstruction condition. Also, the results explain how we can obtain such functions with as much differentiability and as many vanishing moments as we want. We will use these results in the next section to construct interesting wavelets. There we will also cover how we can construct the simplest possible such filters.

There are some details which have been left out in this section: We have not addressed why the wavelet bases we have constructed are linearly independent, and

why they span $L^2(\mathbb{R})$. Dual Riesz bases. These details are quite technical, and we refer to [5] for them. Let us also express what we have found in terms of analog filters.

Observation 7.15. Let

$$f(t) = \sum_n c_{m,n} \phi_{m,n} = \sum_n c_{0,n} \phi_{0,n} + \sum_{m' < m, n} w_{m',n} \psi_{m',n} \in V_m.$$

$c_{m,n}$ and $w_{m,n}$ can be computed by sampling the output of an analog filter. To be more precise,

$$\begin{aligned} c_{m,n} &= \langle f, \tilde{\phi}_{m,n} \rangle = \int_0^N f(t) \tilde{\phi}_{m,n}(t) dt = \int_0^N (-\tilde{\phi}_{m,0}(-t)) f(2^{-m}n - t) dt \\ w_{m,n} &= \langle f, \tilde{\psi}_{m,n} \rangle = \int_0^N f(t) \tilde{\psi}_{m,n}(t) dt = \int_0^N (-\tilde{\psi}_{m,0}(-t)) f(2^{-m}n - t) dt. \end{aligned}$$

In other words, $c_{m,n}$ can be obtained by sampling $s_{-\tilde{\phi}_{m,0}(-t)}(f(t))$ at the points $2^{-m}n$, $w_{m,n}$ by sampling $s_{-\tilde{\psi}_{m,0}(-t)}(f(t))$ at $2^{-m}n$, where the analog filters $s_{-\tilde{\phi}_{m,0}(-t)}$, $s_{-\tilde{\psi}_{m,0}(-t)}$ were defined in Theorem 1.40, i.e.

$$s_{-\tilde{\phi}_{m,0}(-t)}(f(t)) = \int_0^N (-\tilde{\phi}_{m,0}(-s)) f(t-s) ds \quad (7.11)$$

$$s_{-\tilde{\psi}_{m,0}(-t)}(f(t)) = \int_0^N (-\tilde{\psi}_{m,0}(-s)) f(t-s) ds. \quad (7.12)$$

A similar statement can be made for $\tilde{f} \in \tilde{V}_m$. Here the convolution kernels of the filters were as before, with the exception that ϕ, ψ were replaced by $\tilde{\phi}, \tilde{\psi}$. Note also that, if the functions $\tilde{\phi}, \tilde{\psi}$ are symmetric, we can increase the precision in the DWT with the method of symmetric extension also in this more general setting.

Exercises for Section 7.2

- Define the filters h_0, h_1, g_0, g_1 as in Example 6.16, and let us consider the following code, which shows how the cascade algorithm can be used to plot the scaling functions and the mother wavelet of the alternative piecewise linear wavelet and its dual wavelet, as shown in Figure 7.1.

```
m=10;
t=linspace(-2,6,8*2^m);

coodsvm=IDWTImpl(g0,g1,[0; 0; 1; 0; 0; 0; 0; 0; zeros(8*2^m-8,1)],m);
plot(t,coodsvm*2^(m/2)) % phi function

coodsvm=IDWTImpl(g0,g1,[zeros(8,1);...
0; 0; 1; 0; 0; 0; 0; 0; zeros(8*2^m-16,1)],m);
```

```

plot(t,coordsvm*2^(m/2)) % psi function

coordsvm=IDWTImpl(h0,h1,[0; 0; 1; 0; 0; 0; 0; 0; zeros(8*2^m-8,1)],m);
plot(t,coordsvm*2^(m/2)) % tildephi function

coordsvm=IDWTImpl(h0,h1,[zeros(8,1);...
    0; 0; 1; 0; 0; 0; 0; 0; zeros(8*2^m-16,1)],m);
plot(t,coordsvm*2^(m/2)) % tildepsi function

```

- a. Explain that the input to IDWTImpl in the code above are the coordinates of $\phi_{0,2}$, $\psi_{0,2}$, $\tilde{\phi}_{0,2}$, and $\tilde{\psi}_{0,2}$ in the basis $(\phi_0, \psi_0, \psi_1, \psi_2, \dots, \psi_{m-1})$, respectively.
- b. If you changed to the coordinates of $\psi_{0,0}$ instead in a., you would see a different function. Try to explain why this is the case. Hint: This has to do with the fact that the function IDWTImpl performs a symmetric extension of the signal. Try to write down what the symmetric extension of this function would be.
- c. In the code you see that all values are scaled with the factor $2^{m/2}$ before they are plotted. Can you think out an explanation to why this is done?

7.3 Characterization of wavelets w.r.t. number of vanishing moments

We have seen that wavelets are particularly suitable for approximation of functions when the mother wavelet or the dual mother wavelet have vanishing moments. The more vanishing moments they have, the more attractive they are. In this section we will attempt to characterize wavelets which have a given number of vanishing moments. In particular we will characterize the simplest such, those where the filters have few filters coefficients.

There are two particular cases we will look at. First we will consider the case when all filters are symmetric. Then we will look at the case of orthonormal wavelets. It turns out that these two cases are mutually disjoint (except for trivial examples), but that there is a common result which can be used to characterize the solutions to both problems. We will state the results in terms of the multiplicities of the zeros of $\lambda_{H_0}, \lambda_{G_0}$ at π , which we proved are the same as the number of vanishing moments.

7.3.1 Symmetric filters

The main result when the filters are symmetric looks as follows.

Theorem 7.16. Assume that H_0, H_1, G_0, G_1 are the filters of a wavelet, and that

1. the filters are symmetric,

2. λ_{H_0} has a zero of multiplicity N_1 at π ,
3. λ_{G_0} has a zero of multiplicity N_2 at π .

Then N_1 and N_2 are even, and there exist a polynomial Q which satisfies

$$u^{(N_1+N_2)/2} Q(1-u) + (1-u)^{(N_1+N_2)/2} Q(u) = 2. \quad (7.13)$$

so that $\lambda_{H_0}(\omega), \lambda_{G_0}(\omega)$ can be written on the form

$$\lambda_{H_0}(\omega) = \left(\frac{1}{2}(1+\cos\omega)\right)^{N_1/2} Q_1\left(\frac{1}{2}(1-\cos\omega)\right) \quad (7.14)$$

$$\lambda_{G_0}(\omega) = \left(\frac{1}{2}(1+\cos\omega)\right)^{N_2/2} Q_2\left(\frac{1}{2}(1-\cos\omega)\right), \quad (7.15)$$

where $Q = Q_1 Q_2$.

Proof: Since the filters are symmetric, $\lambda_{H_0}(\omega) = \lambda_{H_0}(-\omega)$ and $\lambda_{G_0}(\omega) = \lambda_{G_0}(-\omega)$. Since $e^{in\omega} + e^{-in\omega} = 2\cos(n\omega)$, and since $\cos(n\omega)$ is the real part of $(\cos\omega + i\sin\omega)^n$, which is a polynomial in $\cos^k\omega \sin^l\omega$ with l even, and since $\sin^2\omega = 1 - \cos^2\omega$, λ_{H_0} and λ_{G_0} can both be written on the form $P(\cos\omega)$, with P a real polynomial.

Note that a zero at π in $\lambda_{H_0}, \lambda_{G_0}$ corresponds to a factor of the form $1 + e^{-i\omega}$, so that we can write

$$\lambda_{H_0}(\omega) = \left(\frac{1+e^{-i\omega}}{2}\right)^{N_1} f(e^{i\omega}) = e^{-iN_1\omega/2} \cos^{N_1}(\omega/2) f(e^{i\omega}),$$

where f is a polynomial. In order for this to be real, we must have that $f(e^{i\omega}) = e^{iN_1\omega/2} g(e^{i\omega})$ where g is real-valued, and then we can write $g(e^{i\omega})$ as a real polynomial in $\cos\omega$. This means that $\lambda_{H_0}(\omega) = \cos^{N_1}(\omega/2) P_1(\cos\omega)$, and similarly for $\lambda_{G_0}(\omega)$. Clearly this can be a polynomial in $e^{i\omega}$ only if N_1 is even. Both N_1 and N_2 must then be even, and we can write

$$\begin{aligned} \lambda_{H_0}(\omega) &= \cos^{N_1}(\omega/2) P_1(\cos\omega) = (\cos^2(\omega/2))^{N_1/2} P_1(1-2\sin^2(\omega/2)) \\ &= (\cos^2(\omega/2))^{N_1/2} Q_1(\sin^2(\omega/2)), \end{aligned}$$

where we have used that $\cos\omega = 1 - 2\sin^2(\omega/2)$, and defined Q_1 by the relation $Q_1(x) = P_1(1-2x)$. Similarly we can write $\lambda_{G_0}(\omega) = (\cos^2(\omega/2))^{N_2/2} Q_2(\sin^2(\omega/2))$ for another polynomial Q_2 . Using the identities

$$\cos^2 \frac{\omega}{2} = \frac{1}{2}(1+\cos\omega) \quad \sin^2 \frac{\omega}{2} = \frac{1}{2}(1-\cos\omega),$$

we see that λ_{H_0} and λ_{G_0} satisfy equations (7.14) and (7.15). With $Q = Q_1 Q_2$, Equa-

tion (6.25) can now be rewritten as

$$\begin{aligned} 2 &= \lambda_{G_0}(\omega)\lambda_{H_0}(\omega) + \lambda_{G_0}(\omega+\pi)\lambda_{H_0}(\omega+\pi) \\ &= (\cos^2(\omega/2))^{(N_1+N_2)/2} Q(\sin^2(\omega/2)) + (\cos^2((\omega+\pi)/2))^{(N_1+N_2)/2} Q(\sin^2((\omega+\pi)/2)) \\ &= (\cos^2(\omega/2))^{(N_1+N_2)/2} Q(\sin^2(\omega/2)) + (\sin^2(\omega/2))^{(N_1+N_2)/2} Q(\cos^2(\omega/2)) \\ &= (\cos^2(\omega/2))^{(N_1+N_2)/2} Q(1 - \cos^2(\omega/2)) + (1 - \cos^2(\omega/2))^{(N_1+N_2)/2} Q(\cos^2(\omega/2)) \end{aligned}$$

Setting $u = \cos^2(\omega/2)$ we see that Q must fulfill the equation

$$u^{(N_1+N_2)/2} Q(1-u) + (1-u)^{(N_1+N_2)/2} Q(u) = 2,$$

which is Equation (7.13). This completes the proof. \blacksquare

While this result characterizes all wavelets with a given number of vanishing moments, it does not say which of these have fewest filter coefficients. The polynomial Q decides the length of the filters H_0, G_0 , however, so that what we need to do is to find the polynomial Q of smallest degree. In this direction, note first that the polynomials $u^{N_1+N_2}$ and $(1-u)^{N_1+N_2}$ have no zeros in common. Bezouts theorem, proved in Section 7.3.3, states that the equation

$$u^N q_1(u) + (1-u)^N q_2(u) = 1 \quad (7.16)$$

has unique solutions q_1, q_2 with $\deg(q_1), \deg(q_2) < (N_1 + N_2)/2$. To find these solutions, substituting $1-u$ for u gives the following equations:

$$\begin{aligned} u^N q_1(u) + (1-u)^N q_2(u) &= 1 \\ u^N q_2(1-u) + (1-u)^N q_1(1-u) &= 1, \end{aligned}$$

and uniqueness in Bezouts theorem gives that $q_1(u) = q_2(1-u)$, and $q_2(u) = q_1(1-u)$. Equation (7.16) can thus be stated as

$$u^N q_2(1-u) + (1-u)^N q_2(u) = 1,$$

and comparing with Equation (7.13) (set $N = (N_1 + N_2)/2$) we see that $Q(u) = 2q_2(u)$. $u^N q_1(u) + (1-u)^N q_2(u) = 1$ now gives

$$\begin{aligned} q_2(u) &= (1-u)^{-N}(1-u^N q_1(u)) = (1-u)^{-N}(1-u^N q_2(1-u)) \\ &= \left(\sum_{k=0}^{N-1} \binom{N+k-1}{k} u^k + O(u^N) \right) (1-u^N q_2(1-u)) \\ &= \sum_{k=0}^{N-1} \binom{N+k-1}{k} u^k + O(u^N), \end{aligned}$$

where we have used the first N terms in the Taylor series expansion of $(1-u)^{-N}$ around 0. Since q_2 is a polynomial of degree $N-1$, we must have that

$$Q(u) = 2q_2(u) = 2 \sum_{k=0}^{N-1} \binom{N+k-1}{k} u^k. \quad (7.17)$$

Define $Q^{(N)}(u) = 2 \sum_{k=0}^{N-1} \binom{N+k-1}{k} u^k$. The first $Q^{(N)}$ are

$$\begin{aligned} Q^{(0)}(u) &= 2 \\ Q^{(1)}(u) &= 2 + 4u \\ Q^{(2)}(u) &= 2 + 6u + 12u^2 \\ Q^{(3)}(u) &= 2 + 8u + 20u^2 + 40u^3, \end{aligned}$$

for which we compute

$$\begin{aligned} Q^{(0)}\left(\frac{1}{2}(1 - \cos \omega)\right) &= 2 \\ Q^{(1)}\left(\frac{1}{2}(1 - \cos \omega)\right) &= -e^{-i\omega} + 4 - e^{i\omega} \\ Q^{(2)}\left(\frac{1}{2}(1 - \cos \omega)\right) &= \frac{3}{4}e^{-2i\omega} - \frac{9}{2}e^{-i\omega} + \frac{19}{2} - \frac{9}{2}e^{i\omega} + \frac{3}{4}e^{2i\omega} \\ Q^{(3)}\left(\frac{1}{2}(1 - \cos \omega)\right) &= -\frac{5}{8}e^{-3i\omega} + 5e^{-2i\omega} - \frac{131}{8}e^{-i\omega} + 26 - \frac{131}{8}e^{i\omega} + 5e^{2i\omega} - \frac{5}{8}e^{3i\omega}, \end{aligned}$$

Thus in order to construct wavelets where $\lambda_{H_0}, \lambda_{G_0}$ have as many zeros at π as possible, and where there are as few filter coefficients as possible, we need to compute the polynomials above, factorize them into polynomials Q_1 and Q_2 , and distribute these among λ_{H_0} and λ_{G_0} . Since we need real factorizations, we must in any case pair complex roots. If we do this we obtain the factorizations

$$\begin{aligned} Q^{(0)}\left(\frac{1}{2}(1 - \cos \omega)\right) &= 2 \\ Q^{(1)}\left(\frac{1}{2}(1 - \cos \omega)\right) &= \frac{1}{3.7321}(e^{i\omega} - 3.7321)(e^{-i\omega} - 3.7321) \\ Q^{(2)}\left(\frac{1}{2}(1 - \cos \omega)\right) &= \frac{3}{4} \frac{1}{9.4438}(e^{2i\omega} - 5.4255e^{i\omega} + 9.4438) \\ &\quad \times (e^{-2i\omega} - 5.4255e^{-i\omega} + 9.4438) \\ Q^{(3)}\left(\frac{1}{2}(1 - \cos \omega)\right) &= \frac{5}{8} \frac{1}{3.0407} \frac{1}{7.1495}(e^{i\omega} - 3.0407)(e^{2i\omega} - 4.0623e^{i\omega} + 7.1495) \\ &\quad \times (e^{-i\omega} - 3.0407)(e^{-2i\omega} - 4.0623e^{-i\omega} + 7.1495), \end{aligned} \tag{7.18}$$

The factors in these factorizations can be distributed as factors in the frequency responses of $\lambda_{H_0}(\omega)$, and $\lambda_{G_0}(\omega)$. One possibility is to let one of these frequency responses absorb all the factors, another possibility is to split the factors as evenly as possible across the two. When a frequency response absorbs more factors, the corresponding filter gets more filter coefficients. In the following examples, both factor distribution strategies will be encountered. Note that it is straightforward to use tools such as Matlab to factor Q into a product of polynomials Q_1 and Q_2 . First the `roots` function can be used to find the roots in the polynomials. Then the `conv` function can be used to multiply together factors corresponding to different roots, to obtain the coefficients in the polynomials Q_1 and Q_2 .

7.3.2 Orthonormal wavelets

Now we turn to the case of orthonormal wavelets, i.e. where $G_0 = (H_0)^T$, $G_1 = (H_1)^T$. For simplicity we will assume $d = 0$, $\alpha = -1$ in conditions (6.23) and (6.24) (this corresponded to requiring $\lambda_{H_1}(\omega) = -\overline{\lambda_{H_0}(\omega + \pi)}$ in the definition of alternative QMF filter banks). We will also assume for simplicity that G_0 is causal (the other solutions can be derived from this). We saw that the Haar wavelet was such an orthonormal wavelet. We have the following result:

Theorem 7.17. Assume that H_0, H_1, G_0, G_1 are the filters of an orthonormal wavelet (i.e. $H_0 = (G_0)^T$ and $H_1 = (G_1)^T$) which also is an alternative QMF filter bank (i.e. $\lambda_{H_1}(\omega) = -\overline{\lambda_{H_0}(\omega + \pi)}$). Assume also that $\lambda_{G_0}(\omega)$ has a zero of multiplicity N at π and that G_0 is causal. Then there exists a polynomial Q which satisfies

$$u^N Q(1-u) + (1-u)^N Q(u) = 2, \quad (7.19)$$

so that if f is another polynomial which satisfies $f(e^{i\omega})f(e^{-i\omega}) = Q\left(\frac{1}{2}(1-\cos\omega)\right)$, $\lambda_{G_0}(\omega)$ can be written on the form

$$\lambda_{G_0}(\omega) = \left(\frac{1+e^{-i\omega}}{2}\right)^N f(e^{-i\omega}), \quad (7.20)$$

We avoided stating $\lambda_{H_0}(\omega)$ in this result, since the relation $H_0 = (G_0)^T$ gives that $\lambda_{H_0}(\omega) = \overline{\lambda_{G_0}(\omega)}$. In particular, $\lambda_{H_0}(\omega)$ also has a zero of multiplicity N at π . That G_0 is causal is included to simplify the expression further.

Proof: The proof is very similar to the proof of Theorem 7.16. N vanishing moments means that we can write

$$\lambda_{G_0}(\omega) = \left(\frac{1+e^{-i\omega}}{2}\right)^N f(e^{-i\omega}) = (\cos(\omega/2))^N e^{-iN\omega/2} f(e^{-i\omega}),$$

where f is a real polynomial (we used causality here). Also

$$\lambda_{H_0}(\omega) = \overline{\lambda_{G_0}(\omega)} = (\cos(\omega/2))^N e^{iN\omega/2} f(e^{i\omega}).$$

Condition (6.25) now says that

$$\begin{aligned} 2 &= \lambda_{G_0}(\omega)\lambda_{H_0}(\omega) + \lambda_{G_0}(\omega + \pi)\lambda_{H_0}(\omega + \pi) \\ &= (\cos^2(\omega/2))^N f(e^{i\omega})f(e^{-i\omega}) + (\sin^2(\omega/2))^N f(e^{i(\omega+\pi)})f(e^{-i(\omega+\pi)}). \end{aligned}$$

Now, the function $f(e^{i\omega})f(e^{-i\omega})$ is symmetric around 0, so that it can be written on the form $P(\cos\omega)$ with P a polynomial, so that

$$\begin{aligned} 2 &= (\cos^2(\omega/2))^N P(\cos\omega) + (\sin^2(\omega/2))^N P(\cos(\omega + \pi)) \\ &= (\cos^2(\omega/2))^N P(1 - 2\sin^2(\omega/2)) + (\sin^2(\omega/2))^N P(1 - 2\cos^2(\omega/2)). \end{aligned}$$

If we as in the proof of Theorem 7.16 define Q by $Q(x) = P(1 - 2x)$, we can write this as

$$(\cos^2(\omega/2))^N Q(\sin^2(\omega/2)) + (\sin^2(\omega/2))^N Q(\cos^2(\omega/2)) = 2,$$

which again gives Equation (7.13) for finding Q . What we thus need to do is to compute the polynomial $Q\left(\frac{1}{2}(1 - \cos \omega)\right)$ as before, and consider the different factorizations of this on the form $f(e^{i\omega})f(e^{-i\omega})$. Since this polynomial is symmetric, a is a root if and only $1/a$ is, and if and only if \bar{a} is. If the real roots are

$$b_1, \dots, b_m, 1/b_1, \dots, 1/b_m,$$

and the complex roots are

$$a_1, \dots, a_n, \bar{a}_1, \dots, \bar{a}_n \text{ and } 1/a_1, \dots, 1/a_n, \bar{1/a_1}, \dots, \bar{1/a_n},$$

we can write

$$\begin{aligned} & Q\left(\frac{1}{2}(1 - \cos \omega)\right) \\ &= K(e^{-i\omega} - b_1) \dots (e^{-i\omega} - b_m) \\ &\quad \times (e^{-i\omega} - a_1)(e^{-i\omega} - \bar{a}_1)(e^{-i\omega} - a_2)(e^{-i\omega} - \bar{a}_2) \dots (e^{-i\omega} - a_n)(e^{-i\omega} - \bar{a}_n) \\ &\quad \times (e^{i\omega} - b_1) \dots (e^{i\omega} - b_m) \\ &\quad \times (e^{i\omega} - a_1)(e^{i\omega} - \bar{a}_1)(e^{i\omega} - a_2)(e^{i\omega} - \bar{a}_2) \dots (e^{i\omega} - a_n)(e^{i\omega} - \bar{a}_n) \end{aligned}$$

where K is a constant. We now can define the polynomial f by

$$\begin{aligned} f(e^{i\omega}) &= \sqrt{K}(e^{i\omega} - b_1) \dots (e^{i\omega} - b_m) \\ &\quad \times (e^{i\omega} - a_1)(e^{i\omega} - \bar{a}_1)(e^{i\omega} - a_2)(e^{i\omega} - \bar{a}_2) \dots (e^{i\omega} - a_n)(e^{i\omega} - \bar{a}_n) \end{aligned}$$

in order to obtain a factorization $Q\left(\frac{1}{2}(1 - \cos \omega)\right) = f(e^{i\omega})f(e^{-i\omega})$. This concludes the proof. \blacksquare

In the previous proof we note that the polynomial f is not unique - we could pair the roots in many different ways. The new algorithm is thus as follows:

1. As before, write $Q\left(\frac{1}{2}(1 - \cos \omega)\right)$ as a polynomial in $e^{i\omega}$, and find the roots.
2. Split the roots into the two classes

$$\{b_1, \dots, b_m, a_1, \dots, a_n, \bar{a}_1, \dots, \bar{a}_n\}$$

and

$$\{1/b_1, \dots, 1/b_m, 1/a_1, \dots, 1/a_n, \bar{1/a_1}, \dots, \bar{1/a_n}\},$$

and form the polynomial f as above.

3. Compute $\lambda_{G_0}(\omega) = \left(\frac{1+e^{-i\omega}}{2}\right)^N f(e^{-i\omega})$.

Clearly the filters obtained with this strategy are not symmetric since f is not symmetric. In Section 7.6 we will take a closer look at wavelets constructed in this way.

7.3.3 The proof of Bezouts theorem

Theorem 7.18. If p_1 and p_2 are two polynomials, of degrees n_1 and n_2 respectively, with no common zeros, then there exist unique polynomials q_1, q_2 , of degree less than n_2, n_1 , respectively, so that

$$p_1(x)q_1(x) + p_2(x)q_2(x) = 1. \quad (7.21)$$

Proof: We first establish the existence of q_1, q_2 satisfying Equation (7.21). Denote by $\deg(P)$ the degree of the polynomial P . Renumber the polynomials if necessary, so that $n_1 \geq n_2$. By polynomial division, we can now write

$$p_1(x) = a_2(x)p_2(x) + b_2(x),$$

where $\deg(a_2) = \deg(p_1) - \deg(p_2)$, $\deg(b_2) < \deg(p_2)$. Similarly, we can write

$$p_2(x) = a_3(x)b_2(x) + b_3(x),$$

where $\deg(a_3) = \deg(p_2) - \deg(b_2)$, $\deg(b_3) < \deg(b_2)$. We can repeat this procedure, so that we obtain a sequence of polynomials $a_n(x), b_n(x)$ so that

$$b_{n-1}(x) = a_{n+1}(x)b_n(x) + b_{n+1}(x), \quad (7.22)$$

where $\deg a_{n+1} = \deg(b_{n-1}) - \deg(b_n)$, $\deg(b_{n+1}) < \deg(b_n)$. Since $\deg(b_n)$ is strictly decreasing, we must have that $b_{N+1} = 0$ and $b_N \neq 0$ for some N , i.e. $b_{N-1}(x) = a_{N+1}(x)b_N(x)$. Since $b_{N-2} = a_N b_{N-1} + b_N$, it follows that b_{N-2} can be divided by b_N , and by induction that all b_n can be divided by b_N , in particular p_1 and p_2 can be divided by b_N . Since p_1 and p_2 have no common zeros, b_N must be a nonzero constant.

Using Equation (7.22), we can write recursively

$$\begin{aligned} b_N &= b_{N-2} - a_N b_{N-1} \\ &= b_{N-2} - a_N(b_{N-3} - a_{N-1} b_{N-2}) \\ &= (1 + a_N a_{N-1})b_{N-2} - a_N b_{N-3}. \end{aligned}$$

By induction we can write

$$b_N = a_{N,k}^{(1)} b_{N-k} + a_{N,k}^{(2)} b_{N-k-1}.$$

We see that the leading order term for $a_{N,k}^{(1)}$ is $a_N \cdots a_{N-k+1}$, which has degree

$$(\deg(b_{N-2}) - \deg(b_{N-1}) + \cdots + (\deg(b_{N-k-1}) - \deg(b_{N-k})) = \deg(b_{N-k-1}) - \deg(b_{N-1}),$$

while the leading order term for $a_{N,k}^{(2)}$ is $a_N \cdots a_{N-k+2}$, which similarly has order $\deg(b_{N-k}) - \deg(b_{N-1})$. For $k = N - 1$ we find

$$b_N = a_{N,N-1}^{(1)} b_1 + a_{N,N-1}^{(2)} b_0 = a_{N,N-1}^{(1)} p_2 + a_{N,N-1}^{(2)} p_1, \quad (7.23)$$

with $\deg(a_{N,N-1}^{(1)}) = \deg(p_1) - \deg(b_{N-1}) < \deg(p_1)$ (since by construction $\deg(b_{N-1}) > 0$), and $\deg(a_{N,N-1}^{(2)}) = \deg(p_2) - \deg(b_{N-1}) < \deg(p_2)$. From Equation (7.23) it follows that $q_1 = a_{N,N-1}^{(2)} / b_N$ and $q_2 a_{N,N-1}^{(1)} / b_N$ satisfies Equation (7.21), and that they satisfy the required degree constraints.

Now we turn to uniqueness of solutions q_1, q_2 . Assume that r_1, r_2 are two other solutions to Equation (7.21). Then

$$p_1(q_1 - r_1) + p_2(q_2 - r_2) = 0.$$

Since p_1 and p_2 have no zeros in common this means that every zero of p_2 is a zero of $q_1 - r_1$, with at least the same multiplicity. If $q_1 \neq r_1$, this means that $\deg(q_1 - r_1) \geq \deg(p_2)$, which is impossible since $\deg(q_1) < \deg(p_2)$, $\deg(r_1) < \deg(p_2)$. Hence $q_1 = r_1$. Similarly $q_2 = r_2$, establishing uniqueness. \blacksquare

Exercises for Section 7.3

1. Calculate the number of vanishing moment for the ψ as defined in Lemma ??.

7.4 A design strategy suitable for lossless compression

We choose $Q_1 = Q$, $Q_2 = 1$. In this case there is no need to find factors in Q . The filters in the filter factorization take the form

$$\begin{aligned}\lambda_{H_0}(\omega) &= \left(\frac{1}{2}(1 + \cos \omega)\right)^{N_1/2} Q\left(\frac{1}{2}(1 - \cos \omega)\right) \\ \lambda_{G_0}(\omega) &= \left(\frac{1}{2}(1 + \cos \omega)\right)^{N_2/2}.\end{aligned}$$

Since Q has degree $\frac{N_1+N_2}{2} - 1$, λ_{H_0} has degree $N_1 + N_2 - 2 = 2N_1 + N_2 - 2$, and λ_{G_0} has degree N_2 . These are both even numbers, so that the filters have odd length. The names of these filters are indexed by the filter lengths, and are called *Spline wavelets*, since, as we now will show, the scaling function for this design strategy is the B -spline of order N_2 : we have that

$$\lambda_{G_0}(\omega) = \frac{1}{2^{N_2/2}} (1 + \cos \omega)^{N_2/2} = \cos(\omega/2)^{N_2}.$$

Letting s be the analog filter with convolution kernel ϕ we can as in Equation (??) write

$$\begin{aligned}\lambda_s(f) &= \lambda_s(f/2^k) \prod_{i=1}^k \frac{\lambda_{G_0}(2\pi f/2^i)}{2} \\ &= \lambda_s(f/2^k) \prod_{i=1}^k \frac{\cos^{N_2}(\pi f/2^i)}{2} \\ &= \lambda_s(f/2^k) \prod_{i=1}^k \left(\frac{\sin(2\pi f/2^i)}{2 \sin(\pi f/2^i)} \right)^{N_2} \\ &= \lambda_s(f/2^k) \left(\frac{\sin(\pi f)}{2^k \sin \pi f/2^k} \right)^{N_2},\end{aligned}$$

where we have used the identity $\cos \omega = \frac{\sin(2\omega)}{2 \sin \omega}$. If we here let $k \rightarrow \infty$, and use the identity $\lim_{f \rightarrow 0} \frac{\sin f}{f} = 1$, we get that

$$\lambda_s(f) = \lambda_s(0) \left(\frac{\sin(\pi f)}{\pi f} \right)^{N_2}.$$

On the other hand, the frequency response of $\chi_{[-1/2, 1/2]}(t)$

$$\begin{aligned}&= \int_{-1/2}^{1/2} e^{-2\pi i f t} dt = \left[\frac{1}{-2\pi i f} e^{-2\pi i f t} \right]_{-1/2}^{1/2} \\ &= \frac{1}{-2\pi i f} (e^{-\pi i f} - e^{\pi i f}) = \frac{1}{-2\pi i f} 2i \sin(-\pi f) \\ &= \frac{\sin(\pi f)}{\pi f}.\end{aligned}$$

Due to this $\left(\frac{\sin(\pi f)}{\pi f} \right)^{N_2}$ is the frequency response of $*_{k=1}^{N_2} \chi_{[-1/2, 1/2]}(t)$. By the uniqueness of the frequency response we have that $\phi(t) = \hat{\phi}(0) *_{k=1}^{N_2} \chi_{[-1/2, 1/2]}(t)$. In Exercise 2 you will be asked to show that this scaling function gives rise to the multiresolution analysis of functions which are piecewise polynomials which are differentiable at the borders, also called *splines*. This explains why this type of wavelet is called a spline wavelet. To be more precise, the resolution spaces are as follows

Definition 7.19 (Resolution spaces of piecewise polynomials). We define V_m as the subspace of functions which are $r - 1$ times continuously differentiable and equal to a polynomial of degree r on any interval of the form $[n2^{-m}, (n+1)2^{-m}]$.

Note that the piecewise linear wavelet can be considered as the first Spline wavelet. This is further considered in the following example.

Example 7.20. For the case of $N_1 = N_2 = 2$ when the first design strategy is used,

equations (7.14) and (7.15) take the form

$$\begin{aligned}\lambda_{G_0}(\omega) &= \frac{1}{2}(1 + \cos \omega) = \frac{1}{4}e^{i\omega} + \frac{1}{2} + \frac{1}{4}e^{-i\omega} \\ \lambda_{H_0}(\omega) &= \frac{1}{2}(1 + \cos \omega)Q^{(1)}\left(\frac{1}{2}(1 - \cos \omega)\right) \\ &= \frac{1}{4}(2 + e^{i\omega} + e^{-i\omega})(4 - e^{i\omega} - e^{-i\omega}) \\ &= -\frac{1}{4}e^{2i\omega} + \frac{1}{2}e^{i\omega} + \frac{3}{2} + \frac{1}{2}e^{-i\omega} - \frac{1}{4}e^{-2i\omega}.\end{aligned}$$

The filters G_0, H_0 are thus

$$\begin{aligned}G_0 &= \left\{ \frac{1}{4}, \underline{\frac{1}{2}}, \frac{1}{4} \right\} \\ H_0 &= \left\{ -\frac{1}{4}, \frac{1}{2}, \underline{\frac{3}{2}}, \frac{1}{2}, -\frac{1}{4} \right\}\end{aligned}$$

The length of the filters are 3 and 5 in this case, so that this wavelet is called the *Spline 5/3 wavelet*. Up to a constant, the filters are seen to be the same as those of the alternative piecewise linear wavelet, see Example 6.16. Now, how do we find the filters (G_1, H_1) ? Previously we saw how to find the constant α in Theorem 6.14 when we knew one of the two pairs $(G_0, G_1), (H_0, H_1)$. This was the last part of information we needed in order to construct the other two filters. Here we know (G_0, H_0) instead. In this case it is even easier to find (G_1, H_1) since we can set $\alpha = 1$. This means that (G_1, H_1) can be obtained simply by adding alternating signs to (G_0, H_0) , i.e. they are the corresponding highpass filters. We thus can set

$$\begin{aligned}G_1 &= \left\{ -\frac{1}{4}, -\frac{1}{2}, \underline{\frac{3}{2}}, -\frac{1}{2}, -\frac{1}{4} \right\} \\ H_1 &= \left\{ -\frac{1}{4}, \underline{\frac{1}{2}}, -\frac{1}{4} \right\}.\end{aligned}$$

We have now found all the filters. It is clear that the scaling function and mother wavelet of this wavelet equals those of the alternative piecewise linear wavelet, up to a constant. ♣

The coefficients for the Spline wavelets are always dyadic fractions, and are therefore suitable for lossless compression, as they can be computed using low precision arithmetic and bitshift operations. The particular Spline wavelet from Example 7.20 is used for lossless compression in the JPEG2000 standard.

Exercises for Section 7.4

1. Write code which computes the coefficients in $(\frac{1}{2}(1 \pm \cos \omega))^k$ as a polynomial in $e^{i\omega}$. Since $1 \pm \cos \omega = 1 \pm \frac{1}{2}e^{i\omega} \pm \frac{1}{2}e^{-i\omega}$ you should think of it as the sequence $\{\pm \frac{1}{2}, 1, \pm \frac{1}{2}\}$, and then use the `conv`-function to compute the coefficients in the k 'th power.

2. Show that $B_r(t) = *_{k=1}^r \chi_{[-1/2, 1/2]}(t)$ is $r-2$ times differentiable, and equals a polynomial of degree $r-1$ on subintervals of the form $[n, n+1]$. Explain why these functions can be used as basis for the spaces V_j of functions which are piecewise polynomials of degree $r-1$ on intervals of the form $[n2^{-m}, (n+1)2^{-m}]$, and $r-2$ times differentiable. B_r is also called the B -spline of order r .

7.5 A design strategy suitable for lossy compression

The factors of Q are split evenly among Q_1 and Q_2 . In this case we need to factorize Q into a product of real polynomials. This can be done by finding all roots, and pairing the complex conjugate roots into real second degree polynomials (if Q is real, its roots come in conjugate pairs), and then distribute these as evenly as possible among Q_1 and Q_2 . These filters are called the CDF-wavelets, after Cohen, Daubechies, and Feauveau, who discovered them.

Example 7.21. We choose $N_1 = N_2 = 4$. In Equation (7.18) we pair inverse terms to obtain

$$\begin{aligned} Q^{(3)}\left(\frac{1}{2}(1-\cos\omega)\right) &= \frac{5}{8} \frac{1}{3.0407} \frac{1}{7.1495} (e^{i\omega} - 3.0407)(e^{-i\omega} - 3.0407) \\ &\quad \times (e^{2i\omega} - 4.0623e^{i\omega} + 7.1495)(e^{-2i\omega} - 4.0623e^{-i\omega} + 7.1495) \\ &= \frac{5}{8} \frac{1}{3.0407} \frac{1}{7.1495} (-3.0407e^{i\omega} + 10.2456 - 3.0407e^{-i\omega}) \\ &\quad \times (7.1495e^{2i\omega} - 33.1053e^{i\omega} + 68.6168 - 33.1053e^{-i\omega} + 7.1495e^{-2i\omega}). \end{aligned}$$

We can write this as $Q_1 Q_2$ with $Q_1(0) = Q_2(0)$ when

$$\begin{aligned} Q_1(\omega) &= -1.0326e^{i\omega} + 3.4795 - 1.0326e^{-i\omega} \\ Q_2(\omega) &= 0.6053e^{2i\omega} - 2.8026e^{i\omega} + 5.8089 - 2.8026e^{-i\omega} + 0.6053e^{-2i\omega}, \end{aligned}$$

from which we obtain

$$\begin{aligned} \lambda_{G_0}(\omega) &= \left(\frac{1}{2}(1+\cos\omega)\right)^2 Q_1(\omega) \\ &= -0.0645e^{3i\omega} - 0.0407e^{2i\omega} + 0.4181e^{i\omega} + 0.7885 \\ &\quad + 0.4181e^{-i\omega} - 0.0407e^{-2i\omega} - 0.0645e^{-3i\omega} \\ \lambda_{H_0}(\omega) &= \left(\frac{1}{2}(1+\cos\omega)\right)^2 40Q_2(\omega) \\ &= 0.0378e^{4i\omega} - 0.0238e^{3i\omega} - 0.1106e^{2i\omega} + 0.3774e^{i\omega} + 0.8527 \\ &\quad + 0.3774e^{-i\omega} - 0.1106e^{-2i\omega} - 0.0238e^{-3i\omega} + 0.0378e^{-4i\omega}. \end{aligned}$$

The filters G_0, H_0 are thus

$$\begin{aligned} G_0 &= \{-0.0645, -0.0407, 0.4181, 0.7885, 0.4181, -0.0407, -0.0645\} \\ H_0 &= \{0.0378, -0.0238, -0.1106, 0.3774, 0.8527, 0.3774, -0.1106, -0.0238, 0.0378\}. \end{aligned}$$

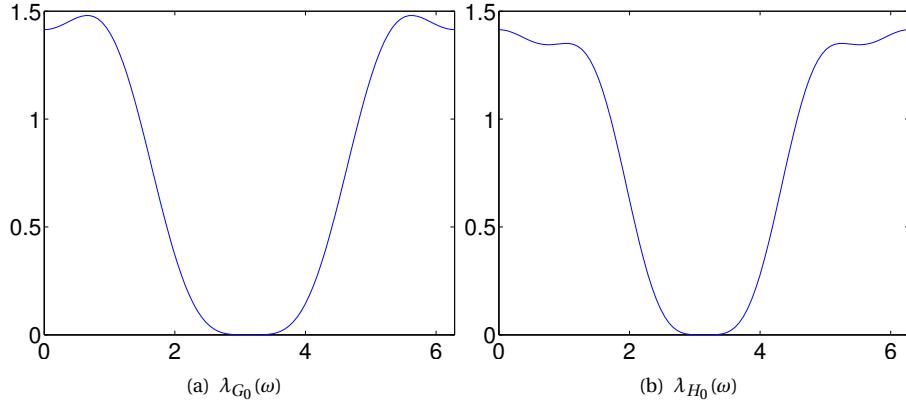


Figure 7.2: The frequency responses for the filters of Example 7.21.

The corresponding frequency responses are plotted in Figure 7.2. It is seen that both filters are lowpass filters also here, and that they are closer to an ideal bandpass filter. Here, the frequency response acts even more like the constant zero function close to π , proving that our construction has worked. We also get

$$G_1 = \{0.0378, 0.0238, -0.1106, -0.3774, \underline{0.8527}, -0.3774, -0.1106, 0.0238, 0.0378\}$$

$$H_1 = \{0.0645, -0.0407, -0.4181, \underline{0.7885}, -0.4181, -0.0407, 0.0645\}.$$

The length of the filters are 9 and 7 in this case, so that this wavelet is called the *CDF 9/7 wavelet*. This wavelet is for instance used for lossy compression with JPEG2000 since it gives a good tradeoff between complexity and compression.

In Example 6.16 we saw that we had analytical expressions for the scaling functions and the mother wavelet, but that we could not obtain this for the dual functions. For the CDF 9/7 wavelet it turns out that none of the four functions have analytical expressions. Let us therefore use the cascade algorithm, as we did in Example 6.16 to plot these functions. The results are shown in Figure 7.3. Again they have irregular shapes, but now at least the functions and dual functions more resemble each other. From Theorem 6.10 it follows that the supports of all these functions should be 6 and 7, which is compatible with these figures. ♣

In the above example there was a unique way of factoring Q into a product of real polynomials. For higher degree polynomials there is no unique way to form to distribute the factors, and we will not go into what strategy can be used for this. In general, the steps we must go through are as follows:

1. Compute the polynomial Q , and find its roots.
2. Pair complex conjugate roots into real second degree polynomials, and form polynomials Q_1, Q_2 .
3. Compute the coefficients in equations (7.14) and (7.15).

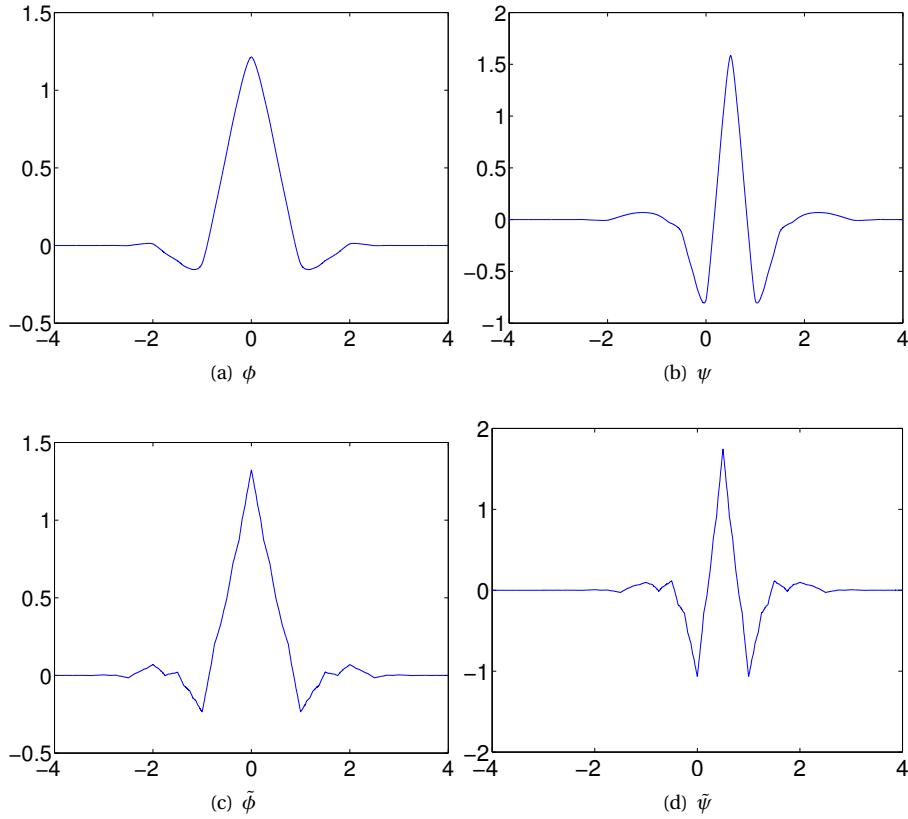


Figure 7.3: Scaling functions and mother wavelets for the CDF 9/7 wavelet.

Exercises for Section 7.5

1. Generate the plots from Figure 7.3 using the cascade algorithm. Reuse the code from Exercise ???.1 in order to achieve this.

7.6 Orthonormal wavelets

Since the filters here are not symmetric, the method of symmetric extension does not work in the same simple way as before. This partially explains why symmetric filters are used more often: They may not be as efficient in representing functions, since the corresponding basis is not orthogonal, but their simple implementation still makes them attractive.

The polynomials $Q^{(0)}$, $Q^{(1)}$, and $Q^{(2)}$ require no further action to obtain the factorization $f(e^{i\omega})f(e^{-i\omega}) = Q\left(\frac{1}{2}(1 - \cos \omega)\right)$. The polynomial $Q^{(3)}$ in Equation (7.18)

can be factored further as

$$Q^{(3)}\left(\frac{1}{2}(1-\cos\omega)\right) = \frac{5}{8} \frac{1}{3.0407} \frac{1}{7.1495} (e^{-3i\omega} - 7.1029e^{-2i\omega} + 19.5014^{-i\omega} - 21.7391) \\ \times (e^{3i\omega} - 7.1029e^{2i\omega} + 19.5014^{i\omega} - 21.7391),$$

which gives that $f(e^{i\omega}) = \sqrt{\frac{5}{8} \frac{1}{3.0407} \frac{1}{7.1495}} (e^{3i\omega} - 7.1029e^{2i\omega} + 19.5014^{i\omega} - 21.7391)$. This factorization is not unique, however. This gives the frequency response $\lambda_{G_0}(\omega) = \left(\frac{1+e^{-i\omega}}{2}\right)^N f(e^{-i\omega})$ as

$$\begin{aligned} & \frac{1}{2}(e^{-i\omega} + 1)\sqrt{2} \\ & \frac{1}{4}(e^{-i\omega} + 1)^2 \sqrt{\frac{1}{3.7321}} (e^{-i\omega} - 3.7321) \\ & \frac{1}{8}(e^{-i\omega} + 1)^3 \sqrt{\frac{3}{4} \frac{1}{9.4438}} (e^{-2i\omega} - 5.4255e^{-i\omega} + 9.4438) \\ & \frac{1}{16}(e^{-i\omega} + 1)^4 \sqrt{\frac{5}{8} \frac{1}{3.0407} \frac{1}{7.1495}} (e^{-3i\omega} - 7.1029e^{-2i\omega} + 19.5014^{-i\omega} - 21.7391), \end{aligned}$$

which gives the filters

$$\begin{aligned} G_0 &= (H_0)^T = (\underline{\sqrt{2}/2}, \sqrt{2}/2) \\ G_0 &= (H_0)^T = (-0.4830, -0.8365, -0.2241, 0.1294) \\ G_0 &= (H_0)^T = (0.3327, 0.8069, 0.4599, -0.1350, -0.0854, 0.0352) \\ G_0 &= (H_0)^T = (-0.2304, -0.7148, -0.6309, 0.0280, 0.1870, -0.0308, -0.0329, 0.0106) \end{aligned}$$

so that we get 2, 4, 6 and 8 filter coefficients in $G_0 = (H_0)^T$. We see that the filter coefficients when $N = 1$ are those of the Haar wavelet. The three next filters we have not seen before. The filter $G_1 = (H_1)^T$ can be obtained from the relation $\lambda_{G_1}(\omega) = -\bar{\lambda}_{G_0}(\omega + \pi)$, i.e. by reversing the elements and adding an alternating sign, plus an extra minus sign, so that

$$\begin{aligned} G_1 &= (H_1)^T = (\sqrt{2}/2, -\underline{\sqrt{2}/2}) \\ G_1 &= (H_1)^T = (0.1294, \underline{0.2241}, -0.8365, 0.4830) \\ G_1 &= (H_1)^T = (0.0352, \underline{0.0854}, -0.1350, -0.4599, 0.8069, -0.3327) \\ G_1 &= (H_1)^T = (0.0106, \underline{0.0329}, -0.0308, -0.1870, 0.0280, 0.6309, -0.7148, 0.2304). \end{aligned}$$

Frequency responses are shown in Figure 7.4. The frequency responses are now complex, so their magnitudes are plotted. Clearly these filters have lowpass characteristic. We also see that the highpass characteristics resemble the lowpass characteristics. We also see that the frequency response gets flatter near the high and low frequencies, as N increases. One can verify that this is the case also when N is increased further. The shapes for higher N are very similar to the frequency responses

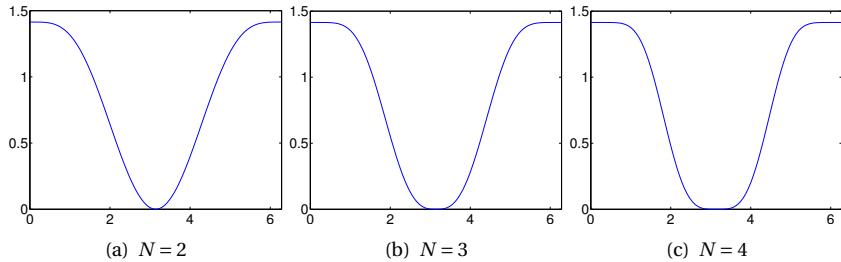


Figure 7.4: The magnitude of the frequency responses $\lambda_{G_0}(\omega) = \lambda_{H_0}(\omega)$ for the first orthonormal wavelets with vanishing moments.

of those filters used in the MP3 standard (see Figure 3.8). One difference is that the support of the latter is concentrated on a smaller set of frequencies.

The way we have defined the filters, the supports for both the scaling function and the mother wavelet must be $[0,3]$ when $N = 2$, $[0,5]$ when $N = 3$, and $[0,7]$ when $N = 4$. These are shown in Figure 7.5. Also in this case we have used the cascade algorithm to approximate the scaling functions and mother wavelets. The code for doing so is a bit different from before, since our IDWTImpl assumes symmetric filters, and the current filters are not symmetric (we thus need to reimplement IDWTImpl so that it does not perform symmetric extension of the input. Otherwise the previous code can be used unchanged). Since the number of nonzero filter coefficients are 4, 6, and 8 for these three wavelets, from Theorem 6.10 it follows that the support sizes of ϕ and ψ should be 3, 5, and 7, respectively, for these three wavelets. This seems to be the case in the figures also.

7.7 Summary

We started the section by showing how filters from filter bank matrices can give rise to scaling functions and mother wavelets. We saw that we obtained dual function pairs in this way, which satisfied a mutual property called biorthogonality. We then saw how differentiable scaling functions or mother wavelets with vanishing moments could be constructed, and we saw how we could construct the simplest such. These could be found in terms of the frequency responses of the involved filters. Finally we studied some examples with applications to image compression.

For the wavelets we constructed in this chapter, we also plotted the corresponding scaling functions and mother wavelets (see figures 7.1, 7.3, 7.5). The importance of these functions are that they are particularly suited for approximation of regular functions, and providing a compact representation of these functions which is localized in time. It seems difficult to guess that these strange shapes are connected to such approximation. Moreover, it may seem strange that, although these functions are useful, we can't write down exact expressions for them, and they are only approximated in terms of the Cascade Algorithm.

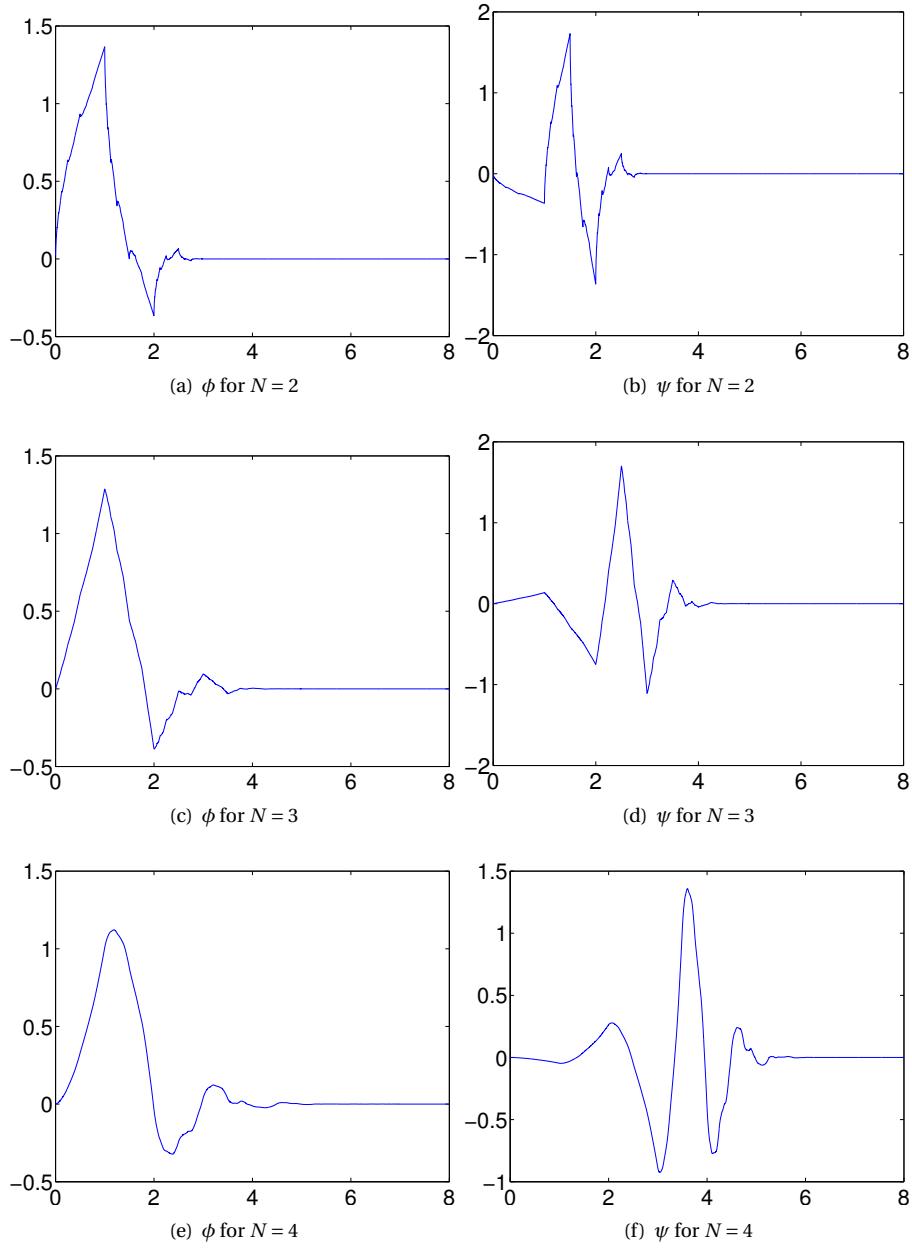


Figure 7.5: The scaling functions and mother wavelets for the first orthonormal wavelets with vanishing moments.

Chapter 8

The polyphase representation and wavelets

In Chapter 6 we saw that we could express wavelet transformations and more general transformations in terms of filters. Through this we obtained intuition for what information the different parts of a wavelet transformation represent, in terms of lowpass and highpass filters. We also obtained some insight into the filters used in the transformation used in the MP3 standard.

In this chapter we will look at another representation of such transformations, which we will call the *polyphase representation*. The polyphase representation will turn out to be useful for different reasons than the filter representation. First of all, with the polyphase representation transformations can be viewed as block matrices where the blocks are filters. This allows us to prove results in a different way than for filter bank transforms, since we can prove results through block matrix manipulation. There will be two major results we will prove in this way.

First, in Section 8.3 we obtain a factorization of a wavelet transformation into sparse matrices, called elementary lifting matrices. We will show that this factorization reduces the number of arithmetic operations. This is important: recall that we previously factored a filter into a product of smaller filters which is useful for efficient hardware implementations. But this did not address the fact that only every second component of the filters needs to be evaluated in the DWT, something any efficient implementation of the DWT should take into account. The factorization into sparse matrices will be called the *lifting factorization*. We will also see how we can use the polyphase representation to prove the remaining parts of Theorem 6.14.

Secondly, in Section 8.5 we will use the polyphase representation to analyze how the forward and reverse filter bank transforms from the MP3 standard can be chosen in order for us to have perfect or near perfect reconstruction. Actually, we will obtain a factorization of the polyphase representation into block matrices also here, and the conditions we need to put on the prototype filters will be clear from this.

8.1 The polyphase representation

Let us start by defining the basic concepts in the polyphase representation.

Definition 8.1 (Polyphase components and representation). Assume that S is a matrix, and that M is a number. By the *polyphase components* of S we mean the matrices $S^{(i,j)}$ defined by $S_{r_1, r_2}^{(i,j)} = S_{i+r_1M, j+r_2M}$, i.e. the matrices obtained by taking every M 'th component of S . When S is a (forward or reverse) filter bank transform, the $S^{(i,j)}$ are filters. By the *polyphase representation* of S we mean the block matrix with entries $S^{(i,j)}$.

In the following we will interchangeably represent a filter bank transform with its polyphase representation, and also write $A \leftrightarrow B$ when A is the polyphase representation arising from the filter bank transform B . Even though the polyphase components of a filter bank transform are filters, they are not the same as the filters G_i, H_i , since they are formed by taking every M 'th element from these in all possible ways.

The polyphase representation applies in particular for vectors. Since a vector \mathbf{x} only has one column, we write $\mathbf{x}^{(p)}$ for its polyphase components.

We have the following result on the polyphase representation. This result is easily proved from manipulation with block matrices, and is therefore left to the reader.

Theorem 8.2. Let A and B be (forward or reverse) filter bank transforms, and denote the corresponding polyphase components by $A^{(i,j)}, B^{(i,j)}$. The following hold

1. $C = AB$ is also a filter bank transform, with polyphase components $C^{(i,j)} = \sum_k A^{(i,k)} B^{(k,j)}$.
2. A^T is also a filter bank transform, with polyphase components $((A^T)^{(i,j)})_{k,l} = (A^{(j,i)})_{l,k}$.

Also, the polyphase components of the identity matrix is the $M \times M$ -block matrix with the identity matrix on the diagonal, and $\mathbf{0}$ elsewhere.

Example 8.3. Polyphase components can be defined for any matrix, regardless of whether it arises from a filter bank trasnfom. As an example, consider the 6×6 MRA-matrix

$$S = \begin{pmatrix} 2 & 3 & 0 & 0 & 0 & 1 \\ 4 & 5 & 6 & 0 & 0 & 0 \\ 0 & 1 & 2 & 3 & 0 & 0 \\ 0 & 0 & 4 & 5 & 6 & 0 \\ 0 & 0 & 0 & 1 & 2 & 3 \\ 6 & 0 & 0 & 0 & 4 & 5 \end{pmatrix}. \quad (8.1)$$

The polyphase components of S are

$$S^{(0,0)} = \begin{pmatrix} 2 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 2 \end{pmatrix} \quad S^{(0,1)} = \begin{pmatrix} 3 & 0 & 1 \\ 1 & 3 & 0 \\ 0 & 1 & 3 \end{pmatrix}$$

$$S^{(1,0)} = \begin{pmatrix} 4 & 6 & 0 \\ 0 & 4 & 6 \\ 6 & 0 & 4 \end{pmatrix} \quad S^{(1,1)} = \begin{pmatrix} 5 & 0 & 0 \\ 0 & 5 & 0 \\ 0 & 0 & 5 \end{pmatrix}$$



We saw that the matrices we encountered for wavelets were filter bank transforms, so that also for them the polyphase representation is a block matrix where the blocks are filters. Since a wavelet transformations can be interpreted as a change of coordinates, let us first interpret the polyphase representation as a change of coordinates also. If we instead consider the change of coordinates from (ϕ_1, ψ_1) to

$$\mathcal{D}_1 = \{\phi_{1,0}, \phi_{1,2}, \phi_{1,4}, \dots, \phi_{1,1}, \phi_{1,3}, \phi_{1,5}, \dots\}, \quad (8.2)$$

i.e. we reorder the basis vectors in ϕ_1 so that the even-indexed ones come first. Clearly $P_{\mathcal{D}_1 \leftarrow (\phi_1, \psi_1)}$ is formed from $P_{\phi_1 \leftarrow \mathcal{C}_1}$ by taking even- and odd-indexed elements in all possible ways from $P_{\phi_1 \leftarrow \mathcal{C}_1}$, and storing these as the four different blocks in the resulting matrix. Clearly then we have the following result.

Theorem 8.4. The polyphase representation of a DWT is the matrix $P_{\mathcal{D}_1 \leftarrow (\phi_1, \psi_1)}$, and the polyphase representation of an IDWT is the matrix $P_{(\phi_1, \psi_1) \leftarrow \mathcal{D}_1}$.

Thus, similarly to how we reordered the basis vectors in (ϕ_0, ψ_1) in Chapter 6 to obtain the filter representation, the polyphase representation can be obtained by reordering the basis vectors in ϕ_1 .

Exercises for Section 8.1

1. Show that S is a filter of length kM if and only if the entries $\{S^{i,j}\}_{i,j=0}^{M-1}$ in the polyphase representation of S satisfy $S^{(i+r) \bmod M, (j+r) \bmod M} = S_{i,j}$. In other words, S is a filter if and only if the polyphase representation of S is a “block-circulant Toeplitz matrix”. This implies a fact that we will use: GH is a filter (and thus provides alias cancellation) if blocks in the polyphase representations repeat cyclically as in a Toeplitz matrix (in particular when the matrix is block-diagonal with the same block repeating on the diagonal).

2. Recall from Definition 6.17 that we defined a classical QMF filter bank as one where $M = 2$, $G_0 = H_0$, $G_1 = H_1$, and $\lambda_{H_1}(\omega) = \lambda_{H_0}(\omega + \pi)$. Show that the forward and reverse filter bank transforms of a classical QMF filter bank take the form

$$H = G = \begin{pmatrix} A & -B \\ B & A \end{pmatrix}$$

3. Recall from Definition 6.18 that we defined an alternative QMF filter bank as one where $M = 2$, $G_0 = (H_0)^T$, $G_1 = (H_1)^T$, and $\lambda_{H_1}(\omega) = \overline{\lambda_{H_0}(\omega + \pi)}$. Show that the forward and reverse filter bank transforms of an alternative QMF filter bank take the form.

$$H = \begin{pmatrix} A^T & B^T \\ -B & A \end{pmatrix} \quad G = \begin{pmatrix} A & -B^T \\ B & A^T \end{pmatrix} = \begin{pmatrix} A^T & B^T \\ -B & A \end{pmatrix}^T.$$

4. Consider alternative QMF filter banks where we take in an additional sign, so that $\lambda_{H_1}(\omega) = -\overline{\lambda_{H_0}(\omega + \pi)}$ (the Haar wavelet was an example of such a filter bank). Show that the forward and reverse filter bank transforms now take the form

$$H = \begin{pmatrix} A^T & B^T \\ B & -A \end{pmatrix} \quad G = \begin{pmatrix} A & B^T \\ B & -A^T \end{pmatrix} = \begin{pmatrix} A^T & B^T \\ B & -A \end{pmatrix}^T.$$

It is straightforward to check that also these satisfy the alias cancellation condition, and that the perfect reconstruction condition also here takes the form $|\lambda_{H_0}(\omega)|^2 + |\lambda_{H_0}(\omega + \pi)|^2 = 2$.

8.2 Motivation for the lifting factorization

Let us again consider the piecewise linear wavelet from Section 5.4, for which we found that the change of coordinate matrix $S = P_{\phi_1 \leftarrow \psi_1}$ was given by Equation (6.3). In the four different polyphase components of S , let us underline the corresponding elements:

$$\frac{1}{\sqrt{2}} \begin{pmatrix} \underline{1} & 0 \\ 1/2 & 1 \\ 0 & 0 \\ \vdots & \vdots \\ 0 & 0 \\ 1/2 & 0 \end{pmatrix}, \frac{1}{\sqrt{2}} \begin{pmatrix} 1 & \underline{0} \\ 1/2 & 1 \\ 0 & 0 \\ \vdots & \vdots \\ 0 & 0 \\ 1/2 & 0 \end{pmatrix}, \frac{1}{\sqrt{2}} \begin{pmatrix} 1 & 0 \\ \underline{1/2} & 1 \\ 0 & 0 \\ \vdots & \vdots \\ 0 & 0 \\ 1/2 & 0 \end{pmatrix}, \frac{1}{\sqrt{2}} \begin{pmatrix} 1 & 0 \\ 1/2 & \underline{1} \\ 0 & 0 \\ \vdots & \vdots \\ 0 & 0 \\ 1/2 & 0 \end{pmatrix}. \quad (8.3)$$

we get the following:

- The upper left corner $S^{(0,0)}$ of $P_{\phi_1 \leftarrow (\phi_1, \psi_1)}$ is $\frac{1}{\sqrt{2}}I$.
- The upper right corner $S^{(0,1)}$ of $P_{\phi_1 \leftarrow (\phi_1, \psi_1)}$ is $\mathbf{0}$.
- The lower left corner $S^{(1,0)}$ of $P_{\phi_1 \leftarrow (\phi_1, \psi_1)}$ is $\frac{1}{\sqrt{2}}S_1$, where S_1 is the filter $\{1/2, \underline{1/2}\}$.
- The lower right corner $S^{(1,1)}$ of $P_{\phi_1 \leftarrow (\phi_1, \psi_1)}$ is $\frac{1}{\sqrt{2}}I$.

In other words, the polyphase representation of S is

$$P_{\phi_1 \leftarrow (\phi_1, \psi_1)} = \frac{1}{\sqrt{2}} \begin{pmatrix} I & \mathbf{0} \\ S_1 & I \end{pmatrix}. \quad (8.4)$$

Note also that this matrix is easily inverted, since

$$\begin{pmatrix} I & \mathbf{0} \\ S_1 & I \end{pmatrix} \begin{pmatrix} I & \mathbf{0} \\ -S_1 & I \end{pmatrix} = (I - \mathbf{0}S_1 - S_1 I) = \begin{pmatrix} I & \mathbf{0} \\ \mathbf{0} & I \end{pmatrix} = I.$$

This means that

$$P_{(\phi_1, \psi_1) \leftarrow \mathcal{D}_1} = \sqrt{2} \begin{pmatrix} I & \mathbf{0} \\ -S_1 & I \end{pmatrix}.$$

We have thus seen two advantages of reordering the wavelet basis ϕ_1 to \mathcal{D}_1 : First of all, we obtain block matrices where all blocks are filters. Secondly, these block matrices are easily inverted.

$P_{\mathcal{D}_1 \leftarrow (\phi_1, \psi_1)}$ is different from Equation (8.4) in general, so that it is not so easily inverted in general. To look further into this, let us consider the alternative piecewise linear wavelet. In this case, Equation (6.5) shows that $P_{\mathcal{D}_1 \leftarrow (\phi_1, \hat{\psi}_1)}$ is not on the form $\begin{pmatrix} I & \mathbf{0} \\ S_1 & I \end{pmatrix}$ for some filter S_1 , since there is more than one element in every column. Recall, however, that the alternative piecewise linear wavelet was obtained by constructing a new mother wavelet $\hat{\psi}$ from the old ψ . $\hat{\psi}$ is defined in Section 5.5 by Equation 5.34. From this equation it is clear that

$$P_{(\phi_1, \psi_1) \leftarrow (\phi_1, \hat{\psi}_1)} = \begin{pmatrix} I & S_2 \\ 0 & I \end{pmatrix},$$

where

$$S_2 = -\frac{1}{4} \begin{pmatrix} 1 & 0 & 0 & \cdots & 0 & 1 \\ 1 & 1 & 0 & \cdots & 0 & 0 \\ 0 & 1 & 1 & \cdots & 0 & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & \cdots & 1 & 0 \\ 0 & 0 & 0 & \cdots & 1 & 1 \end{pmatrix} = -\frac{1}{4} \{1, 1\}. \quad (8.5)$$

We can also write

$$P_{\mathcal{D}_1 \leftarrow (\phi_1, \hat{\psi}_1)} = P_{\mathcal{D}_1 \leftarrow (\phi_1, \psi_1)} P_{(\phi_1, \psi_1) \leftarrow (\phi_1, \hat{\psi}_1)}.$$

Since we already have computed $P_{\mathcal{D}_1 \leftarrow (\phi_1, \psi_1)}$, this means that

$$P_{\mathcal{D}_1 \leftarrow (\phi_1, \hat{\psi}_1)} = \frac{1}{\sqrt{2}} \begin{pmatrix} I & \mathbf{0} \\ S_1 & I \end{pmatrix} \begin{pmatrix} I & S_2 \\ \mathbf{0} & I \end{pmatrix}.$$

In other words, also here the same type of matrix could be used to express the change of coordinates. This matrix is also easily invertible, and

$$P_{(\phi_1, \hat{\psi}_1) \leftarrow \mathcal{D}_1} = \sqrt{2} \begin{pmatrix} I & -S_2 \\ \mathbf{0} & I \end{pmatrix} \begin{pmatrix} I & \mathbf{0} \\ -S_1 & I \end{pmatrix}.$$

It is also clear how the DWT for the Haar wavelet can be written in the same way: If we first defined $\hat{\psi} = -\sqrt{2}\phi_{1,1}$ for the Haar wavelet, and then continued by defining

$\psi = \hat{\psi} + \phi$, ψ is the same function as we originally defined for the Haar wavelet. This gives us that $P_{(\phi_1, \psi_1) - (\phi_1, \hat{\psi}_1)} = \begin{pmatrix} I & S_2 \\ 0 & I \end{pmatrix}$ for some filter S_2 also here. As above, it also follows that $P_{\mathcal{D}_1 \leftarrow (\phi_1, \psi_1)}$ is on the form $(I \quad \mathbf{0} \quad S_1 \quad I)$ for some filter S_1 , so that the same type of factorization can be written down for the Haar wavelet as well.

8.3 The lifting factorization

For the wavelets we have considered we saw in the previous section that their polyphase representation could be factored into a product of matrices on the form $\begin{pmatrix} I & S \\ \mathbf{0} & I \end{pmatrix}$ or $\begin{pmatrix} I & \mathbf{0} \\ S & I \end{pmatrix}$. In this section we will show that these can be used as building blocks for general wavelets, in particular for the more advanced wavelets we constructed in Chapter 7. This will be the basis for what we will define as the lifting factorization. Let us start with the following definition.

Definition 8.5 (Elementary lifting matrices). A matrix on the form $\begin{pmatrix} I & S \\ 0 & I \end{pmatrix}$ where S is a filter is called an *elementary lifting matrix of even type*. A matrix on the form $\begin{pmatrix} I & 0 \\ S & I \end{pmatrix}$ is called an *elementary lifting matrix of odd type*.

The following are the most useful properties of elementary lifting matrices:

Lemma 8.6. The following hold:

1. $\begin{pmatrix} I & S \\ 0 & I \end{pmatrix}^T = \begin{pmatrix} I & 0 \\ S^T & I \end{pmatrix}$, and $\begin{pmatrix} I & 0 \\ S & I \end{pmatrix}^T = \begin{pmatrix} I & S^T \\ 0 & I \end{pmatrix}$,
2. $\begin{pmatrix} I & S_1 \\ 0 & I \end{pmatrix} \begin{pmatrix} I & S_2 \\ 0 & I \end{pmatrix} = \begin{pmatrix} I & S_1 + S_2 \\ 0 & I \end{pmatrix}$,
3. $\begin{pmatrix} I & 0 \\ S_1 & I \end{pmatrix} \begin{pmatrix} I & 0 \\ S_2 & I \end{pmatrix} = \begin{pmatrix} I & 0 \\ S_1 + S_2 & I \end{pmatrix}$,
4. $\begin{pmatrix} I & S \\ 0 & I \end{pmatrix}^{-1} = \begin{pmatrix} I & -S \\ 0 & I \end{pmatrix}^{-1}$
5. $\begin{pmatrix} I & 0 \\ S & I \end{pmatrix}^{-1} = \begin{pmatrix} I & 0 \\ -S & I \end{pmatrix}^{-1}$

These statements follow directly from Theorem 8.2. Due to Property 2, one can assume that odd and even types of lifting matrices appear in alternating order, since

matrices of the same type can be grouped together. The following result says that the elementary lifting matrices can be used as general building blocks:

Theorem 8.7. Any invertible matrix on the form $S = \begin{pmatrix} S^{(0,0)} & S^{(0,1)} \\ S^{(1,0)} & S^{(1,1)} \end{pmatrix}$, where the $S^{(i,j)}$ are filters with a finite number of filter coefficients, can be written on the form

$$\Lambda_1 \cdots \Lambda_n \begin{pmatrix} \alpha_0 E_p & 0 \\ 0 & \alpha_1 E_q \end{pmatrix}, \quad (8.6)$$

where Λ_i are elementary lifting matrices, p, q are integers, α_0, α_1 are nonzero scalars, and E_p, E_q are time delay filters. The inverse is given by

$$\begin{pmatrix} \alpha_0^{-1} E_{-p} & 0 \\ 0 & \alpha_1^{-1} E_{-q} \end{pmatrix} (\Lambda_n)^{-1} \cdots (\Lambda_1)^{-1}. \quad (8.7)$$

Note that $(\Lambda_i)^{-1}$ can be computed with the help of properties 4 and 5 of Lemma 8.6.

Proof: The proof will use the concept of the length of a filter, as defined in Definition 3.24. Let $S = \begin{pmatrix} S^{(0,0)} & S^{(0,1)} \\ S^{(1,0)} & S^{(1,1)} \end{pmatrix}$ be an arbitrary invertible matrix. We will incrementally find an elementary lifting matrix Λ_i with filter S_i in the lower left or upper right corner so that $\Lambda_i S$ has filters of lower length in the first column. Assume first that $l(S^{(0,0)}) \geq l(S^{(1,0)})$, where $l(S)$ is the length of a filter as given by Definition 3.24. If Λ_i is of even type, then the first column in $\Lambda_i S$ is

$$\begin{pmatrix} I & S_i \\ 0 & I \end{pmatrix} \begin{pmatrix} S^{(0,0)} \\ S^{(1,0)} \end{pmatrix} = \begin{pmatrix} S^{(0,0)} + S_i S^{(1,0)} \\ S^{(1,0)} \end{pmatrix}. \quad (8.8)$$

S_i can now be chosen so that $l(S^{(0,0)} + S_i S^{(1,0)}) < l(S^{(1,0)})$. To see how, recall that we in Section 3.4 stated that multiplying filters corresponds to multiplying polynomials. S_i can thus be found from polynomial division with remainder: when we divide $S^{(0,0)}$ by $S^{(1,0)}$, we actually find polynomials S_i and P with $l(P) < l(S^{(1,0)})$ so that $S^{(0,0)} = S_i S^{(1,0)} + P$, so that the length of $P = S^{(0,0)} - S_i S^{(1,0)}$ is less than $l(S^{(1,0)})$. The same can be said if Λ_i is of odd type, in which case the first and second components are simply swapped. This procedure can be continued until we arrive at a product

$$\Lambda_n \cdots \Lambda_1 S$$

where either the first or the second component in the first column is 0. If the first component in the first column is 0, the identity

$$\begin{pmatrix} I & 0 \\ -I & I \end{pmatrix} \begin{pmatrix} I & I \\ 0 & I \end{pmatrix} \begin{pmatrix} 0 & X \\ Y & Z \end{pmatrix} = \begin{pmatrix} Y & X+Z \\ 0 & -X \end{pmatrix}$$

explains that we can bring the matrix to a form where the second element in the first column is zero instead, with the help of the additional lifting matrices $\Lambda_{n+1} =$

$\begin{pmatrix} I & I \\ 0 & I \end{pmatrix}$, and $\Lambda_{n+2} = \begin{pmatrix} I & 0 \\ -I & I \end{pmatrix}$, so that we always can assume that the second element in the first column is 0, i.e.

$$\Lambda_n \cdots \Lambda_1 S = \begin{pmatrix} P & Q \\ 0 & R \end{pmatrix},$$

for some matrices P, Q, R . From the proof of Theorem 6.14 we will see that in order for S to be invertible, we must have that $S^{(0,0)}S^{(1,1)} - S^{(0,1)}S^{(1,0)} = -\alpha^{-1}E_d$ for some nonzero scalar α and integer d . Since $\begin{pmatrix} P & Q \\ 0 & R \end{pmatrix}$ is also invertible, we must thus have that PR must be on the form αE_n . When the filters have a finite number of filter coefficients, the only possibility for this to happen is when $P = \alpha_0 E_p$ and $R = \alpha_1 E_q$ for some p, q, α_0, α_1 . Using this, and also isolating S on one side, we obtain that

$$S = (\Lambda_1)^{-1} \cdots (\Lambda_n)^{-1} \begin{pmatrix} \alpha_0 E_p & Q \\ 0 & \alpha_1 E_q \end{pmatrix}, \quad (8.9)$$

Noting that

$$\begin{pmatrix} \alpha_0 E_p & Q \\ 0 & \alpha_1 E_q \end{pmatrix} = \begin{pmatrix} 1 & \frac{1}{\alpha_1} E_{-q} Q \\ 0 & 1 \end{pmatrix} \begin{pmatrix} \alpha_0 E_p & 0 \\ 0 & \alpha_1 E_q \end{pmatrix},$$

we can rewrite Equation (8.9) as

$$S = (\Lambda_1)^{-1} \cdots (\Lambda_n)^{-1} \begin{pmatrix} 1 & \frac{1}{\alpha_1} E_{-q} Q \\ 0 & 1 \end{pmatrix} \begin{pmatrix} \alpha_0 E_p & 0 \\ 0 & \alpha_1 E_q \end{pmatrix},$$

which is a lifting factorization of the form we wanted to arrive at. The last matrix in the lifting factorization is not really a lifting matrix, but it too can easily be inverted, so that we arrive at Equation (8.7). This completes the proof. \blacksquare

Factorization (8.6) is what we call the *lifting factorization* of S . In practice, one starts with a given wavelet with certain proved properties such as the ones from Chapter 7, and applies an algorithm to obtain a lifting factorization of it. The algorithm can easily be written down from the proof of Theorem 8.7. The lifting factorization is far from unique, and the algorithm only gives one of them.

It is desirable for an implementation to obtain a lifting factorization where the lifting steps are as simple as possible. Let us restrict to the case of wavelets with symmetric filters, since the wavelets used in most applications are symmetric. In particular this means that $S^{(0,0)}$ is a symmetric matrix, and that $S^{(1,0)}$ is symmetric about $-1/2$ (see Exercise 1).

Assume that we in the proof of Theorem 8.7 add an elementary lifting of even type. At this step we then compute $S^{(0,0)} + S_i S^{(1,0)}$ in the first entry of the first column. Since $S^{(0,0)}$ is now assumed symmetric, $S_i S^{(1,0)}$ must also be symmetric in order for the length to be reduced. And since the filter coefficients of $S^{(1,0)}$ are assumed symmetric about $-1/2$, S_i must be chosen with symmetry around $1/2$. If the difference in the lengths of the filters in the first column is 1, we can choose a filter of length 2 to reduce the length by 2, so that the S_i in an even lifting step can be chosen

on the form $S_i = \lambda_i \{\underline{1}, 1\}$. Similarly, for an odd lifting step, S_i can be chosen on the form $S_i = \lambda_i \{1, \underline{1}\}$. Let us summarize this as follows:

Theorem 8.8. When the filters in a wavelet are symmetric and the lengths of the filters in the first column differ by 1 at all steps in the lifting factorization, the lifting steps of even type take the simplified form $\begin{pmatrix} I & \lambda_i \{\underline{1}, 1\} \\ 0 & I \end{pmatrix}$, and the lifting steps of odd type take the simplified form $\begin{pmatrix} I & 0 \\ \lambda_i \{1, \underline{1}\} & I \end{pmatrix}$.

The lifting steps mentioned in this theorem are quickly computed due to their simple structure. Writing them as MRA-matrices we get

$$\begin{aligned} \begin{pmatrix} I & \lambda_i \{\underline{1}, 1\} \\ 0 & I \end{pmatrix} &\leftrightarrow \begin{pmatrix} 1 & \lambda & 0 & 0 & \cdots & 0 & 0 & \lambda \\ 0 & 1 & 0 & 0 & \cdots & 0 & 0 & 0 \\ 0 & \lambda & 1 & \lambda & \cdots & 0 & 0 & 0 \\ \vdots & \vdots \\ 0 & 0 & 0 & 0 & \cdots & \lambda & 1 & \lambda \\ 0 & 0 & 0 & 0 & \cdots & 0 & 0 & 1 \end{pmatrix} \\ \begin{pmatrix} I & 0 \\ \lambda_i \{1, \underline{1}\} & I \end{pmatrix} &\leftrightarrow \begin{pmatrix} 1 & 0 & 0 & 0 & \cdots & 0 & 0 & 0 \\ \lambda & 1 & \lambda & 0 & \cdots & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & \cdots & 0 & 0 & 0 \\ \vdots & \vdots \\ 0 & 0 & 0 & 0 & \cdots & 0 & 1 & 0 \\ \lambda & 0 & 0 & 0 & \cdots & 0 & \lambda & 1 \end{pmatrix}. \end{aligned}$$

These lifting steps are thus also MRA-matrices with symmetric filters. In other words, when the lifting factorization is applied as above, it means that an MRA-matrix with symmetric filters is factored into simpler MRA-matrices which also have symmetric filters, i.e. $S = A_1 \cdots A_n$ where all S, A_i are MRA-matrices with symmetric filters. If we apply the DWT to symmetric extensions, we clearly also have that $S_r =$

$(A_1)_r \cdots (A_n)_r$, where S_r is given by Theorem 5.44. We also have that

$$\begin{pmatrix} I & \lambda_i\{\underline{1}, 1\} \\ 0 & I \end{pmatrix}_r \leftrightarrow \begin{pmatrix} 1 & 2\lambda & 0 & 0 & \cdots & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & \cdots & 0 & 0 & 0 \\ 0 & \lambda & 1 & \lambda & \cdots & 0 & 0 & 0 \\ \vdots & \vdots \\ 0 & 0 & 0 & 0 & \cdots & \lambda & 1 & \lambda \\ 0 & 0 & 0 & 0 & \cdots & 0 & 0 & 1 \end{pmatrix} \quad (8.10)$$

$$\begin{pmatrix} I & 0 \\ \lambda_i\{\underline{1}, 1\} & I \end{pmatrix}_r \leftrightarrow \begin{pmatrix} 1 & 0 & 0 & 0 & \cdots & 0 & 0 & 0 \\ \lambda & 1 & \lambda & 0 & \cdots & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & \cdots & 0 & 0 & 0 \\ \vdots & \vdots \\ 0 & 0 & 0 & 0 & \cdots & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & \cdots & 0 & 2\lambda & 1 \end{pmatrix}. \quad (8.11)$$

These matrices have very simple implementations: every second element is left unchanged, while for the remaining elements we simply add the two neighbours. The number of arithmetic operations needed to apply this matrix is easily computed. An implementation which exploits the symmetry of H_0, H_1 requires $N/2$ multiplications, and $2N/2 = N$ additions, due to Observation 4.21.

For the wavelets we will consider in the following examples it will turn out the filters in the first column differ by 1 at all steps in the lifting factorization, so that these simple lifting matrices can be used. The examples also explain how lifting helps us to obtain some computational improvements in the number of arithmetic operations.

Example 8.9 (Lifting of the Spline 5/3 wavelet). Let us consider the Spline 5/3 wavelet, which we defined in Example 7.20. Let us start by looking at the matrix G^T , where G is the MRA-matrix of the DWT of this wavelet. This has the filter coefficients of G_0 and G_1 in the first two rows, and we recall that

$$G_0 = \left\{ \frac{1}{4}, \underline{\frac{1}{2}}, \frac{1}{4} \right\}$$

$$G_1 = \left\{ -\frac{1}{4}, -\frac{1}{2}, \underline{\frac{3}{2}}, -\frac{1}{2}, -\frac{1}{4} \right\},$$

from which we see that the polyphase components of G^T are

$$\begin{pmatrix} S^{(0,0)} & S^{(0,1)} \\ S^{(1,0)} & S^{(1,1)} \end{pmatrix} = \begin{pmatrix} \frac{1}{2}I & \{\frac{1}{4}, \frac{1}{4}\} \\ \{-\frac{1}{2}, \underline{-\frac{1}{2}}\} & \{-\frac{1}{4}, \underline{\frac{3}{2}}, -\frac{1}{4}\} \end{pmatrix}$$

We see here that the lower filter has biggest length in the first column, so that we must start with an elementary lifting of odd type. We must then find an S_1 so that $S_i \frac{1}{2}I + \{-\frac{1}{2}, \underline{-\frac{1}{2}}\}$ has lower length than $\frac{1}{2}I$. It is clear that we can choose $S_i = \{1, \underline{1}\}$,

and that we then get $\mathbf{0}$. The first lifting step thus gives

$$\begin{aligned}\Lambda_1 G^T &= \begin{pmatrix} I & \mathbf{0} \\ \{1, \underline{1}\} & I \end{pmatrix} \begin{pmatrix} \frac{1}{2}I & \{\frac{1}{4}, \frac{1}{4}\} \\ \{-\frac{1}{2}, \underline{-\frac{1}{2}}\} & \{-\frac{1}{4}, \underline{\frac{3}{2}}, -\frac{1}{4}\} \end{pmatrix} = \begin{pmatrix} \frac{1}{2}I & \{\frac{1}{4}, \frac{1}{4}\} \\ \mathbf{0} & 2I \end{pmatrix} \\ &= \begin{pmatrix} I & \frac{1}{8}\{\underline{1}, 1\} \\ \mathbf{0} & I \end{pmatrix} \begin{pmatrix} \frac{1}{2}I & \mathbf{0} \\ \mathbf{0} & 2I \end{pmatrix},\end{aligned}$$

where we also used the same rewriting as in the proof of Theorem 8.7. This gives

$$G^T = \begin{pmatrix} I & \mathbf{0} \\ \{-1, \underline{-1}\} & I \end{pmatrix} \begin{pmatrix} I & \frac{1}{8}\{\underline{1}, 1\} \\ \mathbf{0} & I \end{pmatrix} \begin{pmatrix} \frac{1}{2}I & \mathbf{0} \\ \mathbf{0} & 2I \end{pmatrix}.$$

Transposing this expression gives

$$G = \begin{pmatrix} \frac{1}{2}I & \mathbf{0} \\ \mathbf{0} & 2I \end{pmatrix} \begin{pmatrix} I & \mathbf{0} \\ \frac{1}{8}\{1, \underline{1}\} & I \end{pmatrix} \begin{pmatrix} I & \{\underline{-1}, -1\} \\ \mathbf{0} & I \end{pmatrix}$$

Taking inverses of this expression gives

$$H = \begin{pmatrix} I & \{\underline{1}, 1\} \\ \mathbf{0} & I \end{pmatrix} \begin{pmatrix} I & \mathbf{0} \\ -\frac{1}{8}\{1, \underline{1}\} & I \end{pmatrix} \begin{pmatrix} 2I & \mathbf{0} \\ \mathbf{0} & \frac{1}{2}I \end{pmatrix}$$

We now have obtained the lifting factorization. Only two lifting steps were required. We also see that the lifting steps involve only dyadic fractions, just as the filter coefficients did. This means that both the lifting factorization also can be used for lossless operations.

Let us also compute the number of additions and multiplications performed when we apply the lifting factorization. We saw that each lifting step uses a number of N additions and $N/2$ multiplications. Since there are two lifting steps, the number of additions and multiplications are $2N$ and $2N$, respectively, where we have taken into account multiplication with the diagonal matrix as well. A direct implementation of DWT in terms of filters would need a total of $3N$ additions and $2.5N$ multiplications, where we have taken the symmetry of the filters into account. We thus see that the lifting factorization gives a small decrease in the number of arithmetic operations. ♣

Example 8.10 (Lifting of the CDF 9/7 wavelet). For the wavelet we considered in Example 7.21, it is more cumbersome to compute the lifting factorization by hand. It is however, straightforward to write an algorithm which computes the lifting steps, as these are performed in the proof of Theorem 8.7. You will be spared the details of this algorithm. Also, when we use these wavelets in implementations later they will use precomputed values of these lifting steps, and you can take these implementations for granted too. If we run the algorithm for computing the lifting factorization

we obtain

$$\begin{aligned}
H &= \begin{pmatrix} I & 0.5861\{1, 1\} \\ 0 & I \end{pmatrix} \begin{pmatrix} I & 0 \\ 0.6681\{1, 1\} & I \end{pmatrix} \begin{pmatrix} I & -0.0700\{1, 1\} \\ 0 & I \end{pmatrix} \\
&\quad \times \begin{pmatrix} I & 0 \\ -1.2002\{1, 1\} & I \end{pmatrix} \begin{pmatrix} 1/0.8699 & 0 \\ 0 & 1/1.1496 \end{pmatrix} \\
G &= \begin{pmatrix} 0.8699 & 0 \\ 0 & 1.1496 \end{pmatrix} \begin{pmatrix} I & 0 \\ 1.2002\{1, 1\} & I \end{pmatrix} \begin{pmatrix} I & 0.0700\{1, 1\} \\ 0 & I \end{pmatrix} \\
&\quad \times \begin{pmatrix} I & 0 \\ -0.6681\{1, 1\} & I \end{pmatrix} \begin{pmatrix} I & -0.5861\{1, 1\} \\ 0 & I \end{pmatrix}.
\end{aligned}$$

In this case four lifting steps were required. Let us also here compute the number of additions and multiplications performed. Since there now are four lifting steps, the number of additions and multiplications are $4N$ and $3N$, respectively. A direct implementation of DWT in terms of filters would need a total of $7N$ additions and $4.5N$ multiplications, where we again have taken the symmetry of the filters into account. Still the decrease in the number of operations is not very big, but at least the decrease is clearer than in the previous example. ♣

Based on the previous examples, we can now show the following:

Theorem 8.11. As the number of lifting steps grow big, the lifting factorization approximately halves the number of additions and multiplications needed when compared to a direct implementation which exploits the symmetry of the filters as in Observation 4.21.

Proof: For simplicity we assume also here that the wavelet filters H_0, H_1 differ in length by two. Assume that their lengths are $2s + 1$ and $2s + 3$. Following the same reasoning as in Observation 4.21, an implementation of the DWT (which filters half the elements with H_0 , half the elements with H_1) which exploits the symmetry of H_0, H_1 (without exploiting the lifting factorization), requires

$$(s+1)\frac{N}{2} + (s+2)\frac{N}{2} = \left(s + \frac{3}{2}\right)N$$

multiplications, and

$$2s\frac{N}{2} + (2s+2)\frac{N}{2} = (2s+1)N$$

additions. Let us compute the same numbers when we apply a lifting factorization of the wavelet instead. It is clear from the two previous examples that each lifting step, when viewed in terms of MRA matrices, reduces the total number of filter coefficients by 4 (except for the last lifting), all the way until what remains is a diagonal matrix. The Spline 5/3 wavelet requires two liftings since $5 + 3 = 8 = 2 \times 4$. The CDF 9/7 wavelet requires four liftings since $9 + 7 = 16 = 4 \times 4$. If the filters have lengths $2s + 1$ and $2s + 3$, we have that $2s + 1 + 2s + 3 = 4s + 4 = 4(s + 1)$, so that a total of $s + 1$ lifting steps are required. This requires a total of $(s + 1)N/2 + N = \frac{1}{2}(s + 3)N$

multiplications, and $(s + 1)N$ additions, where we have taken into account the multiplication of the diagonal matrix as well. Clearly then, in the limit as s increases, the lifting factorization halves the number of arithmetic operations needed. ■

Perhaps more important than the reduction in the number of arithmetic operations is the fact that the lifting factorization splits the DWT and IDWT into simpler components, each very attractive for hardware implementations since a lifting step only requires the additional value λ_i from Theorem 8.8. Lifting actually provides us with a complete implementation strategy for the DWT and IDWT, in which the λ_i are used as precomputed values.

Exercises for Section 8.3

1. Assume that the filters H_0, H_1 of a wavelet are symmetric, and denote by $S^{(i,j)}$ the polyphase components of the corresponding MRA-matrix H . Show that $S^{(0,0)}$ and $S^{(1,1)}$ are symmetric filters, that the filter coefficients of $S^{(1,0)}$ has symmetry about $-1/2$, and that $S^{(0,1)}$ has symmetry about $1/2$. Also show a similar statement for the MRA-matrix G of the inverse DWT.

2. Write functions

```
function x=liftingstepapplyA(lambda,x)
function x=liftingstepapplyB(lambda,x)
```

which applies the elementary lifting matrices (8.10) and (8.11), respectively, to the vector x . You can assume that N has even length. The function should not perform matrix multiplication to achieve this, and the result should be returned in the same vector as the input. The function should also perform as few multiplications as possible. How can you achieve this?

3. Write functions `DWTImpl53` and `IDWTImpl53` which implements the DWT and the IDWT for the Spline 5/3 wavelet, using the lifting factorization obtained in Example 8.9. You should use the functions you implemented in Exercise 2, and again you can make the assumption that N can be written as an appropriate power of 2.

4. Write functions `DWTImpl97` and `IDWTImpl97` which implements the DWT and the IDWT for the CDF 9/7 wavelet, using the lifting factorization obtained in Example 8.10.

8.4 The proof of Theorem 6.14

Let $H^{(i,j)}$ be the polyphase components of the DWT-matrix, $G^{(i,j)}$ the polyphase components of the IDWT matrix. $GH = I$ means that

$$\begin{pmatrix} G^{(0,0)} & G^{(0,1)} \\ G^{(1,0)} & G^{(1,1)} \end{pmatrix} \begin{pmatrix} H^{(0,0)} & H^{(0,1)} \\ H^{(1,0)} & H^{(1,1)} \end{pmatrix} = \begin{pmatrix} I & \mathbf{0} \\ \mathbf{0} & I \end{pmatrix}.$$

If we here multiply with $\begin{pmatrix} G^{(1,1)} & -G^{(0,1)} \\ -G^{(1,0)} & G^{(0,0)} \end{pmatrix}$ on both sides to the left, or with $\begin{pmatrix} H^{(1,1)} & -H^{(0,1)} \\ -H^{(1,0)} & H^{(0,0)} \end{pmatrix}$ on both sides to the right, we get

$$\begin{pmatrix} G^{(1,1)} & -G^{(0,1)} \\ -G^{(1,0)} & G^{(0,0)} \end{pmatrix} = \begin{pmatrix} (G^{(0,0)}G^{(1,1)} - G^{(1,0)}G^{(0,1)})H^{(0,0)} & (G^{(0,0)}G^{(1,1)} - G^{(1,0)}G^{(0,1)})H^{(0,1)} \\ (G^{(0,0)}G^{(1,1)} - G^{(1,0)}G^{(0,1)})H^{(1,0)} & (G^{(0,0)}G^{(1,1)} - G^{(1,0)}G^{(0,1)})H^{(1,1)} \end{pmatrix}$$

$$\begin{pmatrix} H^{(1,1)} & -H^{(0,1)} \\ -H^{(1,0)} & H^{(0,0)} \end{pmatrix} = \begin{pmatrix} (H^{(0,0)}H^{(1,1)} - H^{(1,0)}H^{(0,1)})G^{(0,0)} & (H^{(0,0)}H^{(1,1)} - H^{(1,0)}H^{(0,1)})G^{(0,1)} \\ (H^{(0,0)}H^{(1,1)} - H^{(1,0)}H^{(0,1)})G^{(1,0)} & (H^{(0,0)}H^{(1,1)} - H^{(1,0)}H^{(0,1)})G^{(1,1)} \end{pmatrix}$$

Now since $G^{(0,0)}G^{(1,1)} - G^{(1,0)}G^{(0,1)}$ and $H^{(0,0)}H^{(1,1)} - H^{(1,0)}H^{(0,1)}$ also are circulant Toeplitz matrices, the expressions above give that

$$\begin{aligned} l(H^{(0,0)}) &\leq l(G^{(1,1)}) \leq l(H^{(0,0)}) \\ l(H^{(0,1)}) &\leq l(G^{(0,1)}) \leq l(H^{(0,1)}) \\ l(H^{(1,0)}) &\leq l(G^{(1,0)}) \leq l(H^{(1,0)}) \end{aligned}$$

so that we must have equality here, and with both $G^{(0,0)}G^{(1,1)} - G^{(1,0)}G^{(0,1)}$ and $H^{(0,0)}H^{(1,1)} - H^{(1,0)}H^{(0,1)}$ having only one nonzero diagonal. In particular we can define the diagonal matrix $D = H^{(0,0)}H^{(1,1)} - H^{(0,1)}H^{(1,0)} = \alpha^{-1}E_d$ (for some α, d), and we have that

$$\begin{pmatrix} G^{(0,0)} & G^{(0,1)} \\ G^{(1,0)} & G^{(1,1)} \end{pmatrix} = \begin{pmatrix} \alpha E_{-d} H^{(1,1)} & -\alpha E_{-d} H^{(0,1)} \\ -\alpha E_{-d} H^{(1,0)} & \alpha E_{-d} H^{(0,0)} \end{pmatrix}.$$

Now, the first column in the matrix on the left hand side gives the even and odd filter coefficients of G_0 . On the right hand side, the block $H^{(1,1)}$ gives the even filter coefficients of H_1 , the block $H^{(1,0)}$ gives the odd filter coefficients of H_1 . Together with the delay, the first column on the right hand side thus produces the filter $E_{-2d}H_1$, which has frequency response $e^{2id\omega}\lambda_{H_1}(\omega)$. If we add the alternating sign and the constant, the first column on the right hand side thus produces the filter

$$\alpha e^{2id(\omega+\pi)}\lambda_{H_1}(\omega+\pi) = \alpha e^{2id\omega}\lambda_{H_1}(\omega+\pi)$$

We thus have that $\lambda_{G_0}(\omega) = \alpha e^{2id\omega}\lambda_{H_1}(\omega+\pi)$, so that

$$\lambda_{H_1}(\omega) = \alpha^{-1}e^{-2id\omega}\lambda_{G_0}(\omega+\pi),$$

which is Equation (6.23). Now, the second column in the matrix on the left hand side gives the even and odd filter coefficients of E_1G_1 . On the right hand side, the block $H^{(0,1)}$ gives the odd filter coefficients of H_0 , delayed with 1. The block $H^{(0,0)}$ gives the even filter coefficients of H_0 . Together these produce the output E_1H_0 . Together with the delay, the second column on the right hand side thus produces the filter $E_{-2d+1}H_0$, which has frequency response $e^{i(2d-1)\omega}\lambda_{H_0}(\omega)$. If we add the alternating sign and the constant, the second column on the right hand side this produces the filter

$$-\alpha e^{i(2d-1)(\omega+\pi)}\lambda_{H_0}(\omega+\pi) = \alpha e^{2id\omega}e^{-id\omega}\lambda_{H_0}(\omega+\pi).$$

We thus have that $e^{-id\omega}\lambda_{G_1}(\omega) = \alpha e^{2id\omega}e^{-id\omega}\lambda_{H_0}(\omega+\pi)$, so that

$$\lambda_{G_1}(\omega) = \alpha e^{2id\omega}\lambda_{H_0}(\omega+\pi),$$

which is Equation (6.24).

8.5 Cosine-modulated filter banks and the MP3 standard

Previously we saw that the MP3 standard used a certain filter bank, called a cosine-modulated filter bank. We also illustrated that, surprisingly for a much used international standard, the synthesis system did not exactly invert the analysis system, i.e. we do not have perfect reconstruction, only “near-perfect reconstruction”. In this section we will first explain how this filter bank can be constructed, and why it can not give perfect reconstruction. In particular it will be clear how the prototype filter can be constructed. We will then construct a very similar filter bank, which actually can give perfect reconstruction. It may seem very surprising that the MP3 standard does not use this filter bank instead due to this. The explanation may lie in that the MP3 standard was established at about the same time as these filter banks were developed, so that the standard did not capture this very similar filter bank with perfect reconstruction.

8.5.1 Polyphase representations of the filter bank transforms of the MP3 standard

The main idea is to find the polyphase representations of the forward and reverse filter bank transforms of the MP3 standard. We start with the expression

$$z_{32(s-1)+n} = \sum_{k=0}^{511} \cos((n + 1/2)(k - 16)\pi/32) h_k x_{32s-k-1}, \quad (8.12)$$

which lead to the expression of the forward filter bank transform (Theorem 6.25). Using that any $k < 512$ can be written uniquely on the form $k = m + 64r$, where $0 \leq m < 64$, and $0 \leq r < 8$, we can rewrite this as

$$\begin{aligned} &= \sum_{m=0}^{63} \sum_{r=0}^7 (-1)^r \cos(2\pi(n + 1/2)(m - 16)/64) h_{m+64r} x_{32s-(m+64r)-1} \\ &= \sum_{m=0}^{63} \cos(2\pi(n + 1/2)(m - 16)/64) \sum_{r=0}^7 (-1)^r h_{m+32 \cdot 2r} x_{32(s-2r)-m-1}. \end{aligned}$$

Here we also used Property (6.45). If we write

$$V^{(m)} = \{\underline{(-1)^0 h_m}, 0, \underline{(-1)^1 h_{m+64}}, 0, \underline{(-1)^2 h_{m+128}}, \dots, \underline{(-1)^7 h_{m+7 \cdot 64}}, 0\}, \quad (8.13)$$

for $0 \leq m \leq 63$, and we can write the expression above as

$$\begin{aligned} &\sum_{m=0}^{63} \cos(2\pi(n + 1/2)(m - 16)/64) \sum_{r=0}^{15} V_r^{(m)} x_{32(s-r)-m-1} \\ &= \sum_{m=0}^{63} \cos(2\pi(n + 1/2)(m - 16)/64) \sum_{r=0}^{15} V_r^{(m)} x_{s-1-r}^{(32-m-1)} \\ &= \sum_{m=0}^{63} \cos(2\pi(n + 1/2)(m - 16)/64) (V^{(m)} \mathbf{x}^{(32-m-1)})_{s-1}, \end{aligned}$$

where we recognized $x_{32(s-r)-m-1}$ in terms of the polyphase components of \mathbf{x} , and the inner sum as a convolution. We remark that the inner terms $\{(V^{(m)}\mathbf{x}^{(32-m-1)})_{s-1}\}_{m=0}^{63}$ here are what the standard calls partial calculations (windowing refers to multiplication with the combined set of filter coefficients of the $V^{(m)}$), and that matrixing here represents the multiplication with the cosine entries. Since $\mathbf{z}^{(n)} = \{z_{32(s-1)+n}\}_{s=0}^{\infty}$ is the n 'th polyphase component of \mathbf{z} , this can be written as

$$\mathbf{z}^{(n)} = \sum_{m=0}^{63} \cos(2\pi(n+1/2)(m-16)/64) IV^{(m)}\mathbf{x}^{(32-m-1)}.$$

In terms of matrices this can be written as

$$\begin{aligned} \mathbf{z} &= \begin{pmatrix} \cos(2\pi(0+1/2)\cdot(-16)/64)I & \cdots & \cos(2\pi(0+1/2)\cdot(47)/64)I \\ \vdots & \ddots & \vdots \\ \cos(2\pi(31+1/2)\cdot(-16)/64)I & \cdots & \cos(2\pi(31+1/2)\cdot(47)/64)I \end{pmatrix} \\ &\times \begin{pmatrix} V^{(0)} & \mathbf{0} & \cdots & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & V^{(1)} & \cdots & \mathbf{0} & \mathbf{0} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ \mathbf{0} & \mathbf{0} & \cdots & V^{(62)} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \cdots & \mathbf{0} & V^{(63)} \end{pmatrix} \begin{pmatrix} \mathbf{x}^{(31)} \\ \mathbf{x}^{(30)} \\ \vdots \\ \mathbf{x}^{(-32)} \end{pmatrix}. \end{aligned}$$

If we place the 15 first columns in the cosine matrix last using Property (6.45) (we must then also place the 15 first rows last in the second matrix), we obtain

$$\begin{aligned} \mathbf{z} &= \begin{pmatrix} \cos(2\pi(0+1/2)\cdot(0)/64)I & \cdots & \cos(2\pi(0+1/2)\cdot(63)/64)I \\ \vdots & \ddots & \vdots \\ \cos(2\pi(31+1/2)\cdot(0)/64)I & \cdots & \cos(2\pi(31+1/2)\cdot(63)/64)I \end{pmatrix} \\ &\times \begin{pmatrix} \mathbf{0} & \cdots & \mathbf{0} & V^{(16)} & \cdots & \mathbf{0} \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ \mathbf{0} & \cdots & \mathbf{0} & \mathbf{0} & \cdots & V^{(63)} \\ -V^{(0)} & \cdots & \cdots & \mathbf{0} & \cdots & \mathbf{0} \\ \vdots & \ddots & \vdots & \vdots & \vdots & \mathbf{0} \\ \mathbf{0} & \cdots & -V^{(15)} & \mathbf{0} & \cdots & \mathbf{0} \end{pmatrix} \begin{pmatrix} \mathbf{x}^{(31)} \\ \mathbf{x}^{(30)} \\ \vdots \\ \mathbf{x}^{(-32)} \end{pmatrix}. \end{aligned}$$

Using Property (6.46) to combine column k and $64-k$ in the cosine matrix (as well as row k and $64-k$ in the second matrix), we can write this as

$$\begin{pmatrix} \cos(2\pi(0+1/2)\cdot(0)/64)I & \cdots & \cos(2\pi(0+1/2)\cdot(31)/64)I \\ \vdots & \ddots & \vdots \\ \cos(2\pi(31+1/2)\cdot(0)/64)I & \cdots & \cos(2\pi(31+1/2)\cdot(31)/64)I \end{pmatrix} (A' \quad B') \begin{pmatrix} \mathbf{x}^{(31)} \\ \mathbf{x}^{(30)} \\ \vdots \\ \mathbf{x}^{(-32)} \end{pmatrix}.$$

where

$$A' = \begin{pmatrix} \mathbf{0} & \mathbf{0} & \cdots & \mathbf{0} & V^{(16)} & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \cdots & V^{(15)} & \mathbf{0} & V^{(17)} & \cdots & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots & \vdots & \vdots & \ddots & \mathbf{0} \\ \mathbf{0} & V^{(1)} & \cdots & \mathbf{0} & \mathbf{0} & \mathbf{0} & \cdots & V^{(31)} \\ V^{(0)} & \mathbf{0} & \cdots & \mathbf{0} & \mathbf{0} & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \cdots & \mathbf{0} & \mathbf{0} & \mathbf{0} & \cdots & \mathbf{0} \\ \vdots & \vdots \\ \mathbf{0} & \mathbf{0} & \cdots & \mathbf{0} & \mathbf{0} & \mathbf{0} & \cdots & \mathbf{0} \end{pmatrix}$$

$$B' = \begin{pmatrix} \mathbf{0} & \mathbf{0} & \cdots & \mathbf{0} & \mathbf{0} & \mathbf{0} & \cdots & \mathbf{0} \\ \vdots & \vdots \\ V_{(32)} & \mathbf{0} & \cdots & \mathbf{0} & \mathbf{0} & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{0} & V^{(33)} & \cdots & \mathbf{0} & \mathbf{0} & \mathbf{0} & \cdots & -V^{(63)} \\ \vdots & \vdots & \ddots & \vdots & \vdots & \vdots & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \cdots & V^{(47)} & \mathbf{0} & -V^{(49)} & \cdots & \mathbf{0} \end{pmatrix}.$$

Using Equation (4.3), the cosine matrix here can be written as

$$\sqrt{\frac{M}{2}}(D_M)^T \begin{pmatrix} \sqrt{2} & 0 & \cdots & 0 \\ 0 & 1 & \cdots & 0 \\ \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & \cdots & 1 \end{pmatrix}.$$

The above can thus be written as

$$4(D_{32})^T (A - B) \begin{pmatrix} \mathbf{x}^{(31)} \\ \mathbf{x}^{(30)} \\ \vdots \\ \mathbf{x}^{(-32)} \end{pmatrix},$$

where A and B are the matrices A', B' with the first row multiplied by $\sqrt{2}$ (i.e. replace $V^{(16)}$ with $\sqrt{2}V^{(16)}$ in the matrix A'). Using that $\mathbf{x}^{(-i)} = E_1 \mathbf{x}_i$ for $1 \leq i \leq 32$, we can write this as

$$4(D_{32})^T (A - B) \begin{pmatrix} \mathbf{x}^{(31)} \\ \vdots \\ \mathbf{x}^{(0)} \\ E_1 \mathbf{x}^{(31)} \\ \vdots \\ E_1 \mathbf{x}^{(0)} \end{pmatrix} = 4(D_{32})^T \left(A \begin{pmatrix} \mathbf{x}^{(31)} \\ \vdots \\ \mathbf{x}^{(0)} \end{pmatrix} + B \begin{pmatrix} E_1 \mathbf{x}^{(31)} \\ \vdots \\ E_1 \mathbf{x}^{(0)} \end{pmatrix} \right),$$

which can be written as

$$4(D_{32})^T \begin{pmatrix} \mathbf{0} & \mathbf{0} & \cdots & \mathbf{0} & \sqrt{2}V^{(16)} & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \cdots & V^{(15)} & \mathbf{0} & V^{(17)} & \cdots & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots & \ddots & \vdots \\ \mathbf{0} & V^{(1)} & \cdots & \mathbf{0} & \mathbf{0} & \mathbf{0} & \cdots & V^{(31)} \\ V^{(0)} + E_1 V^{(32)} & \mathbf{0} & \cdots & \mathbf{0} & \mathbf{0} & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{0} & E_1 V^{(33)} & \cdots & \mathbf{0} & \mathbf{0} & \mathbf{0} & \cdots & -E_1 V^{(63)} \\ \vdots & \vdots & \ddots & \vdots & \vdots & \vdots & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \cdots & E_1 V^{(47)} & \mathbf{0} & -E_1 V^{(49)} & \cdots & \mathbf{0} \end{pmatrix} \begin{pmatrix} \mathbf{x}^{(31)} \\ \vdots \\ \mathbf{x}^{(0)} \end{pmatrix},$$

which also can be written as

$$4(D_{32})^T \begin{pmatrix} \mathbf{0} & \cdots & \mathbf{0} & \sqrt{2}V^{(16)} & \mathbf{0} & \cdots & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \cdots & V^{(17)} & \mathbf{0} & V^{(15)} & \cdots & \mathbf{0} & \mathbf{0} \\ \vdots & \ddots & \vdots & \ddots & \vdots & \ddots & \vdots & \vdots \\ V^{(31)} & \cdots & \mathbf{0} & \mathbf{0} & \mathbf{0} & \cdots & V^{(1)} & \mathbf{0} \\ \mathbf{0} & \cdots & \mathbf{0} & \mathbf{0} & \mathbf{0} & \cdots & \mathbf{0} & V^{(0)} + E_1 V^{(32)} \\ -E_1 V^{(63)} & \cdots & \mathbf{0} & \mathbf{0} & \mathbf{0} & \cdots & E_1 V^{(33)} & \mathbf{0} \\ \vdots & \ddots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ \mathbf{0} & \cdots & -E_1 V^{(49)} & \mathbf{0} & E_1 V^{(47)} & \cdots & \mathbf{0} & \mathbf{0} \end{pmatrix} \begin{pmatrix} \mathbf{x}^{(0)} \\ \vdots \\ \mathbf{x}^{(31)} \end{pmatrix}.$$

We have therefore proved the following result.

Theorem 8.12. (Polyphase factorization of a forward filter bank transform based on a prototype filter). The polyphase form of a forward filter bank transform based on a prototype filter can be factored as

$$4(D_{32})^T \begin{pmatrix} \mathbf{0} & \cdots & \mathbf{0} & \sqrt{2}V^{(16)} & \mathbf{0} & \cdots & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \cdots & V^{(17)} & \mathbf{0} & V^{(15)} & \cdots & \mathbf{0} & \mathbf{0} \\ \vdots & \ddots & \vdots & \ddots & \vdots & \ddots & \vdots & \vdots \\ V^{(31)} & \cdots & \mathbf{0} & \mathbf{0} & \mathbf{0} & \cdots & V^{(1)} & \mathbf{0} \\ \mathbf{0} & \cdots & \mathbf{0} & \mathbf{0} & \mathbf{0} & \cdots & \mathbf{0} & V^{(0)} + E_1 V^{(32)} \\ -E_1 V^{(63)} & \cdots & \mathbf{0} & \mathbf{0} & \mathbf{0} & \cdots & E_1 V^{(33)} & \mathbf{0} \\ \vdots & \ddots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ \mathbf{0} & \cdots & -E_1 V^{(49)} & \mathbf{0} & E_1 V^{(47)} & \cdots & \mathbf{0} & \mathbf{0} \end{pmatrix} \quad (8.14)$$

Due to theorem 6.27, it is also very simple to write down the polyphase factorization of the reverse filter bank transform as well. Since $E_{481}G^T$ is a forward filter bank transform where the prototype filter has been reversed, $E_{481}G^T$ can be factored as above, with $V^{(m)}$ replaced by $W^{(m)}$, with $W^{(m)}$ being the filters derived from the synthesis prototype filter in reverse order. This means that the polyphase form of G

can be factored as

$$4 \begin{pmatrix} \mathbf{0} & \mathbf{0} & \cdots & (W^{(31)})^T & \mathbf{0} & -E_{-1}(W^{(63)})^T & \cdots & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots & \ddots & \vdots \\ \mathbf{0} & (W^{(17)})^T & \cdots & \mathbf{0} & \mathbf{0} & \mathbf{0} & \cdots & -E_{-1}(W^{(49)})^T \\ \sqrt{2}(W^{(16)})^T & \mathbf{0} & \cdots & \mathbf{0} & \mathbf{0} & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{0} & (W^{(15)})^T & \cdots & \mathbf{0} & \mathbf{0} & \mathbf{0} & \cdots & E_{-1}(W^{(47)})^T \\ \vdots & \vdots & \ddots & \vdots & \vdots & \vdots & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \cdots & (W^{(1)})^T & \mathbf{0} & E_{-1}(W^{(33)})^T & \cdots & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \cdots & \mathbf{0} & (W^{(0)})^T + E_{-1}(W^{(32)})^T & \mathbf{0} & \cdots & \mathbf{0} \end{pmatrix} \times D_{32}E_{481}. \quad (8.15)$$

Now, if we define $U^{(m)}$ as the filters derived from the synthesis prototype filter itself, we have that

$$(W^{(k)})^T = -E_{-14}V^{(64-k)}, \quad 1 \leq k \leq 15 \quad (W^{(0)})^T = E_{-16}V^{(0)}.$$

Inserting this in Equation (8.15) we get the following result:

Theorem 8.13. (Polyphase factorization of a reverse filter bank transform based on a prototype filter). Assume that G is a reverse filter filter bank transform based on a prototype filter, and that $U^{(m)}$ are the filters derived from this prototype filter. Then the polyphase form of G can be factored as

$$4 \begin{pmatrix} \mathbf{0} & \mathbf{0} & \cdots & -U^{(33)} & \mathbf{0} & E_{-1}U^{(1)} & \cdots & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots & \ddots & \vdots \\ \mathbf{0} & -U^{(47)} & \cdots & \mathbf{0} & \mathbf{0} & \mathbf{0} & \cdots & E_{-1}U^{(15)} \\ -\sqrt{2}U^{(48)} & \mathbf{0} & \cdots & \mathbf{0} & \mathbf{0} & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{0} & -U^{(49)} & \cdots & \mathbf{0} & \mathbf{0} & \mathbf{0} & \cdots & -E_{-1}U^{(17)} \\ \vdots & \vdots & \ddots & \vdots & \vdots & \vdots & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \cdots & -U^{(63)} & \mathbf{0} & -E_{-1}U^{(31)} & \cdots & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \cdots & \mathbf{0} & E_{-2}U^{(0)} - E_{-1}U^{(32)} & \mathbf{0} & \cdots & \mathbf{0} \end{pmatrix} \times D_{32}E_{33}. \quad (8.16)$$

Now, consider the matrices

$$\begin{pmatrix} V^{(32-i)} & V^{(i)} \\ -E_1V^{(64-i)} & E_1V^{(32+i)} \end{pmatrix} \text{ and } \begin{pmatrix} -U^{(32+i)} & E_{-1}U^{(i)} \\ -U^{(64-i)} & -E_{-1}U^{(32-i)} \end{pmatrix}. \quad (8.17)$$

for $1 \leq i \leq 15$. These make out submatrices in the matrices in equations (8.14)

and (8.16). Clearly, only the product of these matrices influence the result. Since

$$\begin{aligned} & \begin{pmatrix} -U^{(32+i)} & E_{-1}U^{(i)} \\ -U^{(64-i)} & -E_{-1}U^{(32-i)} \end{pmatrix} \begin{pmatrix} V^{(32-i)} & V^{(i)} \\ -E_1V^{(64-i)} & E_1V^{(32+i)} \end{pmatrix} \\ &= \begin{pmatrix} -U^{(32+i)} & U^{(i)} \\ -U^{(64-i)} & -U^{(32-i)} \end{pmatrix} \begin{pmatrix} V^{(32-i)} & V^{(i)} \\ -V^{(64-i)} & V^{(32+i)} \end{pmatrix} \end{aligned} \quad (8.18)$$

we have the following result.

Theorem 8.14. Let H, G be forward and reverse filter bank transforms defined from analysis and synthesis prototype filters. Let also $V^{(k)}$ be the prototype filter of H , and $U^{(k)}$ the reverse of the prototype filter of G . If

$$\begin{aligned} & \begin{pmatrix} -U^{(32+i)} & U^{(i)} \\ -U^{(64-i)} & -U^{(32-i)} \end{pmatrix} \begin{pmatrix} V^{(32-i)} & V^{(i)} \\ -V^{(64-i)} & V^{(32+i)} \end{pmatrix} = c \begin{pmatrix} E_d & \mathbf{0} \\ \mathbf{0} & E_d \end{pmatrix} \\ & (\sqrt{2}V^{(16)})(-\sqrt{2}U^{(48)}) = cE_d \\ & (V^{(0)} + E_1V^{(32)})(E_{-2}U^{(0)} - E_{-1}U^{(32)}) = cE_d \end{aligned} \quad (8.19)$$

for $1 \leq i \leq 15$, then $GH = 16cE_{33+32d}$.

This result is the key ingredient we need in order to construct forward and reverse systems which together give perfect reconstruction. In Exercise 3 we go through how we can use lifting in order to express a wide range of possible (U, V) matrix pairs which satisfy Equation (8.19). This turns the problem of constructing cosine-modulated filter banks which are useful for audio coding into an optimization problem: the optimization variables are values λ_i which characterize lifting steps, and the objective function is the deviation of the corresponding prototype filter from an ideal bandpass filter. This optimization problem has been subject to a lot of research, and we will not go into details on this.

8.5.2 The prototype filters chosen in the MP3 standard

Now, let us return to the MP3 standard. We previously observed that in this standard the coefficients in the synthesis prototype filter seemed to equal 32 times the analysis prototype filter. This indicates that $U^{(k)} = 32V^{(k)}$. A closer inspection also yields that there is a symmetry in the values of the prototype filter: We see that $C_i = -C_{512-i}$ (i.e. antisymmetry) for most values of i . The only exception is for $i = 64, 128, \dots, 448$, for which $C_i = C_{512-i}$ (i.e. symmetry). The antisymmetry can be translated to that the filter coefficients of $V^{(k)}$ equal those of $V^{(64-k)}$ in reverse order, with a minus sign. The symmetry can be translated to that $V^{(0)}$ is symmetric. These observations can be rewritten as

$$V^{(64-k)} = -E_{14}(V^{(k)})^T, 1 \leq k \leq 15. \quad (8.20)$$

$$V^{(0)} = E_{16}(V^{(0)})^T. \quad (8.21)$$

Inserting first that $U^{(k)} = 32V^{(k)}$ in Equation (8.18) gives

$$\begin{aligned} & \begin{pmatrix} -U^{(32+i)} & U^{(i)} \\ -U^{(64-i)} & -U^{(32-i)} \end{pmatrix} \begin{pmatrix} V^{(32-i)} & V^{(i)} \\ -V^{(64-i)} & V^{(32+i)} \end{pmatrix} \\ &= 32 \begin{pmatrix} -V^{(32+i)} & V^{(i)} \\ -V^{(64-i)} & -V^{(32-i)} \end{pmatrix} \begin{pmatrix} V^{(32-i)} & V^{(i)} \\ -V^{(64-i)} & V^{(32+i)} \end{pmatrix}. \end{aligned}$$

Substituting for $V^{(32+i)}$ and $V^{(64-i)}$ after what we found by inspection now gives

$$\begin{aligned} & 32 \begin{pmatrix} E_{14}(V^{(32-i)})^T & V^{(i)} \\ E_{14}(V^{(i)})^T & -V^{(32-i)} \end{pmatrix} \begin{pmatrix} V^{(32-i)} & V^{(i)} \\ E_{14}(V^{(i)})^T & -E_{14}(V^{(32-i)})^T \end{pmatrix} \\ &= 32 \begin{pmatrix} E_{14} & \mathbf{0} \\ \mathbf{0} & E_{14} \end{pmatrix} \begin{pmatrix} (V^{(32-i)})^T & V^{(i)} \\ (V^{(i)})^T & -V^{(32-i)} \end{pmatrix} \begin{pmatrix} V^{(32-i)} & V^{(i)} \\ (V^{(i)})^T & -(V^{(32-i)})^T \end{pmatrix} \\ &= 32 \begin{pmatrix} E_{14} & \mathbf{0} \\ \mathbf{0} & E_{14} \end{pmatrix} \begin{pmatrix} V^{(32-i)} & V^{(i)} \\ (V^{(i)})^T & -(V^{(32-i)})^T \end{pmatrix}^T \begin{pmatrix} V^{(32-i)} & V^{(i)} \\ (V^{(i)})^T & -(V^{(32-i)})^T \end{pmatrix} \\ &= 32 \begin{pmatrix} E_{14} & \mathbf{0} \\ \mathbf{0} & E_{14} \end{pmatrix} \begin{pmatrix} V^{(i)}(V^{(i)})^T + V^{(32-i)}(V^{(32-i)})^T & \mathbf{0} \\ \mathbf{0} & V^{(i)}(V^{(i)})^T + V^{(32-i)}(V^{(32-i)})^T \end{pmatrix}. \end{aligned} \tag{8.22}$$

Due to Exercise 8.2.4 (set $A = (V^{(32-i)})^T, B = (V^{(i)})^T$), with

$$H = \begin{pmatrix} V^{(32-i)} & V^{(i)} \\ (V^{(i)})^T & -(V^{(32-i)})^T \end{pmatrix} \quad G = \begin{pmatrix} (V^{(32-i)})^T & V^{(i)} \\ (V^{(i)})^T & -V^{(32-i)} \end{pmatrix}$$

we recognize an alternative QMF filter bank. We thus have alias cancellation, with perfect reconstruction only if $|\lambda_{H_0}(\omega)|^2 + |\lambda_{H_0}(\omega + \pi)|^2$. For the two remaining filters we compute

$$\begin{aligned} & (\sqrt{2}V^{(16)})(-\sqrt{2}U^{(48)}) \\ &= -64V^{(16)}V^{(48)} = 64E_{14}V^{(16)}(V^{(16)})^T = 32E_{14}(V^{(16)}(V^{(16)})^T + V^{(16)}(V^{(16)})^T) \end{aligned} \tag{8.23}$$

and

$$\begin{aligned} & (V^{(0)} + E_1 V^{(32)})(E_{-2}U^{(0)} - E_{-1}U^{(32)}) \\ &= 32(V^{(0)} + E_1 V^{(32)})(E_{-2}V^{(0)} - E_{-1}V^{(32)}) = 32E_{-2}(V^{(0)} + E_1 V^{(32)})(V^{(0)} - E_1 V^{(32)}) \\ &= 32E_{-2}(V^{(0)})^2 - (V^{(32)})^2 = 32E_{14}((V^{(0)}(V^{(0)})^T + V^{(32)}(V^{(32)})^T)). \end{aligned} \tag{8.24}$$

We see that the filters from equations (8.22)-(8.24) are similar, and that we thus can combine them into

$$\{V^{(i)}(V^{(i)})^T + V^{(32-i)}(V^{(32-i)})^T\}_{i=0}^{16}. \tag{8.25}$$

All of these can be the identity, except for $1024V^{(16)}(V^{(16)})^T$, since we know that the product of two FIR filters is never the identity, except when both are delays (And all $V^{(m)}$ are FIR, since the prototype filters defined by the MP3 standard are FIR). This single filter is thus what spoils for perfect reconstruction, so that we can only hope

for alias cancellation, and this happens when the filters from Equation (8.25) all are equal. Ideally this is close to cI for some scalar c , and we then have that

$$GH = 16 \cdot 32cE_{33+448} = 512cE_{481}I.$$

This explains the observation from the MP3 standard that GH seems to be close to E_{481} . Since all the filters $V^{(i)}(V^{(i)})^T + V^{(32-i)}(V^{(32-i)})^T$ are symmetric, GH is also a symmetric filter due to Theorem 8.2, so that its frequency response is real, so that we have no phase distortion. We can thus summarize our findings as follows.

Observation 8.15. The prototype filters from the MP3 standard do not give perfect reconstruction. They are found by choosing 17 filters $\{V^{(k)}\}_{k=0}^{16}$ so that the filters from Equation 8.25 are equal, and so that their combination into a prototype filter using equations (8.13) and (8.20) is as close to an ideal bandpass filter as possible. When we have equality the alias cancellation condition is satisfied, and we also have no phase distortion. When the common value is close to $\frac{1}{512}I$, GH is close to E_{481} , so that we have near-perfect reconstruction.

This states clearly the optimization problem which the values stated in the MP3 standard solves.

8.5.3 How can we obtain perfect reconstruction?

How can we overcome the problem that $1024V^{(16)}(V^{(16)})^T \neq I$, which spoiled for perfect reconstruction in the MP3 standard? It turns out that we can address this a simple change in our procedure. In Equation (8.12) we replace with

$$z_{32(s-1)+n} = \sum_{k=0}^{511} \cos((n+1/2)(k+1/2-16)\pi/32) h_k x_{32s-k-1}, \quad (8.26)$$

i.e. $1/2$ is added inside the cosine. We now have the properties

$$\cos(2\pi(n+1/2)(k+64r+1/2)/(2N)) = (-1)^r \cos(2\pi(n+1/2)(k+1/2)/(2N)) \quad (8.27)$$

$$\cos(2\pi(n+1/2)(2N-k-1+1/2)/(2N)) = -\cos(2\pi(n+1/2)(k+1/2)/(2N)). \quad (8.28)$$

Due to the first property, we can deduce as before that

$$\mathbf{z}^{(n)} = \sum_{m=0}^{63} \cos(2\pi(n+1/2)(m+1/2-16)/64) IV^{(m)} \mathbf{x}^{(32-m-1)},$$

where the filters $V^{(m)}$ are defined as before. As before placing the 15 first columns of the cosine-matrix last, but instead using Property (8.28) to combine columns k and $64-k-1$ of the cosine-matrix, we can write this as

$$\begin{pmatrix} \cos(2\pi(0+1/2) \cdot (0+1/2)/64) I & \cdots & \cos(2\pi(0+1/2) \cdot (31+1/2)/64) I \\ \vdots & \ddots & \vdots \\ \cos(2\pi(31+1/2) \cdot (0+1/2)/64) I & \cdots & \cos(2\pi(31+1/2) \cdot (31+1/2)/64) I \end{pmatrix} (A \quad B) \begin{pmatrix} \mathbf{x}^{(31)} \\ \vdots \\ \mathbf{x}^{(-32)} \end{pmatrix}$$

where

$$A = \begin{pmatrix} \mathbf{0} & \mathbf{0} & \cdots & V^{(15)} & V^{(16)} & \cdots & \cdots & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots & \vdots \\ \mathbf{0} & V^{(1)} & \cdots & \mathbf{0} & \mathbf{0} & \cdots & \cdots & V^{(30)} & \mathbf{0} \\ V^{(0)} & \mathbf{0} & \cdots & \mathbf{0} & \mathbf{0} & \cdots & \cdots & \cdots & V^{(31)} \\ \mathbf{0} & \mathbf{0} & \cdots & \mathbf{0} & \mathbf{0} & \mathbf{0} & \cdots & \mathbf{0} \\ \vdots & \vdots \\ \mathbf{0} & \mathbf{0} & \cdots & \mathbf{0} & \mathbf{0} & \mathbf{0} & \cdots & \mathbf{0} \end{pmatrix}$$

$$B = \begin{pmatrix} \mathbf{0} & \mathbf{0} & \cdots & \mathbf{0} & \mathbf{0} & \mathbf{0} & \cdots & \mathbf{0} \\ \vdots & \vdots \\ V_{(32)} & \mathbf{0} & \cdots & \mathbf{0} & \mathbf{0} & \mathbf{0} & \cdots & -V^{(63)} \\ V^{(33)} & \cdots & \mathbf{0} & \mathbf{0} & \mathbf{0} & \cdots & -V^{(62)} & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots & \vdots \\ \mathbf{0} & \mathbf{0} & \cdots & V^{(47)} & -V^{(48)} & \cdots & \cdots & \mathbf{0} \end{pmatrix}.$$

Since the cosine matrix can be written as $\sqrt{\frac{M}{2}} D_M^{(iv)}$, the above can be written as

$$4D_M^{(iv)} (A - B) \begin{pmatrix} \mathbf{x}^{(31)} \\ \vdots \\ \mathbf{x}^{(-32)} \end{pmatrix}.$$

As before we can rewrite this as

$$4D_M^{(iv)} (A - B) \begin{pmatrix} \mathbf{x}^{(31)} \\ \vdots \\ \mathbf{x}^{(0)} \\ E_1 \mathbf{x}^{(31)} \\ \vdots \\ E_1 \mathbf{x}^{(0)} \end{pmatrix} = 4D_M^{(iv)} \left(A \begin{pmatrix} \mathbf{x}^{(31)} \\ \vdots \\ \mathbf{x}^{(0)} \end{pmatrix} + B \begin{pmatrix} E_1 \mathbf{x}^{(31)} \\ \vdots \\ E_1 \mathbf{x}^{(0)} \end{pmatrix} \right),$$

which can be written as

$$4D_M^{(iv)} \begin{pmatrix} \mathbf{0} & \mathbf{0} & \cdots & V^{(15)} & V^{(16)} & \cdots & \cdots & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots & \vdots \\ \mathbf{0} & V^{(1)} & \cdots & \mathbf{0} & \mathbf{0} & \cdots & V^{(30)} & \mathbf{0} \\ V^{(0)} & \mathbf{0} & \cdots & \mathbf{0} & \mathbf{0} & \cdots & \cdots & V^{(31)} \\ E_1 V_{(32)} & \mathbf{0} & \cdots & \mathbf{0} & \mathbf{0} & \cdots & \cdots & -E_1 V^{(63)} \\ \mathbf{0} & E_1 V^{(33)} & \cdots & \mathbf{0} & \mathbf{0} & \cdots & -E_1 V^{(62)} & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots & \vdots \\ \mathbf{0} & \mathbf{0} & \cdots & E_1 V^{(47)} & -E_1 V^{(48)} & \cdots & \cdots & \mathbf{0} \end{pmatrix} \begin{pmatrix} \mathbf{x}^{(31)} \\ \vdots \\ \mathbf{x}^{(0)} \end{pmatrix},$$

which also can be written as

$$4D_M^{(iv)} \begin{pmatrix} \mathbf{0} & \mathbf{0} & \cdots & V^{(16)} & V^{(15)} & \cdots & \cdots & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots & \vdots \\ \mathbf{0} & V^{(30)} & \cdots & \mathbf{0} & \mathbf{0} & \cdots & V^{(1)} & \mathbf{0} \\ V^{(31)} & \mathbf{0} & \cdots & \mathbf{0} & \mathbf{0} & \cdots & \cdots & V^{(0)} \\ -E_1 V_{(63)} & \mathbf{0} & \cdots & \mathbf{0} & \mathbf{0} & \cdots & \cdots & E_1 V^{(32)} \\ \mathbf{0} & -E_1 V^{(62)} & \cdots & \mathbf{0} & \mathbf{0} & \cdots & E_1 V^{(33)} & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots & \vdots \\ \mathbf{0} & \mathbf{0} & \cdots & -E_1 V^{(48)} & E_1 V^{(47)} & \cdots & \cdots & \mathbf{0} \end{pmatrix} \begin{pmatrix} \mathbf{x}^{(0)} \\ \vdots \\ \mathbf{x}^{(31)} \end{pmatrix}.$$

We therefore have the following result

Theorem 8.16. (Polyphase factorization of a forward filter bank transform based on a prototype filter, modified version). The modified version of the polyphase form of a forward filter bank transform based on a prototype filter can be factored as

$$4D_M^{(iv)} \begin{pmatrix} \mathbf{0} & \mathbf{0} & \cdots & V^{(16)} & V^{(15)} & \cdots & \cdots & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots & \vdots \\ \mathbf{0} & V^{(30)} & \cdots & \mathbf{0} & \mathbf{0} & \cdots & V^{(1)} & \mathbf{0} \\ V^{(31)} & \mathbf{0} & \cdots & \mathbf{0} & \mathbf{0} & \cdots & \cdots & V^{(0)} \\ -E_1 V_{(63)} & \mathbf{0} & \cdots & \mathbf{0} & \mathbf{0} & \cdots & \cdots & E_1 V^{(32)} \\ \mathbf{0} & -E_1 V^{(62)} & \cdots & \mathbf{0} & \mathbf{0} & \cdots & E_1 V^{(33)} & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots & \vdots \\ \mathbf{0} & \mathbf{0} & \cdots & -E_1 V^{(48)} & E_1 V^{(47)} & \cdots & \cdots & \mathbf{0} \end{pmatrix} \quad (8.29)$$

Clearly this factorization avoids having two blocks of filters: There are now 16×2 -polyphase matrices, and as we know, each of them can be invertible, so that the full matrix can be inverted in a similar fashion as before. It is therefore now possible to obtain perfect reconstruction. Although we do not state recipes for implementing this, one has just as efficient implementations as in the MP3 standard.

Since we ended up with the 2×2 polyphase matrices M_k , we can apply the lifting factorization in order to halve the number of multiplications/additions. This is not done in practice, since a lifting factorization requires that we compute all outputs at once. In audio coding it is required that we compute the output progressively, due to the large size of the input vector. The procedure above is therefore mostly useful for providing the requirements for the filters, while the preceding comments can be used for the implementation.

Exercises for Section 8.5

- Run the forward and then the reverse transform from Exercise 6.4.2 on the vector

```
x=(1:8192)';
```

Verify that there seems to be a delay on 481 elements, as promised by Therorem 8.15.
Do you get the exact same result back?

2. Use Matlab to verify the symmetries we have stated for the symmetries in the prototype filters, i.e. that

$$C_i = \begin{cases} -C_{512-i} & i \neq 64, 128, \dots, 448 \\ C_{512-i} & i = 64, 128, \dots, 448. \end{cases}$$

Explain also that this implies that $h_i = h_{512-i}$ for $i = 1, \dots, 511$. In other words, the prototype filter has symmetry around $(511 + 1)/2 = 256$, so that it has linear phase.

3. We mentioned that we could use the lifting factorization to construct filters on the form stated in Equation (8.13), so that the matrices on the form given by Equation (8.17), i.e.

$$\begin{pmatrix} V^{(32-i)} & V^{(i)} \\ -V^{(64-i)} & V^{(32+i)} \end{pmatrix},$$

are invertible. Let us see what kind of lifting steps produce such matrices.

a. Show that the lifting steps $\begin{pmatrix} I & \lambda E_2 \\ \mathbf{0} & I \end{pmatrix}$ and $\begin{pmatrix} I & \mathbf{0} \\ \lambda I & I \end{pmatrix}$ applied in alternating order to a matrix on the form given by Equation (8.17), where the filters are on the from given by Equation (8.13), again produces matrices and filters on these forms. This explains how we can parametrize a larger number of such matrices with the help of lifting steps. It also explain why the inverse matrix is on the form stated in Equation (8.17) with filters on the same form, since the inverse lifting steps are of the same type.

b. Explain that 16 numbers $\{\lambda_i\}_{i=1}^{16}$ are needed (together with what we start with on the diagonal in the lifting construction), in order to construct filters so that the prototype filter has 512 coefficients. Since there are 15 submatrices, this gives 240 optimization variables.

Lifting gives the following strategy for finding a corresponding synthesis prototype filter which gives perfect reconstruction: First compute matrices V, W which are inverses of one another using lifting (using the lifting steps of this exercise ensures that all filters will be on the form stated in Equation (8.13)), and write

$$\begin{aligned} VW &= \begin{pmatrix} V^{(1)} & V^{(2)} \\ -V^{(3)} & V^{(4)} \end{pmatrix} \begin{pmatrix} W^{(1)} & -W^{(3)} \\ W^{(2)} & W^{(4)} \end{pmatrix} = \begin{pmatrix} V^{(1)} & V^{(2)} \\ -V^{(3)} & V^{(4)} \end{pmatrix} \begin{pmatrix} (W^{(1)})^T & (W^{(2)})^T \\ -(W^{(3)})^T & (W^{(4)})^T \end{pmatrix}^T \\ &= \begin{pmatrix} V^{(1)} & V^{(2)} \\ -V^{(3)} & V^{(4)} \end{pmatrix} \begin{pmatrix} E_{15}(W^{(1)})^T & E_{15}(W^{(2)})^T \\ -E_{15}(W^{(3)})^T & E_{15}(W^{(4)})^T \end{pmatrix}^T \begin{pmatrix} E_{15} & \mathbf{0} \\ \mathbf{0} & E_{15} \end{pmatrix} = I. \end{aligned}$$

Now, the matrices $U^{(i)} = E_{15}(W^{(i)})^T$ are on the form stated in Equation (8.13), and we have that

$$\begin{pmatrix} V^{(1)} & V^{(2)} \\ -V^{(3)} & V^{(4)} \end{pmatrix} \begin{pmatrix} U^{(1)} & U^{(2)} \\ -U^{(3)} & U^{(4)} \end{pmatrix} = \begin{pmatrix} E_{-15} & \mathbf{0} \\ \mathbf{0} & E_{-15} \end{pmatrix}$$

We can now conclude from Theorem 8.14 that if we define the synthesis prototype filter as therein, and set $c = 1$, $d = -15$, we have that $GH = 16E_{481-32\cdot 15} = 16E_1$.

8.6 Summary

We defined the polyphase representation of a matrix, and proved some useful properties. For filter bank transforms, the polyphase representation was a block matrix where the blocks are filters, and these blocks/filters were called polyphase components. In particular, the filter bank transforms of wavelets were 2×2 -block matrices of filters. We saw that, for wavelets, the polyphase representation could be realized through a rearrangement of the wavelet bases, and thus paralleled the development in Chapter 6 for expressing the DWT in terms of filters, where we instead rearranged the target base of the DWT.

We showed with two examples that factoring the polyphase representation into simpler matrices (also referred to as a polyphase factorization) could be a useful technique. First, for wavelets ($M = 2$), we established the lifting factorization. This is useful not only since it factorizes the DWT and the IDWT into simpler operations, but also since it reduces the number of arithmetic operations in these. The lifting factorization is therefore also used in practical implementations, and we applied it to some of the wavelets we constructed in Chapter 7. The JPEG2000 standard document [1] explains a procedure for implementing some of these wavelet transforms using lifting, and the values of the lifting steps used in the standard thus also appear here.

The polyphase representation was also useful for proving the characterization of wavelets we encountered in Chapter 7, which we used to find expressions for many useful wavelets.

The polyphase representation was also useful to explain how the prototype filters of the MP3 standard should be chosen, in order for the reverse filter bank transform to invert the forward filter bank transform. Again this was attacked by factoring the polyphase representation of the forward and reverse filter bank transforms. The parts of the factorization which represented the prototype filters were represented by a sparse matrix, and it was clear from this matrix what properties we needed to put on the prototype filter, in order to have alias cancellation, and no phase distortion. In fact, we proved that the MP3 standard could not possibly give perfect reconstruction, but it was very clear from our construction how the filter bank could be modified in order for the overall system to provide perfect reconstruction.

This development concludes the one-dimensional aspect of wavelets in this book. In the following we will extend our theory to also apply for images. Images will be presented in Chapter 9. After that we will define the tensor product concept, which will be the key ingredient to apply wavelets to two-dimensional objects such as images.

Digital images

Upto now we have presented wavelets in a one-dimensional setting. Images, however, are two-dimensional by nature. This poses another challenge, which we did not encounter in the case of sound signals. In this chapter we will establish the mathematics to handle this, but first we will present some basics on images, as well as how they can be represented and manipulated with simple mathematics. Images are a very important type of digital media, and this material is thus useful, general knowledge for anyone with a digital camera and a computer. For many scientists this material is also an essential tool. As an example, in astrophysics data from both satellites and distant stars and galaxies is collected in the form of images, and information is extracted from the images with advanced image processing techniques. As another example, medical imaging makes it possible to gather different kinds of information in the form of images, even from the inside of the body. By analysing these images it is possible to discover tumours and other disorders.

We will see how filter-based operations extend naturally to the two-dimensional setting of images. Smoothing and edge detections are the two main examples of filter-based operations we will consider for images. The key mathematical concept in this extension is the *tensor product*, which can be thought of as a general tool for constructing two-dimensional objects from one-dimensional counterparts. We will also see that the tensor product allows us to establish an efficient implementation of filtering for images, efficient meaning a complexity substantially less than what is required by general linear transformations.

We will finally consider useful coordinate changes for images. Recall that the DFT, the DCT, and the wavelet transform were all defined as changes of coordinates for vectors or functions of one variable, and therefore cannot be directly applied to two-dimensional data like images. It turns out that the tensor product can also be used to extend changes of coordinates to a two-dimensional setting.

9.1 What is an image?

Before we do computations with images, it is helpful to be clear about what an image really is. Images cannot be perceived unless there is some light present, so we first review superficially what light is.

9.1.1 Light

Fact 9.1 (Light). Light is electromagnetic radiation with wavelengths in the range 400–700 nm (1 nm is 10^{-9} m): Violet has wavelength 400 nm and red has wavelength 700 nm. White light contains roughly equal amounts of all wave lengths.

Other examples of electromagnetic radiation are gamma radiation, ultraviolet and infrared radiation and radio waves, and all electromagnetic radiation travel at the speed of light ($\approx 3 \times 10^8$ m/s). Electromagnetic radiation consists of waves and may be reflected and refracted, just like sound waves (but sound waves are not electromagnetic waves).

We can only see objects that emit light, and there are two ways that this can happen. The object can emit light itself, like a lamp or a computer monitor, or it reflects light that falls on it. An object that reflects light usually absorbs light as well. If we perceive the object as red it means that the object absorbs all light except red, which is reflected. An object that emits light is different; if it is to be perceived as being red it must emit only red light.

9.1.2 Digital output media

Our focus will be on objects that emit light, for example a computer display. A computer monitor consists of a matrix of small dots which emit light. In most technologies, each dot is really three smaller dots, and each of these smaller dots emit red, green and blue light. If the amounts of red, green and blue is varied, our brain merges the light from the three small light sources and perceives light of different colours. In this way the colour at each set of three dots can be controlled, and a colour image can be built from the total number of dots.

It is important to realise that it is possible to generate most, but not all, colours by mixing red, green and blue. In addition, different computer monitors use slightly different red, green and blue colours, and unless this is taken into consideration, colours will look different on the two monitors. This also means that some colours that can be displayed on one monitor may not be displayable on a different monitor.

Printers use the same principle of building an image from small dots. On most printers however, the small dots do not consist of smaller dots of different colours. Instead as many as 7–8 different inks (or similar substances) are mixed to the right colour. This makes it possible to produce a wide range of colours, but not all, and the problem of matching a colour from another device like a monitor is at least as difficult as matching different colours across different monitors.

Video projectors builds an image that is projected onto a wall. The final image is therefore a reflected image and it is important that the surface is white so that it reflects all colours equally.

The quality of a device is closely linked to the density of the dots.

Fact 9.2 (Resolution). The resolution of a medium is the number of dots per inch (dpi). The number of dots per inch for monitors is usually in the range 70–120, while for printers it is in the range 150–4800 dpi. The horizontal and vertical densities may be different. On a monitor the dots are usually referred to as *pixels* (picture elements).

9.1.3 Digital input media

The two most common ways to acquire digital images is with a digital camera or a scanner. A scanner essentially takes a photo of a document in the form of a matrix of (possibly coloured) dots. As for printers, an important measure of quality is the number of dots per inch.

Fact 9.3. The resolution of a scanner usually varies in the range 75 dpi to 9600 dpi, and the colour is represented with up to 48 bits per dot.

For digital cameras it does not make sense to measure the resolution in dots per inch, as this depends on how the image is printed (its size). Instead the resolution is measured in the number of dots recorded.

Fact 9.4. The number of pixels recorded by a digital camera usually varies in the range 320×240 to 6000×4000 with 24 bits of colour information per pixel. The total number of pixels varies in the range 76 800 to 24 000 000 (0.077 megapixels to 24 megapixels).

For scanners and cameras it is easy to think that the more dots (pixels), the better the quality. Although there is some truth to this, there are many other factors that influence the quality. The main problem is that the measured colour information is very easily polluted by noise. And of course high resolution also means that the resulting files become very big; an uncompressed 6000×4000 image produces a 72 MB file. The advantage of high resolution is that you can magnify the image considerably and still maintain reasonable quality.

9.1.4 Definition of digital image

We have already talked about digital images, but we have not yet been precise about what they are. From a mathematical point of view, an image is quite simple.



Figure 9.1: Different versions of the same image; black and white (a), grey-level (b), and colour (c).

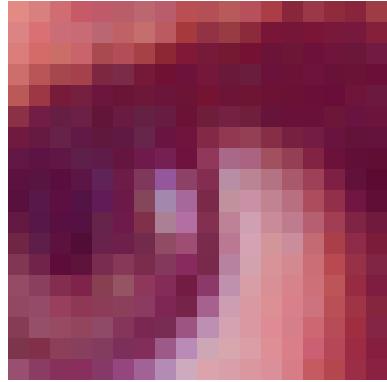
Fact 9.5 (Digital image). A digital image P is a matrix of *intensity values* $\{p_{i,j}\}_{i,j=1}^{M,N}$. For grey-level images, the value $p_{i,j}$ is a single number, while for colour images each $p_{i,j}$ is a vector of three or more values. If the image is recorded in the rgb-model, each $p_{i,j}$ is a vector of three values,

$$p_{i,j} = (r_{i,j}, g_{i,j}, b_{i,j}),$$

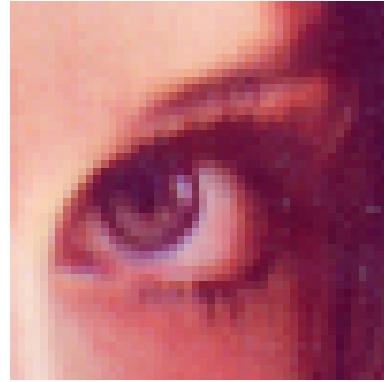
that denote the amount of red, green and blue at the point (i, j) .

Note that, when referring to the coordinates (i, j) in an image, i will refer to row index, j to column index, in the same way as for matrices. In particular, the top row in the image have coordinates $\{(0, j)\}_{j=0}^{N-1}$, while the left column in the image has coordinates $\{(i, 0)\}_{i=0}^{M-1}$. With this notation, the dimension of the image is $M \times N$. The value $p_{i,j}$ gives the colour information at the point (i, j) . It is important to remember that there are many formats for this. The simplest case is plain black and white images in which case $p_{i,j}$ is either 0 or 1. For grey-level images the intensities are usually integers in the range 0–255. However, we will assume that the intensities vary in the interval [0, 1], as this sometimes simplifies the form of some mathematical functions. For colour images there are many different formats, but we will just consider the rgb-format mentioned in the fact box. Usually the three components are given as integers in the range 0–255, but as for grey-level images, we will assume that they are real numbers in the interval [0, 1] (the conversion between the two ranges is straightforward, see Example 9.10 below). In Figure 9.1(c) we have shown a test image which we will use a lot, called the *Lena image*, named after the girl in the image. This image is also used as a test image in many image processing textbooks. In (a) and (b) the corresponding black and white, and grey-level versions of this image are shown.

Fact 9.6. In these notes the intensity values $p_{i,j}$ are assumed to be real numbers in the interval [0, 1]. For colour images, each of the red, green, and blue intensity



(a) 18×18 pixels



(b) 50×50 pixels

Figure 9.2: Two excerpts of the colour image in figure 9.1. The grid indicates the borders between the pixels.

values are assumed to be real numbers in $[0, 1]$.

If we magnify the part of the colour image in figure 9.1 around one of the eyes, we obtain the images in figure 9.2. As we can see, the pixels have been magnified to big squares. This is a standard representation used by many programs — the actual shape of the pixels will depend on the output medium. Nevertheless, we will consider the pixels to be square, with integer coordinates at their centres, as indicated by the grids in figure 9.2.

Fact 9.7 (Shape of pixel). The pixels of an image are assumed to be square with sides of length one, with the pixel with value $p_{i,j}$ centred at the point (i, j) .

9.2 Some simple operations on images

Images are two-dimensional matrices of numbers, contrary to the sound signals we considered in the previous section. In this respect it is quite obvious that we can manipulate an image by performing mathematical operations on the numbers. In this section we will consider some of the simpler operations. In later sections we will go through more advanced operations, and explain how the theory for these can be generalized from the corresponding theory for one-dimensional (sound) signals (which we will go through first).

In order to perform these operations, we need to be able to use images with a programming environment.

9.2.1 Images and MATLAB

An image can also be thought of as a matrix, by associating each pixel with an element in a matrix. The matrix indices thus correspond to positions in the pixel grid. Black and white images correspond to matrices where the elements are natural numbers between 0 and 255. To store a colour image, we need 3 matrices, one for each colour component. This enables us to use linear algebra packages, such as MATLAB, on order to work with images. After we now have made the connection with matrices, we can create images from mathematical formulas, just as we could with sound in the previous sections. But what we also need before we go through operations on images, is, as in the sections on sound, means of reading an image from a file so that its contents are accessible as a matrix, and write images represented by a matrix which we have constructed ourselves to file. Reading a function from file can be done with help of the function `imread`. If we write

```
X = double(imread('filename(fmt','fmt')));
```

the image with the given path and format is read, and stored in the matrix which we call `X`. `'fmt'` can be `'jpg'`, `'tif'`, `'png'`... You should consult the MATLAB help pages to see which formats are supported. After the call to `imread`, we have a matrix where the entries represent the pixel values, and of integer data type (more precisely, the data type `uint8` in Matlab). To perform operations on the image, we must first convert the entries to the data type `double`. This is done with a call to the Matlab function `double`. Similarly, the function `imwrite` can be used to write the image represented by a matrix to file. If we write

```
imwrite(uint8(X), 'filename(fmt','fmt')
```

the image represented by the matrix is written to the given path, in the given format. Before the image is written to file, you see that we have converted the matrix values back to the integer data type with the help of the function `uint8`. In other words: `imread` and `imwrite` both assume integer matrix entries, while operations on matrices assume double matrix entries. If you want to print images you have created yourself, you can use this function first to write the image to a file, and then send that file to the printer using another program. Finally, we need an alternative to playing a sound, namely displaying an image. The function `imshow(uint8(X))` displays the matrix `X` as an image in a separate window. Also here we needed to convert the samples using the function `uint8`.

The following examples go through some much used operations on images.

Example 9.8 (Normalising the intensities). We have assumed that the intensities all lie in the interval $[0, 1]$, but as we noted, many formats in fact use integer values in the range $[0, 255]$. And as we perform computations with the intensities, we quickly end up with intensities outside $[0, 1]$ even if we start out with intensities within this interval. We therefore need to be able to *normalise* the intensities. This we can do with the simple linear function

$$g(x) = \frac{x - a}{b - a}, \quad a < b,$$

which maps the interval $[a, b]$ to $[0, 1]$. A simple case is mapping $[0, 255]$ to $[0, 1]$ which we accomplish with the scaling $g(x) = x/255$. More generally, we typically perform computations that result in intensities outside the interval $[0, 1]$. We can then compute the minimum and maximum intensities p_{\min} and p_{\max} and map the interval $[p_{\min}, p_{\max}]$ back to $[0, 1]$. Below we have shown a function `mapto01` which achieves this task.

```
function Z=mapto01(X)
    minval = min(min(X));
    maxval = max(max(X));
    Z = (X - minval)/(maxval-minval);
```

Several examples of using this function will be shown below. A good question here is why the functions `min` and `max` are called three times in succession. The reason is that there is a third “dimension” in play, besides the spatial x - and y -directions. This dimension describes the color components in each pixel, which are usually the red-, green-, and blue colour components. ♣

Example 9.9 (Extracting the different colours). If we have a colour image

$$P = (r_{i,j}, g_{i,j}, b_{i,j})_{i,j=1}^{m,n},$$

it is often useful to manipulate the three colour components separately as the three images

$$P_r = (r_{i,j})_{i,j=1}^{m,n}, \quad P_g = (g_{i,j})_{i,j=1}^{m,n}, \quad P_b = (b_{i,j})_{i,j=1}^{m,n}.$$

As an example, let us first see how we can produce three separate images, showing the R, G, and B colour components, respectively. Let us take the image `lena.png` used in Figure 9.1. When the image is read (first line below), the returned object has three dimensions. The first two dimensions represent the spatial directions (the row-index and column-index). The third dimension represents the colour component. One can therefore view images representing the different colour components with the help of the following code:

```
X=double(imread('lena.png','png'));

Z=zeros(size(X));
Z(:,:,1)=X(:,:,1);
imshow(uint8(Z))

Z=zeros(size(X));
Z(:,:,2)=X(:,:,2);
imshow(uint8(Z))

Z=zeros(size(X));
Z(:,:,3)=X(:,:,3);
imshow(uint8(Z))
```

The resulting image files are shown in Figure 9.3. ♣

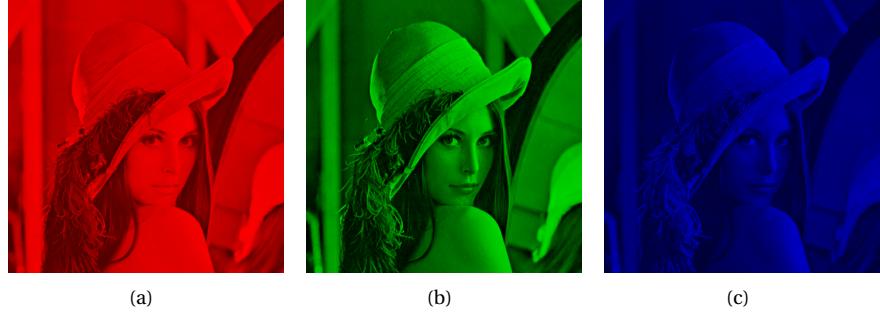


Figure 9.3: The red (a), green (b), and blue (c) components of the colour image in Figure 9.1.

Example 9.10 (Converting from colour to grey-level). If we have a colour image we can convert it to a grey-level image. This means that at each point in the image we have to replace the three colour values (r, g, b) by a single value p that will represent the grey level. If we want the grey-level image to be a reasonable representation of the colour image, the value p should somehow reflect the intensity of the image at the point. There are several ways to do this.

It is not unreasonable to use the largest of the three colour components as a measure of the intensity, i.e., to set $p = \max(r, g, b)$. The result of this can be seen in Figure 9.4(a).

An alternative is to use the sum of the three values as a measure of the total intensity at the point. This corresponds to setting $p = r + g + b$. Here we have to be a bit careful with a subtle point. We have required each of the r , g and b values to lie in the range $[0, 1]$, but their sum may of course become as large as 3. We also require our grey-level values to lie in the range $[0, 1]$ so after having computed all the sums we must normalise as explained above. The result can be seen in Figure 9.4(b).

A third possibility is to think of the intensity of (r, g, b) as the length of the colour vector, in analogy with points in space, and set $p = \sqrt{r^2 + g^2 + b^2}$. Again, we may end up with values in the range $[0, \sqrt{3}]$ so we have to normalise like we did in the second case. The result is shown in Figure 9.4(c).

Let us sum this up as follows.

Observation 9.11 (Conversion from colour to grey level). A colour image $P = (r_{i,j}, g_{i,j}, b_{i,j})_{i,j=1}^{m,n}$ can be converted to a grey level image $Q = (q_{i,j})_{i,j=1}^{m,n}$ by one of the following three operations:

1. Set $q_{i,j} = \max(r_{i,j}, g_{i,j}, b_{i,j})$ for all i and j .
2. (a) Compute $\hat{q}_{i,j} = r_{i,j} + g_{i,j} + b_{i,j}$ for all i and j .
(b) Transform all the values to the interval $[0, 1]$ by setting

$$q_{i,j} = \frac{\hat{q}_{i,j}}{\max_{k,l} \hat{q}_{k,l}}.$$



(a) Each colour triple has been replaced by its maximum

(b) Each colour triple has been replaced by its sum and the result mapped to $(0, 1)$

(c) Each triple has been replaced by its length and the result mapped to $(0, 1)$

Figure 9.4: Alternative ways to convert the colour image in Figure 9.1 to a grey level image.

3. (a) Compute $\hat{q}_{i,j} = \sqrt{r_{i,j}^2 + g_{i,j}^2 + b_{i,j}^2}$ for all i and j .
(b) Transform all the values to the interval $[0, 1]$ by setting

$$q_{i,j} = \frac{\hat{q}_{i,j}}{\max_{k,l} \hat{q}_{k,l}}.$$

This can be implemented by using most of the code from the previous example, and replacing with the lines

```
Z1=max(X, [], 3);
newvals=X(:,:,1)+X(:,:,2)+X(:,:,3);
Z2=newvals/max(max(newvals))*255;
newvals=sqrt(X(:,:,1).^2+X(:,:,2).^2+X(:,:,3).^2);
Z3=newvals/max(max(newvals))*255;
```

respectively. In practice one of the last two methods are usually preferred, perhaps with a preference for the last method, but the actual choice depends on the application.



Example 9.12 (Computing the negative image). In film-based photography a negative image was obtained when the film was developed, and then a positive image was created from the negative. We can easily simulate this and compute a negative digital image.

Suppose we have a grey-level image $P = (p_{i,j})_{i,j=1}^{m,n}$ with intensity values in the interval $[0, 1]$. Here intensity value 0 corresponds to black and 1 corresponds to white. To obtain the negative image we just have to replace an intensity p by its 'mirror value' $1 - p$.

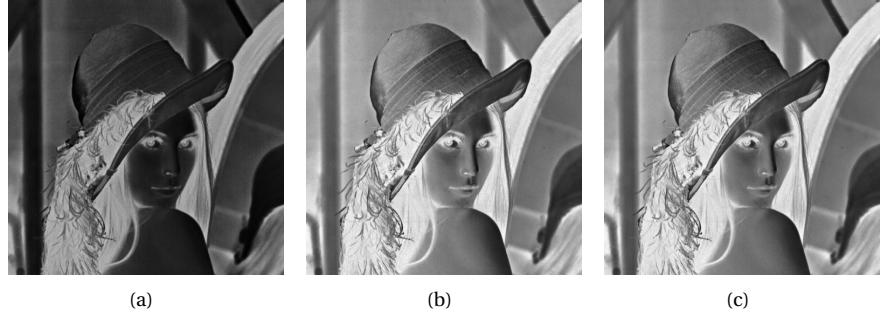


Figure 9.5: The negative versions of the corresponding images in figure 9.4.

Fact 9.13 (Negative image). Suppose the grey-level image $P = (p_{i,j})_{i,j=1}^{m,n}$ is given, with intensity values in the interval $[0, 1]$. The negative image $Q = (q_{i,j})_{i,j=1}^{m,n}$ has intensity values given by $q_{i,j} = 1 - p_{i,j}$ for all i and j .

This is also easily translated to code as above. The resulting image is shown in Figure 9.5. ♣

Example 9.14 (Increasing the contrast). A common problem with images is that the contrast often is not good enough. This typically means that a large proportion of the grey values are concentrated in a rather small subinterval of $[0, 1]$. The obvious solution to this problem is to somehow spread out the values. This can be accomplished by applying a function f to the intensity values, i.e., new intensity values are computed by the formula

$$\hat{p}_{i,j} = f(p_{i,j})$$

for all i and j . If we choose f so that its derivative is large in the area where many intensity values are concentrated, we obtain the desired effect.

Figure 9.6 shows some examples. The functions in the left plot have quite large derivatives near $x = 0.5$ and will therefore increase the contrast in images with a concentration of intensities with value around 0.5. The functions are all on the form

$$f_n(x) = \frac{\arctan(n(x - 1/2))}{2 \arctan(n/2)} + \frac{1}{2}. \quad (9.1)$$

For any $n \neq 0$ these functions satisfy the conditions $f_n(0) = 0$ and $f_n(1) = 1$. The three functions in figure 9.6(a) correspond to $n = 4, 10$, and 100 .

Functions of the kind shown in figure 9.6(b) have a large derivative near $x = 0$ and will therefore increase the contrast in an image with a large proportion of small intensity values, i.e., very dark images. The functions are given by

$$g_\epsilon(x) = \frac{\ln(x + \epsilon) - \ln \epsilon}{\ln(1 + \epsilon) - \ln \epsilon}, \quad (9.2)$$

and the ones shown in the plot correspond to $\epsilon = 0.1, 0.01$, and 0.001 .

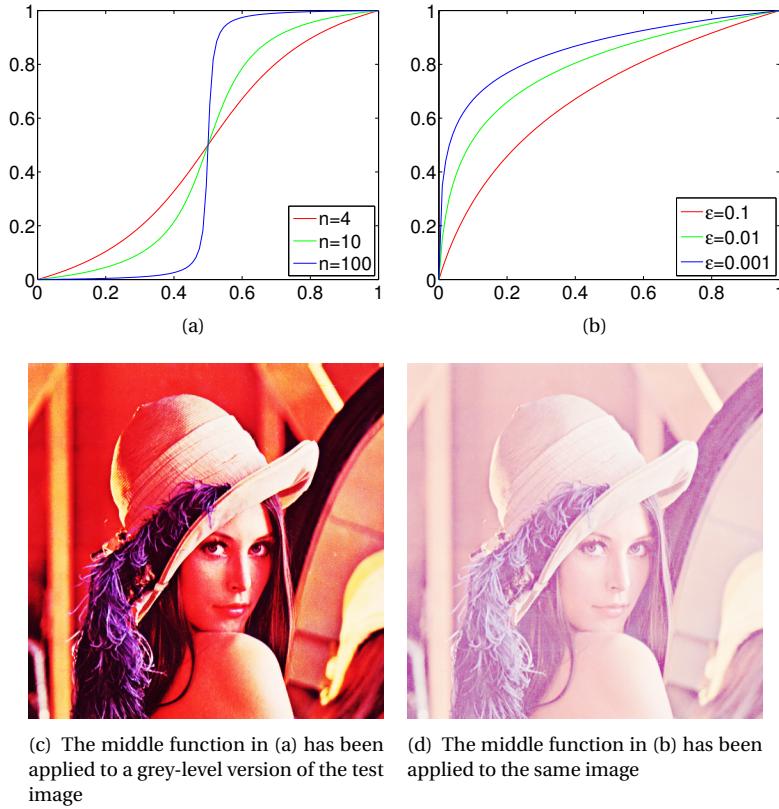


Figure 9.6: The plots in (a) and (b) show some functions that can be used to improve the contrast of an image.

In figure 9.6(c) the middle function in (a) has been applied to the image in figure 9.4(c). Since the image was quite well balanced, this has made the dark areas too dark and the bright areas too bright. In figure 9.6(d) the function in (b) has been applied to the same image. This has made the image as a whole too bright, but has brought out the details of the road which was very dark in the original.

Observation 9.15. Suppose a large proportion of the intensity values $p_{i,j}$ of a grey-level image P lie in a subinterval I of $[0, 1]$. Then the contrast of the image can be improved by computing new intensities $\hat{p}_{i,j} = f(p_{i,j})$ where f is a function with a large derivative in the interval I .

Increasing the contrast is easy to implement. The following function uses the contrast adjusting function from Equation (9.2), with ϵ as in that equation as parameter

```

function Z=contrastadjust(X,epsilon)
    Z = X/255; % Maps the pixel values to [0,1]
    Z = (log(Z+epsilon) - log(epsilon))/...
        (log(1+epsilon)-log(epsilon));
    Z = Z*255; % Maps the values back to [0,255]

```

This has been used to generate the image in Figure 9.6(d). ♣

What you should have learnt in this section

How to read, write, and show images with Matlab. How to extract different colour components, convert from colour to grey-level images, and use functions for adjusting the contrast.

Exercises for Section 9.2

1. Black and white images can be generated from greyscale images (with values between 0 and 255) by replacing each pixel value with the one of 0 and 255 which is closest. Use this strategy to generate the black and white image shown in Figure 9.1(a).
2. Generate the image in Figure 9.6(d) on your own by writing code which uses the function `contrastadjust`.
3. Let us also consider the second way we mentioned for increasing the contrast.
 - a. Write a function `contrastadjust0` which instead uses the function from Equation (9.1) to increase the contrast. n should be a parameter to the function.
 - b. Generate the image in Figure 9.6(c) on your own by using your code from Exercise 2, and instead calling the function `contrastadjust0`.
4. In this exercise we will look at another function for increasing the contrast of a picture.
 - a. Show that the function $f : \mathbb{R} \rightarrow \mathbb{R}$ given by

$$f_n(x) = x^n,$$

for all n maps the interval $[0, 1] \rightarrow [0, 1]$, and that $f'(1) \rightarrow \infty$ as $n \rightarrow \infty$.

- b. The color image `secret.jpg`, shown in Figure 9.7, contains some information that is nearly invisible to the naked eye on most computer monitors. Use the function $f(x)$, to reveal the secret message.
Hint: You will first need to convert the image to a greyscale image. You can then use the function `contrastadjust` as a starting point for your own program.



Figure 9.7: Secret message

9.3 Filter-based operations on images

The next examples of operations on images we consider will use filters. These examples define what it means to apply a filter to two-dimensional data. We start with the following definition of a computational molecule. This term stems from image processing, and seems at the outset to be unrelated to filters.

Definition 9.16 (Computational molecules). We say that an operation S on an image X is given by the *computational molecule*

$$A = \begin{pmatrix} \vdots & \vdots & \vdots & \vdots & \vdots \\ \cdots & a_{-1,-1} & a_{-1,0} & a_{-1,1} & \cdots \\ \cdots & a_{0,-1} & \underline{a_{0,0}} & a_{0,1} & \cdots \\ \cdots & a_{1,-1} & a_{1,0} & a_{1,1} & \cdots \\ \vdots & \vdots & \vdots & \vdots & \vdots \end{pmatrix}$$

if we have that

$$(SX)_{i,j} = \sum_{k_1, k_2} a_{k_1, k_2} X_{i+k_1, j+k_2}. \quad (9.3)$$

In the molecule, indices are allowed to be both positive and negative, we underline the element with index $(0,0)$ (the centre of the molecule), and assume that $a_{i,j}$ with indices falling outside those listed in the molecule are zero (as for compact filter notation).

The interpretation of a computational molecule is that we place the centre of the molecule on a pixel, multiply the pixel and its neighbours by the corresponding weights $a_{i,j}$, and finally sum up in order to produce the resulting value. This type of operation will turn out to be particularly useful for images. The following result expresses how computational molecules and filters are related. It states that, if we apply one filter to all the columns, and then another filter to all the rows, the end

result can be expressed with the help of a computational molecule.

Theorem 9.17 (Filtering and computational molecules). Let S_1 and S_2 be filters with compact filter notation \mathbf{t}_1 and \mathbf{t}_2 , respectively, and consider the operation S where S_1 is first applied to the columns in the images, and then S_2 is applied to the rows in the image. Then S is an operation which can be expressed in terms of the computational molecule $a_{ij} = (\text{rev}(\mathbf{t}_1))_i (\text{rev}(\mathbf{t}_2))_j$, where $\text{rev}(\mathbf{t})$ denotes the vector with the elements of \mathbf{t} reversed.

Proof: Let $X_{i,j}$ be the pixels in the image. When we apply S_1 to the columns of X we get the image Y defined by

$$Y_{i,j} = \sum_{k_1} \sum_{k_2} (t_1)_{k_1} X_{i-k_1, j}.$$

When we apply S_2 to the rows of Y we get the image Z defined by

$$\begin{aligned} Z_{i,j} &= \sum_{k_2} (t_2)_{k_2} Y_{i, j-k_2} = \sum_{k_2} (t_2)_{k_2} \sum_{k_1} (t_1)_{k_1} X_{i-k_1, j-k_2} \\ &= \sum_{k_1} \sum_{k_2} (t_1)_{k_1} (t_2)_{k_2} X_{i-k_1, j-k_2} = \sum_{k_1} \sum_{k_2} (\text{rev}(t_1))_{-k_1} (\text{rev}(t_2))_{-k_2} X_{i-k_1, j-k_2} \\ &= \sum_{k_1} \sum_{k_2} (\text{rev}(t_1))_{k_1} (\text{rev}(t_2))_{k_2} X_{i+k_1, j+k_2}, \end{aligned}$$

where we finally substituted $-k_1, -k_2$ with k_1, k_2 . Comparing with Equation (9.3) we see that S is given by the computational molecule with entries $a_{ij} = (\text{rev}(\mathbf{t}_1))_i (\text{rev}(\mathbf{t}_2))_j$. ■

Note that, when we filter an image with S_1 and S_2 in this way, the order does not matter: since computing $S_1 X$ is the same as applying S_1 to all columns of X , and computing $Y(S_2)^T$ is the same as applying S_2 to all rows of Y , the combined filtering operation takes the form

$$SX = S_1 X (S_2)^T, \quad (9.4)$$

and the fact that the order does not matter simply boils down to the fact that it does not matter which of the left or right multiplications we perform first. Applying S_1 to the columns of X is what we call a *vertical filtering operation*, while applying S_2 to the rows of X is what we call a *horizontal filtering operation*. We can thus state the following.

Observation 9.18. The order of vertical and horizontal filtering of an image does not matter.

Most computational molecules we will consider in the following can be expressed in terms of filters as in this theorem, but clearly there exist also computational molecules which are not on this form, since the matrix A with entries $a_{ij} = (\text{rev}(\mathbf{t}_1))_i (\text{rev}(\mathbf{t}_2))_j$ has rank one, and a general computational molecule can have any rank.

In the previous theorem, note that the filter coefficients need to be reversed when we apply the computational molecule. This may seem strange. The difference lies in that we define the application of a computational molecule in a different way than application of a filter: a filter is “placed” over the samples in reverse order, contrary to a computational molecule. To be more precise, this has to do with the difference in the sign in front of k and k_1, k_2 , respectively, in the equations

$$(Sx)_n = \sum_k s_k x_{n-k} \quad (SX)_{i,j} = \sum_{k_1, k_2} a_{k_1, k_2} X_{i+k_1, j+k_2}.$$

In most of the examples the filters are symmetric. In this case the reversal of the filter coefficients causes no confusion.

It is straightforward to write a function which filters the rows and the columns. Assume that the image is stored as the matrix X , and that the computational molecule is obtained by applying the filter S_1 to the columns, and the filter S_2 to the rows. Then the following function will apply the computational molecule to the image (we have assumed that the filter lengths are odd, and that the middle filter coefficient has index 0):

```
function Z=filterimage(X,S1,S2)
    Z=X;
    [m,n]=size(X);
    S1skip=(length(S1)-1)/2; S2skip=(length(S2)-1)/2;
    for col=1:n
        res=conv(Z(:,col),S1);
        Z(:,col)=res((1+S1skip):(length(res)-S1skip));
    end
    for row=1:m
        res=conv(Z(row,:),S2);
        Z(row,:)=res((1+S2skip):(length(res)-S2skip));
    end
```

We have here adapted the simple convention that all pixels outside the image have intensity 0, i.e. we do not make a periodic extension of the image. Instead of using the `conv`-function, we could here also have used any efficient filter implementation.

9.3.1 Tensor product notation for operations on images

Filter-based operations on images can be written compactly using what we will call *tensor product notation*. This is part of a very general tensor product framework, and we will review parts of this framework for the sake of completeness. Let us first define the tensor product of vectors.

Definition 9.19 (Tensor product of vectors). If \mathbf{x}, \mathbf{y} are vectors of length M and N , respectively, their tensor product $\mathbf{x} \otimes \mathbf{y}$ is defined as the $M \times N$ -matrix defined by $(\mathbf{x} \otimes \mathbf{y})_{ij} = x_i y_j$. In other words, $\mathbf{x} \otimes \mathbf{y} = \mathbf{x}\mathbf{y}^T$.

The tensor product $\mathbf{x}\mathbf{y}^T$ is also called the *outer product* of \mathbf{x} and \mathbf{y} (contrary to the inner product $\langle \mathbf{x}, \mathbf{y} \rangle = \mathbf{x}^T \mathbf{y}$). In particular $\mathbf{x} \otimes \mathbf{y}$ is a matrix of rank 1, which means that most matrices cannot be written as a tensor product of two vectors. The special case $\mathbf{e}_i \otimes \mathbf{e}_j$ is the matrix which is 1 at (i, j) and 0 elsewhere, and the set of all such matrices forms a basis for the set of $M \times N$ -matrices.

Observation 9.20. Let $\mathcal{E}_M = \{\mathbf{e}_i\}_{i=0}^{M-1}$, $\mathcal{E}_N = \{\mathbf{e}_i\}_{i=0}^{N-1}$ be the standard bases for \mathbb{R}^M and \mathbb{R}^N . Then

$$\mathcal{E}_{M,N} = \{\mathbf{e}_i \otimes \mathbf{e}_j\}_{(i,j)=(0,0)}^{(M-1,N-1)}$$

is a basis for $L_{M,N}(\mathbb{R})$, the set of $M \times N$ -matrices. This basis is often referred to as the standard basis for $L_{M,N}(\mathbb{R})$.

The standard basis thus consists of rank 1-matrices. An image can simply be thought of as a matrix in $L_{M,N}(\mathbb{R})$, and a computational molecule is simply a special type of linear transformation from $L_{M,N}(\mathbb{R})$ to itself. Let us also define the tensor product of matrices.

Definition 9.21 (Tensor product of matrices). If $S_1 : \mathbb{R}^M \rightarrow \mathbb{R}^M$ and $S_2 : \mathbb{R}^N \rightarrow \mathbb{R}^N$ are matrices, we define the linear mapping $S_1 \otimes S_2 : L_{M,N}(\mathbb{R}) \rightarrow L_{M,N}(\mathbb{R})$ by linear extension of $(S_1 \otimes S_2)(\mathbf{e}_i \otimes \mathbf{e}_j) = (S_1 \mathbf{e}_i) \otimes (S_2 \mathbf{e}_j)$. The linear mapping $S_1 \otimes S_2$ is called the tensor product of the matrices S_1 and S_2 .

A couple of remarks are in order. First, from linear algebra we know that, when S is linear mapping from V and $S(\mathbf{v}_i)$ is known for a basis $\{\mathbf{v}_i\}_i$ of V , S is uniquely determined. In particular, since the $\{\mathbf{e}_i \otimes \mathbf{e}_j\}_{i,j}$ form a basis, there exists a unique linear transformation $S_1 \otimes S_2$ so that $(S_1 \otimes S_2)(\mathbf{e}_i \otimes \mathbf{e}_j) = (S_1 \mathbf{e}_i) \otimes (S_2 \mathbf{e}_j)$. This unique linear transformation is what we call the linear extension from the values in the given basis. Clearly, by linearity, also $(S_1 \otimes S_2)(\mathbf{x} \otimes \mathbf{y}) = (S_1 \mathbf{x}) \otimes (S_2 \mathbf{y})$, since

$$\begin{aligned} (S_1 \otimes S_2)(\mathbf{x} \otimes \mathbf{y}) &= (S_1 \otimes S_2)((\sum_i x_i \mathbf{e}_i) \otimes (\sum_j y_j \mathbf{e}_j)) = (S_1 \otimes S_2)(\sum_{i,j} x_i y_j (\mathbf{e}_i \otimes \mathbf{e}_j)) \\ &= \sum_{i,j} x_i y_j (S_1 \otimes S_2)(\mathbf{e}_i \otimes \mathbf{e}_j) = \sum_{i,j} x_i y_j (S_1 \mathbf{e}_i) \otimes (S_2 \mathbf{e}_j) \\ &= \sum_{i,j} x_i y_j S_1 \mathbf{e}_i ((S_2 \mathbf{e}_j))^T = S_1 (\sum_i x_i \mathbf{e}_i) (S_2 (\sum_j y_j \mathbf{e}_j))^T \\ &= S_1 \mathbf{x} (S_2 \mathbf{y})^T = (S_1 \mathbf{x}) \otimes (S_2 \mathbf{y}). \end{aligned}$$

Here we used the result from Exercise 5. We can now prove the following.

Theorem 9.22. If $S_1 : \mathbb{R}^M \rightarrow \mathbb{R}^M$ and $S_2 : \mathbb{R}^N \rightarrow \mathbb{R}^N$ are matrices of linear transformations, then $(S_1 \otimes S_2)X = S_1 X (S_2)^T$ for any $X \in L_{M,N}(\mathbb{R})$. In particular $S_1 \otimes S_2$ is the operation which applies S_1 to the columns of X , and S_2 to the resulting rows. In other words, if S_1, S_2 have compact filter notations \mathbf{t}_1 and \mathbf{t}_2 , respectively, then $S_1 \otimes S_2$ has computational molecule $\text{rev}(\mathbf{t}_1) \otimes \text{rev}(\mathbf{t}_2)$.

We have not formally defined the tensor product of compact filter notations. This is a straightforward extension of the usual tensor product of vectors, where we additionally mark the element at index $(0,0)$.

Proof: We have that

$$\begin{aligned}(S_1 \otimes S_2)(\mathbf{e}_i \otimes \mathbf{e}_j) &= (S_1 \mathbf{e}_i) \otimes (S_2 \mathbf{e}_j) \\ &= (\text{col}_i(S_1)) \otimes (\text{col}_j(S_2)) = \text{col}_i(S_1)(\text{col}_j(S_2))^T \\ &= \text{col}_i(S_1)\text{row}_j((S_2)^T) = S_1(\mathbf{e}_i \otimes \mathbf{e}_j)(S_2)^T.\end{aligned}$$

This means that $(S_1 \otimes S_2)X = S_1 X (S_2)^T$ for any $X \in L_{M,N}(\mathbb{R})$ also, since equality holds on the basis vectors $\mathbf{e}_i \otimes \mathbf{e}_j$. and since the matrix A with entries $a_{ij} = (\text{rev}(\mathbf{t}_1))_i (\text{rev}(\mathbf{t}_2))_j$ can also be written as $\text{rev}(\mathbf{t}_1) \otimes \text{rev}(\mathbf{t}_2)$, the result follows. ■

We have thus shown that we alternatively can write $S_1 \otimes S_2$ for the operations we have considered. This notation also makes it easy to combine several two-dimensional filtering operations:

Corollary 9.23. We have that $(S_1 \otimes T_1)(S_2 \otimes T_2) = (S_1 S_2) \otimes (T_1 T_2)$.

Proof: By Theorem 9.22 we have that

$$(S_1 \otimes T_1)(S_2 \otimes T_2)X = S_1(S_2 X T_2^T)T_1^T = (S_1 S_2)X(T_1 T_2)^T = ((S_1 S_2) \otimes (T_1 T_2))X.$$

for any $X \in L_{M,N}(\mathbb{R})$. This proves the result. ■

Suppose that we want to apply the operation $S_1 \otimes S_2$ to an image. We can factorize $S_1 \otimes S_2$ as

$$S_1 \otimes S_2 = (S_1 \otimes I)(I \otimes S_2) = (I \otimes S_2)(S_1 \otimes I). \quad (9.5)$$

Moreover, since

$$(S_1 \otimes I)X = S_1 X \quad (I \otimes S_2)X = X(S_2)^T = (S_2 X^T)^T,$$

$S_1 \otimes I$ is a vertical filtering operation, and $I \otimes S_2$ is a horizontal filtering operation in this factorization. For filters we have an even stronger result: If S_1, S_2, S_3, S_4 all are filters, we have from Corollary 9.23 that $(S_1 \otimes S_2)(S_3 \otimes S_4) = (S_3 \otimes S_4)(S_1 \otimes S_2)$, since all filters commute. This does not hold in general since general matrices do not commute.

We will now consider two important examples of filtering operations on images: smoothing and edge detection/computing partial derivatives. For all examples we will use the tensor product notation for these operations.

Example 9.24 (Smoothing an image). When we considered filtering of digital sound, we observed that replacing each sample of a sound by an average of the sample and its neighbours damped the high frequencies of the sound. Let us consider the computational molecules where such a filter is applied to both the rows and the columns. For the one-dimensional case on sound, we argued that filter coefficients taken from Pascal's triangle give good smoothing effects. The same can be argued for

images. If we use the filter $S = \frac{1}{4}\{1, 2, 1\}$ (row 2 from Pascal's triangle), Theorem 9.17 says that we obtain the computational molecule

$$A = \frac{1}{16} \begin{pmatrix} 1 & 2 & 1 \\ 2 & 4 & 2 \\ 1 & 2 & 1 \end{pmatrix}. \quad (9.6)$$

This means that we compute the new pixels by

$$\hat{p}_{i,j} = \frac{1}{16} (4p_{i,j} + 2(p_{i,j-1} + p_{i-1,j} + p_{i+1,j} + p_{i,j+1}) + p_{i-1,j-1} + p_{i+1,j-1} + p_{i-1,j+1} + p_{i+1,j+1}).$$

If we instead use the filter $S = \frac{1}{64}\{1, 6, 15, 20, 15, 6, 1\}$ (row 6 from Pascal's triangle), we get the computational molecule

$$\frac{1}{4096} \begin{pmatrix} 1 & 6 & 15 & 20 & 15 & 6 & 1 \\ 6 & 36 & 90 & 120 & 90 & 36 & 6 \\ 15 & 90 & 225 & 300 & 225 & 90 & 15 \\ 20 & 120 & 300 & 400 & 300 & 120 & 20 \\ 15 & 90 & 225 & 300 & 225 & 90 & 15 \\ 6 & 36 & 90 & 120 & 90 & 36 & 6 \\ 1 & 6 & 15 & 20 & 15 & 6 & 1 \end{pmatrix}. \quad (9.7)$$

For both molecules the weights sum to one, so that the new intensity values $\hat{p}_{i,j}$ are weighted averages of the intensity values on the right. We anticipate that both molecules give a smoothing effect, but that the second molecules provides more smoothing. The result of applying the two molecules in (9.6) and (9.7) to our greyscale-image is shown in Figure 9.8(b) and - (c) respectively. With the help of the function `filterimage`, smoothing with the first molecule (9.6) above can be obtained by writing

```
filterimage(X, (1/4)*[ 1 2 1 ], (1/4)*[ 1 2 1 ]);
```

To make the smoothing effect visible, we have zoomed in on the face in the image. The smoothing effect is clearly best visible in the second image. Smoothing effects are perhaps more visible if we use a simple image, as the one in Figure 9.24(a). Again we have used the filter $S = \frac{1}{4}\{1, 2, 1\}$. Here we also have shown in (b) and (c) what happens if we only smooth the image in one of the directions. The smoothing effects are then only seen in one of the vertical or horizontal directions. In (d) we have smoothed in both directions. We clearly see the union of the two one-dimensional smoothing operations then. Let us summarize as follows.

Observation 9.25. An image P can be smoothed by applying a smoothing filter to the rows, and then to the columns.



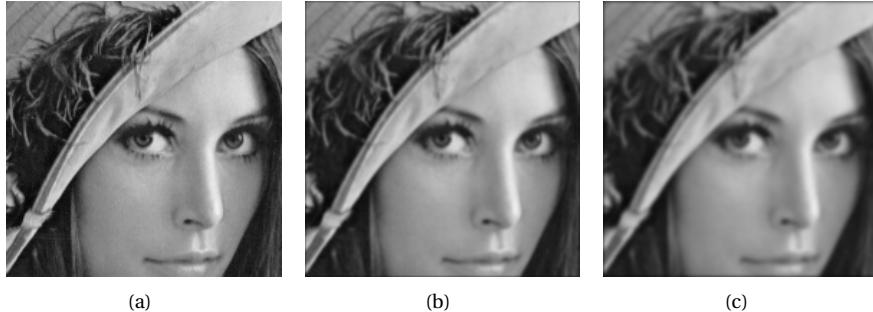


Figure 9.8: The images in (b) and (c) show the effect of smoothing the image in (a).

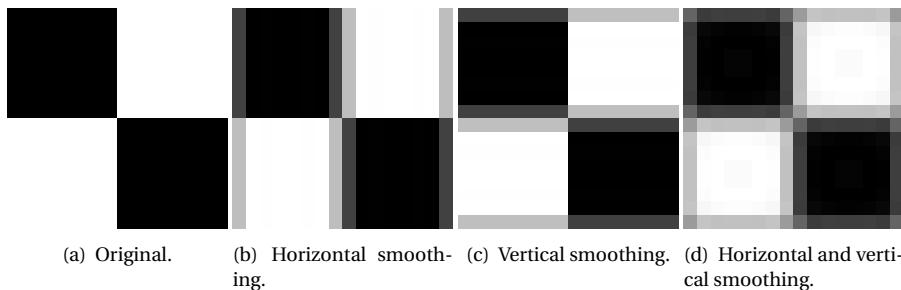


Figure 9.9: The results of smoothing a simple image with the filter $\frac{1}{4}\{1, 2, 1\}$ horizontally, vertically, and both.

Example 9.26 (Computing partial derivatives and detecting edges). Another operation on images which can be expressed in terms of computational molecules is edge detection. An edge in an image is characterised by a large change in intensity values over a small distance in the image. For a continuous function this corresponds to a large derivative. An image is only defined at isolated points, so we cannot compute derivatives, but we have a perfect situation for applying numerical differentiation. Since a grey-level image is a scalar function of two variables, numerical differentiation techniques can be applied.

Partial derivative in x -direction. Let us first consider computation of the partial derivative $\partial P / \partial x$ at all points in the image. Note first that it is the second coordinate in an image which refers to the x -direction used when plotting functions. This means that the familiar symmetric Newton quotient approximation for the partial derivative [23] takes the form

$$\frac{\partial P}{\partial x}(i, j) \approx \frac{p_{i,j+1} - p_{i,j-1}}{2}, \quad (9.8)$$

where we have used the convention $h = 1$ which means that the derivative is measured in terms of 'intensity per pixel'. This corresponds to applying the bass-reducing

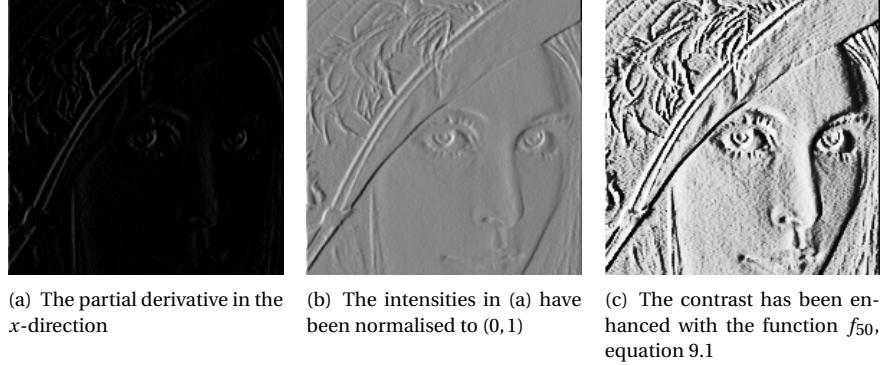


Figure 9.10: Experimenting with the partial derivative for the image in 9.4.

filter $S = \frac{1}{2}\{1, 0, -1\}$ to all the rows (alternatively, applying the tensor product $I \otimes S$ to the image). We can thus express this in terms of computational molecules as follows.

Observation 9.27. Let $P = (p_{i,j})_{i,j=1}^{m,n}$ be a given image. The partial derivative $\partial P / \partial x$ of the image can be computed with the computational molecule

$$\frac{1}{2} \begin{pmatrix} 0 & 0 & 0 \\ -1 & 0 & 1 \\ 0 & 0 & 0 \end{pmatrix}. \quad (9.9)$$

We have included the two rows of 0s just to make it clear how the computational molecule is to be interpreted when we place it over the pixels. If we apply the smooth-function to the same excerpt of the Lena image with this molecule, we obtain figure 9.10(a). This image is not very helpful since it is almost completely black. The reason is that many of the intensities are in fact negative, and these are just displayed as black. More specifically, the intensities turn out to vary in the interval $[-0.424, 0.418]$. We therefore normalise and map all intensities to $[0, 1]$. The result of this is shown in (b). The predominant colour of this image is an average grey, i.e., an intensity of about 0.5. To get more detail in the image we therefore try to increase the contrast by applying the function f_{50} in equation 9.1 to each intensity value. The result is shown in figure 9.10(c) which does indeed show more detail.

It is important to understand the colours in these images. We have computed the derivative in the x -direction, and we recall that the computed values varied in the interval $[-0.424, 0.418]$. The negative value corresponds to the largest average decrease in intensity from a pixel $p_{i-1,j}$ to a pixel $p_{i+1,j}$. The positive value on the other hand corresponds to the largest average increase in intensity. A value of 0 in figure 9.10(a) corresponds to no change in intensity between the two pixels.

When the values are mapped to the interval $[0, 1]$ in figure 9.10(b), the small values are mapped to something close to 0 (almost black), the maximal values are

mapped to something close to 1 (almost white), and the values near 0 are mapped to something close to 0.5 (grey). In figure 9.10(c) these values have just been emphasised even more.

Figure 9.10(c) tells us that in large parts of the image there is very little variation in the intensity. However, there are some small areas where the intensity changes quite abruptly, and if you look carefully you will notice that in these areas there is typically both black and white pixels close together, like down the vertical front corner of the bus. This will happen when there is a stripe of bright or dark pixels that cut through an area of otherwise quite uniform intensity.

Partial derivative in y -direction. The partial derivative $\partial P / \partial y$ can be computed analogously to $\partial P / \partial x$, i.e. we apply the filter $-S = \frac{1}{2}\{-1, 0, 1\}$ to all columns of the image (alternatively, apply the tensor product $-S \otimes I$ to the image), where S is the filter which we used for edge detection in the x -direction. Note that the positive direction of this axis in an image is opposite to the direction of the y -axis we use when plotting functions.

Observation 9.28. Let $P = (p_{i,j})_{i,j=1}^{m,n}$ be a given image. The partial derivative $\partial P / \partial y$ of the image can be computed with the computational molecule

$$\frac{1}{2} \begin{pmatrix} 0 & -1 & 0 \\ 0 & 0 & 0 \\ 0 & 1 & 0 \end{pmatrix}. \quad (9.10)$$

The result is shown in figure 9.12(b). The intensities have been normalised and the contrast enhanced by the function f_{50} from Equation (9.1).

The gradient. The gradient of a scalar function is often used as a measure of the size of the first derivative. The gradient is defined by the vector

$$\nabla P = \left(\frac{\partial P}{\partial x}, \frac{\partial P}{\partial y} \right),$$

so its length is given by

$$|\nabla P| = \sqrt{\left(\frac{\partial P}{\partial x} \right)^2 + \left(\frac{\partial P}{\partial y} \right)^2}.$$

When the two first derivatives have been computed it is a simple matter to compute the gradient vector and its length; the resulting is shown as an image in figure 9.11c.

The image of the gradient looks quite different from the images of the two partial derivatives. The reason is that the numbers that represent the length of the gradient are (square roots of) sums of squares of numbers. This means that the parts of the image that have virtually constant intensity (partial derivatives close to 0) are coloured black. In the images of the partial derivatives these values ended up in the

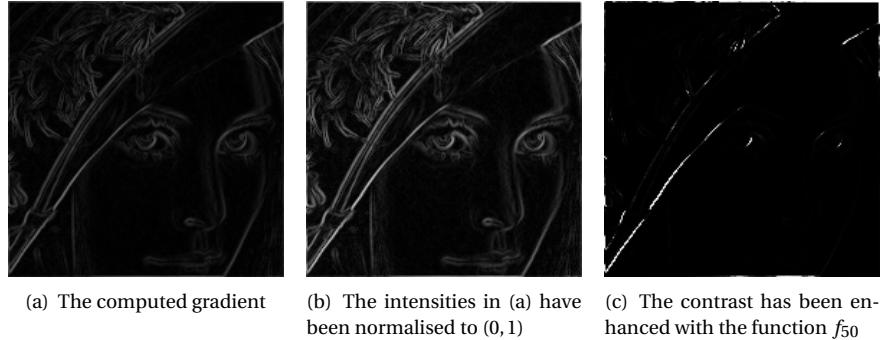


Figure 9.11: Computing the gradient.

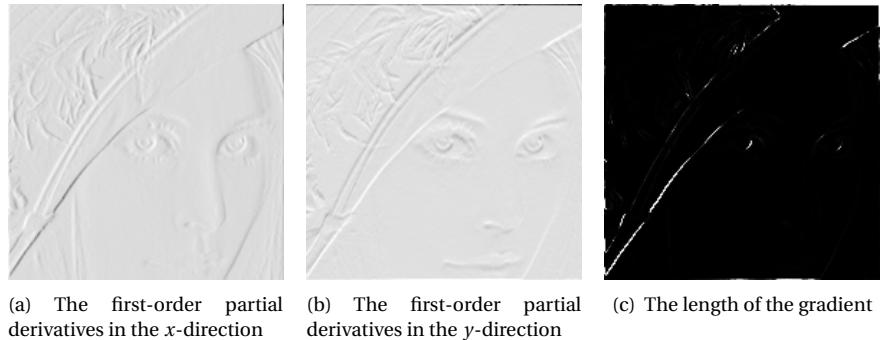


Figure 9.12: In all images, the computed numbers have been normalised and the contrast enhanced.

middle of the range of intensity values, with a final colour of grey, since there were both positive and negative values.

Figure 9.11(a) shows the computed values of the gradient. Although it is possible that the length of the gradient could become larger than 1, the maximum value in this case is about 0.876. By normalising the intensities we therefore increase the contrast slightly and obtain the image in figure 9.11(b).

To enhance the contrast further we have to do something different from what was done in the other images since we now have a large number of intensities near 0. The solution is to apply a function like the ones shown in figure 9.6(b) to the intensities. If we use the function $g_{0.01}$ defined in equation(9.2) we obtain the image in figure 9.11(c).



9.3.2 Comparing the first derivatives

Figure 9.12 shows the two first-order partial derivatives and the gradient. If we compare the two partial derivatives we see that the x -derivative seems to emphasise vertical edges while the y -derivative seems to emphasise horizontal edges. This is precisely what we must expect. The x -derivative is large when the difference between neighbouring pixels in the x -direction is large, which is the case across a vertical edge. The y -derivative enhances horizontal edges for a similar reason.

The gradient contains information about both derivatives and therefore emphasises edges in all directions. It also gives a simpler image since the sign of the derivatives has been removed.

9.3.3 Second-order derivatives

To compute the three second order derivatives we can combine the two computational molecules which we already have described. For the mixed second order derivative we get $(I \otimes S)((-S) \otimes I) = -S \otimes S$. For the last two second order derivative $\frac{\partial^2 P}{\partial x^2}, \frac{\partial^2 P}{\partial y^2}$, we can also use the three point approximation to the second derivative [23]

$$\frac{\partial P}{\partial x^2}(i, j) \approx p_{i,j+1} - 2p_{i,j} + p_{i,j-1} \quad (9.11)$$

to the second derivative (again we have set $h = 1$). This gives a smaller molecule than if we combine the two molecules for order one differentiation (i.e. $(I \otimes S)(I \otimes S) = (I \otimes S^2)$ and $((-S) \otimes I)((-S) \otimes I) = (S^2 \otimes I)$), since $S^2 = \frac{1}{2}\{1, 0, -1\}\frac{1}{2}\{1, 0, -1\} = \frac{1}{4}\{1, 0, -2, 0, 1\}$.

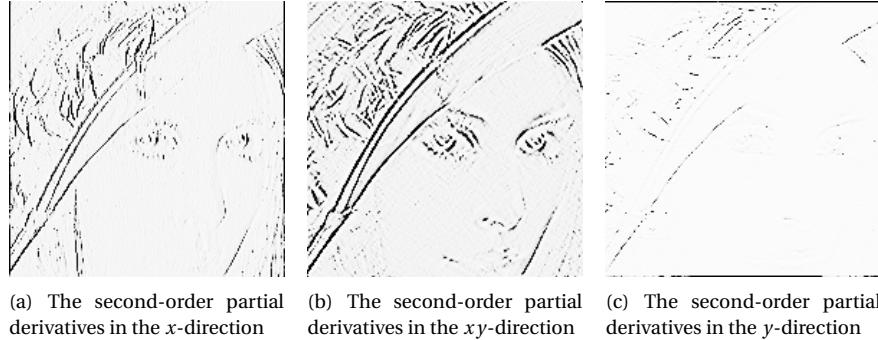
Observation 9.29 (Second order derivatives of an image). The second order derivatives of an image P can be computed by applying the computational molecules

$$\frac{\partial^2 P}{\partial x^2} : \quad \begin{pmatrix} 0 & 0 & 0 \\ 1 & -2 & 1 \\ 0 & 0 & 0 \end{pmatrix}, \quad (9.12)$$

$$\frac{\partial^2 P}{\partial y \partial x} : \quad \frac{1}{4} \begin{pmatrix} -1 & 0 & 1 \\ 0 & 0 & 0 \\ 1 & 0 & -1 \end{pmatrix}, \quad (9.13)$$

$$\frac{\partial^2 P}{\partial y^2} : \quad \begin{pmatrix} 0 & 1 & 0 \\ 0 & -2 & 0 \\ 0 & 1 & 0 \end{pmatrix}. \quad (9.14)$$

With the information in observation 9.29 it is quite easy to compute the second-order derivatives, and the results are shown in figure 9.13. The computed derivatives were first normalised and then the contrast enhanced with the function f_{100} in each image, see equation 9.1.



(a) The second-order partial derivatives in the x -direction

(b) The second-order partial derivatives in the xy -direction

(c) The second-order partial derivatives in the y -direction

Figure 9.13: In all images, the computed numbers have been normalised and the contrast enhanced.

As for the first derivatives, the xx -derivative seems to emphasise vertical edges and the yy -derivative horizontal edges. However, we also see that the second derivatives are more sensitive to noise in the image (the areas of grey are less uniform). The mixed derivative behaves a bit differently from the other two, and not surprisingly it seems to pick up both horizontal and vertical edges.

This procedure can be generalized to higher order derivatives also. To apply $\frac{\partial^{k+l} P}{\partial x^k \partial y^l}$ to an image we can compute $S_l \otimes S_k$ where S_r corresponds to any point method for computing the r 'th order derivative. We can also compute $(S^l) \otimes (S^k)$, where we iterate the filter $S = \frac{1}{2}\{1, 0, -1\}$ for the first derivative, but this gives longer filters.

Let us also apply the molecules for differentiation to a chess pattern test image. In Figure 9.14 we have applied $S \otimes I$, $I \otimes S$, and $S \otimes S$ to the image shown in (a). These images make it clear that $S \otimes I$ detects all horizontal edges, that $I \otimes S$ detects all vertical edges, and that $S \otimes S$ detects all points where abrupt changes appear in both directions. We also see that the second order partial derivative detects exactly the same edges which the first order partial derivative found. Note that the edges detected with $I \otimes S^2$ are wider than the ones detected with $I \otimes S$. The reason is that the filter S^2 has more filter coefficients than S . Also, edges are detected with different colours. This reflects whether the difference between the neighbouring pixels is positive or negative. The values after we have applied the tensor product may thus not lie in the legal range of pixel values (since they may be negative). The figures have taken this into account by mapping the values back to a legal range of values, as we did in Chapter 9. Finally, we also see additional edges at the first and last rows/edges in the images. The reason is that the filter S is defined by assuming that the pixels repeat periodically (i.e. it is a circulant Toeplitz matrix). Due to this, we have additional edges at the first/last rows/edges. This effect can also be seen in Chapter 9, although there we did not assume that the pixels repeat periodically.

Defining a two-dimensional filter by filtering columns and then rows is not the only way we can define a two-dimensional filter. Another possible way is to let the $MN \times MN$ -matrix itself be a filter. Unfortunately, this is a bad way to define filter-

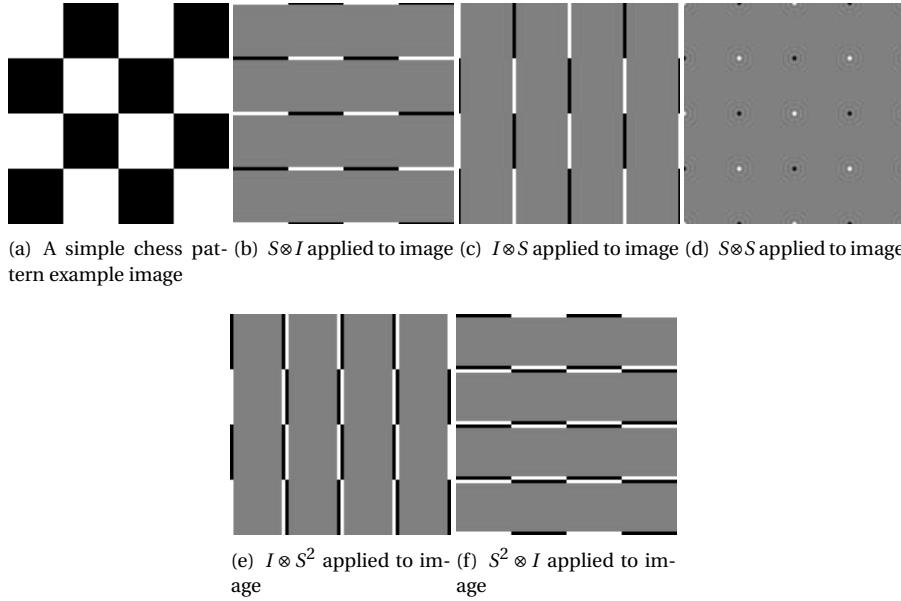


Figure 9.14: Different combinations of tensor products applied to a simple image

ing of an image, since there are some undesirable effects near the boundaries between rows: in the vector we form, the last element of one row is followed by the first element of the next row. These boundary effects are unfortunate when a filter is applied.

What you should have learnt in this section

The operation $X \rightarrow S_1 X (S_2)^T$ can be used to define operations on images, based on one-dimensional operations S_1 and S_2 . This amounts to applying S_1 to all columns in the image, and then S_2 to all rows in the result. You should know how this operation can be conveniently expressed with tensor product notation, and that in the typical case when S_1 and S_2 are filters, this can equivalently be expressed in terms of computational molecules. The case when the S_i are smoothing filters gives rise to smoothing operations on images. A simple highpass filter, corresponding to taking the derivative, gives rise to edge-detection operations on images.

Exercises for Section 9.3

1. Generate the image in Figure 9.8(b) and -(c) by writing code which calls the function `filterimage` with appropriate filters.
2. Generate the image in Figure 9.10(c) by writing code in the same way. Also generate the images in figures 9.11, 9.12, and 9.13.
3. Let the filter S be defined by $S = \{-1, 1\}$.

- a.** Let X be a matrix which represents the pixel values in an image. What can you say about how the new images $(S \otimes I)X$ og $(I \otimes S)X$ look? What are the interpretations of these operations?
- b.** Write down the 4×4 -matrix $X = (1, 1, 1, 1) \otimes (0, 0, 1, 1)$. Compute $(S \otimes I)X$ by applying the filters to the corresponding rows/columns of X as we have learnt, and interpret the result. Do the same for $(I \otimes S)X$.
- 4.** Let S be the moving average filter of length $2L + 1$, i.e. $T = \frac{1}{L} \underbrace{\{1, \dots, 1, \underline{1}, 1, \dots, 1\}}_{2L+1 \text{ times}}$. What is the computational molecule of $S \otimes S$?
- 5.** Show that the mapping $F(\mathbf{x}, \mathbf{y}) = \mathbf{x} \otimes \mathbf{y}$ is bi-linear, i.e. that $F(\alpha\mathbf{x}_1 + \beta\mathbf{x}_2, \mathbf{y}) = \alpha F(\mathbf{x}_1, \mathbf{y}) + \beta F(\mathbf{x}_2, \mathbf{y})$, and $F(\mathbf{x}, \alpha\mathbf{y}_1 + \beta\mathbf{y}_2) = \alpha F(\mathbf{x}, \mathbf{y}_1) + \beta F(\mathbf{x}, \mathbf{y}_2)$.
- 6.** Attempt to find matrices $S_1 : \mathbb{R}^M \rightarrow \mathbb{R}^M$ and $S_2 : \mathbb{R}^N \rightarrow \mathbb{R}^N$ so that the following mappings from $L_{M,N}(\mathbb{R})$ to $L_{M,N}(\mathbb{R})$ can be written on the form $X \rightarrow S_1 X (S_2)^T = (S_1 \otimes S_2)X$. In all the cases, it may be that no such S_1, S_2 can be found. If this is the case, prove it.
- a.** The mapping which reverses the order of the rows in a matrix.
- b.** The mapping which reverses the order of the columns in a matrix.
- c.** The mapping which transposes a matrix.
- 7.** Let the filter S be defined by $S = \{1, \underline{2}, 1\}$.
- a.** Write down the computational molecule of $S \otimes S$.
- b.** Let us define $\mathbf{x} = (1, 2, 3)$, $\mathbf{y} = (3, 2, 1)$, $\mathbf{z} = (2, 2, 2)$, and $\mathbf{w} = (1, 4, 2)$. Compute the matrix $A = \mathbf{x} \otimes \mathbf{y} + \mathbf{z} \otimes \mathbf{w}$.
- c.** Compute $(S \otimes S)A$ by applying the filter S to every row and column in the matrix the way we have learnt. If the matrix A was more generally an image, what can you say about how the new image will look?
- 8.** Let $S = \frac{1}{4}\{1, \underline{2}, 1\}$ be a filter.
- a.** What is the effect of applying the tensor products $S \otimes I$, $I \otimes S$, and $S \otimes S$ on an image represented by the matrix X ?
- b.** Compute $(S \otimes S)(\mathbf{x} \otimes \mathbf{y})$, where $\mathbf{x} = (4, 8, 8, 4)$, $\mathbf{y} = (8, 4, 8, 4)$ (i.e. both \mathbf{x} and \mathbf{y} are column vectors).
- 9.** Suppose that we have an image given by the $N \times N$ -matrix X , and consider the following Matlab code:

```

for k=1:2
    for s=1:N
        X(1,s)=0.25*X(N,s)+0.5*X(1,s)+0.25*X(2,s);
        X(2:(N-1),s)=0.25*X(1:(N-2),s)+0.5*X(2:(N-1),s)+0.25*X(3:N,s);
        X(N,s)=0.25*X(N-1,s)+0.5*X(N,s)+0.25*X(1,s);
    end
    X=X';
end

```

Which tensor product is applied to the image? Comment what the code does, in particular the first and third line in the inner **for**-loop. What effect does the code have on the image?

10. Let \mathbf{v}_A be an eigenvector of A with eigenvalue λ_A , and \mathbf{v}_B an eigenvector of B with eigenvalue λ_B . Show that $\mathbf{v}_A \otimes \mathbf{v}_B$ is an eigenvector of $A \otimes B$ with eigenvalue $\lambda_A \lambda_B$. Explain from this why $\|A \otimes B\| = \|A\| \|B\|$, where $\|\cdot\|$ denotes the operator norm of a matrix.

11. The *Kronecker tensor product* of two matrices A and B , written $A \otimes^k B$, is defined as

$$A \otimes^k B = \begin{pmatrix} a_{11}B & a_{12}B & \cdots & a_{1M}B \\ a_{21}B & a_{22}B & \cdots & a_{2M}B \\ \vdots & \vdots & \ddots & \vdots \\ a_{p1}B & a_{p2}B & \cdots & a_{pM}B \end{pmatrix},$$

where the entries of A are a_{ij} . The tensor product of a $p \times M$ -matrix, and a $q \times N$ -matrix is thus a $(pq) \times (MN)$ -matrix. Note that this tensor product in particular gives meaning for vectors: if $\mathbf{x} \in \mathbb{R}^M$, $\mathbf{y} \in \mathbb{R}^N$ are column vectors, then $\mathbf{x} \otimes^k \mathbf{y} \in \mathbb{R}^{MN}$ is also a column vector. In this exercise we will investigate how the Kronecker tensor product is related to tensor products as we have defined them in this section.

a. Explain that, if $\mathbf{x} \in \mathbb{R}^M$, $\mathbf{y} \in \mathbb{R}^N$ are column vectors, then $\mathbf{x} \otimes^k \mathbf{y}$ is the column vector where the rows of $\mathbf{x} \otimes \mathbf{y}$ have first been stacked into one large row vector, and this vector transposed. The linear extension of the operation defined by

$$\mathbf{x} \otimes \mathbf{y} \in \mathbb{R}^{M,N} \rightarrow \mathbf{x} \otimes^k \mathbf{y} \in \mathbb{R}^{MN}$$

thus stacks the rows of the input matrix into one large row vector, and transposes the result.

b. Show that $(A \otimes^k B)(\mathbf{x} \otimes^k \mathbf{y}) = (A\mathbf{x}) \otimes^k (B\mathbf{y})$. We can thus use any of the defined tensor products \otimes , \otimes_k to produce the same result, i.e. we have the following commutative diagram,

$$\begin{array}{ccc} \mathbf{x} \otimes \mathbf{y} & \xrightarrow{A \otimes B} & (A\mathbf{x}) \otimes (B\mathbf{y}) \\ \downarrow & & \downarrow \\ \mathbf{x} \otimes^k \mathbf{y} & \xrightarrow{A \otimes^k B} & (A\mathbf{x}) \otimes^k (B\mathbf{y}), \end{array}$$

where the vertical arrows represent stacking the rows in the matrix, and transposing, and the horizontal arrows represent the two tensor product linear transformations we have defined. In particular, we can compute the tensor product in terms of vectors, or in terms of matrices, and it is clear that the Kronecker tensor product gives the matrix of tensor product operations.

c. Using the Euclidean inner product on $L(M, N) = \mathbb{R}^{MN}$, i.e.

$$\langle X, Y \rangle = \sum_{i=0}^{M-1} \sum_{j=0}^{N-1} X_{i,j} \overline{Y_{i,j}}.$$

and the correspondence in a. we can define the inner product of $\mathbf{x}_1 \otimes \mathbf{y}_1$ and $\mathbf{x}_2 \otimes \mathbf{y}_2$ by

$$\langle \mathbf{x}_1 \otimes \mathbf{y}_1, \mathbf{x}_2 \otimes \mathbf{y}_2 \rangle = \langle \mathbf{x}_1 \otimes^k \mathbf{y}_1, \mathbf{x}_2 \otimes^k \mathbf{y}_2 \rangle.$$

Show that

$$\langle \mathbf{x}_1 \otimes \mathbf{y}_1, \mathbf{x}_2 \otimes \mathbf{y}_2 \rangle = \langle \mathbf{x}_1, \mathbf{x}_2 \rangle \langle \mathbf{y}_1, \mathbf{y}_2 \rangle.$$

Clearly this extends linearly to an inner product on $L_{M,N}$.

d. Show that the FFT factorization can be written as

$$\begin{pmatrix} F_{N/2} & D_{N/2}F_{N/2} \\ F_{N/2} & -D_{N/2}F_{N/2} \end{pmatrix} = \begin{pmatrix} I_{N/2} & D_{N/2} \\ I_{N/2} & -D_{N/2} \end{pmatrix} (I_2 \otimes_k F_{N/2}).$$

Also rewrite the sparse matrix factorization for the FFT from Equation (2.19) in terms of tensor products.

9.4 Change of coordinates for images

Filter-based operations were not the only operations we considered for sound. We also considered the DFT, the DCT, and the wavelet transform, which were changes of coordinates which gave us useful frequency- or time-frequency information. We would like to define similar changes of coordinates for images, which also give useful such information. Tensor product notation will also be useful in this respect, and we start with the following result.

Theorem 9.30. If $\mathcal{B}_1 = \{\mathbf{v}_i\}_{i=0}^{M-1}$ is a basis for \mathbb{R}^M , and $\mathcal{B}_2 = \{\mathbf{w}_j\}_{j=0}^{N-1}$ is a basis for \mathbb{R}^N , then $\{\mathbf{v}_i \otimes \mathbf{w}_j\}_{(i,j)=(0,0)}^{(M-1,N-1)}$ is a basis for $L_{M,N}(\mathbb{R})$. We denote this basis by $\mathcal{B}_1 \otimes \mathcal{B}_2$.

Proof: Suppose that $\sum_{(i,j)=(0,0)}^{(M-1,N-1)} \alpha_{i,j} (\mathbf{v}_i \otimes \mathbf{w}_j) = 0$. Setting $\mathbf{h}_i = \sum_{j=0}^{N-1} \alpha_{i,j} \mathbf{w}_j$ we get

$$\sum_{j=0}^{N-1} \alpha_{i,j} (\mathbf{v}_i \otimes \mathbf{w}_j) = \mathbf{v}_i \otimes \left(\sum_{j=0}^{N-1} \alpha_{i,j} \mathbf{w}_j \right) = \mathbf{v}_i \otimes \mathbf{h}_i.$$

where we have used the bi-linearity of the tensor product mapping $(\mathbf{x}, \mathbf{y}) \rightarrow \mathbf{x} \otimes \mathbf{y}$ (Exercise 9.3.5). This means that

$$\mathbf{0} = \sum_{(i,j)=(0,0)}^{(M-1,N-1)} \alpha_{i,j} (\mathbf{v}_i \otimes \mathbf{w}_j) = \sum_{i=0}^{M-1} \mathbf{v}_i \otimes \mathbf{h}_i = \sum_{i=0}^{M-1} \mathbf{v}_i \mathbf{h}_i^T.$$

Column k in this matrix equation says $\mathbf{0} = \sum_{i=0}^{M-1} h_{i,k} \mathbf{v}_i$, where $h_{i,k}$ are the components in \mathbf{h}_i . By linear independence of the \mathbf{v}_i we must have that $h_{0,k} = h_{1,k} = \dots = h_{M-1,k} = 0$. Since this applies for all k , we must have that all $\mathbf{h}_i = 0$. This means that $\sum_{j=0}^{N-1} \alpha_{i,j} \mathbf{w}_j = \mathbf{0}$ for all i , from which it follows by linear independence of the \mathbf{w}_j that $\alpha_{i,j} = 0$ for all j , and for all i . This means that $\mathcal{B}_1 \otimes \mathcal{B}_2$ is a basis. ■

In particular, as we have already seen, the standard basis for $L_{M,N}(\mathbb{R})$ can be written $\mathcal{E}_{M,N} = \mathcal{E}_M \otimes \mathcal{E}_N$. This is the basis for a useful convention: For a tensor product the bases are most naturally indexed in two dimensions, rather than the usual sequential indexing. This difference translates also to the meaning of coordinate vectors, which now are more naturally thought of as coordinate matrices:

Definition 9.31 (Coordinate matrix). Let $\{\mathbf{v}_i\}_{i=0}^{M-1}, \{\mathbf{w}_j\}_{j=0}^{N-1}$ be bases for \mathbb{R}^M and \mathbb{R}^N . By the coordinate matrix of $\sum_{k,l} \alpha_{k,l} (\mathbf{v}_k \otimes \mathbf{w}_l)$ we will mean the $M \times N$ -matrix X with entries $X_{kl} = \alpha_{k,l}$.

We will have use for the following theorem, which shows how change of coordinates in \mathbb{R}^M and \mathbb{R}^N translate to a change of coordinates in the tensor product:

Theorem 9.32 (Change of coordinates in tensor products). Assume that

1. $\mathcal{B}_1, \mathcal{C}_1$ are bases for \mathbb{R}^M , and that S_1 is the change of coordinates matrix from \mathcal{B}_1 to \mathcal{C}_1 ,
2. $\mathcal{B}_2, \mathcal{C}_2$ are bases for \mathbb{R}^N , and that S_2 is the change of coordinates matrix from \mathcal{B}_2 to \mathcal{C}_2 .

Both $\mathcal{B}_1 \otimes \mathcal{B}_2$ and $\mathcal{C}_1 \otimes \mathcal{C}_2$ are bases for $L_{M,N}(\mathbb{R})$, and if X is the coordinate matrix in $\mathcal{B}_1 \otimes \mathcal{B}_2$, and Y the coordinate matrix in $\mathcal{C}_1 \otimes \mathcal{C}_2$, then the change of coordinates from $\mathcal{B}_1 \otimes \mathcal{B}_2$ to $\mathcal{C}_1 \otimes \mathcal{C}_2$ can be computed as

$$Y = S_1 X (S_2)^T. \quad (9.15)$$

Proof: Let \mathbf{c}_{ki} be the i 'th basis vector in \mathcal{C}_k , \mathbf{b}_{ki} the i 'th basis vector in \mathcal{B}_k , $k = 1, 2$. Since any change of coordinates is linear, it is enough to show that it coincides with $X \rightarrow S_1 X (S_2)^T$ on the basis $\mathcal{C}_1 \otimes \mathcal{C}_2$. The basis vector $\mathbf{c}_{1i} \otimes \mathbf{c}_{2j}$ has coordinate vector $X = \mathbf{e}_i \otimes \mathbf{e}_j$ in $\mathcal{C}_1 \otimes \mathcal{C}_2$. With the mapping $X \rightarrow S_1 X (S_2)^T$ this is sent to

$$S_1 X (S_2)^T = S_1 (\mathbf{e}_i \otimes \mathbf{e}_j) (S_2)^T = \text{col}_i(S_1) \text{row}_j((S_2)^T).$$

On the other hand, since column i in S_1 is the coordinates of \mathbf{c}_{1i} in the basis \mathcal{B}_1 , and column j in S_2 is the coordinates of \mathbf{c}_{2j} in the basis \mathcal{B}_2 , we can write

$$\begin{aligned}\mathbf{c}_{1i} \otimes \mathbf{c}_{2j} &= \left(\sum_k (S_1)_{k,i} \mathbf{b}_{1k} \right) \otimes \left(\sum_l (S_2)_{l,j} \mathbf{b}_{2l} \right) = \sum_{k,l} (S_1)_{k,i} (S_2)_{l,j} (\mathbf{b}_{1k} \otimes \mathbf{b}_{2l}) \\ &= \sum_{k,l} (S_1)_{k,i} ((S_2)^T)_{j,l} (\mathbf{b}_{1k} \otimes \mathbf{b}_{2l}) = \sum_{k,l} (\text{col}_i(S_1) \text{row}_j((S_2)^T)_{k,l}) (\mathbf{b}_{1k} \otimes \mathbf{b}_{2l})\end{aligned}$$

we see that the coordinate vector of $\mathbf{c}_{1i} \otimes \mathbf{c}_{2j}$ in the basis $\mathcal{B}_1 \otimes \mathcal{B}_2$ is $\text{col}_i(S_1) \text{row}_j((S_2)^T)$. In other words, change of coordinates coincides with $X \rightarrow S_1 X (S_2)^T$, and the proof is done. \blacksquare

In both cases of filtering and change of coordinates in tensor products, we see that we need to compute the mapping $X \rightarrow S_1 X (S_2)^T$. As we have seen, this amounts to a row/column-wise operation, which we restate as follows:

Observation 9.33. The change of coordinates from $\mathcal{B}_1 \otimes \mathcal{B}_2$ to $\mathcal{C}_1 \otimes \mathcal{C}_2$ can be implemented as follows:

1. For every column in the coordinate matrix in $\mathcal{B}_1 \otimes \mathcal{B}_2$, perform a change of coordinates from \mathcal{B}_1 to \mathcal{C}_1 .
2. For every row in the resulting matrix, perform a change of coordinates from \mathcal{B}_2 to \mathcal{C}_2 .

If we let S_1 and S_2 be functions which implement the changes of coordinates in question, we can modify the algorithm for `filterimage` as follows, in order to apply change of coordinates to a tensor product.

```
[M,N]=size(X);
for col=1:N
    X(:,col)=S1(X(:,col));
end
for row=1:M
    X(row,:)=S2((X(row,:))');
end
```

The operation $X \rightarrow (S_1)X(S_2)^T$, which we now have encountered in two different ways, is one particular type of linear transformation from \mathbb{R}^{N^2} to itself (see Exercise 9.3.11 for how the matrix of this linear transformation can be constructed). While a general such linear transformation requires N^4 multiplications (i.e. when we perform a full matrix multiplication), $X \rightarrow (S_1)X(S_2)^T$ can be implemented generally with only $2N^3$ multiplications (since multiplication of two $N \times N$ -matrices require N^3 multiplications in general). The operation $X \rightarrow (S_1)X(S_2)^T$ is thus computationally simpler than linear transformations in general. In practice the operations S_1 and S_2 are also computationally simpler, since they can be filters, FFT's, or wavelet transformations, so that the complexity in $X \rightarrow (S_1)X(S_2)^T$ can be even lower.

In the following examples, we will interpret the pixel values in an image as coordinates in the standard basis, and perform a change of coordinates.

Example 9.34. (Change of coordinates with the DFT). The DFT is one particular change of coordinates which we have considered. It was the change of coordinates from the standard basis to the Fourier basis. A corresponding change of coordinates in a tensor product is obtained by substituting the DFT as the functions S_1, S_2 for implementing the changes of coordinates above. The change of coordinates in the opposite direction is obtained by using the IDFT instead of the DFT.

Modern image standards do typically not apply a change of coordinates to the entire image. Rather the image is split into smaller squares of appropriate size, called blocks, and a change of coordinates is performed independently for each block. In this example we have split the image into blocks of size 8×8 .

Recall that the DFT values express frequency components. The same applies for the two-dimensional DFT and thus for images, but frequencies are now represented in two different directions. Let us introduce a neglection threshold in the same way as in Example 2.30, to view the image after we set certain frequencies to zero. As for sound, this has little effect on the human perception of the image, if we use a suitable neglection threshold. After we have performed the two-dimensional DFT on an image, we can neglect DFT-coefficients below a threshold on the resulting matrix X with the following code:

```
X=X.*(abs(X)>=threshold);
```

`abs(X)>=threshold` now instead returns a *threshold matrix* with 1 and 0 of the same size as X .

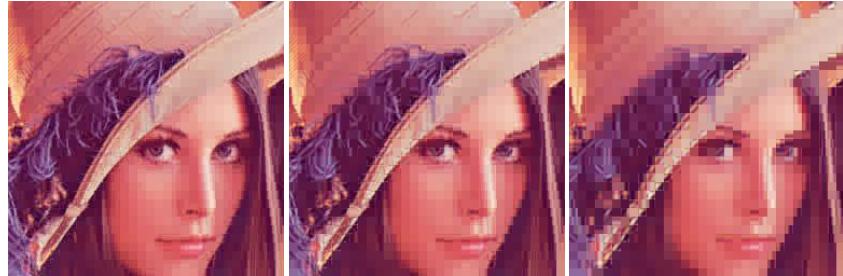
In Figure 9.15 we have applied the two-dimensional DFT to our test image. We have then neglected DFT coefficients which are below certain thresholds, and transformed the samples back to reconstruct the image. When increasing the threshold, the image becomes more and more unclear, but the image is quite clear in the first case, where as much as more than 90% of the samples have been neglected. A blocking effect at the block boundaries is clearly visible. ♣

Example 9.35. (Change of coordinates with the DCT). Similarly to the DFT, the DCT was the change of coordinates from the standard basis to what we called the DCT basis. Change of coordinates in tensor products between the standard basis and the DCT basis is obtained by substituting with the DCT and the IDCT for the changes of coordinates S_1, S_2 above.

The DCT is used more than the DFT in image processing. In particular, the JPEG standard applies a two-dimensional DCT, rather than a two-dimensional DFT. With the JPEG standard, the blocks are always 8×8 , as in the previous example. It is of course not a coincidence that a power of 2 is chosen here, since the DCT, as the DFT, has an efficient implementation for powers of 2.

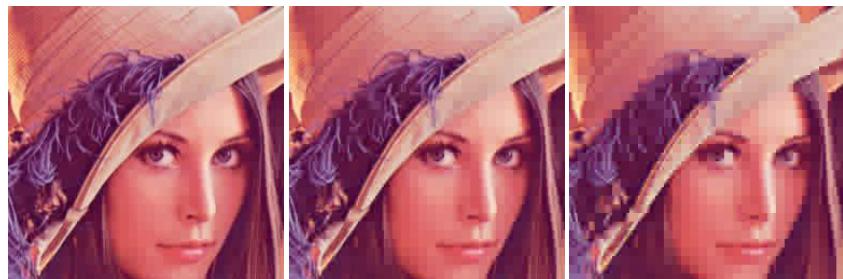
If we follow the same strategy for the DCT as for the DFT example, so that we neglect DCT-coefficients which are below a given threshold,¹ and use the same block sizes, we get the images shown in Figure 9.16. We see similar effects as with the DFT,

¹The JPEG standard does not do exactly the kind of thresholding described here. Rather it performs what is called a quantization.



(a) Threshold 30. 91.4% of the DFT-values were neglected (b) Threshold 50. 95.3% of the DFT-values were neglected (c) Threshold 100. 97.7% of the DFT-values were neglected

Figure 9.15: The effect on an image when it is transformed with the DFT, and the DFT-coefficients below a certain threshold were neglected.



(a) Threshold 30. 93.1% of the DCT-values were neglected (b) Threshold 50. 96.6% of the DCT-values were neglected (c) Threshold 100. 98.8% of the DCT-values were neglected

Figure 9.16: The effect on an image when it is transformed with the DCT, and the DCT-coefficients below a certain threshold were neglected.

but it seems that the latter images are a bit clearer, verifying that the DCT is a better choice than the DFT. It is also interesting to compare with what happens when we drop splitting the image into blocks. Of course, when we neglect many of the DCT-coefficients, we should see some artifacts, but there is no reason to believe that these should be at the old block boundaries. The new artifacts can be seen in Figure 9.17, where the same thresholds as before have been used. Clearly, the new artifacts take a completely different shape. ♣

In the exercises you will be asked to implement functions which generate the images shown in these examples.



(a) Threshold 30. 93.2% of the DCT-values were neglected (b) Threshold 50. 95.8% of the DCT-values were neglected (c) Threshold 100. 97.7% of the DCT-values were neglected

Figure 9.17: The effect on an image when it is transformed with the DCT, and the DCT-coefficients below a certain threshold were neglected. The image has not been split into blocks here.

What you should have learnt in this section

The operation $X \rightarrow S_1 X (S_2)^T$ can also be used to facilitate change of coordinates in images, in addition to filtering images. In other words, change of coordinates is done first column by column, then row by row. The DCT and the DFT are particular changes of coordinates used for images.

Exercises for Section 9.4

1. Implement functions

```
function newx=FFT2Impl(x)
function x=IFF2Impl(newx)
function newx=DCT2Impl(x)
function x=IDCT2Impl(newx)
```

which implement the two-dimensional DCT, FFT, and their inverses as described in this section. Base your code on the code from this section, but where you set the operations S_1 and S_2 accordingly.

2. Implement functions

```
function samples=transform2jpeg(x)
function samples=transform2invjpeg(x)
```

which splits the image into blocks of size 8×8 , and performs the DCT2/IDCT2 on each block. Finally run the code

```
function showDCThigher(threshold)
img=double(imread('lena.png','png'));
zeroedout=0;
```

```

for k=1:3
    newimg=transform2jpeg(img(:,:,k));
    thresholdmatr=(abs(newimg)>=threshold);
    zeroedout=zeroedout+size(img,1)*size(img,2)-sum(sum(thresholdmatr));
    img(:,:,:,k)=transform2invjpeg(newimg.*thresholdmatr);
end
imshow(uint8(255*mapto01(img)));
fprintf('%i percent of samples zeroed out\n',...
    100*zeroedout/(3*size(img,1)*size(img,2)));

```

for different threshold parameters, and check that this reproduces the test images of this section, and prints the correct numbers of values which have been neglected (i.e. which are below the threshold) on screen.

3. Suppose that FFTImpl and IFFTImpl are implementations of the FFT and IFFT algorithms. Suppose also that we have given an image by the matrix X. Consider the following Matlab code:

```

threshold=30;
[m,n]=siz(X);
for col=1:n
    X(:,col)=FFTImpl(X(:,col));
end
for row=1:m
    X(row,:)=FFTImpl((X(row,:))');
end

X=X.*(abs(X)>=threshold);

for col=1:n
    X(:,col)=IFFTImpl(X(:,col));
end
for row=1:m
    X(row,:)=IFFTImpl((X(row,:))');
end

```

Comment what the code does. Comment in particular on the meaning of the parameter `threshold`, and what effect this has on the image.

9.5 Summary

We started by discussing the basic question what an image is, and took a closer look at digital images. We then went through several operations which give meaning for digital images. Many of these operations could be described in terms of a row/column-wise application of filters, and more generally in term of what we called computational molecules. We defined the tensor product, and saw how our

operations could be expressed within this framework. The tensor product framework could also be used to state change of coordinates for images, so that we could consider changes of coordinates such as the DFT and the DCT also for images. The algorithm for computing filtering operations or changes of coordinates for images turned out to be similar, in the sense that the one-dimensional counterparts were simply applied to the rows and the columns in the image.

In introductory image processing textbooks, many other image processing methods are presented. We have limited to the techniques presented here, since our interest in images is mainly for transformation operations which are useful for compression. An excellent textbook on image processing which also uses Matlab is [16]. This contains important topics such as image restoration and reconstruction, geometric transformations, morphology, and object recognition. None of these are considered in this book.

In much literature, one only mentions that filtering can be extended to images by performing one-dimensional filtering for the rows, followed by one-dimensional filtering for the columns, without properly explaining why this is the natural thing to do. The tensor product may be the most natural concept to explain this, and a concept which is firmly established in mathematical literature. Tensor products are usually not part of beginning courses in linear algebra. We have limited the focus here to an introduction to tensor products, and the theory needed to explain filtering an image, and computing the two-dimensional wavelet transform. Some linear algebra books (such as [22]) present tensor products in exercise form only, and often only mentions the Kronecker tensor product, as we defined it.

Many international standards exist for compression of images, and we will take a closer look at two of them in this book. The JPEG standard, perhaps the most popular format for images on the Internet, applies a change of coordinates with a two-dimensional DCT, as described in this chapter. The compression level in JPEG images is selected by the user and may result in conspicuous artefacts if set too high. JPEG is especially prone to artefacts in areas where the intensity changes quickly from pixel to pixel. JPEG is usually lossy, but may also be lossless and has become . The standard defines both the algorithms for encoding and decoding and the storage format. The extension of a JPEG-file is .jpg or .jpeg. JPEG is short for *Joint Photographic Experts Group*, and was approved as an international standard in 1994. A more detailed description of the standard can be found in [27].

The second standard we will consider is JPEG2000. It was developed to address some of the shortcomings of JPEG, and is based on wavelets. The standard document for this [1] does not focus on explaining the theory behind the standard. As the MP3 standard document, it rather states step-by-step procedures for implementing the standard.

The theory we present related to these image standards concentrate on transforming the image (either with a DWT or a DCT) to obtain something which is more suitable for (lossless or lossy) compression. However, many other steps are also needed in order to obtain a full image compression system. One of these is *quantization*. In the simplest form of quantization, every resulting sample from the transformation is rounded to a fixed number of bits. Quantization can also be done in more advanced ways than this: We have already mentioned that the MP3 standard

may use different number of bits for values in the different subbands, depending on the importance of the samples for the human perception. The JPEG2000 standard quantizes in such a way that there is bigger interval around 0 which is quantized to 0, i.e. the rounding error is allowed to be bigger in an interval around 0. Standards which are lossless do not apply quantization, since this always leads to loss.

Somewhere in the image processing or sound processing pipeline, we also need a step which actually achieves compression of the data. Different standards use different lossless coding techniques for this. JPEG2000 uses an advances type of arithmetic coding for this. JPEG can also use arithmetic coding, but also Huffman coding.

Besides transformation, quantization, and coding, many other steps are used, which have different tasks. Many standards preprocess the pixel values before a transform is applied. Preprocessing may mean to center the pixel values around a certain value (JPEG2000 does this), or extracting the different image components before they are processed separately. Also, the image is often split into smaller parts (often called tiles), which are processed separately. For big images this is very important, since it allows users to zoom in on a small part of the image, without processing larger uninteresting parts of the image. Independent processing of the separate tiles makes the image compression what we call *error-resilient*, to errors such as transmission errors, since errors in one tile does not propagate to errors in the other tiles. It is also much more memory-friendly to process the image in several smaller parts, since it is not required to have the entire image in memory at any time. It also gives possibilities for parallel computing. For standards such as JPEG and JPEG2000, tiles are split into even smaller parts, called blocks, where parts of the processing within each block also is performed independently. This makes the possibilities for parallel computing even bigger.

An image standard also defines how to store metadata about an image, and what metadata is accepted, like resolution, time when the image was taken, where the image was taken (such as GPS coordinates), and similar information. Metadata can also tell us how the colour in the image are represented. As we have already seen, in most colour images the colour of a pixel is represented in terms of the amount of red, green and blue or (r, g, b) . But there are other possibilities as well: Instead of storing all 24 bits of colour information in cases where each of the three colour components needs 8 bits, it is common to create a table of up to 256 colours with which a given image could be represented quite well. Instead of storing the 24 bits, one then just stores a colour table in the metadata, and at each pixel, the eight bits corresponding to the correct entry in the table. This is usually referred to as *eight-bit colour*, and the table is called a *look-up table* or *palette*. For large photographs, however, 256 colours is far from sufficient to obtain reasonable colour reproduction. Metadata is usually stored in the beginning of the file, formatted in a very specific way.

Using tensor products to apply wavelets to images

Previously we have used the theory of wavelets to analyze sound. We would also like to use wavelets in a similar way to analyze images. Since the tensor product concept constructs two dimensional objects (matrices) from one-dimensional objects (vectors), we are lead to believe that tensor products can also be used to apply wavelets to images. In this chapter we will see that this can indeed be done. The vector spaces we V_m encountered for wavelets were function spaces, however. What we therefore need first is to establish a general definition of tensor products of function spaces. This will be done in the first section of this chapter. In the second section we will then specialize the function spaces to the spaces V_m we use for wavelets, and interpret the tensor product of these and the wavelet transform applied to images more carefully. Finally we will look at some examples on this theory applied to some example images.

10.1 Tensor product of function spaces

In the setting of functions, it will turn out that the tensor product of two univariate functions can be most intuitively defined as a function in two variables. This seems somewhat different from the strategy of Chapter ??, but we will see that the results we obtain will be very similar.

Definition 10.1 (Tensor product of function spaces). Let U_1 and U_2 be vector spaces of functions, defined on the intervals $[0, M)$ and $[0, N)$, respectively, and suppose that $f_1 \in U_1$ and $f_2 \in U_2$. The tensor product of f_1 and f_2 , denoted $f_1 \otimes f_2$, is the function in two variables defined on $[0, M) \times [0, N)$ by

$$(f_1 \otimes f_2)(t_1, t_2) = f_1(t_1)f_2(t_2).$$

$f_1 \otimes f_2$ is also called the separable extension of f_1 and f_2 to two variables. The tensor product of the spaces $U_1 \otimes U_2$ is the vector space spanned by the two-variable functions $\{f_1 \otimes f_2\}_{f_1 \in U_1, f_2 \in U_2}$.

We will always assume that the spaces U_1 and U_2 consist of functions which are at least integrable. In this case $U_1 \otimes U_2$ is also an inner product space, with the inner product given by a double integral,

$$\langle f, g \rangle = \int_0^N \int_0^M f(t_1, t_2) g(t_1, t_2) dt_1 dt_2. \quad (10.1)$$

In particular, this says that

$$\begin{aligned} \langle f_1 \otimes f_2, g_1 \otimes g_2 \rangle &= \int_0^N \int_0^M f_1(t_1) f_2(t_2) g_1(t_1) g_2(t_2) dt_1 dt_2 \\ &= \int_0^M f_1(t_1) g_1(t_1) dt_1 \int_0^N f_2(t_2) g_2(t_2) dt_2 = \langle f_1, g_1 \rangle \langle f_2, g_2 \rangle. \end{aligned} \quad (10.2)$$

This means that for tensor products, a double integral can be computed as the product of two one-dimensional integrals. This formula also ensures that inner products of tensor products of functions obey the same rule as we found for tensor products of vectors in Exercise ??11.

The tensor product space defined in Definition 10.1 is useful for approximation of functions of two variables if each of the two spaces of univariate functions have good approximation properties.

Idea 10.2. If the spaces U_1 and U_2 can be used to approximate functions in one variable, then $U_1 \otimes U_2$ can be used to approximate functions in two variables.

We will not state this precisely, but just consider some important examples.

Example 10.3. Let $U_1 = U_2$ be the space of all polynomials of finite degree. We know that U_1 can be used for approximating many kinds of functions, such as continuous functions, for example by Taylor series. The tensor product $U_1 \otimes U_1$ consists of all functions on the form $\sum_{i,j} \alpha_{i,j} t_1^i t_2^j$. It turns out that polynomials in several variables have approximation properties analogous to univariate polynomials. ♣

Example 10.4. Let $U_1 = U_2 = V_{N,T}$ be the N th order Fourier space which is spanned by the functions

$$e^{-2\pi i Nt/T}, \dots, e^{-2\pi it/T}, 1, e^{2\pi it/T}, \dots, e^{2\pi i Nt/T}$$

The tensor product space $U_1 \otimes U_1$ now consists of all functions on the form

$$\sum_{k,l=-N}^N \alpha_{k,l} e^{2\pi i k t_1 / T} e^{2\pi i l t_2 / T}.$$

One can show that this space has approximation properties similar to $V_{N,T}$ for functions in two variables. This is the basis for the theory of Fourier series in two variables. ♣

In the following we think of $U_1 \otimes U_2$ as a space which can be used for approximating a general class of functions. By associating a function with the vector of coordinates relative to some basis, and a matrix with a function in two variables, we have the following parallel to Theorem 9.30:

Theorem 10.5. If $\{f_i\}_{i=0}^{M-1}$ is a basis for U_1 and $\{g_j\}_{j=0}^{N-1}$ is a basis for U_2 , then $\{f_i \otimes g_j\}_{(i,j)=(0,0)}^{(M-1,N-1)}$ is a basis for $U_1 \otimes U_2$. Moreover, if the bases for U_1 and U_2 are orthogonal/orthonormal, then the basis for $U_1 \otimes U_2$ is orthogonal/orthonormal.

Proof. The proof is similar to that of Theorem 9.30: if

$$\sum_{(i,j)=(0,0)}^{(M-1,N-1)} \alpha_{i,j} (f_i \otimes g_j) = 0,$$

we define $h_i(t_2) = \sum_{j=0}^{N-1} \alpha_{i,j} g_j(t_2)$. It follows as before that $\sum_{i=0}^{M-1} h_i(t_2) f_i = 0$ for any t_2 , so that $h_i(t_2) = 0$ for any t_2 due to linear independence of the f_i . But then $\alpha_{i,j} = 0$ also, due to linear independence of the g_j . The statement about orthogonality follows from Equation 10.2. \blacksquare

We can now define the tensor product of two bases of functions as before, and coordinate matrices as before:

Definition 10.6. if $\mathcal{B} = \{f_i\}_{i=0}^{M-1}$ and $\mathcal{C} = \{g_j\}_{j=0}^{N-1}$, we define $\mathcal{B} \otimes \mathcal{C}$ as the basis $\{f_i \otimes g_j\}_{(i,j)=(0,0)}^{(M-1,N-1)}$ for $U_1 \otimes U_2$. We say that X is the coordinate matrix of f if $f(t_1, t_2) = \sum X_{i,j} (f_i \otimes g_j)(t_1, t_2)$, where $X_{i,j}$ are the elements of X .

Theorem 9.32 can also be proved in the same way in the context of function spaces. We state this as follows:

Theorem 10.7 (Change of coordinates in tensor products of function spaces).
Assume that U_1 and U_2 are function spaces, and that

1. $\mathcal{B}_1, \mathcal{C}_1$ are bases for U_1 , and that S_1 is the change of coordinates matrix from \mathcal{B}_1 to \mathcal{C}_1 ,
2. $\mathcal{B}_2, \mathcal{C}_2$ are bases for U_2 , and that S_2 is the change of coordinates matrix from \mathcal{B}_2 to \mathcal{C}_2 .

Both $\mathcal{B}_1 \otimes \mathcal{B}_2$ and $\mathcal{C}_1 \otimes \mathcal{C}_2$ are bases for $U_1 \otimes U_2$, and if X is the coordinate matrix in $\mathcal{B}_1 \otimes \mathcal{B}_2$, Y the coordinate matrix in $\mathcal{C}_1 \otimes \mathcal{C}_2$, then the change of coordinates from $\mathcal{B}_1 \otimes \mathcal{B}_2$ to $\mathcal{C}_1 \otimes \mathcal{C}_2$ can be computed as

$$Y = S_1 X (S_2)^T. \quad (10.3)$$

10.2 Tensor product of function spaces in a wavelet setting

We will now specialize the spaces U_1, U_2 from Definition 10.1 to the resolution spaces V_m and the detail spaces W_m , arising from a given wavelet. We can in particular form the tensor products $\phi_{0,n_1} \otimes \phi_{0,n_2}$. We will assume that

1. the first component ϕ_{0,n_1} has period M (so that $\{\phi_{0,n_1}\}_{n_1=0}^{M-1}$ is a basis for the first component space),
2. the second component ϕ_{0,n_2} has period N (so that $\{\phi_{0,n_2}\}_{n_2=0}^{N-1}$ is a basis for the second component space).

When we speak of $V_0 \otimes V_0$ we thus mean an MN -dimensional space with basis $\{\phi_{0,n_1} \otimes \phi_{0,n_2}\}_{(n_1,n_2)=(0,0)}^{(M-1,N-1)}$, where the coordinate matrices are $M \times N$. This difference in the dimension of the two components is done to allow for images where the number of rows and columns may be different. In the following we will implicitly assume that the component spaces have dimension M and N , to ease notation. If we use that (ϕ_{m-1}, ψ_{m-1}) also is a basis for V_m , we get the following corollary to Theorem 10.5:

Corollary 10.8. Let ϕ, ψ be a scaling function and a mother wavelet. Then the two sets of tensor products given by

$$\boldsymbol{\phi}_m \otimes \boldsymbol{\phi}_m = \{\phi_{m,n_1} \otimes \phi_{m,n_2}\}_{n_1,n_2}$$

and

$$\begin{aligned} & (\boldsymbol{\phi}_{m-1}, \boldsymbol{\psi}_{m-1}) \otimes (\boldsymbol{\phi}_{m-1}, \boldsymbol{\psi}_{m-1}) \\ &= \{\phi_{m-1,n_1} \otimes \phi_{m-1,n_2}, \\ & \quad \phi_{m-1,n_1} \otimes \psi_{m-1,n_2}, \\ & \quad \psi_{m-1,n_1} \otimes \phi_{m-1,n_2}, \\ & \quad \psi_{m-1,n_1} \otimes \psi_{m-1,n_2}\}_{n_1,n_2} \end{aligned}$$

are both bases for $V_m \otimes V_m$. This second basis is orthogonal/orthonormal whenever the first basis is.

From this we observe that, while the one-dimensional wavelet decomposition splits V_m into a direct sum of the two vector spaces V_{m-1} and W_{m-1} , the corresponding two-dimensional decomposition splits $V_m \otimes V_m$ into a direct sum of four tensor product vector spaces. These vector spaces deserve individual names:

Definition 10.9. We define the following tensor product spaces:

1. The space $W_m^{(0,1)}$ spanned by $\{\phi_{m,n_1} \otimes \psi_{m,n_2}\}_{n_1,n_2}$,
2. The space $W_m^{(1,0)}$ spanned by $\{\psi_{m,n_1} \otimes \phi_{m,n_2}\}_{n_1,n_2}$,

3. The space $W_m^{(1,1)}$ spanned by $\{\psi_{m,n_1} \otimes \psi_{m,n_2}\}_{n_1, n_2}$.

Since these spaces are linearly independent, we can write

$$V_m \otimes V_m = (V_{m-1} \otimes V_{m-1}) \oplus W_{m-1}^{(0,1)} \oplus W_{m-1}^{(1,0)} \oplus W_{m-1}^{(1,1)}. \quad (10.4)$$

Also in the setting of tensor products we refer to $V_{m-1} \otimes V_{m-1}$ as the space of low-resolution approximations. The remaining parts, $W_{m-1}^{(0,1)}$, $W_{m-1}^{(1,0)}$, and $W_{m-1}^{(1,1)}$, are referred to as detail spaces. The coordinate matrix of

$$\sum_{n_1, n_2=0}^{2^{m-1}N} (c_{m-1, n_1, n_2} (\phi_{m-1, n_1} \otimes \phi_{m-1, n_2}) + w_{m-1, n_1, n_2}^{(0,1)} (\phi_{m-1, n_1} \otimes \psi_{m-1, n_2}) + w_{m-1, n_1, n_2}^{(1,0)} (\psi_{m-1, n_1} \otimes \phi_{m-1, n_2}) + w_{m-1, n_1, n_2}^{(1,1)} (\psi_{m-1, n_1} \otimes \psi_{m-1, n_2})) \quad (10.5)$$

in the basis $(\phi_{m-1}, \psi_{m-1}) \otimes (\phi_{m-1}, \psi_{m-1})$ is

$$\left(\begin{array}{cc|cc} c_{m-1, 0, 0} & \cdots & w_{m-1, 0, 0}^{(0,1)} & \cdots \\ \vdots & \vdots & \vdots & \vdots \\ \hline w_{m-1, 0, 0}^{(1,0)} & \cdots & w_{m-1, 0, 0}^{(1,1)} & \cdots \\ \vdots & \vdots & \vdots & \vdots \end{array} \right). \quad (10.6)$$

The coordinate matrix is thus split into four submatrices:

- The c_{m-1} -values, i.e. the coordinates for $V_{m-1} \oplus V_{m-1}$. This is the upper left corner in Equation (10.6).
- The $w_{m-1}^{(0,1)}$ -values, i.e. the coordinates for $W_{m-1}^{(0,1)}$. This is the upper right corner in Equation (10.6).
- The $w_{m-1}^{(1,0)}$ -values, i.e. the coordinates for $W_{m-1}^{(1,0)}$. This is the lower left corner in Equation (10.6).
- The $w_{m-1}^{(1,1)}$ -values, i.e. the coordinates for $W_{m-1}^{(1,1)}$. This is the lower right corner in Equation (10.6).

The $w_{m-1}^{(i,j)}$ -values are as in the one-dimensional situation often referred to as wavelet coefficients. Let us consider the Haar wavelet as an example.

Example 10.10. If V_m is the vector space of piecewise constant functions on any interval of the form $[k2^{-m}, (k+1)2^{-m}]$ (as in the piecewise constant wavelet), $V_m \otimes V_m$ is the vector space of functions in two variables which are constant on any square of the form $[k_1 2^{-m}, (k_1+1)2^{-m}] \times [k_2 2^{-m}, (k_2+1)2^{-m}]$. Clearly $\phi_{m,k_1} \otimes \phi_{m,k_2}$ is constant on such a square and 0 elsewhere, and these functions are a basis for $V_m \otimes V_m$.

Let us compute the orthogonal projection of $\phi_{1,k_1} \otimes \phi_{1,k_2}$ onto $V_0 \otimes V_0$. Since the Haar wavelet is orthonormal, the basis functions in (10.4) are orthonormal, and we

can thus use the orthogonal decomposition formula to find this projection. Clearly $\phi_{1,k_1} \otimes \phi_{1,k_2}$ has different support from all except one of $\phi_{0,n_1} \otimes \phi_{0,n_2}$. Since

$$\langle \phi_{1,k_1} \otimes \phi_{1,k_2}, \phi_{0,n_1} \otimes \phi_{0,n_2} \rangle = \langle \phi_{1,k_1}, \phi_{0,n_1} \rangle \langle \phi_{1,k_2}, \phi_{0,n_2} \rangle = \frac{\sqrt{2}}{2} \frac{\sqrt{2}}{2} = \frac{1}{2}$$

when the supports intersect, we obtain

$$\text{proj}_{V_0 \otimes V_0}(\phi_{1,k_1} \otimes \phi_{1,k_2}) = \begin{cases} \frac{1}{2}(\phi_{0,k_1/2} \otimes \phi_{0,k_2/2}) & \text{when } k_1, k_2 \text{ are even} \\ \frac{1}{2}(\phi_{0,k_1/2} \otimes \phi_{0,(k_2-1)/2}) & \text{when } k_1 \text{ is even, } k_2 \text{ is odd} \\ \frac{1}{2}(\phi_{0,(k_1-1)/2} \otimes \phi_{0,k_2/2}) & \text{when } k_1 \text{ is odd, } k_2 \text{ is even} \\ \frac{1}{2}(\phi_{0,(k_1-1)/2} \otimes \phi_{0,(k_2-1)/2}) & \text{when } k_1, k_2 \text{ are odd} \end{cases}$$

So, in this case there were 4 different formulas, since there were 4 different combinations of even/odd. Let us also compute the projection onto the orthogonal complement of $V_0 \otimes V_0$ in $V_1 \otimes V_1$, and let us express this in terms of the $\phi_{0,n}, \psi_{0,n}$, like we did in the one-variable case. Also here there are 4 different formulas. When k_1, k_2 are both even we obtain

$$\begin{aligned} & \phi_{1,k_1} \otimes \phi_{1,k_2} - \text{proj}_{V_0 \otimes V_0}(\phi_{1,k_1} \otimes \phi_{1,k_2}) \\ &= \phi_{1,k_1} \otimes \phi_{1,k_2} - \frac{1}{2}(\phi_{0,k_1/2} \otimes \phi_{0,k_2/2}) \\ &= \left(\frac{1}{\sqrt{2}}(\phi_{0,k_1/2} + \psi_{0,k_1/2}) \right) \otimes \left(\frac{1}{\sqrt{2}}(\phi_{0,k_2/2} + \psi_{0,k_2/2}) \right) - \frac{1}{2}(\phi_{0,k_1/2} \otimes \phi_{0,k_2/2}) \\ &= \frac{1}{2}(\phi_{0,k_1/2} \otimes \phi_{0,k_2/2}) + \frac{1}{2}(\phi_{0,k_1/2} \otimes \psi_{0,k_2/2}) \\ &\quad + \frac{1}{2}(\psi_{0,k_1/2} \otimes \phi_{0,k_2/2}) + \frac{1}{2}(\psi_{0,k_1/2} \otimes \psi_{0,k_2/2}) - \frac{1}{2}(\phi_{0,k_1/2} \otimes \phi_{0,k_2/2}) \\ &= \frac{1}{2}(\phi_{0,k_1/2} \otimes \psi_{0,k_2/2}) + \frac{1}{2}(\psi_{0,k_1/2} \otimes \phi_{0,k_2/2}) + \frac{1}{2}(\psi_{0,k_1/2} \otimes \psi_{0,k_2/2}). \end{aligned}$$

Here we have used the relation $\phi_{1,k_i} = \frac{1}{\sqrt{2}}(\phi_{0,k_i/2} + \psi_{0,k_i/2})$, which we have from our first analysis of the Haar wavelet. Checking the other possibilities we find similar formulas for the projection onto the orthogonal complement of $V_0 \otimes V_0$ in $V_1 \otimes V_1$ when either k_1 or k_2 is odd. In all cases, the formulas use the basis functions for $W_0^{(0,1)}, W_0^{(1,0)}, W_0^{(1,1)}$. These functions are shown in Figure 10.1, together with the function $\phi \otimes \phi \in V_0 \otimes V_0$. ♣

Example 10.11. If we instead use any of the wavelets for piecewise linear functions, the wavelet basis functions are not orthogonal anymore, just as in the one-dimensional case. The new basis functions are shown in Figure 10.2 for the alternative piecewise linear wavelet. ♣

An immediate corollary of Theorem 10.7 is the following:

Corollary 10.12. Let

$$\begin{aligned} A_m &= P_{(\phi_{m-1}, \psi_{m-1}) \leftarrow \phi_m} \\ B_m &= P_{\phi_m \leftarrow (\phi_{m-1}, \psi_{m-1})} \end{aligned}$$

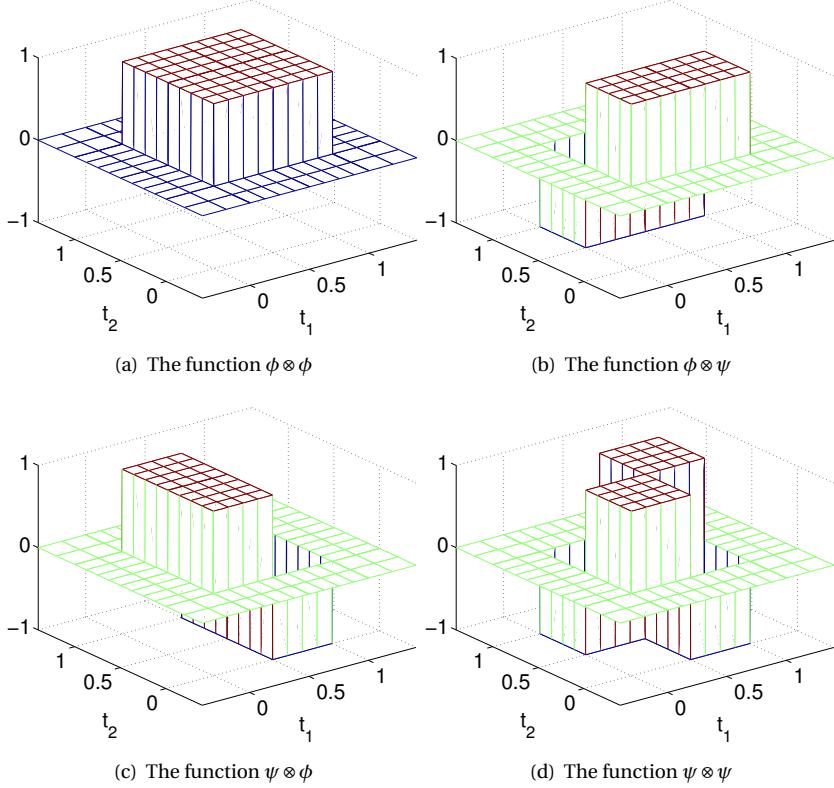


Figure 10.1: The basis functions for $(V_0 \otimes V_0) \oplus W_0^{(0,1)} \oplus W_0^{(1,0)} \oplus W_0^{(1,1)}$ for the Haar wavelet.

be the stages in the DWT and the IDWT, and let

$$X = (c_{m,i,j})_{i,j} \quad Y = \begin{pmatrix} (c_{m-1,i,j})_{i,j} & (w_{m-1,i,j}^{(0,1)})_{i,j} \\ (w_{m-1,i,j}^{(1,0)})_{i,j} & (w_{m-1,i,j}^{(1,1)})_{i,j} \end{pmatrix} \quad (10.7)$$

be the coordinate matrices in $\phi_m \otimes \phi_m$, and $(\phi_{m-1}, \psi_{m-1}) \otimes (\phi_{m-1}, \psi_{m-1})$, respectively. Then

$$Y = A_m X A_m^T \quad (10.8)$$

$$X = B_m Y B_m^T \quad (10.9)$$

By the m -level two-dimensional DWT/IDWT (or DWT2/IDWT2) we mean the change of coordinates where this is repeated m times as in a DWT/IDWT.

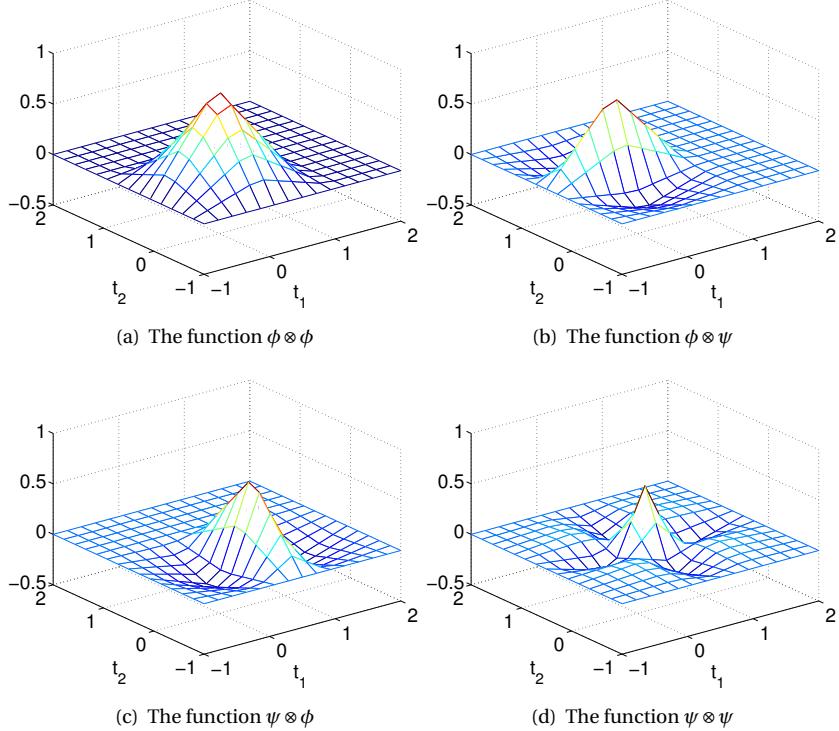


Figure 10.2: The basis functions for $(V_0 \otimes V_0) \oplus W_0^{(0,1)} \oplus W_0^{(1,0)} \oplus W_0^{(1,1)}$ for the alternative piecewise linear wavelet.

Each stage in DWT2 and IDWT2 can now be implemented by substituting the matrices A_m, B_m above into Theorem ???. This implementation can reuse an efficient implementation of the one-dimensional DWT/IDWT. When using many levels of the DWT2, the next stage is applied only to the upper left corner of the matrix, just as the DWT at the next stage only is applied to the first part of the coordinates. At each stage, the upper left corner of the coordinate matrix (which gets smaller at each iteration), is split into four equally big parts. This is illustrated in Figure 10.3, where the different types of coordinates which appear in the first two stages in a DWT2 are indicated.

It is instructive to see what information the different types of coordinates in an image represent. In the following examples we will discard some types of coordinates, and view the resulting image. Discarding a type of coordinates will be illustrated by coloring the corresponding regions from Figure 10.3 black. As an example, if we perform a two-level DWT2 (i.e. we start with a coordinate matrix in the basis $\phi_2 \otimes \phi_2$), Figure 10.4 illustrates first the collection of all coordinates, and then the resulting collection of coordinates after removing subbands at the first level suc-

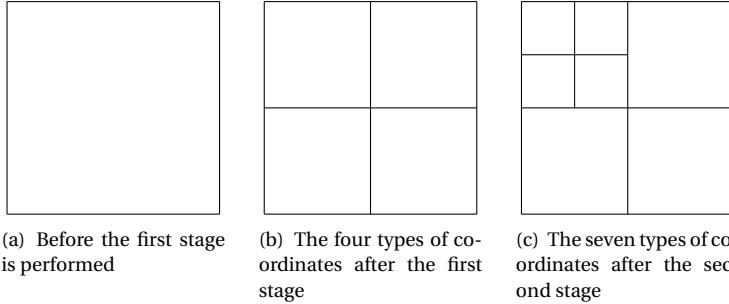


Figure 10.3: Illustration of the different coordinates in a two level DWT2.

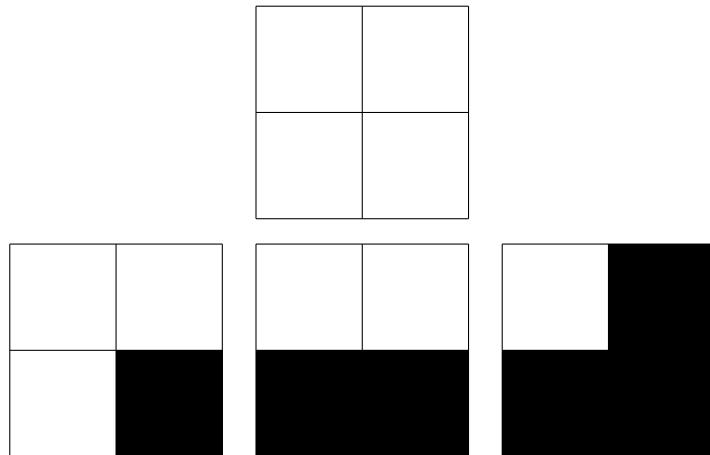


Figure 10.4: Graphical representation of neglecting the wavelet coefficients at the first level. After applying DWT2, the wavelet coefficients are split into four parts, as shown in the first figure. In the following figures we have removed coefficients from $W_1^{(1,1)}$, $W_1^{(1,0)}$, and $W_1^{(0,1)}$, in that order.

sively. Figure 10.5 illustrates in the same way incremental removal of the subbands at the second level.

Before we turn to experiments on images using wavelets, we would like to make another interpretation on the corners in the matrices after the DWT2, which correspond to the different coordinates $(c_{m-1,i,j})_{i,j}$, $(w^{(0,1)})_{m-1,i,j}$, $(w^{(1,0)})_{m-1,i,j}$, and $(w^{(1,1)})_{m-1,i,j}$. It turns out that these corners have natural interpretations in terms of the filter characterization of wavelets, as given in Chapter 6. Recall again that in a DWT2, the DWT is first applied to the columns in the image, then to the rows in the image. Recall first that the DWT2 applies first the DWT to all columns, and then to all rows in the resulting matrix.

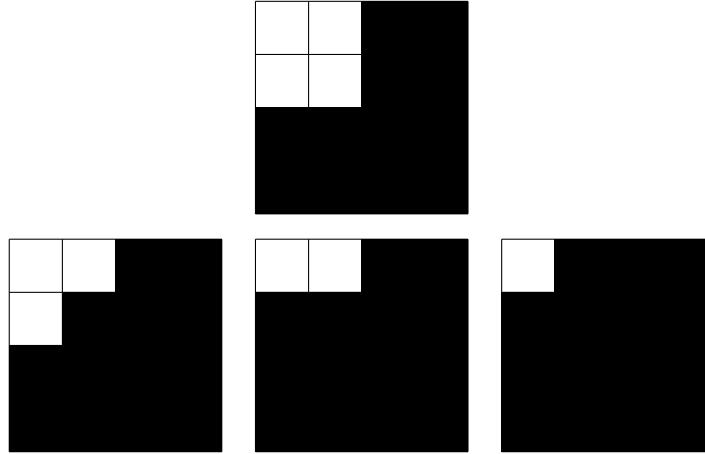


Figure 10.5: Graphical representation of neglecting the wavelet coefficients at the second level. After applying the second stage in DWT2, the wavelet coefficients from the upper left corner are also split into four parts, as shown in the first figure. In the following figures we have removed coefficients from $W_2^{(1,1)}$, $W_2^{(1,0)}$, and $W_2^{(0,1)}$, in that order.

First the DWT is applied to all columns in the image. Since the first half of the coordinates in a DWT are outputs from a lowpass filter H_0 (Theorem 6.3), the upper half after the DWT has now been subject to a lowpass filter to the columns. Similarly, the second half of the coordinates in a DWT are outputs from a highpass filter H_1 (Theorem 6.3 again), so that the bottom half after the DWT has been subject to a highpass filter to the columns.

Then the DWT is applied to all rows in the image. Similarly as when we applied the DWT to the columns, the left half after the DWT has been subject to the same lowpass filter to the rows, and the right half after the DWT has been subject to the same highpass filter to the rows.

These observations split the resulting matrix after DWT2 into four blocks, with each block corresponding to a combination of lowpass and highpass filters. The following names are thus given to these blocks:

- The upper left corner is called the LL-subband,
- The upper right corner is called the LH-subband,
- The lower left corner is called the HL-subband,
- The lower right corner is called the HH-subband.

The two letters indicate the type of filters which have been applied (L=lowpass, H=highpass). The first letter indicates the type of filter which is applied to the columns, the second indicates which is applied to the rows. The order is therefore important. The name

subband comes from the interpretation of these filters as being selective on a certain frequency band. In conclusion, a block in the matrix after the DWT2 corresponds to applying a combination of lowpass/highpass filters to the rows of the columns of the image. Due to this, and since lowpass filters extract slow variations, highpass filters abrupt changes, the following holds:

Observation 10.13. After the DWT2 has been applied to an image, we expect to see the following:

- In the upper left corner, slow variations in both the vertical and horizontal directions are captured, i.e. this is a low-resolution version of the image.
- In the upper right corner, slow variations in the vertical direction are captured, together with abrupt changes in the horizontal direction.
- In the lower left corner, slow variations in the horizontal direction are captured, together with abrupt changes in the vertical direction.
- In the lower right corner, abrupt changes in both directions appear are captured.

These effects will be studied through examples in the next section.

10.3 Experiments with images using wavelets

In this section we will make some experiments with images using the wavelets we have considered¹. The wavelet theory is applied to images in the following way: We first visualize the pixels in the image as coordinates in the basis $\phi_m \otimes \phi_m$ (so that the image has size $(2^m M) \times (2^m N)$). As in the case for sound, this will represent a good approximation when m is large. We then perform a change of coordinates with the DWT2. As we did for sound, we can then either set the detail components from the $W_k^{(i,j)}$ -spaces to zero, or the low-resolution approximation from $V_0 \otimes V_0$ to zero, depending on whether we want to inspect the detail components or the low-resolution approximation. Finally we apply the IDWT2 to end up with coordinates in $\phi_m \otimes \phi_m$ again, and display the new image with pixel values equal to these coordinates.

Example 10.14 (Applying the Haar wavelet to a very simple example image). Let us apply the Haar wavelet to the sample chess pattern example image from Figure 9.14. The lowpass filter of the Haar wavelet was essentially a smoothing filter with two elements. Also, as we have seen, the highpass filter essentially computes an approximation to the partial derivative. Clearly, abrupt changes in the vertical and horizontal directions appear here only at the edges in the chess pattern, and abrupt changes in both directions appear only at the grid points in the chess pattern. Due to Observation 10.13, after a DWT2 we expect to see the following:

¹Note also that Matlab has a wavelet toolbox which could be used for these purposes. We will however not go into the usage of this, since we implement the DWT from scratch.

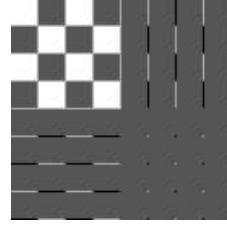


Figure 10.6: The chess pattern example image after application of the DWT2. The Haar wavelet was used.

- In the upper left corner, we should see a low-resolution version of the image.
- In the upper right corner, only the vertical edges in the chess pattern should be visible.
- In the lower left corner, only the horizontal edges in the chess pattern should be visible.
- In the lower right corner, only the grid points in the chess pattern should be visible.

In Figure 10.6 we have applied one level of the DWT2 to the chess pattern example image, and all these effects are seen clearly here. ♣

Example 10.15 (Creating thumbnail images). Let us apply the Haar wavelet to our sample image. In Exercise 3 you will be asked to implement a function which computes DWT2 for the Haar wavelet. After the DWT2, the upper left submatrices represent the low-resolution approximations from $V_{m-1} \otimes V_{m-1}$, $V_{m-2} \otimes V_{m-2}$, and so on. We can now use the following code to store the low-resolution approximation for $m = 1$:

```
X = double(imread('lena.png','png'));
[11,12]=size(X);
Y=DWT2HaarImpl(X,1);
Y=Y(1:(11/2),1:(12/2));
imshow(uint8(Y));
```

In Figure 10.7 the results are shown up to 4 resolutions. In Figure 10.8 we have also shown the entire result after a 1- and 2-stage DWT2 on the image. The first two thumbnail images can be seen as the the upper left corners of the first two images. The other corners represent detail. ♣

Example 10.16 (Detail and low-resolution approximations with the Haar wavelet). In Exercise 4 you will be asked to implement a function `showDWTlowerHaar` which displays the low-resolution approximations to our test image for the Haar wavelet, using functions we implement in the exercises. Let us take a closer look at the images generated. Above we viewed the low-resolution approximation as a smaller

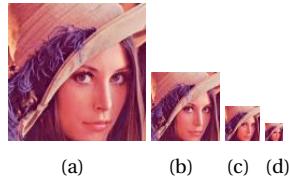


Figure 10.7: The corresponding thumbnail images for the Image of Lena, obtained with a DWT of 1, 2, 3, and 4 levels.

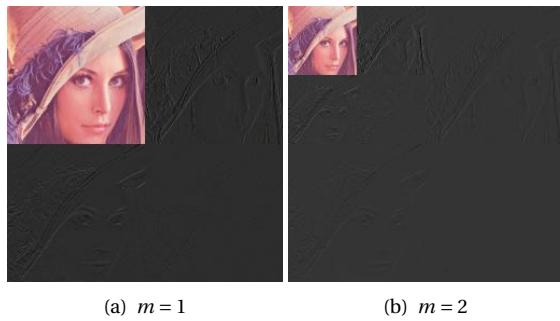


Figure 10.8: The corresponding image resulting from a wavelet transform with the Haar-wavelet.

image. Let us compare with the image resulting from setting the wavelet detail coefficients to zero, and viewing the result as an image of the same size. In particular, let us neglect the wavelet coefficients as pictured in Figure 10.4 and Figure 10.5. Since the Haar wavelet has few vanishing moments, we should expect that the lower order resolution approximations from V_0 are worse when m increase. Figure 10.9 confirms this for the lower order resolution approximations. Alternatively, we should see that the higher order detail spaces contain more information. In Exercise 5 you will be asked to implement a function `showDWTlowerdifferenceHaar` which displays the detail components in the image for a given resolution m for the Haar wavelet. The new images when this function is used are shown in Figure 10.10. The black colour indicates values which are close to 0. In other words, most of the coefficients are close to 0, which reflects one of the properties of the wavelet. ♣

Example 10.17 (The Spline 5/3 wavelet and removing bands in the detail spaces).

In Exercise 8 you will be asked to implement functions `showDWTlower53` and `showDWTlower97` which display the low-resolution approximations to our image test file `lena.png`, for the Spline 5/3 and CDF 9/7 wavelets (these call functions we also implement in the exercises). With these functions we can display the result for all the wavelets we have considered up to now, in succession, and at a given resolution, with the following code:

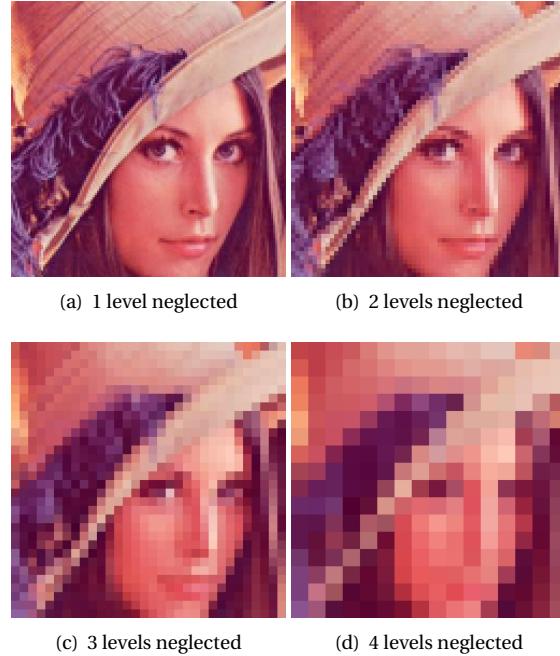


Figure 10.9: Image of Lena, with higher levels of detail neglected for the Haar wavelet.

```
function showDWTall(m)
    disp('Haar wavelet');
    showDWTlowerHaar(m);
    disp('5/3 wavelet');
    showDWTlower53(m);
    disp('9/7 wavelet');
    showDWTlower97(m);
```

The call to `showDWTlowerHaar` first displays the result, using the Haar wavelet. The code then moves on to the function `showDWTlower53` which uses the Spline 5/3 wavelet, and the function `showDWTlower97` which uses the CDF 9/7 wavelet. In the `show`-functions, the image is first read from file, and then code of the following form is called to neglect m levels of detail:

```
Y=DWT2Impl53(X,m);
tokeep=Y(1:(11/2^m),1:(12/2^m));
Y=zeros(size(Y));
Y(1:(11/2^m),1:(12/2^m))=tokeep(1:(11/2^m),1:(12/2^m));
X=IDWT2Impl53(Y,m);
```

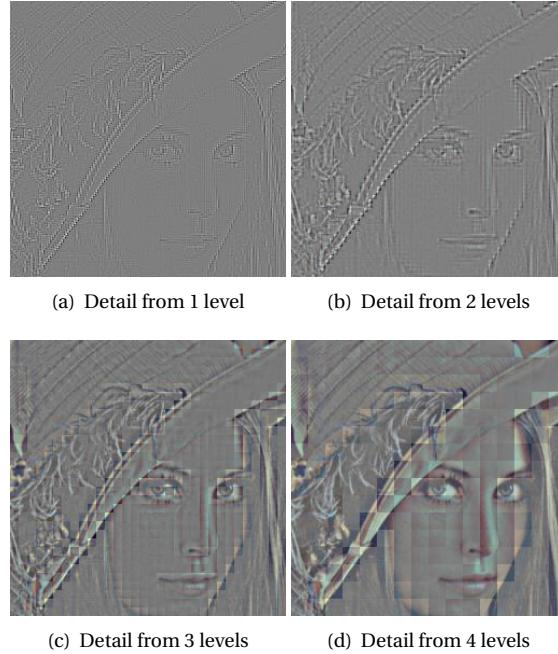


Figure 10.10: The corresponding detail for the images in Figure 10.9, with the Haar wavelet.

Here the Spline 5/3 wavelet was used, and the image had size $l_1 \times l_2$. We can repeat this for various number of levels m , and compare the different images. We can also neglect only parts of the detail, since it at each level is grouped into three bands ($W_m^{(1,1)}$, $W_m^{(1,0)}$, $W_m^{(0,1)}$), contrary to the one-dimensional case. Let us use the Spline 5/3 wavelet. The resulting images when the bands on the first level indicated in Figure 10.4 are removed are shown in Figure 10.11. The resulting images when the bands on the second level indicated in Figure 10.5 are removed are shown in Figure 10.12. The image is seen still to resemble the original one, even after two levels of wavelets coefficients have been neglected. This in itself is good for compression purposes, since we may achieve compression simply by dropping the given coefficients. However, if we continue to neglect more levels of coefficients, the result will look poorer. In Figure 10.13 we have also shown the resulting image after the third and fourth level of detail have been neglected. Although we still can see details in the image, the quality in the image is definitely poorer. Although the quality is poorer when we neglect levels of wavelet coefficients, all information is kept if we additionally include the detail/bands. In Figure 10.14, we have shown the corresponding detail for Figure 10.11(d), Figure 10.12(c), and Figure 10.13. Clearly, more detail can be seen in the image when more of the detail is included. ♣

Example 10.18. Let us repeat the previous example for the CDF 9/7 wavelet, using



(a) The image unaltered



(b) Resulting image after neglecting detail in $W_1^{(1,1)}$, as illustrated in Figure 10.4(b).



(c) Resulting image after neglecting also detail in $W_1^{(1,0)}$, as illustrated in Figure 10.4(c).

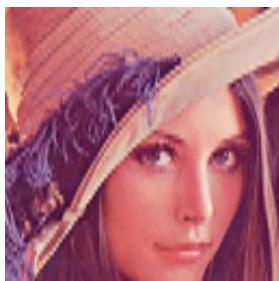


(d) Resulting image after neglecting also detail in $W_1^{(0,1)}$, as illustrated in Figure 10.4(d).

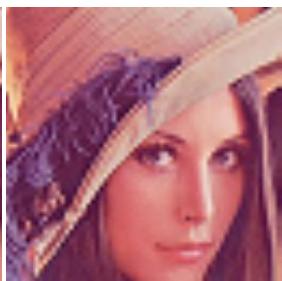
Figure 10.11: Image of Lena, with various bands of detail at the first level neglected. The Spline 5/3 wavelet was used.



(a) Resulting image after also neglecting detail in $W_2^{(1,1)}$, as illustrated in Figure 10.12(a).



(b) Resulting image after also neglecting detail in $W_2^{(1,0)}$, as illustrated in Figure 10.12(b).



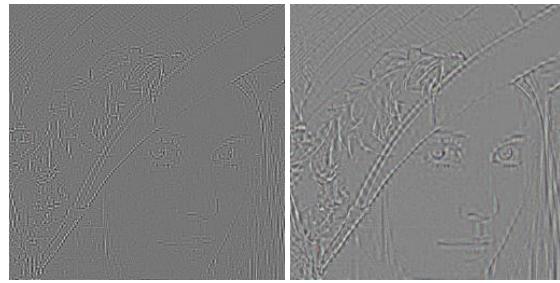
(c) Resulting image after also neglecting detail in $W_2^{(0,1)}$, as illustrated in Figure 10.12(c).

Figure 10.12: Image of Lena, with various bands of detail at the second level neglected. The Spline 5/3 wavelet was used.



(a) 3 levels neglected (b) 4 levels neglected

Figure 10.13: Image of Lena, with higher levels of detail neglected. The Spline 5/3 wavelet was used.



(a) Detail from 1 level (b) Detail from 2 levels



(c) Detail from 3 levels (d) Detail from 4 levels

Figure 10.14: The corresponding detail for the image of Lena. The Spline 5/3 wavelet was used.

the function `showDWTlower97` you implemented in Exercise 8. We should now see improved images when we discard the detail in the images. Figure 10.15 confirms this for the lower resolution spaces, while Figure 10.16 confirms this for the higher order detail spaces. ♣

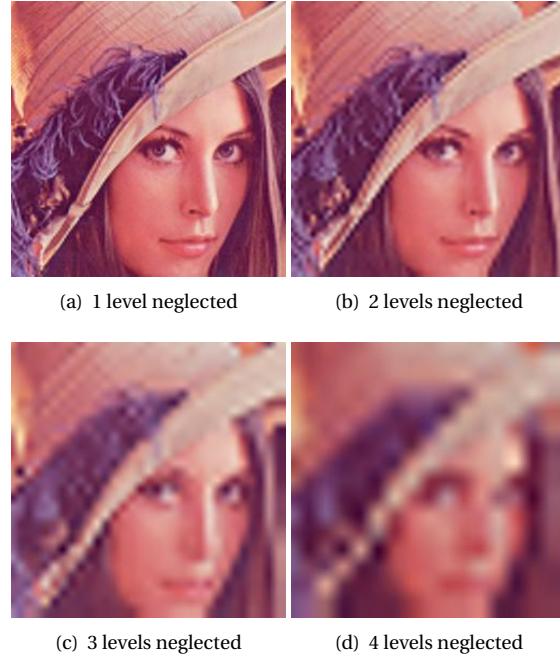


Figure 10.15: Image of Lena, with higher levels of detail neglected. The CDF 9/7 wavelet was used.

As mentioned, the procedure developed in this section for applying a wavelet transform to an image with the help of the tensor product construction, is adopted in the JPEG2000 standard. This lossy (can also be used as lossless) image format was developed by the Joint Photographic Experts Group and published in 2000. After significant processing of the wavelet coefficients, the final coding with JPEG2000 uses an advanced version of arithmetic coding. At the cost of increased encoding and decoding times, JPEG2000 leads to as much as 20 % improvement in compression ratios for medium compression rates, possibly more for high or low compression rates. The artefacts are less visible than in JPEG and appear at higher compression rates. Although a number of components in JPEG2000 are patented, the patent holders have agreed that the core software should be available free of charge, and JPEG2000 is part of most Linux distributions. However, there appear to be some further, rather obscure, patents that have not been licensed, and this may be the reason why JPEG2000 is not used more. The extension of JPEG2000 files is .jp2.

Exercises for Section 10.3

Note that there are three colour components in the test image 'lena.png'. In the following exercises you must therefore run the DWT and IDWT2 on all three com-

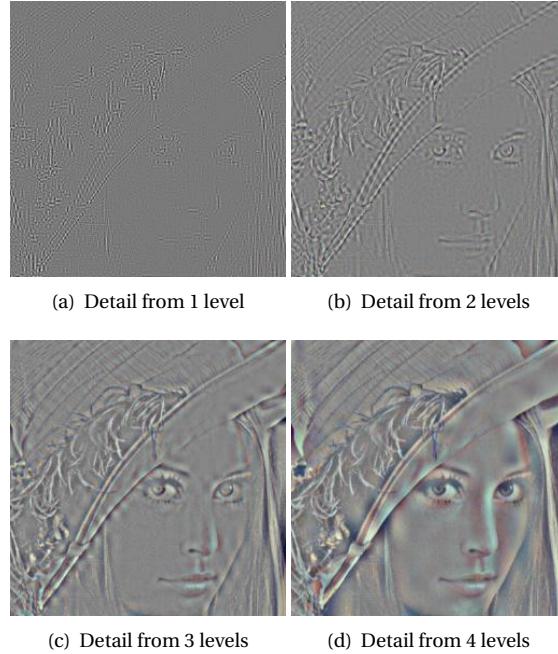


Figure 10.16: The corresponding detail for the image of Lena. The CDF 9/7 wavelet was used.

ponents.

1. Assume that we have an image represented by the $N \times N$ -matrix X , and consider the following Matlab code:

```

for k=1:2
    for s=1:N
        c=(X(1:2:(N-1),s)+X(2:2:N,s))/sqrt(2);
        w=(X(1:2:(N-1),s)-X(2:2:N,s))/sqrt(2);
        X(1:N,s)=[c; w];
    end
    X=X';
end

```

a. Comment what the code does, and explain what you will see if you display X as an image after the code has run.

b. The code above has an inverse transformation, which reproduce the original image from the transformed values which we obtained. Assume that you zero out the values in the lower left and the upper right corner of the matrix

X after the code above has run, and that you then reproduce the image by applying this inverse transformation. What changes can you then expect in the image?

2. In this exercise we will use the filters $G_0 = \{1, 1\}$, $G_1 = \{1, -1\}$.
 - a. Let X be a matrix which represents the pixel values in an image. Define $\mathbf{x} = (1, 0, 1, 0)$ and $\mathbf{y} = (0, 1, 0, 1)$. Compute $(G_0 \otimes G_0)(\mathbf{x} \otimes \mathbf{y})$.
 - b. For a general image X , describe how the images $(G_0 \otimes G_0)X$, $(G_0 \otimes G_1)X$, $(G_1 \otimes G_0)X$, and $(G_1 \otimes G_1)X$ may look.
 - c. Assume that we run the following code on an image represented by the matrix X :

```
[11,12]=size(X);
for s=1:12
    c=(X(1:2:(11-1),s)+X(1:11,s));
    w=(X(1:2:(11-1),s)-X(1:11,s));
    X(1:11,s)=[c w];
end
X=X';

for s=1:11
    c=(X(1:2:(12-1),s)+X(1:12,s));
    w=(X(1:2:(12-1),s)-X(1:12,s));
    X(1:12,s)=[c w];
end
X=X';
```

Comment the code. Describe what will be shown in the upper left corner of X after the code has run. Do the same for the lower left corner of the matrix. What is the connection with the images $(G_0 \otimes G_0)X$, $(G_0 \otimes G_1)X$, $(G_1 \otimes G_0)X$, and $(G_1 \otimes G_1)X$?

3. Implement functions

```
function Y=DWT2HaarImpl(X,m)
function X=IDWT2HaarImpl(Y,m)
```

which implements DWT2 and IDWT2 for the Haar-wavelet. Your functions should call the functions `DWTHaarImpl` and `IDWTHaarImpl` from Section 5.3.

4. In this exercise we will experiment with applying the m -level DWT2 to an image.

- a. Write a function

```
function showDWTlowerHaar(m)
```

which

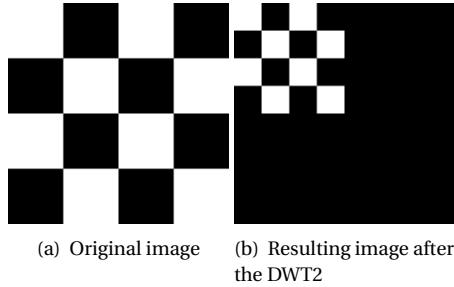


Figure 10.17: A simple image before and after one level of the DWT2. The Haar wavelet was used.

1. reads the image file `lena.png`,
 2. performs an m -level DWT2 to the image samples using the function `DWT2HaarImpl`,
 3. sets all wavelet coefficients representing detail to zero (i.e. keep only wavelet coefficients from $V_0 \otimes V_0$),
 4. performs an IDWT2 on the resulting coefficients using the function `IDWT2HaarImpl`,
 5. displays the resulting image.
- b.** Run the function `showDWTlowerHaar` for different values of m . Describe what you see for different m . degraded? Compare with what you saw with the function `showDCThigher` in Exercise 2, where you performed a DCT on the image samples instead, and set DCT coefficients below a given threshold to zero.
- c.** Do the image samples returned by `showDWTlowerHaar` lie in $[0, 255]$?
5. Repeat Exercise 4, but this time instead keep only wavelet coefficients from the detail spaces. Call the new function `showDWTlowerdifferenceHaar`. What kind of image do you see? Can you recognize the original image in what you see?
 6. In Figure 10.17 we have applied the DWT2 with the Haar wavelet to an image very similar to the one you see in Figure 10.6. You see here, however, that there seems to be no detail components, which is very different from Figure 10.6, even though the images are very similar. Attempt to explain what causes this to happen.
Hint: Compare with Exercise 5.3.5.
 7. In this exercise we will implement DTW2 for the CDF 9/7 wavelet.

a. Implement functions

```
function Y=DWT2Impl153(X,m)
function X=IDWT2Impl153(Y,m)
```

where $DWT2Impl153$ performs the m -level DWT2 on the image given by X , and $IDWT2Impl153$ performs the m -level IDWT2, when the Spline 5/3 wavelet is used. The functions should at each stage call $DWTImpl153$ and $IDWTImpl153$ with $m = 1$, which you implemented in Exercise 8.3.3, and each call to these functions should alter the appropriate upper left submatrix in the coordinate matrix, in the way we have described in this chapter. For simplicity, assume that the image dimensions can be written as appropriate powers of 2.

b. Implement functions

```
function Y=DWT2Impl197(X,m)
function X=IDWT2Impl197(Y,m)
```

which does for the CDF 9/7 wavelet, what the functions in (a) did for the Spline 5/3 wavelet. The same comments apply here, but you should instead call the functions $DWTImpl197$ and $IDWTImpl197$, which you implemented in Exercise 8.3.4.

8. Write functions

```
function showDWTlower53(m)
function showDWTlower97(m)
```

which reimplements the function $showDWTlowerHaar$ from Exercise 4 when the Spline 5/3 wavelet and the CDF 9/7 wavelet are used instead. Look at the result using the different wavelets we have encountered and for different m , using the code from Example 10.16. Can you see any difference from the Haar wavelet? If so, which wavelet gives the best image quality?

9. In this exercise we will change the code in Example 10.16 so that it instead only shows the contribution from the detail spaces.

a. Reimplement the functions you made in Exercise 8 so that they instead show the contribution from the detail spaces. Call the new functions $showDWTlowerdifference53$ and $showDWTlowerdifference97$.

b. In Exercise 5 we implemented a function $showDWTlowerdifferenceHaar$ for looking at the detail/error when the Haar wavelet is used. In the function $showDWTall$ from Example 10.16, replace $showDWTlowerHaar$, $showDWTlower53$, and $showDWTlower97$, with $showDWTlowerdifferenceHaar$, $showDWTlowerdifference53$, and $showDWTlowerdifference97$. Describe the images you see for different m . Try to explain why the images seem to get clearer when you increase m .

10.4 An application to the FBI standard for compression of fingerprint images

In the beginning of the 1990s, the FBI realized that it had a major problem when it came to their archive of fingerprint images. With more than 200 million finger-

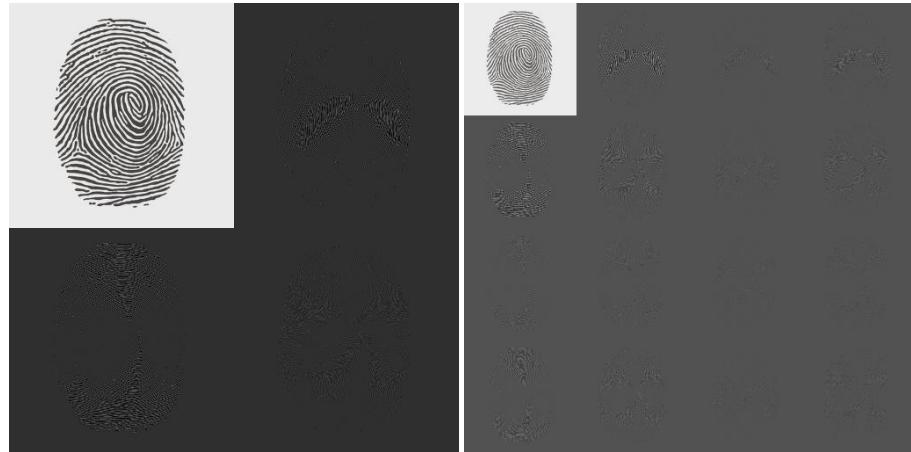


Figure 10.18: A typical fingerprint image.

print records, their digital storage exploded in size, so that some compression strategy needed to be employed. Several strategies were tried, for instance the widely adopted JPEG standard. The problem with JPEG had to do with the blocking artefacts, which we saw in Section 9.4. Among other strategies, FBI chose a wavelet-based strategy due to its nice properties. The particular way wavelets are applied in this strategy is called *Wavelet transform/scalar quantization* (WSQ).

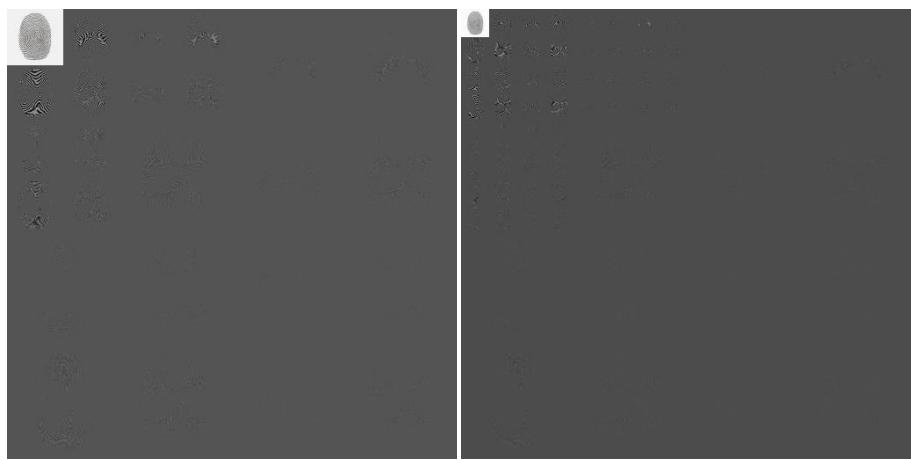
Fingerprint images are a very specific type of images, as seen in Figure 10.18. They differ from natural images by having a large number of abrupt changes. One may ask whether other wavelets than the ones we have used up to now are more suitable for compressing such images. After all, the technique of vanishing moments we have used for constructing wavelets are most suitable when the images display some degree of regularity (such as most natural images do). Extensive tests were undertaken to compare different wavelets, and the CDF 9/7 wavelet used by JPEG2000 turned out to perform very well, also for fingerprint images. One advantage with the choice of this wavelet for the FBI standard is that one then can exploit existing wavelet transformations from the JPEG2000 standard.

Besides the choice of wavelet, one can also ask other questions in the quest to compress fingerprint images: What number of levels is optimal in the application of the DWT2? And, while the levels in a DWT2 (see Figure 10.3) have an interpretation as change of coordinates, one can apply a DWT2 to the other subbands as well. This can not be interpreted as a change of coordinates, but if we assume that these subbands have the same characteristics as the original image, the DWT2 will also help us with compression when applied to them. Let us illustrate how the FBI standard applies the DWT2 to the different subbands. After one stage of the DWT2, we get the image shown in Figure 10.19(a). If we after this also apply a DWT2 to the bands where highpass filters have been applied, we get the following illustrations of the new subbands:



(a) After one stage of the DWT2

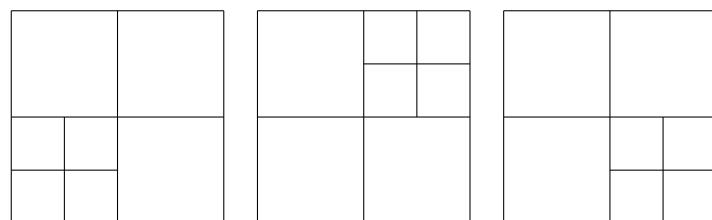
(b) After applying another stage of DWT2 to all bands



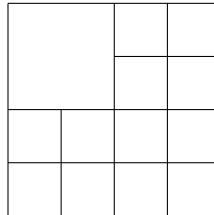
(c) After one stage of the DWT

(d) After one stage of the DWT

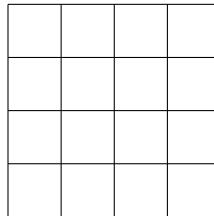
Figure 10.19: The fingerprint image after several DWT's



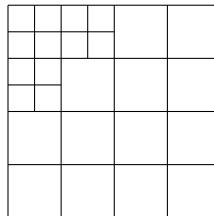
If we apply all these, we get the following subband structure:



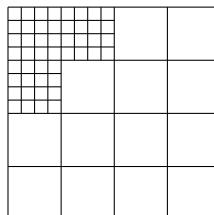
If we also apply the second stage in a DWT2 we arrive at



The resulting image is shown in Figure 10.19(b). Now the FBI standard applies a DWT2 in three of the four resulting subbands, to arrive at the subband structure



The resulting image is shown in Figure 10.19(c). In all the remaining subbands, the DWT2 is now again applied:



The resulting image is shown in Figure 10.19(d). Finally, a DWT2 is again applied, but this time only to the upper left corner. In Figure 10.20 the resulting subband decomposition is shown, together with the resulting image. When establishing the standard for compression of fingerprint images, the FBI chose this subband decomposition. In Figure 10.21 we also show the corresponding low resolution approximation and detail. As can be seen from the subband decomposition, the low-resolution approximation is simply the approximation after a five stage DWT2.

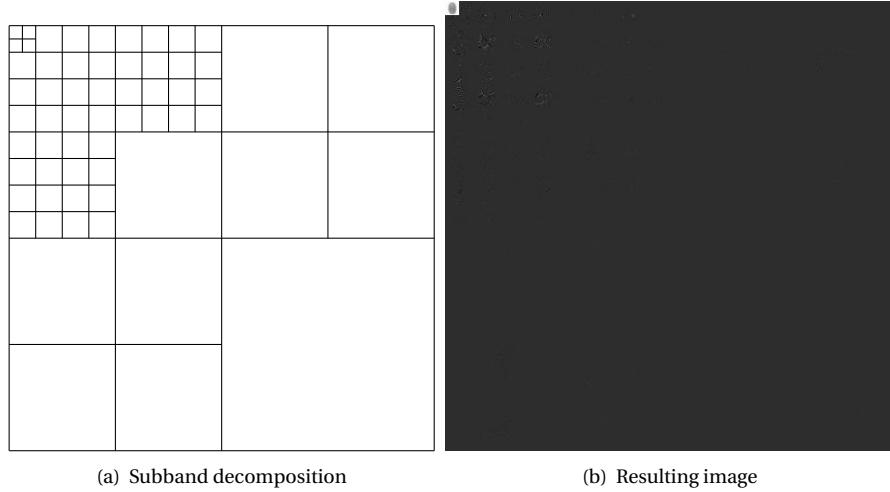


Figure 10.20: The wavelet subband decomposition with the resulting image, as employed by the FBI

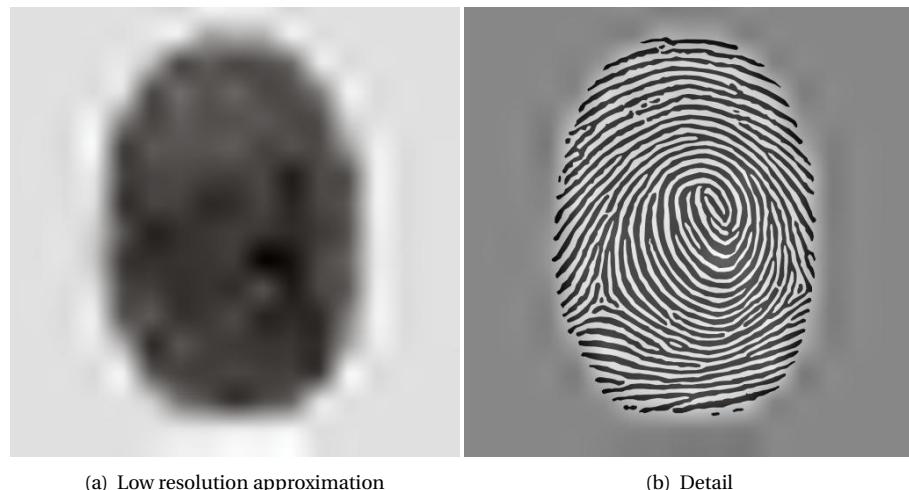


Figure 10.21: The low-resolution approximation and the detail obtained by the FBI standard for compression of fingerprint images, when applied to our sample fingerprint image.

The original JPEG2000 standard did not give the possibility for this type of subband decomposition. This has been added to a later extension of the standard, which makes the two standards more compatible. IN FBI's system, there are also other important parts besides the actual compression strategy, such as *fingerprint*

pattern matching: In order to match a fingerprint quickly with the records in the database, several characteristics of the fingerprints are stored, such as the number of lines in the fingerprint, and points where the lines split or join. When the database is indexed with this information, one may not need to decompress all images in the database to perform matching. We will not go into details on this here.

Exercises for Section 10.4

1. Write code which generates the images shown in figures 10.19, 10.20, and 10.21. Use the functions `DWT2Impl197` and `IDWT2Impl197` from Exercise 10.2. 7 to achieve this.

10.5 Summary

We extended the tensor product construction to functions by defining the tensor product of functions as a function in two variables. We explained with some examples that this made the tensor product formalism useful for approximation of functions in several variables. We extended the wavelet transform to the tensor product setting, so that it too could be applied to images. We also performed several experiments on our test image, such as creating low-resolution images and neglecting wavelet coefficients. We also used different wavelets, such as the Haar wavelet, the Spline 5/3 wavelet, and the CDF 9/7 wavelet. The experiments confirmed what we previously have proved, that wavelets with many vanishing moments are better suited for compression purposes.

The specification of the JPGE2000 standard can be found in [1]. In [33], most details of this theory is covered, in particular details on how the wavelet coefficients are coded (which is not covered here).

One particular application of wavelets in image processing is the compression of fingerprint images. The standard which describes how this should be performed can be found in [11]. In [3], the theory is described. The book [14] uses the application to compression of fingerprint images as an example of the usefulness of recent developments in wavelet theory.

Basic Linear Algebra

This book assumes that the student has taken a beginning course in linear algebra at university level. In this appendix we summarize the most important concepts one needs to know from linear algebra. Note that what is listed here should not be considered as a substitute for such a course: It is important for the student to go through a full course in linear algebra, in order to get good intuition for these concepts through extensive exercises. Such exercises are omitted here.

A.1 Matrices

An $m \times n$ -matrix is simply a set of mn numbers, stored in m rows and n columns. We write a_{kn} for the entry in row k and column n of the matrix A . The zero matrix, denoted $\mathbf{0}$ is the matrix with all zeroes. A square matrix (i.e. where $m = n$) is said to be diagonal if $a_{kn} = 0$ whenever $k \neq n$. The identity matrix, denoted I , or I_n to make the dimension of the matrix clear, is the diagonal matrix where the entries on the diagonal are 1, the rest zeroes. If A is a matrix we will denote the transpose of A by A^T . If A is invertible we denote its inverse by A^{-1} . We say that a matrix A is *orthogonal* if $A^T A = AA^T = I$. A matrix is called sparse if most of the entries in the matrix are zero.

A.2 Vector spaces

A set of vectors V is called a vector space if... We say that the vectors $\{\mathbf{v}_0, \mathbf{v}_1, \dots, \mathbf{v}_{n-1}\}$ are *linearly independent* if, whenever $\sum_{i=0}^{n-1} c_i \mathbf{v}_i = \mathbf{0}$, we must have that all $c_i = 0$. We will say that a set of vectors $\mathcal{B} = \{\mathbf{v}_0, \mathbf{v}_1, \dots, \mathbf{v}_{n-1}\}$ from V is a *basis* for V if the vectors are linearly independent, and span V .

Subspaces of \mathbb{R}^N , and function spaces.

A.3 Inner products and orthogonality

Most vector spaces in this book are inner product spaces. A (real) inner product on a vector space is a binary operation, written as $(\mathbf{u}, \mathbf{v}) \rightarrow \langle \mathbf{u}, \mathbf{v} \rangle$, which fulfills the following properties for any vectors \mathbf{u} , \mathbf{v} , and \mathbf{w} :

1. $\langle \mathbf{u}, \mathbf{v} \rangle = \langle \mathbf{v}, \mathbf{u} \rangle$
2. $\langle \mathbf{u} + \mathbf{v}, \mathbf{w} \rangle = \langle \mathbf{u}, \mathbf{w} \rangle + \langle \mathbf{v}, \mathbf{w} \rangle$
3. $\langle c\mathbf{u}, \mathbf{v} \rangle = c\langle \mathbf{u}, \mathbf{v} \rangle$ for any scalar c
4. $\langle \mathbf{u}, \mathbf{u} \rangle \geq 0$, and $\langle \mathbf{u}, \mathbf{u} \rangle = 0$ if and only if $\mathbf{u} = 0$.

\mathbf{u} and \mathbf{v} are said to be *orthogonal* if $\langle \mathbf{u}, \mathbf{v} \rangle = \mathbf{0}$. In this book we have seen two important examples of inner product spaces. First of all the Euclidean inner product, which is defined by

$$\langle \mathbf{u}, \mathbf{v} \rangle = \sum_{i=0}^{n-1} u_i v_i \quad (\text{A.1})$$

for any \mathbf{u} , \mathbf{v} in \mathbb{R}^n . For functions we have seen examples which are variants of the following form:

$$\langle f, g \rangle = \int f(t)g(t)dt. \quad (\text{A.2})$$

Any set of mutually orthogonal elements are also linearly independent. A basis where all basis vectors are mutually orthogonal is called an *orthogonal basis*. If additionally the vectors all have length 1, we say that the basis is *orthonormal*. If \mathbf{x} is in a vector space with an orthogonal basis $\mathcal{B} = \{\mathbf{v}_k\}_{k=0}^{n-1}$, we can express \mathbf{x} as

$$\frac{\langle \mathbf{x}, \mathbf{v}_0 \rangle}{\langle \mathbf{v}_0, \mathbf{v}_0 \rangle} \mathbf{v}_0 + \frac{\langle \mathbf{x}, \mathbf{v}_1 \rangle}{\langle \mathbf{v}_1, \mathbf{v}_1 \rangle} \mathbf{v}_1 + \cdots + \frac{\langle \mathbf{x}, \mathbf{v}_{n-1} \rangle}{\langle \mathbf{v}_{n-1}, \mathbf{v}_{n-1} \rangle} \mathbf{v}_{n-1}. \quad (\text{A.3})$$

In other words, the weights in linear combinations are easily found when the basis is orthogonal. This is also called the *orthogonal decomposition theorem*.

By the *projection* of a vector \mathbf{x} onto a subspace U we mean the vector $\mathbf{y} = \text{proj}_U \mathbf{x}$ which minimizes the distance $\|\mathbf{y} - \mathbf{x}\|$. If \mathbf{v}_i is an orthogonal basis for U , we have that $\text{proj}_U \mathbf{x}$ can be written by Equation (A.3).

A.4 Coordinates and change of coordinates

If $\mathcal{B} = \{\mathbf{v}_0, \mathbf{v}_1, \dots, \mathbf{v}_{n-1}\}$ is a basis for a vector space, and $\mathbf{x} = \sum_{i=0}^{n-1} x_i \mathbf{v}_i$, we say that $(x_0, x_1, \dots, x_{n-1})$ is the *coordinate vector* of \mathbf{x} w.r.t. the basis \mathcal{B} . We also write $[\mathbf{x}]_{\mathcal{B}}$ for this coordinate vector.

If \mathcal{B} and \mathcal{C} are two different bases for the same vector space, we can write down the two coordinate vectors $[\mathbf{x}]_{\mathcal{B}}$ and $[\mathbf{x}]_{\mathcal{C}}$. A useful operation is to transform the coordinates in \mathcal{B} to those in \mathcal{C} , i.e. apply the transformation which sends $[\mathbf{x}]_{\mathcal{B}}$ to $[\mathbf{x}]_{\mathcal{C}}$. This is a linear transformation, and we will denote the $n \times n$ -matrix of this

linear transformation by $P_{\mathcal{C} \leftarrow \mathcal{B}}$, and call this the *change of coordinate matrix* from \mathcal{B} to \mathcal{C} . In other words, the change of coordinate matrix is defined by requiring that

$$P_{\mathcal{C} \leftarrow \mathcal{B}}[\mathbf{x}]_{\mathcal{B}} = [\mathbf{x}]_{\mathcal{C}}. \quad (\text{A.4})$$

It is straightforward to show that $P_{\mathcal{C} \leftarrow \mathcal{B}} = (P_{\mathcal{B} \leftarrow \mathcal{C}})^{-1}$, so that matrix inversion can be used to compute the change of coordinate matrix the opposite way. It is also straightforward to show that the columns in the change of coordinate matrix can be obtained by expressing the old basis in terms of the new basis, i.e. finding the vectors $[P_{\mathcal{B} \leftarrow \mathcal{C}}(\mathbf{v}_i)]_{\mathcal{C}}$.

If L is a linear transformation between the spaces V and W , and \mathcal{B} is a basis for V , \mathcal{C} a basis for W , we can consider the operation which sends the coordinates of $\mathbf{v} \in V$ in the basis \mathcal{B} to the coordinates of $L\mathbf{v} \in W$ in the basis \mathcal{C} . This is represented by a matrix, called *the matrix of L relative to the bases \mathcal{B} and \mathcal{C}* . Similarly to change of coordinate matrices, the columns of the matrix of L relative to the bases \mathcal{B} and \mathcal{C} are given by $[L(\mathbf{v}_i)]_{\mathcal{C}}$.

A.5 Eigenvectors and eigenvalues

If A is a linear transformation from a vector space to itself, a vector \mathbf{v} is called an *eigenvector* if there exists a scalar λ so that $A\mathbf{v} = \lambda\mathbf{v}$. λ is called the corresponding eigenvalue.

If the matrix A is symmetric, the following hold:

1. The eigenvalues of A are real,
2. the eigenspaces of A are orthonormal,
3. any vector can be decomposed as a sum of eigenvectors from A .

For non-symmetric matrices, these results do not hold in general. But for filters, clearly the second and third property always hold, regardless of whether the filter is symmetric or not.

A.6 Diagonalization

One can show that, for a symmetric matrix, $A = PDP^T$ where D is a diagonal matrix and the eigenvalues of A are the values on the diagonal of D , and P is a matrix where the columns are the eigenvectors of A , with corresponding eigenvalue appearing in the same column in D .

Appendix **B**

Signal processing and linear algebra: a translation guide

This book should not be considered as a standard signal processing textbook. There are several reasons for this. First of all, much signal processing literature is written for people with an engineering background. This book is written for people with a basic linear algebra background. Secondly, the book does not give a comprehensive treatment of all basic signal processing concepts. Signal processing concepts are introduced whenever they are needed to encompass the mathematical exposition. In order to learn more about the different signal processing concepts, the reader can consult many excellent textbooks, such as [28, 2, 25, 32]. The translation guide of this chapter may be of some help in this respect, when one tries to unify material presented here with material from these signal processing textbooks. The translation guide handles both differences in notation between this book and signal processing literature, and topical differences. Most topical differences are also elaborated further in the summaries of the different chapters. The book has adopted most of its notation and concepts from mathematical literature.

B.1 Complex numbers

There are several differences between engineering literature and mathematics. In mathematics literature, i is used for the imaginary complex number which satisfies $i^2 = -1$. In engineering literature, the name j is used instead.

B.2 Functions

What in signal processing are referred to as continuous-time signals, are here referred to as functions. Usually we refer to a function by the letter f , according to the mathematical tradition. The variable is mostly time, represented by the symbol t .

In signal processing, one often uses capital letters to denote a function which is the Fourier transform of another function, so that the Fourier transform of x would be denoted by X . Here we simply denote a periodic function by its Fourier coefficients y_n , and we avoid the CTFT. We use analog filters, however, which also work in continuous time. Analog filters preserve frequencies, and we have used v to denote frequency (variations per second), and not used angular frequency ω . In signal processing literature it is common to jump between the two.

B.3 Vectors

Discrete-time signals, as they are used in signal processing, are here mostly referred to as vectors. To as big extent as possible, we have attempted to keep vectors finite-dimensional. Vectors are in boldface (i.e. \mathbf{x}), but its elements are not in boldface, and with subscripts (i.e. x_n). Superscripts are also used to differ between vectors with the same base name (i.e. $\mathbf{x}^{(1)}$, $\mathbf{x}^{(2)}$ etc.), so that this does not interfere with the vector indices. In signal processing literature the corresponding notation would be x for the signal, and $x[n]$ for its elements, and signals with equal base names could be named like $x_1[n], x_2[n]$.

We have sometimes denoted the Fourier transform of x by $\hat{\mathbf{x}}$, according to the mathematical tradition. More often we have distinguished between a vector and its Discrete Fourier transform by using \mathbf{x} for the first, and \mathbf{y} for the latter. This also makes us distinguish between the input and output to a filter, where we instead use \mathbf{z} for the latter. Much signal processing literature write (capital) X for the DFT of the vector x .

B.4 Inner products and orthogonality

Throughout the book we have defined inner products for functions (for Fourier analysis and wavelets), and we have also used the standard inner product of \mathbb{R}^N . From this we have deduced the orthogonality of several basis functions used in signal processing theory. That the functions are orthogonal, as well as the inner product itself are, however, often not commented on in signal processing literature. As an unfortunate consequence, one has to explain the expression for the Fourier series using other means than the orthogonal decomposition formula and the least squares method. Also, one does not mention that the DFT is a unitary transformation.

B.5 Matrices and filters

Boldface notation is not used for matrices, according to the mathematical tradition. In signal processing, it is not common to formulate matrix equations, such as for the DFT and DCT, or matrix factorizations. Instead one typically writes down each equation, one equation for each row in $\mathbf{y} = A\mathbf{x}$, i.e. not recognizing matrix/vector multiplication. We have stucked to the name filtering operations, but made it clear that this is nothing but a linear transformation with a Toeplitz matrix as its matrix. In

particular, we alternately use the terms filtering and multiplication with a Toeplitz matrix. The characterization of filters as circulant Toeplitz matrices is usually not done in signal processing literature (but see [14]). In this text we allow for matrices also to be of infinite dimensions, expanding on the common use in linear algebra. When infinite dimensions are assumed, infinite in both directions is assumed. Matrices are scaled if necessary to make them unitary, in particular the DCT and the DFT. This scaling is usually not done in signal processing literature.

Representing a filter in terms of a finite matrix and restriction of a filter to a finite signal. This is usually omitted in signal processing literature.

One of the most important statements in signal processing is that convolution in time is equivalent to multiplication in frequency. We have presented a compelling interpretation of this in linear algebra terms. Since the frequency response simply are eigenvalues of the filter, and convolution simply is matrix factorization, multiplication in frequency simply means to multiply two diagonal matrices to obtain the frequency response of the product. Moreover, the Fourier basis vectors can be interpreted as eigenvectors.

B.6 Convolution

While we have defined the concept of convolution, readers familiar with signal processing may have noticed that this concept has not been used much. The reason is that we have wanted to present convolution as a matrix multiplication (to adapt to mathematical tradition), and that we have used the concept of filtering often instead. In signal processing literature one defines convolution in terms of vectors of infinite length. We have avoided this, since in practice vectors always need to be truncated to finite lengths. Due to this, we also have analyzed how a finite vector may be turned into a periodic vector (periodic or symmetric extension), and how this affects our analysis. Also we have concentrated on FIR-filters, and this makes us avoid convergence issues. Note that we do not present matrix multiplication as a method of implementing filtering, due to the special structure of this operation. We do not suggest other methods for implementation than applying the convolution formula in a brute-force way, or factoring the filter in simpler components.

B.7 Polyphase factorizations and lifting

In signal processing literature, it is not common to associate polyphase components with matrices, but rather with Laurent polynomials generated from the corresponding filter. The Laurent polynomial is nothing else than the Z -transform of the associated filter. Associating polyphase components with blocks in a block matrix makes this book fit with block matrix methods in linear algebra textbooks.

The polyphase factorization serves two purposes in this book. Firstly, the lifting factorization (as used for wavelets) is derived from it, and put in a linear algebra framework as a factorization into sparse matrices, similarly to the FFT factorization. Thereby it fits together with many of the matrix factorization results from classical

linear algebra, where also sparsity is what makes the factorization good for computation.

Secondly, the polyphase factorization of the filter bank transforms in the MP3 standard are derived (also as a sparse matrix factorization), and from this it is apparent what properties to put on the prototype filters in order to obtain useful transforms. In fact, from this factorization it became apparent that the MP3 filter bank transforms could be expressed in terms of alternative QMF filter banks (i.e. $M = 2$).

These two topics (lifting and the MP3 filter bank transform polyphase factorization) are usually not presented in a unified way in textbooks. we see here that there is a big advantage of doing this, since the second can build on theory from the first.

B.8 Transforms in general

In signal processing, one often refers to the forward and reverse filter bank transforms as analysis and synthesis, respectively, and for obvious reasons. In mathematical literature, one instead often use the term change of coordinates in a wavelet setting. These terms are not normally used in mathematical literature, where the term basis vectors/change of coordinate matrices would be used instead. Also, the output from a forward filter bank transform is often referred to as the *transformed vector*, and the result we get when we apply the reverse filter bank transform to this is called the *reconstructed vector*.

This exposition takes extra care in presenting how the DCT is derived naturally from the DFT. In particular both the DFT and the DCT are derived as matrices of eigenvectors for finite-dimensional filters. The DCT is derived from the DFT in that one restricts to a certain subset of vectors. The orthogonality of these matrices follows from the orthogonality of distinct eigenspaces.

B.9 Perfect reconstruction systems

The term biorthogonality is not used to describe a mutual property of the filters of wavelets. Biorthogonality corresponds simply to two matrices being inverses of one another. For the same reason, the term perfect reconstruction is not used much. Much wavelet theory refer to a property called delay normalization. This term has been avoided by mostly considering wavelets with symmetric filters, for which delay-normalization is automatic. There are, however, many examples of wavelets where this term is important.

B.10 Z-transform and frequency response

The Z -transform and the frequency response are much used in signal processing literature, and are important concepts for filter design. We have deliberately dropped the Z -transform. Due to this, much signal processing has of course been left out, since placements of poles and zeroes are not performed outside or inside the unit

circle, since the frequency response only captures the values on the unit circle. Placement of poles and circles is perhaps the most-used design feature in filter design. The focus here is on implementing filters, not designing them, however.

In signal processing literature, the DTFT and the Z-transform is used, assuming that the inputs and outputs are vectors of infinite length. In practice of course, some truncation is needed, since only finite-dimensional arithmetic is performed by the computer. How this truncation is to be done without affecting the computations is thus never mentioned in signal processing, although it is always performed somehow. This exposition shows that this truncation can be taken as part of the theory, without seriously affecting the results.

Bibliography

- [1] ISO/IEC FDIS 15444-1. *JPEG2000 Part 1 final draft international standard*, 2000.
- [2] A. Ambardar. *Digital Signal Processing: A Modern Introduction*. Cengage Learning, 2006.
- [3] C. M. Brislawn. Fingerprints go digital. *Notices of the AMS*, 42(11):1278–1283, 1995.
- [4] B. A. Cipra. The best of the 20th century: Editors name top 10 algorithms. *SIAM News*, 33(4). <http://www.uta.edu/faculty/rcli/TopTen/topten.pdf>.
- [5] A. Cohen, I. Daubechies, and J-C. Feauveau. Biorthogonal bases of compactly supported wavelets. *Communications on Pure and Appl. Math.*, 45(5):485–560, June 1992.
- [6] J. W. Cooley and J. W. Tukey. An algorithm for the machine calculation of complex fourier series. *Math. Comput.*, 19:297–301, 1965.
- [7] A. Croisier, D. Esteban, and C. Galand. Perfect channel splitting by use of interpolation/decimation/tree decomposition techniques. *Int. Conf. on Information Sciences and Systems*, pages 443–446, August 1976.
- [8] I. Daubechies. Orthonormal bases of compactly supported wavelets. *Communications on Pure and Applied Mathematics*, 41(7):909–996, October 1988.
- [9] I. Daubechies. *Ten lectures on wavelets*. CBMS-NSF conference series in applied mathematics. SIAM Ed., 1992.
- [10] P. Duhamel and H. Hollmann. 'split-radix' FFT-algorithm. *Electronic letters*, 20(1):14–16, 1984.
- [11] FBI. WSQ gray-scale fingerprint image compression specification. Technical report, IAFIS-IC, 1993.
- [12] G. B. Folland. *Real analysis. Modern techniques and their applications*. John Wiley and sons, 1984.

- [13] M. W. Frazier. *An Introduction to Wavelets Through Linear Algebra*. Springer, 1999.
- [14] M. W. Frazier. *An Introduction to Wavelets Through Linear Algebra*. Springer, 1999.
- [15] M. Frigo and S. G. Johnson. The design and implementation of FFTW3. *Proceedings of the IEEE*, 93(2):216–231, 2005.
- [16] R. C. Gonzalez, R. E. Woods, and S. L. Eddins. *Digital Image Processing Using MATLAB*. Gatesmark publishing, 2009.
- [17] ISI/IEC. Information technology - coding of moving pictures and associated audio for digital storage media at up to about 1.5 mbit/s. Technical report, ISO/IEC, 1993.
- [18] S.G Johnson and M. Frigo. A modified split-radix FFT with fewer arithmetic operations. *IEEE Transactions on Signal Processing*, 54, 2006.
- [19] J.D. Johnston. A filter family designed for use in quadrature mirror filter banks. *Proc. Int. Conf. Acoust. Speech and Sig. Proc.*, pages 291–294, 1980.
- [20] D. C. Lay. *linear algebra and its applications (4th edition)*. Addison Wesley, 2011.
- [21] S. Mallat. *A wavelet tour of signal processing*. Academic Press, 1998.
- [22] C. D. Meyer. *Matrix Analysis and Applied Linear Algebra*. siam, 2000.
- [23] Knut Mørken. UIO, 2013.
- [24] P. Noll. MPEG digital audio coding. *IEEE Signal processing magazine*, pages 59–81, September 1997.
- [25] A. V Oppenheim and R. W. Schafer. *Discrete-Time Signal Processing*. Prentice-Hall, 1989.
- [26] D. Pan. A tutorial on MPEG/audio compression. *IEEE Multimedia*, pages 60–74, Summer 1995.
- [27] W. B. Pennebaker and J. L. Mitchell. *JPEG Still Image Data Compression Standard*. Van Nostrand Reihnold, 1993.
- [28] J. G. Proakis and D. G. Manolakis. *Digital Signal Processing. Principles, algorithms, and applications. Fourth edition*. Pearson, 2007.
- [29] C. M. Rader. Discrete Fourier transforms when the number of data samples is prime. *Proceedings of the IEEE*, 56:1107–1108, June 1968.
- [30] T. A. Ramstad, S. O. Aase, and J. H. Husøy. *Subband Compression of Images: Principles and Examples: Principles and Examples*, volume 6. Elsevier, 1995.

- [31] C. E. Shannon. Communication in the presence of noise. *Proc. Institute of Radio Engineers*, 37(1):10–21, Jan. 1949.
- [32] P. Stoica and R. Moses. *Spectral Analysis of Signals*. Prentice Hall, 2005.
- [33] D. S. Taubman and M. W. Marcellin. *JPEG2000. Image compression. Fundamentals, standards and practice*. Kluwer Academic Publishers, 2002.
- [34] M. Vetterli and J. Kovacevic. *Wavelets and subband coding*. Prentice Hall, 1995.
- [35] M. Vetterli and H. J. Nussbaumer. Simple FFT and DCT algorithms with reduced number of operations. *Signal processing*, 6:267–278, 1984.
- [36] S. Winograd. On computing the discrete Fourier transform. *Math. Comp.*, 32:175–199, 1978.
- [37] R. Yavne. An economical method for calculating the discrete Fourier transform. *Proc. AFIPS Fall Joint Computer Conf.*, 33:115–125, 1968.

Nomenclature

$[x]_{\mathcal{B}}$	Coordinate vector of x relative to the basis \mathcal{B} .
\check{x}	Symmetric extension of a vector
\check{f}	Symmetric extension of the function f
\hat{x}	DFT of the vector x
$\lambda_s(v)$	Frequency response of an analog filter
$\lambda_S(\omega)$	Continuous frequency response of a digital filter
$\lambda_{S,n}$	Vector frequency response of a digital filter
$\langle u, v \rangle$	Inner product
v	Frequency
ω	Angular frequency
\otimes	Tensor product
ϕ	Scaling function
ψ	Mother wavelet
$\tilde{\phi}$	Dual scaling function
$\tilde{\psi}$	Dual mother wavelet
$c_{m,n}$	Wavelet low resolution coefficients
$w_{m,n}$	Wavelet detail coefficients
$x * y$	Convolution of vectors
$x^{(e)}$	Vector of even samples

$\mathbf{x}^{(o)}$	Vector of odd samples
$\boldsymbol{\phi}_m$	Basis for V_m
$\boldsymbol{\psi}_m$	Basis for W_m
A^H	Conjugate transpose of a matrix
A^T	Transpose of a matrix
A^{-1}	Inverse of a matrix
D_N	$N \times N$ -DCT matrix
E_d	Filter which delays with d samples
F_N	$N \times N$ -Fourier matrix
f_N	N 'th order Fourier series of f
f_s	Sampling frequency
$l(S)$	Length of a filter
N	Number of points in a DFT/DCT
$O(f(\mathbf{x}))$	Order of a function
$O(N)$	Order of an algorithm
$P_{\mathcal{C} \leftarrow \mathcal{B}}$	Change of coordinate matrix from \mathcal{B} to \mathcal{C} .
S^f	Matrix with the columns reversed
T	Period of a function
T_s	Sampling period
$U \oplus V$	Direct sum of vector spaces
V_m	Resolution space
$V_{N,T}$	N 'th order Fourier space
W_m	Detail space
$W_m^{(0,1)}$	Resolution m Complementary wavelet space, LH
$W_m^{(1,0)}$	Resolution m Complementary wavelet space, HL
$W_m^{(1,1)}$	Resolution m Complementary wavelet space, HH
\mathcal{C}_m	Reordering of $(\boldsymbol{\phi}_{m-1}, \boldsymbol{\psi}_{m-1})$
\mathcal{D}_m	Reordering of $\boldsymbol{\phi}_m$

$\mathcal{D}_N = \{\mathbf{d}_0, \mathbf{d}_1, \dots, \mathbf{d}_{N-1}\}$ N -point DCT basis for \mathbb{R}^N

$\mathcal{D}_{N,T}$ Order N real Fourier basis for $V_{N,T}$

$\mathcal{E}_N = \{\mathbf{e}_0, \mathbf{e}_1, \dots, \mathbf{e}_{N-1}\}$ Standard basis for \mathbb{R}^N

$\mathcal{F}_N = \{\boldsymbol{\phi}_0, \boldsymbol{\phi}_1, \dots, \boldsymbol{\phi}_{N-1}\}$ Fourier basis for \mathbb{R}^N

$\mathcal{F}_{N,T}$ Order N complex Fourier basis for $V_{N,T}$

Mathematics index

- AD conversion, 42
- algebra, 87
- Alias cancellation, 221, 240
- Aliasing, 220, 239
- analysis, 15
 - equations, 15
- Analysis filter components of a forward filter bank transform, 239
- Angular frequency, 94
- Arithmetic operation count
 - DCT, 147
 - DFT direct implementation, 54
 - FFT, 76
 - Lifting, 278
 - revised DCT, 152
 - revised FFT, 152
 - symmetric filters, 142
 - with tensor products, 321
- bandpass filter, 106
- Basis
 - reordering of ϕ_m , 269
 - DCT, 133
 - for V_m , 164
 - for $V_{N,T}$, complex, 22
 - for $V_{N,T}$, real, 15
 - for W_m , 174
 - Fourier, 50
 - reordering of $\phi_{m-1} \oplus \psi_{m-1}$, 212
- basis, 361
- bit rate, 42
- Bit-reversal, 72
- block diagonal matrices, 212
- block matrix, 71
- Blocks, 314
- cascade algorithm, 227
- Change of coordinate matrix, 363
- Change of coordinates, 363
 - in tensor product, 326
- Channel, 239
- Complex Fourier coefficients, 23
- Condition
 - alias cancellation, 222
 - perfect reconstruction, 222
- Conjugate transpose, 52
- continuous sound, 3
- convolution
 - vectors, 101
- coordinate matrix, 326
- Coordinate vector, 362
- Cosine matrices, 134
- Cosine matrix inverse
 - type I, 232
 - type II, 134
 - type III, 134, 138
- critical sampling, 214
- DCT
 - I, 232
- DCT basis, 133
- DCT coefficients, 133
- DCT matrix, 133
- DCT-I factorization, 232
- DCT-II factorization, 133
- DCT-III factorization, 133
- DCT-IV factorization, 138
- Detail space, 173
- detail space, 168
- DFT coefficients, 51
- Diagonalization
 - with F_N , 87
- digital
 - sound, 3, 41
- digital filter, 87
- Direct sum
 - canonical basis, 175
 - linear transformations, 181
 - vector spaces, 168
- Dirichlet conditions, 15
- Discrete Cosine transform, 133
- Discrete Wavelet Transform, 175
- downsampling, 214
- Dual mother wavelet, 201
- Dual scaling function, 200

eigenvalue, 363
 eigenvector, 363
 elementary lifting matrix
 even type, 272
 odd type, 272
 error-resilient, 314

 FFT
 twiddle factors, 77
 FFT factorization, 72
 Filter
 bandpass, 106
 echo, 104
 highpass, 106
 ideal highpass, 107
 ideal lowpass, 107
 length, 101
 linear phase, 131
 lowpass, 106
 moving average, 105
 MP3 standard, 108
 time delay, 104
 Filter bank, 238
 Cosine-modulated, 242
 Filter bank transform, 238
 Filter coefficients, 85
 flop count, 81
 Forward filter bank transform, 238
 in a wavelet setting, 216
 Fourier analysis, 11
 Fourier coefficients, 13
 Fourier domain, 15
 Fourier matrix, 51
 Fourier series, 13
 square wave, 16
 triangle wave, 18
 Fourier space, 13
 Fourier transform
 discrete, 51
 Frequency domain, 15
 Frequency response
 analog filter, 35
 continuous, 94
 vector, 87

 Haar wavelet, 177

 highpass filter, 106

 ideal highpass filter, 107
 ideal lowpass filter, 107
 IDFT, 53
 IMDCT, 139
 impulse response, 88
 Inner product
 of functions in a Fourier setting, 12
 of functions in a tensor product setting, 334
 of functions in a wavelet setting, 162
 of vectors, 50
 interpolating polynomial, 59
 interpolation formula, 68
 ideal
 periodic functions, 67
 Inverse Discrete Wavelet Transform, 175

 JPEG
 standard, 313
 JPEG2000
 lossless compression, 259
 lossy compression, 261
 standard, 259

 least square error, 12
 length of a filter, 101
 lifting factorization, 274
 Linear phase filter, 131
 linearly independent, 361
 lowpass filter, 106
 LTI filters, 119

 matrix of a linear transformation relative to bases, 363
 MDCT, 138
 mother wavelets, 171
 MP3
 and the DCT, 154
 FFT, 63
 filters, 108
 standard, 36
 window, 98

MP3 standard
partial calculation, 241

MP3 standard
matrixing, 241
windowing, 241

MRA-matrix, 213

multiresolution analysis, 184

multiresolution model, 160

Near-perfect reconstruction, 221

Order N complex Fourier basis for $V_{N,T}$,
22

Order of a function, 75

Order of an algorithm, 75

Orthogonal
basis, 362
matrix, 361
vectors, 362

Orthogonal decomposition theorem,
362

Orthonormal
basis, 362
MRA, 184

Orthonormal wavelets, 228

Parallel computing
with tensor products, 321
with the DCT, 147
with the DWT, 314
with the FFT, 76
with the revised FFT, 152

Perfect reconstruction, 221

Perfect reconstruction filter bank, 240

Phase distortion, 221

Polyphase
component of a vector, 76

Polyphase components, 268

Polyphase representation, 268

projection, 362

psycho-acoustic model, 37

pure digital tone, 50

pure tone, 7

QMF filter banks, 227
Alternative definition, 228

Classical definition, 227

Resolution space, 163

Reverse filter bank transform
in a wavelet setting, 216

Reverse filter bank transforms, 239

samples, 42

sampling, 42
frequency, 42
period, 42
rate, 42

scaling function, 164, 184

separable extension, 333

sound channel, 44

Sparse matrix, 361

square wave, 9

Standard
JPEG, 313
JPEG2000, 259
MP3, 36

subband
HH, 342
HL, 342
LH, 341
LL, 341

Subband coding, 239

Subband samples of a filter bank trans-
form, 239

Symmetric
vector, 127

Symmetric extension
of function, 29
used by the DCT, 126
used by wavelets, 187

synthesis, 15
equation, 15
vectors, 51

Synthesis filter components of a reverse
filter bank transform, 239

tensor product, 315
of function spaces, 334
of functions, 333
of matrices, 317
of vectors, 316

Tiles, 314
time domain, 15
time-invariant, 118
Toeplitz matrix, 84
 circulant, 84
triangle wave, 10

Unitary matrix, 52
upsampling, 215

Vector space
 of symmetric vectors, 126

wavelet coefficients, 175
Wavelet detail coefficients, 174
Wavelet low resolution coefficients, 174
Wavelets
 Alternative piecewise linear, 204
 CDF 9/7, 260
 Haar, 171
 Orthonormal, 262
 Piecewise constant, 171
 Piecewise linear, 196
 Spline, 257
 Spline 5/3, 259
window, 97
 Hamming, 98
 Hanning, 100
 in the MP3 standard, 98
 rectangular, 97

Index for MATLAB commands

audioplayer, 44

conv, 101, 116

dct, 135

double, 299

fft, 74

idct, 135

ifft, 74

imread, 298

imshow, 299

imwrite, 299

play, 43

playblocking, 43

rand, 47

uint8, 299

wavread, 44

wavwrite, 44