

Wavelet-based numerical methods for the solution of the Nonuniform Multiconductor Transmission Lines

*Original*

Wavelet-based numerical methods for the solution of the Nonuniform Multiconductor Transmission Lines / GRIVET TALOCIA, Stefano. - (1998). [10.6092/polito/porto/2497973]

*Availability:*

This version is available at: 11583/2497973 since:

*Publisher:*

Politecnico di Torino

*Published*

DOI:10.6092/polito/porto/2497973

*Terms of use:*

Altro tipo di accesso

This article is made available under terms and conditions as specified in the corresponding bibliographic description in the repository

*Publisher copyright*

(Article begins on next page)

POLITECNICO DI TORINO  
Facoltà di Ingegneria

Corso di Dottorato in  
Ingegneria Elettronica e delle Comunicazioni

Tesi di Dottorato

**Wavelet-based numerical methods for the  
solution of the Nonuniform Multiconductor  
Transmission Lines**

Metodi numerici basati sulle onde per la soluzione  
delle linee multiconduttore non uniformi

**Stefano Grivet Talocia**

Coordinatore  
Prof. I. Montrosset

Tutore  
Prof. F. Canavero

X ciclo



# Acknowledgments

I wish to thank Prof. Canavero, the advisor of this thesis, for its encouragement and support throughout the development of this work. I am also particularly grateful to Prof. Canuto and Prof. Tabacco. They gave me a strong mathematical support in understanding wavelet theory and a continuous help in overcoming the difficulties encountered during my work. Prof. Tabacco helped also to improve the quality of this document by carefully reading and correcting the manuscript. Many thanks also to Dott. Levaggi, who provided much help in developing some of the software related to the construction of wavelets on bounded domains.



# Contents

<b>Sommario</b>	<b>iii</b>
<b>Introduction</b>	<b>vii</b>
<b>Mathematical notations</b>	<b>xi</b>
<b>1 The Time Domain Space Expansion method</b>	<b>1</b>
1.1 Mathematical formulation . . . . .	2
1.1.1 Frequency domain analysis . . . . .	6
1.1.2 Incident field excitation . . . . .	7
1.2 Piecewise linear approximation . . . . .	8
1.3 Numerical examples . . . . .	11
1.3.1 The exponential line . . . . .	11
1.3.2 Three-conductor PCB . . . . .	14
1.3.3 Nonparallel wires above a ground plane . . . . .	16
<b>2 An introduction to multilevel decompositions</b>	<b>21</b>
2.1 Multilevel representation of functions: a primer . . . . .	22
2.2 Multilevel decompositions: the abstract setting . . . . .	32
<b>3 Biorthogonal decomposition of <math>L^2(\mathbb{R})</math></b>	<b>35</b>
3.1 The basic axioms . . . . .	35
3.2 Scaling function spaces in $\mathbb{R}$ . . . . .	37
3.3 Wavelet spaces in $\mathbb{R}$ . . . . .	40
3.3.1 Vanishing moments and polynomials reproduction . . . .	43
3.3.2 Compactly supported wavelets . . . . .	44
3.4 Examples of wavelets . . . . .	45
3.4.1 The Haar wavelets . . . . .	45
3.4.2 The Daubechies wavelets . . . . .	46
3.4.3 Biorthogonal spline wavelets . . . . .	47
<b>4 Biorthogonal decomposition on bounded domains</b>	<b>53</b>
4.1 Scaling function spaces in $\mathbb{R}^+$ . . . . .	54
4.1.1 The refinement equation . . . . .	57
4.1.2 The Gramian matrix $X$ . . . . .	61

4.1.3	Boundary value preserving biorthogonalization . . . . .	63
4.1.4	Other polynomial bases . . . . .	65
4.1.5	Approximation properties . . . . .	68
4.2	Wavelet spaces in $\mathbb{R}^+$ . . . . .	69
4.2.1	Projection operators . . . . .	72
4.2.2	Wavelet filters . . . . .	74
4.2.3	The Gramian matrix $Y$ . . . . .	76
4.2.4	Boundary adaption of wavelets . . . . .	77
4.3	An example . . . . .	77
4.3.1	The case $L = 3, \tilde{L} = 5$ . . . . .	78
4.4	The unit interval . . . . .	89
4.4.1	Scaling function spaces . . . . .	89
4.4.2	Wavelet spaces . . . . .	95
4.4.3	Wavelet analysis and synthesis . . . . .	97
4.4.4	An example (continued) . . . . .	100
4.5	Wavelet approximations . . . . .	103
4.5.1	Linear approximations . . . . .	103
4.5.2	Nonlinear approximations . . . . .	104
<b>5</b>	<b>Integrals of refinable functions</b>	<b>115</b>
5.1	Unbounded domains . . . . .	115
5.1.1	Refinement levels reduction . . . . .	117
5.1.2	The eigenvector equation . . . . .	118
5.1.3	An example: the scalar case . . . . .	121
5.1.4	Inner products . . . . .	122
5.2	Bounded domains . . . . .	124
5.2.1	Inner products . . . . .	125
5.2.2	Expansion of discontinuous functions . . . . .	128
<b>6</b>	<b>TDSE with scaling functions and wavelets</b>	<b>131</b>
6.1	TDSE with scaling functions . . . . .	131
6.1.1	The exponential line . . . . .	133
6.2	TDSE with wavelets . . . . .	137
6.2.1	The uniform line . . . . .	139
6.2.2	The exponential line . . . . .	145
6.3	Adaptive TDSE . . . . .	147
6.3.1	The general formulation . . . . .	147
6.3.2	Time-explicit TDSE with wavelets . . . . .	152
6.3.3	Applications . . . . .	158
	<b>Conclusions</b>	<b>163</b>
<b>A</b>	<b>Index of symbols</b>	<b>165</b>

# Sommario

La simulazione delle interconnessioni elettriche costituisce un passo molto importante per l'analisi ed il progetto di sistemi elettronici. Infatti, effetti parassiti quali diafonia ed accoppiamenti elettromagnetici non possono essere trascurati poichè influiscono sensibilmente sul funzionamento del sistema. L'uso di segnali a frequenza sempre più elevata da un lato accentua questi fenomeni, e dall'altro rende necessario l'utilizzo di modelli circuitali a parametri distribuiti.

Il modello delle linee multiconduttore (MTL) è utilizzato comunemente per la simulazione delle interconnessioni elettriche. Il modello si basa sull'ipotesi che la sezione trasversale della struttura sia piccola rispetto alla lunghezza d'onda del campo elettromagnetico circostante ed invariante per traslazioni nella direzione di propagazione dei segnali. In queste condizioni il modo di propagazione fondamentale della struttura è il modo quasi-TEM.

Molte interconnessioni di interesse pratico sono caratterizzate da una sezione non uniforme. Tipici esempi sono le linee per adattamento d'impedenza o i fasci di cavi in strutture complesse quali automobili o aeroplani. In questi casi il modello MTL non può essere applicato direttamente. Quando però la sezione è elettricamente piccola, i campi elettrici e magnetici hanno una componente dominante nel piano trasversale, e quindi soddisfano le ipotesi di modo quasi-TEM. Per questo motivo, è possibile utilizzare il modello delle linee di trasmissione multiconduttore non uniformi (NMTL) per simulare il comportamento elettrico della struttura. Questo modello introduce una variazione longitudinale nei parametri per unità di lunghezza, preservando la struttura delle equazioni. Non è quindi necessario ricorrere a simulazioni di tipo full-wave, che richiedono potenzialità di calcolo molto elevate.

Un certo numero di tecniche per la soluzione delle equazioni NMTL nel dominio del tempo o della frequenza è stato presentato nella letteratura scientifica. Tecniche nel dominio della frequenza portano a soluzioni analitiche in alcuni casi molto semplici, ma possono anche essere utilizzate per risolvere strutture arbitrarie mediante approssimazioni uniformi a tratti. La risposta nel tempo è calcolabile in un secondo tempo mediante FFT inversa. In ogni caso, ciò che si ottiene è una risposta periodicizzata e non una vera risposta in transitorio. È pertanto necessario utilizzare un numero molto elevato di frequenze per ottenere una soluzione nel tempo in cui tutti i fenomeni transitori, quali riflessioni multiple alle terminazioni, siano estinti. Questo numero può diventare mol-





dettaglio il caso di funzioni locali lineari a tratti. Il metodo è successivamente validato risolvendo alcuni problemi di cui si conosce una soluzione analitica oppure una soluzione approssimata.

Il capitolo 2 introduce le ondivine mediante due approcci differenti. Da un lato le proprietà più importanti delle ondivine vengono illustrate mediante il semplice esempio della decomposizione di Haar. Dall'altro viene descritto un quadro astratto molto più generale di cui la decomposizione di Haar costituisce un caso particolare. In questo quadro astratto si collocano anche le particolarizzazioni dei due capitoli successivi.

Il capitolo 3 riassume le proprietà fondamentali delle ondivine biortogonali definite sulla retta reale. Esse permettono di caratterizzare approssimazioni di ordine arbitrario tramite espansioni in termini di funzioni di base (*funzioni di scala*) ottenute mediante semplice traslazione di una singola funzione di partenza. Vengono poi definiti gli spazi di ondivine, che permettono di raffinare una data approssimazione aggiungendo dettagli via via più fini. Anche questi dettagli possono essere caratterizzati mediante sovrapposizione di funzioni di base (*ondivine*) ottenute per traslazione di una funzione di partenza. Vengono mostrati alcuni esempi di funzioni di scala ed ondivine biortogonali, fra cui le cosiddette B-splines biortogonali.

Le proprietà elencate nel capitolo 3 costituiscono il punto di partenza per la costruzione di basi di ondivine definite sull'intervallo unitario. Non è infatti possibile utilizzare funzioni di base invarianti per traslazioni nell'approssimazione di funzioni aventi dominio su intervalli limitati, quali la soluzione delle equazioni NMTL. È quindi necessario ridefinire spazi di approssimazioni e di dettagli mediante funzioni di base modificate e definite su domini limitati. Questo è l'argomento principale del capitolo 4.

Nel capitolo 4 viene dettagliata la procedura che porta alla definizione di ondivine biortogonali su domini limitati ottenute dalle funzioni B-splines. Questo capitolo ha una natura tecnica nella prima parte, in quanto è necessaria una definizione matematica rigorosa delle funzioni di scala e delle ondivine modificate per garantire l'ordine di approssimazione voluto e per preservare la struttura degli spazi funzionali corrispondenti. Un esempio viene illustrato nel dettaglio per riassumere le caratteristiche delle funzioni di base costruite. Successivamente le basi di funzioni di scala e di ondivine vengono utilizzate per definire approssimazioni adattate. Viene introdotto un criterio automatico per definire il minimo numero di coefficienti necessari per rappresentare in modo adattato una funzione anche singolare, garantendo il controllo dell'errore di approssimazione. Questo criterio si basa sull'eliminazione selettiva del contributo di alcune funzioni di base quando i loro coefficienti sono piccoli. La struttura delle funzioni di base definite nella prima parte del capitolo garantisce, da un punto di vista teorico, l'efficienza di queste rappresentazioni.

Il capitolo 5 affronta il problema del calcolo di integrali con funzioni di scala ed ondivine. Questo è un passo necessario per poter utilizzare queste funzioni per

la soluzione delle equazioni NMTL mediante il metodo TDSE. Viene illustrata la procedura che permette il calcolo degli integrali senza dover ricorrere a formule di quadratura numerica. Questo calcolo è ridotto alla determinazione di un autovettore di una matrice costruita a partire dai filtri che caratterizzano le funzioni di scala e le onde prescelte. Viene anche illustrata una procedura per il calcolo di prodotti interni fra una generica funzione e una funzione di scala.

Il capitolo 6 utilizza le basi di onde B-splines biortogonali insieme al metodo TDSE per la soluzione delle equazioni NMTL. Viene definito un metodo adattativo che permette di determinare la soluzione utilizzando il minimo numero di coefficienti strettamente necessari per controllare l'errore di approssimazione. Ciò porta alla definizione di algoritmi ottimizzati in grado di gestire automaticamente la presenza di eventuali singolarità nella soluzione. Viene mostrato mediante alcuni esempi che è possibile risolvere le equazioni utilizzando un sottoinsieme molto piccolo dei coefficienti delle onde senza perdere precisione nella rappresentazione della soluzione.

In conclusione, in questo lavoro si affrontano due problematiche distinte. Da un lato viene introdotta una classe di metodi numerici denominati TDSE per la soluzione delle equazioni delle linee di trasmissione multiconduttore non uniformi. Dall'altro viene definita una classe di spazi di onde su domini limitati. L'uso di queste onde come funzioni di espansione e di test nel metodo TDSE porta alla definizione di algoritmi ad alta adattatività che permettono il calcolo della soluzione ad un basso costo computazionale e senza perdita di precisione.

# Introduction

The simulation of electrical interconnects has become an extremely important step for the analysis and design of electronic systems. In fact, as the clock frequencies of digital systems increase, structures usually modeled with lumped elements are no longer electrically small, and must be treated as distributed circuits. Parasitic effects like crosstalk and electromagnetic coupling cannot be disregarded anymore, because they can seriously affect the overall performance of the system. This is especially relevant for electrical interconnects, which provide the basic link between different devices and parts of a system, or even different systems.

The Multiconductor Transmission Lines (MTL) model [7] is commonly used for the simulation of practical interconnects. This model assumes a small cross-section with respect to the largest wavelength in the system and quasi-TEM fields in the surrounding of the structure. This is true when the cross-section is translation-invariant in the direction of propagation of the signals.

Many interconnections of practical interest are characterized by cross-sections which are not translation-invariant. Examples can be impedance matching networks or cables in complex structures, like automobiles or aeroplanes. In these cases the MTL model is not appropriate. However, as long as the cross-section remains electrically small, the electric and magnetic fields can be assumed to have a dominant transversal component, i.e., satisfy the quasi-TEM mode of propagation. In this cases, the Nonuniform Multiconductor Transmission Lines (NMTL) model can be used to predict the electrical behavior of the interconnect. This model introduces a longitudinal variation in the per-unit-length parameters, by leaving the structure of the equations unchanged. Consequently, the simulation of NMTL equations does not require a full-wave transient simulation through complex three-dimensional electromagnetic solvers, which are extremely heavy under a computational standpoint.

Several techniques have been presented for the simulation of the NMTLs. These techniques can be subdivided in two main classes, performing simulation in the frequency domain or in the time domain, respectively. The former can obtain closed-form solutions [1] in some cases, but can be used also to analyze more general structures through a piecewise constant discretization of the line [2]. If the transient response is wanted, inverse FFT can be used. However, this technique does not allow a true transient simulation, because FFT

can only be used to obtain the steady state solution. The total simulation time must be long enough for the transients to be extinguished. Therefore, when signals with complex waveforms are applied to unmatched lines and long transients are generated, the number of points for the evaluation of the FFT can be very large. This is the reason why numerical schemes performing the simulation directly in the time domain have been recently proposed. Among these we can cite the methods based on the scattering representation [5], the method of characteristics [8], and the waveform relaxation analysis [4].

This work presents a new Time-Domain Space Expansion (TDSE) method for the numerical solution of the NMTL's. This method is based on a weak formulation of the NMTL equations, which leads to a class of numerical schemes of different approximation order according to the particular choice of some trial and test functions. The core of this work is devoted to the definition of trial and test functions that can be used to produce accurate representations of the solution by keeping the computational effort as small as possible. We will see that bases of wavelets are a good choice.

The mathematical theory of wavelets is relatively recent. After the pioneering work of Haar [26], dated back to 1910, a long time has passed before mathematicians renewed their interest in this subject. In the mid-eighties Morlet and his collaborators [29, 24] discovered that efficient representation of seismic signals could be obtained with functions obtained through dilation and translation of a single *wavelet*. Since then, a theoretical background has been developed, and many books are already available in the literature (see e.g. [22, 19, 28, 55, 56, 30, 52, 57]).

In the mean time, the application of wavelets to the numerical solution of differential and integral equations has been pursued with high activity. A far from complete list of publications in this field can be found in Refs. [39]-[51]. The main reason for this interest is due to the new possibilities that wavelets offer with respect to more standard representations. Wavelets generate spaces of functions with any fixed approximation order. In addition, the hierarchical structure of wavelet spaces can be used to adaptively represent even singular functions with a small number of coefficients. This leads to the possibility of constructing numerical schemes for the solution of a given problem with a high accuracy and a small computational time. We will see that the application of wavelets proves quite efficient for the transient solution of the NMTL equations through the TDSE method.

This work is divided into six chapters. Chapter 1 introduces the TDSE method in a general setting, without reference to any specific choice of trial and test functions. The method is then validated through a few simple examples by using piecewise linear functions. The next three chapters are dedicated to the description and definition of the wavelet functions that will be applied to improve the approximation features of the TDSE method. In particular, Chapter 2 introduces the multilevel decompositions of functional spaces through the Haar

system, and summarizes under an abstract standpoint the main concepts underlying general multilevel decompositions. Chapter 3 describes the biorthogonal decompositions of  $L^2(\mathbb{R})$  and introduces scaling function and wavelet spaces on the real line. These are the building blocks for the construction of scaling function and wavelet spaces on the unit interval, developed in Chapter 4. This is a crucial point, because the intrinsic translation-invariance of wavelets is destroyed when working on bounded domains. However, the use of spaces of functions defined on bounded domains is necessary for the solution of the NMTL equations, because the length of any transmission line is finite. This chapter also introduces the concept of nonlinear approximations based on wavelet thresholding, which is used in this work to obtain adapted representations of voltage and current along the transmission line with a small approximation error and a small number of expansion coefficients. Chapter 5 is dedicated to the evaluation of integrals of products of scaling functions and wavelets and their derivatives. Indeed, the TDSE discretization of the NMTL equations requires the evaluation of inner products of trial and test functions. Finally, Chapter 6 merges the partial results of the foregoing chapters into an improved form of the TDSE method, based on the use of wavelets as trial and test functions. We will see that this method is capable of solving the NMTL equations with any fixed approximation order when the voltage generators produce regular waveforms, and with high adaptivity when the waveforms are singular. In both cases, the approximation error and the computational effort remain small.



# Mathematical notations

This section is devoted to the description of the mathematical notations that will be used throughout this work. Let us recall the definition of the Lebesgue spaces  $L^p(\mathbb{R})$ , with  $1 < p < \infty$  [15]. A function  $f : \mathbb{R} \rightarrow \mathbb{R}$  belongs to  $L^p(\mathbb{R})$  if

$$\|f\|_p = \left( \int_{\mathbb{R}} |f(x)|^p dx \right)^{1/p} < \infty. \quad (1)$$

The quantity  $\|f\|_p$  can be interpreted as a seminorm, and becomes a norm when the set of null functions (i.e. all the functions that vanish almost everywhere) is removed and replaced with a single element, the null function  $f = 0$ . The spaces  $L^p(\mathbb{R})$  are Banach spaces. In particular, for  $p = 2$  we have a Hilbert space, equipped with the inner product

$$\langle f, g \rangle = \int_{\mathbb{R}} f(x)g^*(x)dx. \quad (2)$$

In the following, we will drop the suffix  $p$  when it is clear from the context in which space we are working. When the domain of interest is not  $\mathbb{R}$  but  $\Omega \subseteq \mathbb{R}^n$ , we will use the notation  $L^p(\Omega)$  and the definitions modify in an obvious way.

The discrete space corresponding to  $L^2(\mathbb{R})$  will be denoted as  $\ell^2$ , and consists of all square summable sequences,

$$\{\alpha_k\} \in \ell^2 \iff \|\{\alpha_k\}\| = \left( \sum_{k \in \mathbb{Z}} |\alpha_k|^2 \right)^{1/2} < \infty. \quad (3)$$

We will indicate with  $\ell^2(\mathbb{N})$  the subspace of  $\ell^2$  of sequences  $\{\{\alpha_k\} \in \ell^2 \mid \alpha_k = 0, \forall k < 0\}$ .

Let us now recall the definition of the Fourier transform. A function will be denoted  $f$  in the natural domain and  $\hat{f}$  in the Fourier domain. The direct and inverse Fourier transforms are respectively defined as

$$\hat{f}(\omega) = K_1 \int_{-\infty}^{\infty} f(t) e^{-j\omega t} dt \quad (4)$$

$$f(t) = K_2 \int_{-\infty}^{\infty} \hat{f}(\omega) e^{j\omega t} dt, \quad (5)$$

where the two constants  $K_1$  and  $K_2$  must satisfy

$$K_1 K_2 = \frac{1}{2\pi}. \quad (6)$$



Two different conventions will be used in the following. The first, quite common in the electrical engineering literature, is

$$K_1 = 1, \quad K_2 = \frac{1}{2\pi}. \quad (7)$$

In this case, the Parseval and Plancherel identities read

$$\begin{aligned} \langle \hat{f}, \hat{g} \rangle &= 2\pi \langle f, g \rangle \\ \|\hat{f}\|^2 &= 2\pi \|f\|^2. \end{aligned}$$

The other choice for the normalization constants, which is popular in the mathematics literature, is

$$K_1 = \frac{1}{\sqrt{2\pi}}, \quad K_2 = \frac{1}{\sqrt{2\pi}}. \quad (8)$$

This convention leads to the unitarity of the Fourier operator, which is then norm-preserving according to

$$\begin{aligned} \langle \hat{f}, \hat{g} \rangle &= \langle f, g \rangle \\ \|\hat{f}\|^2 &= \|f\|^2. \end{aligned}$$

We do not list here standard properties of the Fourier operator, like linearity, translation, dilation, etc..

We will need to use spaces of functions with a certain degree of regularity. It is natural then to use the Sobolev spaces  $H^s(\mathbb{R})$  [12]. We recall that a function  $f : \mathbb{R} \rightarrow \mathbb{R}$  is in  $H^s(\mathbb{R})$  if it has all the weak derivatives of order at most  $s$  for  $s \geq 1$  and integer. We will set  $H^0(\mathbb{R}) = L^2(\mathbb{R})$ ; for  $s > 0$  the corresponding spaces can be obtained through interpolation. The Sobolev spaces can be characterized in terms of weighted summability in the Fourier domain. More precisely,

$$f \in H^s(\mathbb{R}) \iff \|f\|_s = \left( \int_{\mathbb{R}} |\hat{f}(\xi)|^2 (1 + \xi^2)^s d\xi \right)^{1/2} < \infty. \quad (9)$$

Finally, we list other notations that will be commonly used.

- With  $f(\cdot)$  we indicate the function  $f$  specifying its argument. For example, the writing  $f(\cdot - y)$  means the function that assumes the value  $f(x - y)$  for almost all  $x$ .
- With  $C$  we will denote a positive constant without specifying its value, which may change from time to time.
- Given two functions  $N_i : V \rightarrow \mathbb{R}_+$  ( $i = 1, 2$ ) defined on a set  $V$ , we will use the notation  $N_1(v) \lesssim N_2(v)$  when there is a constant  $C > 0$  such that  $N_1(v) \leq CN_2(v)$ , for each  $v \in V$ . In addition, we will write  $N_1(v) \asymp N_2(v)$ , if  $N_1(v) \lesssim N_2(v)$  and  $N_2(v) \lesssim N_1(v)$ .

- A set of basis functions  $\{f_k\}$  of  $L^2(\mathbb{R})$  will be said 2-stable when the norm of any function  $v \in L^2(\mathbb{R})$  is equivalent to the discrete norm in  $\ell^2$  of its coefficients. Using the notation in the item above, we can express the 2-stability as

$$\|\{\alpha_k\}\|_{\ell^2} \asymp \left\| \sum_k \alpha_k f_k \right\|_{L^2}, \quad \forall \{\alpha_k\} \in \ell^2.$$

- For any  $x \in \mathbb{R}$  we will indicate with  $\lceil x \rceil$  (or  $\lfloor x \rfloor$ ) the smallest (largest) integer larger (smaller) than  $x$ .

A summary of the symbols that are recurrently used throughout this work can be found in Appendix A.



# Chapter 1

## The Time Domain Space Expansion method

This chapter presents the Time-Domain Space Expansion (TDSE) method for the transient simulation of the Nonuniform Multiconductor Transmission Lines (NMTL) equations. The aim is the development of accurate numerical schemes for the electrical simulation of interconnects that are characterized by non-translation invariant cross-sectional parameters. Many structures fall into this class, e.g. high speed packaging interconnects, impedance matching networks or cables in complex structures, like automobiles or airplanes. It is therefore of crucial importance to be able to simulate the electrical behavior of such structures including the effects of longitudinal nonuniformities, in order to predict crosstalk noise, spurious electromagnetic couplings, and derive appropriate design rules.

The presented method is based on the spatial expansion of solution and per-unit-length parameters into two sets of approximating functions, and on testing of the NMTL equations with a third set of functions. The mathematical formulation of the TDSE method is presented in Section 1.1 in a general setting. This formulation is then particularized in Section 1.2 to piecewise linear basis functions (linear finite elements bases). Section 1.3 presents some numerical examples with validations of the TDSE method and applications to the transient solution for some structures of practical interest. The results will show that quite accurate solutions can be obtained with linear approximating functions, provided that the dimension of the approximation space is large enough. Some of the structures analyzed in this chapter will be solved again in Chapter 6 by using basis functions (developed throughout this work) with higher regularity and better approximation properties. This will lead to much higher accuracy in the numerical solutions given the dimension of the approximation space, or equivalently to a smaller number of unknowns needed to obtain a solution at a fixed accuracy. It should be noted that the formulation in Section 1.1 is independent of the particular choice of the basis functions.

## 1.1 Mathematical formulation

Let us consider the Nonuniform Multiconductor Transmission Lines (NMTL) equations,

$$\frac{\partial}{\partial z} \mathbf{V}(z, t) = -\mathbf{L}(z) \frac{\partial}{\partial t} \mathbf{I}(z, t) - \mathbf{R}(z) \mathbf{I}(z, t), \quad (1.1)$$

$$\frac{\partial}{\partial z} \mathbf{I}(z, t) = -\mathbf{C}(z) \frac{\partial}{\partial t} \mathbf{V}(z, t) - \mathbf{G}(z) \mathbf{V}(z, t), \quad (1.2)$$

with  $\mathbf{V}(z, t)$  and  $\mathbf{I}(z, t)$  indicating the voltage and current vectors at location  $z$  and time  $t$ . The line is assumed to have  $P + 1$  conductors, labelled with  $i = 0, \dots, P$ , with the zeroth taken as the reference for voltages and the return for currents. The per-unit-length parameters  $\mathbf{L}(z)$ ,  $\mathbf{C}(z)$ ,  $\mathbf{R}(z)$ , and  $\mathbf{G}(z)$  are  $P \times P$  matrices whose entries are arbitrary functions of the space variable  $z$ . We will consider the length of the line to be normalized. The change of variable  $z = \mathcal{L}\varsigma$  can be used for lines of length  $\mathcal{L}$ , with  $\varsigma \in [0, 1]$ . However, hereafter we proceed using  $z$  without loss of generality. The line will be terminated by Thévenin loads, i.e.

$$\mathbf{V}(0, t) = \mathbf{V}_\mathbf{S}(t) - \mathbf{R}_\mathbf{S} \mathbf{I}(0, t), \quad (1.3)$$

$$\mathbf{V}(1, t) = \mathbf{V}_\mathbf{L}(t) + \mathbf{R}_\mathbf{L} \mathbf{I}(1, t), \quad (1.4)$$

where  $\mathbf{R}_\mathbf{S}$  and  $\mathbf{R}_\mathbf{L}$  are arbitrary  $P \times P$  real matrices and  $\mathbf{V}_\mathbf{S}(t)$ ,  $\mathbf{V}_\mathbf{L}(t)$  are voltage source vectors. The whole derivation can be easily modified to account for Norton or mixed terminations, therefore only the Thévenin case will be described in detail.

The true solution of equations (1.1)-(1.2) loaded with (1.3)-(1.4) lies in some functional space  $\mathcal{S}$ . The basic hypothesis underlying the method is that a sequence of approximation spaces  $\mathcal{S}_h \subset \mathcal{S}$  can be constructed such that

$$\mathcal{S}_h \rightarrow \mathcal{S}, \quad h \rightarrow 0. \quad (1.5)$$

This condition insures the consistence of the discretization as the parameter  $h$  vanishes. As for the solution vectors, also the per unit length parameters will be assumed to belong to some functional space  $\mathcal{P}$ , which can be approximated by some spaces  $\mathcal{P}_h \subset \mathcal{P}$  with the same convergence properties,

$$\mathcal{P}_h \rightarrow \mathcal{P}, \quad h \rightarrow 0. \quad (1.6)$$

An approximate solution for the NMTL equations will be sought for in the spaces  $\mathcal{S}_h$ . It is clear that different choices for the approximation spaces  $\mathcal{S}_h$  and  $\mathcal{P}_h$  will lead to different numerical schemes with different convergence properties. A careful choice of these spaces upon knowledge of the forcing waveforms  $\mathbf{V}_\mathbf{S}(t)$  and  $\mathbf{V}_\mathbf{L}(t)$  will be essential for a good behavior of the method.

We introduce now two sets of basis functions for the approximation spaces,

$$\mathcal{S}_h = \text{span} \{ \zeta_n, n = 1, \dots, N_\zeta \}, \quad (1.7)$$

$$\mathcal{P}_h = \text{span} \{ \phi_k, k = 1, \dots, N_\phi \}, \quad (1.8)$$

where  $N_\zeta$  and  $N_\phi$  must be finite and are dependent on the discretization parameter  $h$ . The voltage and current vectors can be expanded into these basis functions,

$$\mathbf{V}(z, t) = \sum_{n=1}^{N_\zeta} \zeta_n(z) \mathbf{V}_n(t), \quad (1.9)$$

$$\mathbf{I}(z, t) = \sum_{n=1}^{N_\zeta} \zeta_n(z) \mathbf{I}_n(t), \quad (1.10)$$

as well as the per unit length parameters,

$$\begin{aligned} \mathbf{L} &= \sum_{k=1}^{N_\phi} \phi_k(z) \mathbf{L}_k, \\ \mathbf{C} &= \sum_{k=1}^{N_\phi} \phi_k(z) \mathbf{C}_k, \\ \mathbf{R} &= \sum_{k=1}^{N_\phi} \phi_k(z) \mathbf{R}_k, \\ \mathbf{G} &= \sum_{k=1}^{N_\phi} \phi_k(z) \mathbf{G}_k. \end{aligned} \quad (1.11)$$

The voltage and current coefficients  $\mathbf{V}_n(t)$  and  $\mathbf{I}_n(t)$  are real vectors of dimension  $P$ , while the coefficients  $\mathbf{L}_k$ ,  $\mathbf{C}_k$ ,  $\mathbf{R}_k$ , and  $\mathbf{G}_k$  are  $P \times P$  real matrices. If we substitute the above expansions in the NMTL equations (1.1) and (1.2), we get

$$\begin{aligned} \sum_{n=1}^{N_\zeta} \frac{\partial}{\partial z} \zeta_n(z) \mathbf{V}_n(t) + \sum_{k=1}^{N_\phi} \phi_k(z) \mathbf{L}_k \sum_{n=1}^{N_\zeta} \zeta_n(z) \frac{d}{dt} \mathbf{I}_n(t) + \\ \sum_{k=1}^{N_\phi} \phi_k(z) \mathbf{R}_k \sum_{n=1}^{N_\zeta} \zeta_n(z) \mathbf{I}_n(t) = 0, \end{aligned} \quad (1.12)$$

$$\begin{aligned} \sum_{n=1}^{N_\zeta} \frac{\partial}{\partial z} \zeta_n(z) \mathbf{I}_n(t) + \sum_{k=1}^{N_\phi} \phi_k(z) \mathbf{C}_k \sum_{n=1}^{N_\zeta} \zeta_n(z) \frac{d}{dt} \mathbf{V}_n(t) + \\ \sum_{k=1}^{N_\phi} \phi_k(z) \mathbf{G}_k \sum_{n=1}^{N_\zeta} \zeta_n(z) \mathbf{V}_n(t) = 0. \end{aligned} \quad (1.13)$$

We introduce now a third set of functions, which will be taken as test functions for the derivation of a weak form of the NMTL equations. The only

restriction on these functions, denoted as  $\{\eta_m, m = 1, \dots, N_\zeta\}$ , is that they are linearly independent. Taking the inner product of (1.12) and (1.13) with each  $\eta_m$ , we get

$$\sum_{n=1}^{N_\zeta} \Lambda_{mn} \mathbf{V}_n(t) + \sum_{n=1}^{N_\zeta} \hat{\mathbf{L}}_{mn} \frac{d}{dt} \mathbf{I}_n(t) + \sum_{n=1}^{N_\zeta} \widehat{\mathbf{R}}_{mn} \mathbf{I}_n(t) = 0, \quad (1.14)$$

$$\sum_{n=1}^{N_\zeta} \Lambda_{mn} \mathbf{I}_n(t) + \sum_{n=1}^{N_\zeta} \hat{\mathbf{C}}_{mn} \frac{d}{dt} \mathbf{V}_n(t) + \sum_{n=1}^{N_\zeta} \widehat{\mathbf{G}}_{mn} \mathbf{V}_n(t) = 0, \quad (1.15)$$

valid  $\forall m = 1, \dots, N_\zeta$ . The matrices used in the above expressions are

$$\Lambda_{mn} = \left\langle \frac{d}{dz} \zeta_n, \eta_m \right\rangle \mathcal{I}_P, \quad (1.16)$$

where  $\mathcal{I}_P$  is the  $P \times P$  identity matrix, and

$$\begin{aligned} \hat{\mathbf{L}}_{mn} &= \sum_{k=1}^{N_\phi} \mathbf{L}_k B_{mn}^{(k)}, \\ \hat{\mathbf{C}}_{mn} &= \sum_{k=1}^{N_\phi} \mathbf{C}_k B_{mn}^{(k)}, \\ \widehat{\mathbf{R}}_{mn} &= \sum_{k=1}^{N_\phi} \mathbf{R}_k B_{mn}^{(k)}, \\ \widehat{\mathbf{G}}_{mn} &= \sum_{k=1}^{N_\phi} \mathbf{G}_k B_{mn}^{(k)}, \end{aligned} \quad (1.17)$$

where

$$B_{mn}^{(k)} = \langle \zeta_n \phi_k, \eta_m \rangle. \quad (1.18)$$

The two sets of equations (1.14) and (1.15) describe the behavior of the non terminated line. We consider now the inclusion of the loads (1.3) and (1.4). For the subsequent derivation, it is convenient to choose the trial and test functions such that only one is nonzero at the boundaries, i.e.,

$$\begin{aligned} \zeta_n(0) &= 0, \quad \forall n = 2, \dots, N_\zeta, \\ \zeta_n(1) &= 0, \quad \forall n = 1, \dots, N_\zeta - 1, \\ \eta_m(0) &= 0, \quad \forall m = 2, \dots, N_\zeta, \\ \eta_m(1) &= 0, \quad \forall m = 1, \dots, N_\zeta - 1. \end{aligned} \quad (1.19)$$

This is not a real restriction because whatever be the initial choice of basis functions, a change of basis can always be performed to obtain only one nonzero function at both edges. The two edge trial functions will also be normalized so that

$$\zeta_1(0) = \zeta_N(1) = 1.$$

In addition, we will consider the case of non-overlapping border functions, i.e.,

$$\begin{aligned}\text{supp } \zeta_1 \cap \text{supp } \zeta_{N_\zeta} &= \emptyset, \\ \text{supp } \eta_1 \cap \text{supp } \eta_{N_\zeta} &= \emptyset.\end{aligned}\tag{1.20}$$

This will simplify the form of the border equations in the final system, because the terms referring to the two edges will not interact with each other. No similar restrictions need to be enforced on the per unit length coefficients expansion functions  $\phi_k$ . Substituting now the expansions (1.9) and (1.10) into the load equations (1.3) and (1.4) and using the conditions (1.19), we get

$$\mathbf{V}_1(t) = \mathbf{V}_S(t) - \mathbf{R}_S \mathbf{I}_1(t) \tag{1.21}$$

$$\mathbf{V}_{N_\zeta}(t) = \mathbf{V}_L(t) + \mathbf{R}_L \mathbf{I}_{N_\zeta}(t) \tag{1.22}$$

These expressions for the loads can be used to eliminate the two unknowns  $\mathbf{V}_1(t)$  and  $\mathbf{V}_{N_\zeta}(t)$  from the system (1.14)-(1.15). It should be noted, however, that the number of scalar unknowns in the system is  $2NP$ , which matches the number of scalar equations. If we eliminate the voltages at the two edges, i.e.  $2P$  scalar unknowns, also  $2P$  equations must be suppressed in order to keep the balance even. These equations are obviously the ones involving the projection onto the border test functions  $\eta_1$  and  $\eta_{N_\zeta}$ . The following derivation shows how the final system of ODE's can be derived.

We begin with the equations (1.14)-(1.15) with  $m = 2, \dots, N_\zeta - 1$ . These are the projections onto the “internal” test functions, and can be rewritten by substituting the load equations (1.21) and (1.22), obtaining

$$\begin{aligned}& \sum_{n=2}^{N_\zeta-1} \mathbf{\Lambda}_{mn} \mathbf{V}_n(t) + \sum_{n=1}^{N_\zeta} \widehat{\mathbf{L}}_{mn} \frac{d}{dt} \mathbf{I}_n(t) + \sum_{n=1}^{N_\zeta} \widehat{\mathbf{R}}_{mn} \mathbf{I}_n(t) \\ & \quad - \mathbf{\Lambda}_{m1} \mathbf{R}_S \mathbf{I}_1(t) + \mathbf{\Lambda}_{mN_\zeta} \mathbf{R}_L \mathbf{I}_{N_\zeta}(t) \\ &= -\mathbf{\Lambda}_{m1} \mathbf{V}_S(t) - \mathbf{\Lambda}_{mN_\zeta} \mathbf{V}_L(t),\end{aligned}\tag{1.23}$$

$$\begin{aligned}& \sum_{n=1}^{N_\zeta} \mathbf{\Lambda}_{mn} \mathbf{I}_n(t) + \sum_{n=2}^{N_\zeta-1} \widehat{\mathbf{C}}_{mn} \frac{d}{dt} \mathbf{V}_n(t) + \sum_{n=2}^{N_\zeta-1} \widehat{\mathbf{G}}_{mn} \mathbf{V}_n(t) \\ & \quad - \widehat{\mathbf{C}}_{m1} \mathbf{R}_S \frac{d}{dt} \mathbf{I}_1(t) + \widehat{\mathbf{C}}_{mN_\zeta} \mathbf{R}_L \frac{d}{dt} \mathbf{I}_{N_\zeta}(t) - \widehat{\mathbf{G}}_{m1} \mathbf{R}_S \mathbf{I}_1(t) + \widehat{\mathbf{G}}_{mN_\zeta} \mathbf{R}_L \mathbf{I}_{N_\zeta}(t) \\ &= -\widehat{\mathbf{G}}_{m1} \mathbf{V}_S(t) - \widehat{\mathbf{G}}_{mN_\zeta} \mathbf{V}_L(t) - \widehat{\mathbf{C}}_{m1} \frac{d}{dt} \mathbf{V}_S(t) - \widehat{\mathbf{C}}_{mN_\zeta} \frac{d}{dt} \mathbf{V}_L(t).\end{aligned}\tag{1.24}$$

Instead of suppressing two of the remaining border equations, we take the linear combination with coefficients  $\alpha_m$  and  $\beta_m$ ,

$$\alpha_m(1.14) + \beta_m(1.15) = 0, \quad m \in \{1, N_\zeta\}.\tag{1.25}$$

After few straightforward steps, the resulting two border equations read

$$\sum_{n=1}^{N_\zeta} \alpha_1 \widehat{\mathbf{L}}_{1n} \frac{d}{dt} \mathbf{I}_n(t) - \beta_1 \widehat{\mathbf{C}}_{11} \mathbf{R}_S \frac{d}{dt} \mathbf{I}_1(t) + \sum_{n=2}^{N_\zeta-1} \beta_1 \widehat{\mathbf{C}}_{1n} \frac{d}{dt} \mathbf{V}_n(t) +$$



$$\begin{aligned}
& \sum_{n=1}^{N_\zeta} [\beta_1 \mathbf{\Lambda}_{1n} + \alpha_1 \widehat{\mathbf{R}}_{1n}] \mathbf{I}_n(t) - [\alpha_1 \mathbf{\Lambda}_{11} + \beta_1 \widehat{\mathbf{G}}_{11}] \mathbf{R}_S \mathbf{I}_1(t) + \\
& \sum_{n=2}^{N_\zeta-1} [\alpha_1 \mathbf{\Lambda}_{1n} + \beta_1 \widehat{\mathbf{G}}_{1n}] \mathbf{V}_n(t) \\
& = -[\alpha_1 \mathbf{\Lambda}_{11} + \beta_1 \widehat{\mathbf{G}}_{11}] \mathbf{V}_S(t) - \beta_1 \widehat{\mathbf{C}}_{11} \frac{d}{dt} \mathbf{V}_S(t)
\end{aligned} \tag{1.26}$$

$$\begin{aligned}
& \sum_{n=1}^{N_\zeta} \alpha_{N_\zeta} \widehat{\mathbf{L}}_{N_\zeta n} \frac{d}{dt} \mathbf{I}_n(t) + \beta_{N_\zeta} \widehat{\mathbf{C}}_{N_\zeta N_\zeta} \mathbf{R}_L \frac{d}{dt} \mathbf{I}_{N_\zeta}(t) + \sum_{n=2}^{N_\zeta-1} \beta_{N_\zeta} \widehat{\mathbf{C}}_{N_\zeta n} \frac{d}{dt} \mathbf{V}_n(t) + \\
& \sum_{n=1}^{N_\zeta} [\beta_{N_\zeta} \mathbf{\Lambda}_{N_\zeta n} + \alpha_{N_\zeta} \widehat{\mathbf{R}}_{N_\zeta n}] \mathbf{I}_n(t) + [\alpha_{N_\zeta} \mathbf{\Lambda}_{N_\zeta N_\zeta} + \beta_{N_\zeta} \widehat{\mathbf{G}}_{N_\zeta N_\zeta}] \mathbf{R}_L \mathbf{I}_{N_\zeta}(t) + \\
& \sum_{n=2}^{N_\zeta-1} [\alpha_{N_\zeta} \mathbf{\Lambda}_{N_\zeta n} + \beta_{N_\zeta} \widehat{\mathbf{G}}_{N_\zeta n}] \mathbf{V}_n(t) \\
& = -[\alpha_{N_\zeta} \mathbf{\Lambda}_{N_\zeta N_\zeta} + \beta_{N_\zeta} \widehat{\mathbf{G}}_{N_\zeta N_\zeta}] \mathbf{V}_L(t) - \beta_{N_\zeta} \widehat{\mathbf{C}}_{N_\zeta N_\zeta} \frac{d}{dt} \mathbf{V}_L(t).
\end{aligned} \tag{1.27}$$

Putting all the equations together we get a system of  $P(2N_\zeta - 2)$  ODE's, which can be solved with a suitable integration method such as Runge-Kutta or Adams-Moulton [11]. All the simulations produced in this work were obtained with a 5<sup>th</sup> – 6<sup>th</sup> order Runge-Kutta scheme [10]. This system can be formally written as

$$\mathbf{\Psi} \frac{d}{dt} \mathbf{x}(t) + \mathbf{\Phi} \mathbf{x}(t) = \mathbf{\Delta}_S \mathbf{V}_S(t) + \mathbf{\Delta}_{SD} \frac{d}{dt} \mathbf{V}_S(t) + \mathbf{\Delta}_L \mathbf{V}_L(t) + \mathbf{\Delta}_{LD} \frac{d}{dt} \mathbf{V}_L(t), \tag{1.28}$$

where  $\mathbf{\Psi}$  is nonsingular if the trial and test functions are linearly independent and  $\mathbf{\Delta}_S$ ,  $\mathbf{\Delta}_{SD}$ ,  $\mathbf{\Delta}_L$ ,  $\mathbf{\Delta}_{LD}$  are  $P(2N_\zeta - 2) \times P$  real matrices. The vector of unknowns  $\mathbf{x}$  collects the voltage and current coefficients vectors,

$$\mathbf{x} = [\mathbf{I}_1^T, \dots, \mathbf{I}_{N_\zeta}^T, \mathbf{V}_2^T, \dots, \mathbf{V}_{N_\zeta-1}^T]^T. \tag{1.29}$$

It should be noted that due to the weak formulation, the forcing terms in the system (1.28) include also the time derivatives of the source vectors  $\mathbf{V}_S(t)$  and  $\mathbf{V}_L(t)$ . Therefore, singular waveforms like delta functions or step functions cannot be handled by this method.

### 1.1.1 Frequency domain analysis

The analysis of the foregoing section has been conducted in the time domain since we aim at the derivation of accurate numerical schemes for the transient simulation of the NMTL equations. However, the spatial discretization method described in Section 1.1 can also be applied to derive the sinusoidal steady state

voltage and current distributions along the conductors of the line. The transformation of the system (1.28) into the frequency domain is straightforward, i.e.,

$$(j\omega\mathbf{\Psi} + \mathbf{\Phi})\mathbf{X}(\omega) = (\mathbf{\Delta_S} + j\omega\mathbf{\Delta_{SD}})\widehat{\mathbf{V_S}}(\omega) + (\mathbf{\Delta_L} + j\omega\mathbf{\Delta_{LD}})\widehat{\mathbf{V_L}}(\omega). \quad (1.30)$$

The vector  $\mathbf{X}(\omega)$  includes the phasors of the unknown voltage and current expansion coefficients according to (1.29), and the terms  $\widehat{\mathbf{V_S}}(\omega)$ ,  $\widehat{\mathbf{V_L}}(\omega)$  are the phasors associated to the sinusoidal voltage sources  $\mathbf{V_S}(t)$  and  $\mathbf{V_L}(t)$ , respectively. It should be noted that the matrix  $(j\omega\mathbf{\Psi} + \mathbf{\Phi})$ , in the case of locally supported trial and test functions, has a sparse structure. Therefore, the use of an efficient solver for sparse complex matrices could reduce to  $O(N)$  the number of operations involved in the solution of (1.30).

### 1.1.2 Incident field excitation

The voltage and current distributions excited along the line by external fields can also be handled by this method. The NMTL equations (1.1)-(1.2) are modified to account for incident fields by simply adding equivalent distributed voltage and current sources [7],

$$\frac{\partial}{\partial z}\mathbf{V}(z, t) = -\mathbf{L}(z)\frac{\partial}{\partial t}\mathbf{I}(z, t) - \mathbf{R}(z)\mathbf{I}(z, t) + \mathbf{V_F}(z, t), \quad (1.31)$$

$$\frac{\partial}{\partial z}\mathbf{I}(z, t) = -\mathbf{C}(z)\frac{\partial}{\partial t}\mathbf{V}(z, t) - \mathbf{G}(z)\mathbf{V}(z, t) + \mathbf{I_F}(z, t). \quad (1.32)$$

The particular form of the incident fields, e.g. plane waves, determines the dependence of the equivalent sources on  $z$  and  $t$ . In the following we will not assume a particular incident field distribution, but we will suppose that  $\mathbf{V_F}$  and  $\mathbf{I_F}$  are known explicitly at any point  $(z, t)$ .

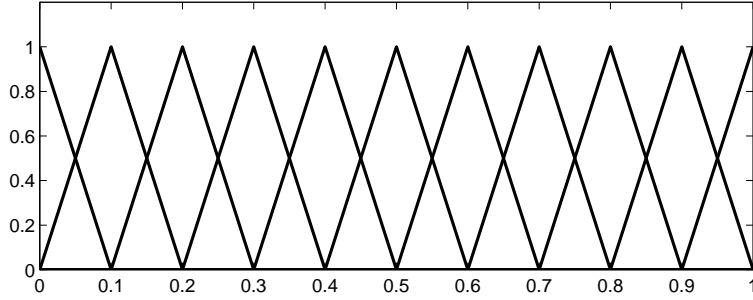
The source terms can be expanded in terms of the trial functions, obtaining

$$\begin{aligned} \mathbf{V_F}(z, t) &= \sum_{n=1}^{N_\zeta} \zeta_n(z) \mathbf{V_{F_n}}(t), \\ \mathbf{I_F}(z, t) &= \sum_{n=1}^{N_\zeta} \zeta_n(z) \mathbf{I_{F_n}}(t). \end{aligned}$$

Projecting now onto the test functions  $\eta_m(z)$  we get

$$\langle \mathbf{V_F}(z, t), \eta_m(z) \rangle = \sum_{n=1}^{N_\zeta} \mathbf{E_{mn}} \mathbf{V_{F_n}}(t), \quad (1.33)$$

$$\langle \mathbf{I_F}(z, t), \eta_m(z) \rangle = \sum_{n=1}^{N_\zeta} \mathbf{E_{mn}} \mathbf{I_{F_n}}(t), \quad (1.34)$$



**Figure 1.1:** Piecewise linear basis functions in the case  $N = 11$ .

where the matrix  $\mathbf{E}$  is defined as

$$\mathbf{E}_{mn} = \langle \zeta_n, \eta_m \rangle \mathcal{I}_P.$$

If we replace now the right end sides of Eqs. (1.14) and (1.15) with the expressions (1.33) and (1.34), respectively, the external field contribution will be automatically accounted for. The treatment of the boundary conditions is obviously not affected by these additional source terms. We skip here the details because the system of ODE's is simply obtained by adding the appropriate source terms to Eq. (1.28).

## 1.2 Piecewise linear approximation

This section will particularize the derivation of section 1.1 to the case of piecewise linear approximations of both solution and per unit length parameters. We will show the details for the simplest case of coinciding approximation spaces  $\mathcal{S}_h = \mathcal{P}_h$ , with the same dimension  $N = N_\zeta = N_\phi$  and basis sets  $\zeta_n = \phi_n = \eta_n$ . We subdivide the unit interval into  $N-1$  equal intervals of length  $h = 1/(N-1)$ , and define the basis functions for  $n = 1, \dots, N$  as

$$\zeta_n(z) = \begin{cases} [z - (n-2)h]/h, & z \in [(n-2)h, (n-1)h] \cap [0, 1], \\ [-z + nh]/h, & z \in [(n-1)h, nh] \cap [0, 1], \\ 0 & \text{otherwise.} \end{cases}$$

These functions are depicted in Fig. 1.1. Note that these are interpolating functions, therefore the computation of the expansion coefficients of any function reduces to its evaluation at the points  $(n-1)h$ . For this choice of basis functions the approximation error is expected to behave as  $O(h^2)$  when  $h$  approaches zero. Of course, this holds only when the solution has a continuous first derivative in the integration domain.

The entries in the matrices forming the building blocks of the system of ODE's (1.28) can be evaluated in closed form with straightforward integrations.

We have

$$\begin{aligned}\mathbf{\Lambda}_{11} &= -(1/2)\mathcal{I}_P, \\ \mathbf{\Lambda}_{mm} &= 0, & m = 2, \dots, N-1, \\ \mathbf{\Lambda}_{NN} &= (1/2)\mathcal{I}_P, \\ \mathbf{\Lambda}_{m,m+1} &= (1/2)\mathcal{I}_P, & m = 1, \dots, N-1, \\ \mathbf{\Lambda}_{m,m-1} &= -(1/2)\mathcal{I}_P, & m = 2, \dots, N,\end{aligned}$$

and

$$\begin{aligned}B_{mm}^{(m)} &= B_{mm}^{(m)} = h/4 & m \in \{1, N\}, \\ B_{mm}^{(m)} &= h/2, & m = 2, \dots, N-1, \\ B_{mm}^{(m+1)} &= B_{mm}^{(m-1)} = h/12, & m = 2, \dots, N-1, \\ B_{m,m+1}^{(m)} &= B_{m,m+1}^{(m+1)} = h/12, & m = 1, \dots, N-1, \\ B_{m,m-1}^{(m)} &= B_{m,m-1}^{(m-1)} = h/12, & m = 2, \dots, N, \\ B_{mn}^{(k)} &= 0 & \text{otherwise.}\end{aligned}$$

These coefficients are such that the matrices  $\mathbf{\Lambda}$ ,  $\widehat{\mathbf{L}}$ ,  $\widehat{\mathbf{C}}$ ,  $\widehat{\mathbf{R}}$ ,  $\widehat{\mathbf{G}}$  have a banded structure, with only one upper and lower codiagonal made of blocks of size  $P$ . For example, the matrix  $\widehat{\mathbf{L}}$  results

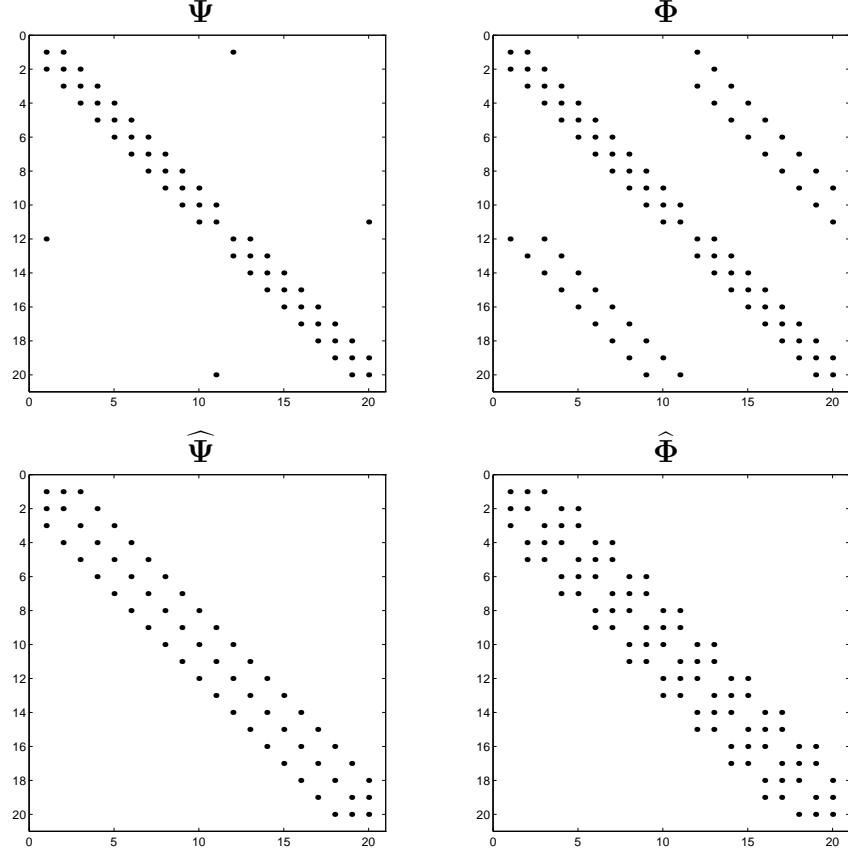
$$\widehat{\mathbf{L}} = \frac{h}{12} \begin{bmatrix} 3\mathbf{L}_1 + \mathbf{L}_2 & \mathbf{L}_1 + \mathbf{L}_2 & & \\ & \mathbf{L}_{m-1} + \mathbf{L}_m & \mathbf{L}_{m-1} + 6\mathbf{L}_m + \mathbf{L}_{m+1} & \mathbf{L}_m + \mathbf{L}_{m+1} \\ & & & \\ & & \mathbf{L}_{N-1} + \mathbf{L}_N & \mathbf{L}_{N-1} + 3\mathbf{L}_N \end{bmatrix}.$$

Note that the main diagonal, apart from the first and last rows, can also be obtained by convolving the coefficients  $\mathbf{L}_k$  with a FIR filter with mask  $\{h/12, h/2, h/12\}$ , and the codiagonals with another FIR filter with mask  $\{h/12, h/12\}$ . Therefore, the dependence of the per unit length parameters on the space variable  $z$  is reproduced along the diagonals. Similar results hold for the other matrices  $\widehat{\mathbf{C}}$ ,  $\widehat{\mathbf{R}}$ ,  $\widehat{\mathbf{G}}$ . As a consequence, also the system matrices  $\mathbf{\Psi}$  and  $\mathbf{\Phi}$  have a block-banded structure. If the unknowns are arranged according to Eq. (1.29) both  $\mathbf{\Psi}$  and  $\mathbf{\Phi}$  have a block-tridiagonal structure, as shown in the top panels of Figure 1.2. However, when the basis sets are formed by functions with local support, such as in the present case, it is convenient to rearrange the unknowns as

$$\hat{\mathbf{x}} = [\mathbf{I}_1^T, \mathbf{I}_2^T, \mathbf{V}_2^T, \dots, \mathbf{I}_n^T, \mathbf{V}_n^T, \dots, \mathbf{V}_{N-1}^T, \mathbf{I}_N^T]^T. \quad (1.35)$$

This can be accomplished by applying a permutation matrix  $\mathbf{T}$  such that  $\mathbf{T}^T = \mathbf{T}^{-1}$ , obtaining the new system

$$\widehat{\mathbf{\Psi}} \frac{d}{dt} \hat{\mathbf{x}}(t) + \widehat{\mathbf{\Phi}} \hat{\mathbf{x}}(t) = \widehat{\mathbf{\Delta}}_{\mathbf{S}} \mathbf{V}_{\mathbf{S}}(t) + \widehat{\mathbf{\Delta}}_{\mathbf{SD}} \frac{d}{dt} \mathbf{V}_{\mathbf{S}}(t) + \widehat{\mathbf{\Delta}}_{\mathbf{L}} \mathbf{V}_{\mathbf{L}}(t) + \widehat{\mathbf{\Delta}}_{\mathbf{LD}} \frac{d}{dt} \mathbf{V}_{\mathbf{L}}(t), \quad (1.36)$$



**Figure 1.2:** Structure of the system matrices of Eqs. (1.28) and (1.36) in the piecewise linear case with  $N = 11$ . Each dot represents a nonzero  $P \times P$  block.

where

$$\widehat{\Psi} = \mathbf{T}^T \Psi \mathbf{T}, \quad \widehat{\Delta}_S = \mathbf{T}^T \Delta_S,$$

and similarly for the other matrices. This results in a banded structure for both  $\widehat{\Psi}$  and  $\widehat{\Phi}$ , as shown in the bottom panels of Fig. 1.2. This is convenient because the numerical schemes for the solution of systems of ODE's require to explicit the time derivative of the state vector. Therefore, the matrix  $\Psi$  must be inverted and a full matrix  $\Psi^{-1}\Phi$  is obtained. This results in  $O(N^2)$  operations to evaluate the matrix-vector product  $\Psi^{-1}\Phi x$ . However, if the inversion is accomplished through LU decomposition, both the lower and upper triangular matrices in the decomposition are still banded. Therefore, the total number of operations involved in the computation of  $\widehat{\Psi}^{-1}\widehat{\Phi}\hat{x}$  is only  $O(N)$ . This allows to increase the accuracy of the method at a low computational cost.

## 1.3 Numerical examples

This section will apply the method outlined in Sections 1.1 and 1.2 to some practical examples. The lossless scalar exponential line will be analyzed in Section 1.3.1. The crosstalk on a three-conductor printed circuit board will be studied in Section 1.3.2. Finally, the crosstalk on a nonuniform MTL made of two nonparallel wires above a ground plane will be determined in Section 1.3.3.

### 1.3.1 The exponential line

We chose the scalar exponential line as a test case for our numerical scheme because the analytical solution in the frequency domain is well known and understood. We will report here the main results. The details can be found in [3].

Let us consider a scalar nonuniform line ( $P = 1$ ) of unitary length characterized by

$$\begin{aligned} L(z) &= L^0 e^{\delta z}, & R(z) &= 0, \\ C(z) &= C^0 e^{-\delta z}, & G(z) &= 0, \end{aligned}$$

where the parameter  $\delta$  controls the rate of taper and  $L^0$ ,  $C^0$  are the nominal per unit length inductance and capacitance at the edge  $z = 0$ . The nominal characteristic impedance of the line is therefore

$$Z(z) = \sqrt{\frac{L(z)}{C(z)}} = Z_0 e^{\delta z},$$

where  $Z_0$  is the nominal characteristic impedance at the edge  $z = 0$ . At a fixed frequency  $\omega$  we can define the propagation constant  $\gamma$  and the transfer constant  $\Gamma$  of the line as

$$\gamma = j\omega\sqrt{L(z)C(z)} = j\omega\sqrt{L^0 C^0}; \quad \Gamma = \sqrt{\gamma^2 + \delta^2/4} = \alpha + j\beta, \quad \alpha, \beta \geq 0.$$

The voltage and current along the line are expressed by

$$\begin{aligned} V(z, \omega) &= A(\omega)e^{-(\Gamma-\delta/2)z} + B(\omega)e^{(\Gamma+\delta/2)z} \\ I(z, \omega) &= \frac{A(\omega)}{Z_0} \frac{\Gamma - \delta/2}{\gamma} e^{-(\Gamma+\delta/2)z} - \frac{B(\omega)}{Z_0} \frac{\Gamma + \delta/2}{\gamma} e^{(\Gamma-\delta/2)z}, \end{aligned}$$

where  $A(\omega)$  and  $B(\omega)$  are frequency-dependent constants that are determined by imposing the load equations, in this case

$$\begin{aligned} V(0, \omega) &= V_S(\omega) - R_S I(0, \omega), \\ V(1, \omega) &= R_L I(1, \omega). \end{aligned}$$

This solution can be interpreted as usual as a superposition of travelling voltage and current waves with positive and negative velocity. As the impedance level

increases at rate  $\delta$  with  $z$ , the positive voltage wave increases in magnitude and the corresponding current wave decreases at rate  $\delta/2$ . The converse holds for negative voltage and current waves.

The parameters of the line that will be investigated here are normalized. More precisely,

$$L^0 = 1\text{H/m}, \quad C^0 = 1\text{F/m}, \quad \delta = \log 4.$$

This corresponds to a 1:4 impedance stepping line. The waveform of the voltage source is set here to a gaussian pulse,

$$v_s(t) = V_0 e^{-\frac{(t-T_s)^2}{2\Delta_s^2}}, \quad (1.37)$$

with amplitude  $V_0 = 1\text{V}$ , center  $T_s = 2\text{ s}$  and width  $\Delta_s = 0.2\text{ s}$ . The source resistance will always be set in our simulations to be

$$R_S = Z_0 = 1\Omega,$$

while different values of the load resistance  $R_L$  will be investigated.

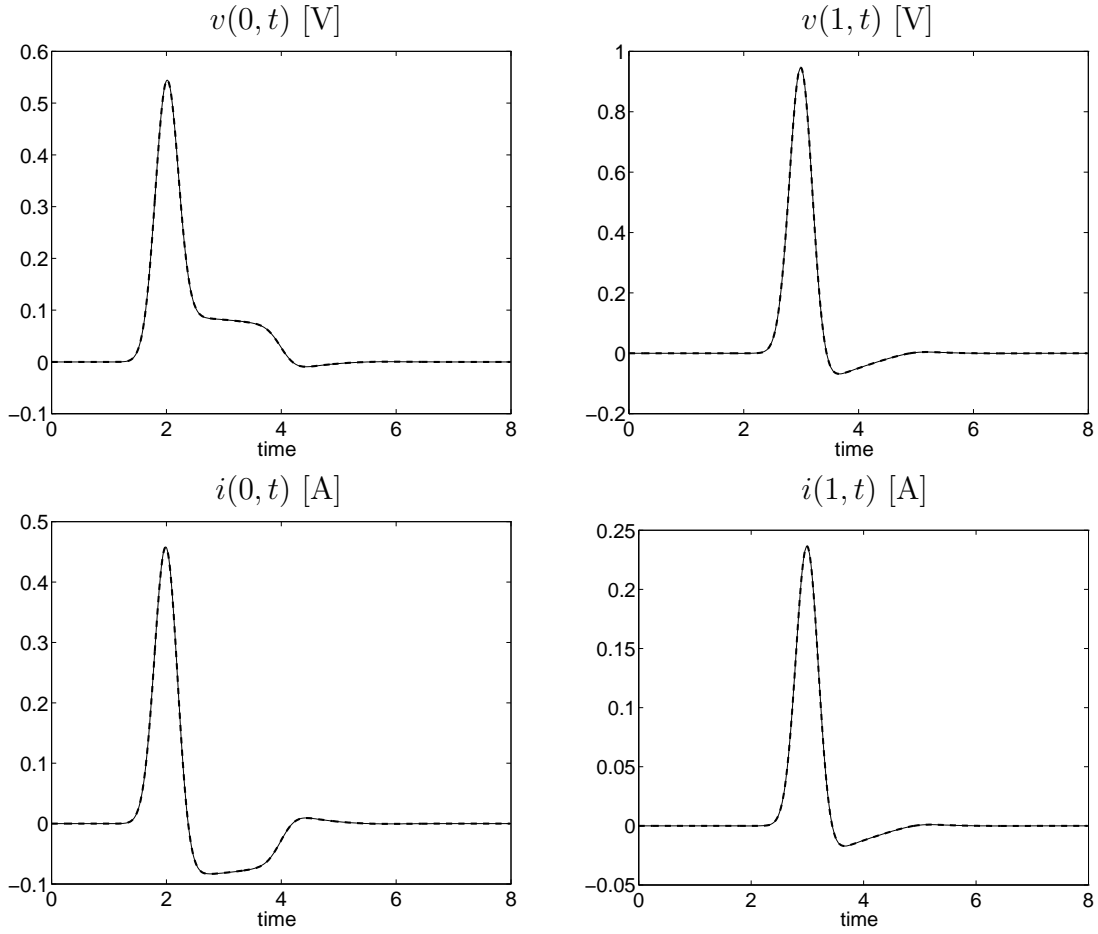
A reference solution in the time domain is obtained from the frequency domain analytical solution reported above through inverse FFT. The total simulation time is set here to  $T_{max} = 8\text{s}$ , which means that the input signal is considered as a periodic pulse train with period  $T_{max}$ . As the one-way delay is 1 s, there are no interactions between two adjacent pulses, because the transient associated to one pulse due to the nonuniformity of the line is already extinguished when the next pulse comes through. Of course, this holds only when at least one of the two ends of the line is matched.

The solution obtained with the weak formulation is plotted and compared to the reference solution in Figures 1.3, 1.4, and 1.5 for matched load, low impedance load and high impedance load, respectively. All these figures have been obtained with  $N = 65$ . The four panels report the voltages and currents at the two terminations. The figures show clearly that the weak solution is undistinguishable from the reference solution.

The numerical method was also tested for the matched exponential line with a trapezoidal pulse voltage source with amplitude 1V, time offset 1 s, rise and fall times 0.4 s and a duration at the 1V level of 3 s. The voltage and current at the two terminations of the matched 1:4 exponential line are depicted in Fig. 1.6 for both our method and the FFT reference solution. Also in this case the two curves are barely distinguishable.

The convergence properties of the method as the dimension  $N$  of the approximation spaces increases is now investigated. We fix the load resistance to  $R_L = 4\Omega$ , i.e. the line is matched at both edges. The approximation error on voltage and current is computed for each  $N$  according to

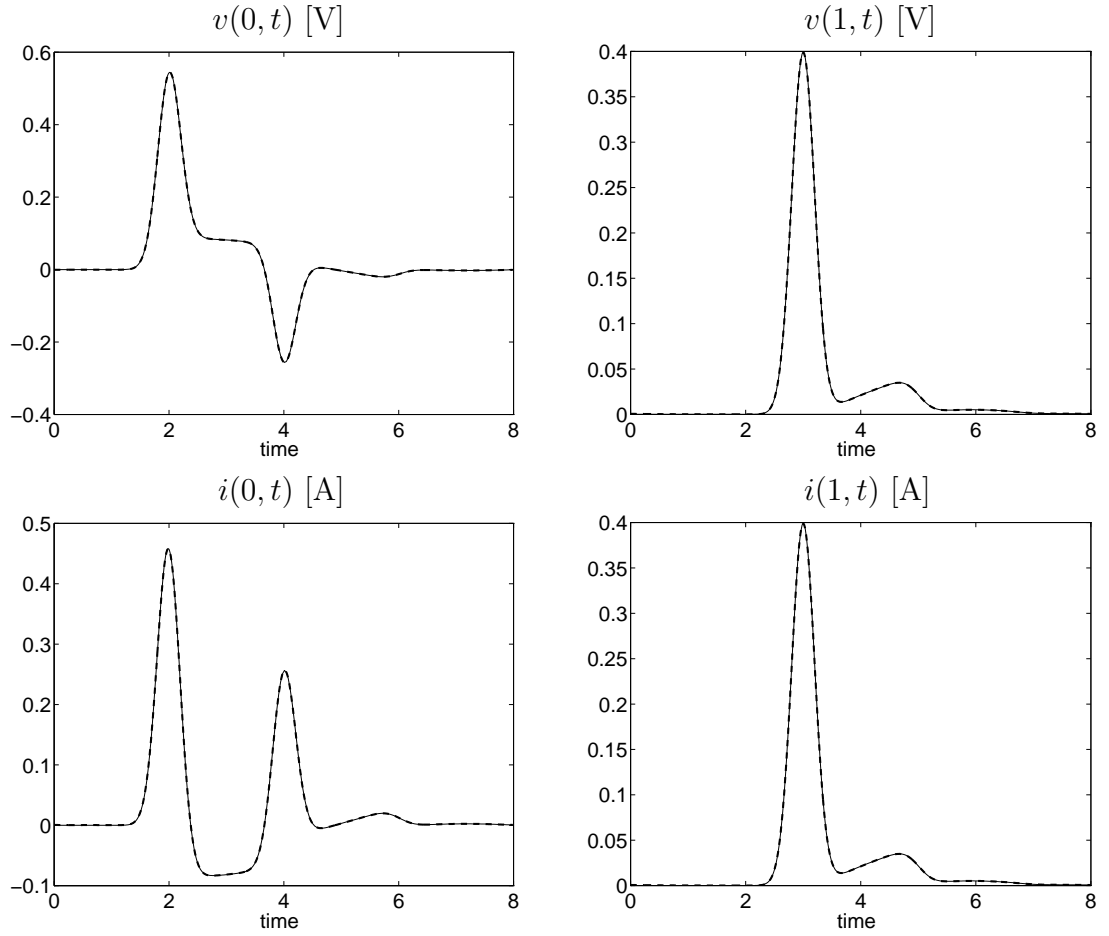
$$\begin{aligned} E_v(N) &= \max_t \max_z |v_N(z, t) - v_{ref}(z, t)|, \\ E_i(N) &= \max_t \max_z |i_N(z, t) - i_{ref}(z, t)|, \end{aligned} \quad (1.38)$$



**Figure 1.3:** Comparison between the weak solution (thin continuous line) and FFT reference solution (thick dashed line) for the 1:4 exponential line with matched load ( $R_L = 4\Omega$ ) and  $N = 65$ .

where  $v_{ref}(z, t)$ ,  $i_{ref}(z, t)$  represent the reference voltage and current while  $v_N(z, t)$ ,  $i_N(z, t)$  are the voltage and current obtained with our method. The approximation errors are reported in Figure 1.7 as functions of  $N$  for the gaussian pulse (left panel) and the trapezoidal pulse (right panel). As expected, the error decreases as  $O(N^{-2})$  as  $N$  increases for the gaussian pulse. The behavior of the error is instead of the type  $O(N^{-1})$  for the trapezoidal pulse due to the singularity in the first derivative of the source waveform.





**Figure 1.4:** As in Figure 1.3, but with  $R_L = 1 \Omega$  (low impedance load).

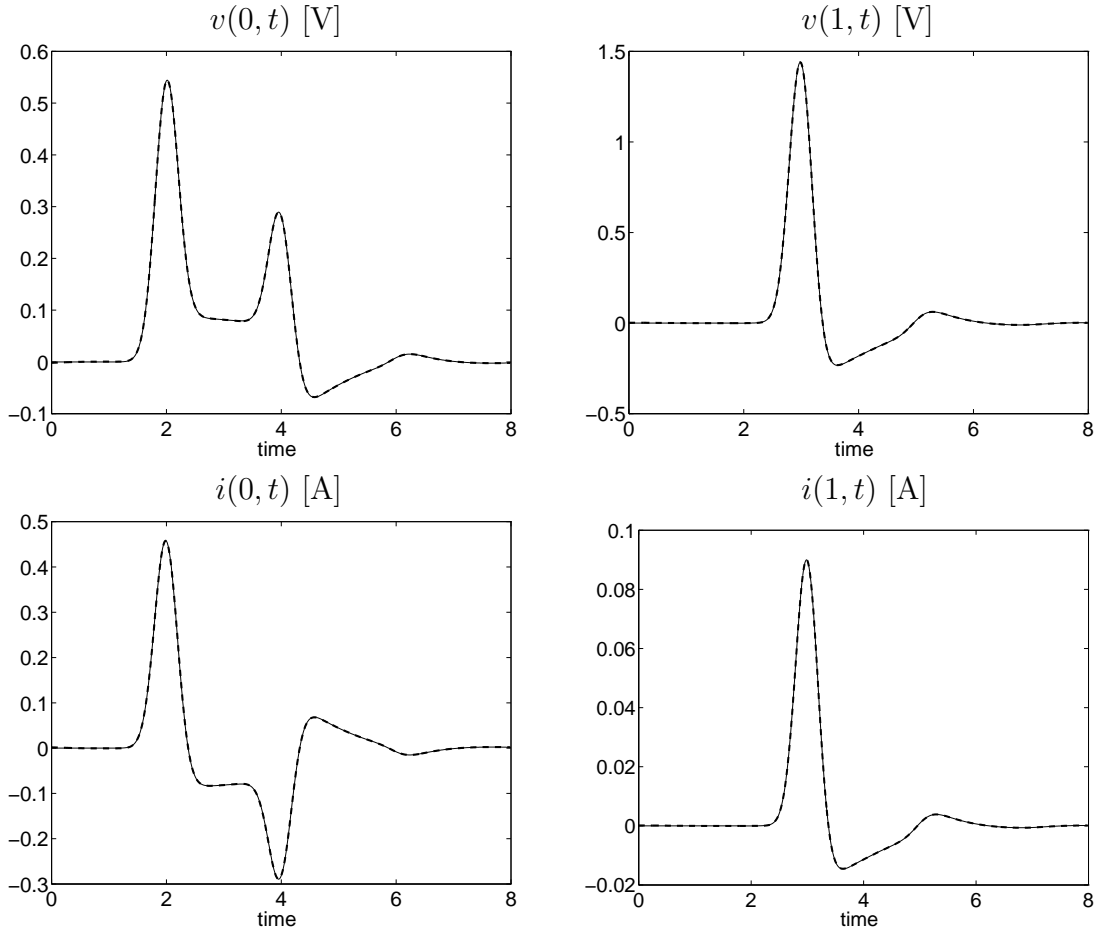
### 1.3.2 Three-conductor PCB

We will analyze in this section the coupled microstrip described in [7] and whose cross-section is shown below.



This structure consists of three rectangular conductors (of width 15 mils, height 1.38 mils, and separation 45 mils) placed above a glass epoxy ( $\epsilon_r = 4.7$ ) substrate 47 mils high. The length of the line is  $\mathcal{L} = 10$  inches. This structure is characterized by the per-unit-length matrices

$$\mathbf{L} = \begin{bmatrix} 1.10418 & 0.690094 \\ 0.690094 & 1.38019 \end{bmatrix} \mu\text{H/m}, \quad \mathbf{C} = \begin{bmatrix} 40.6280 & -20.3140 \\ -20.3140 & 29.7632 \end{bmatrix} \text{pF/m}.$$



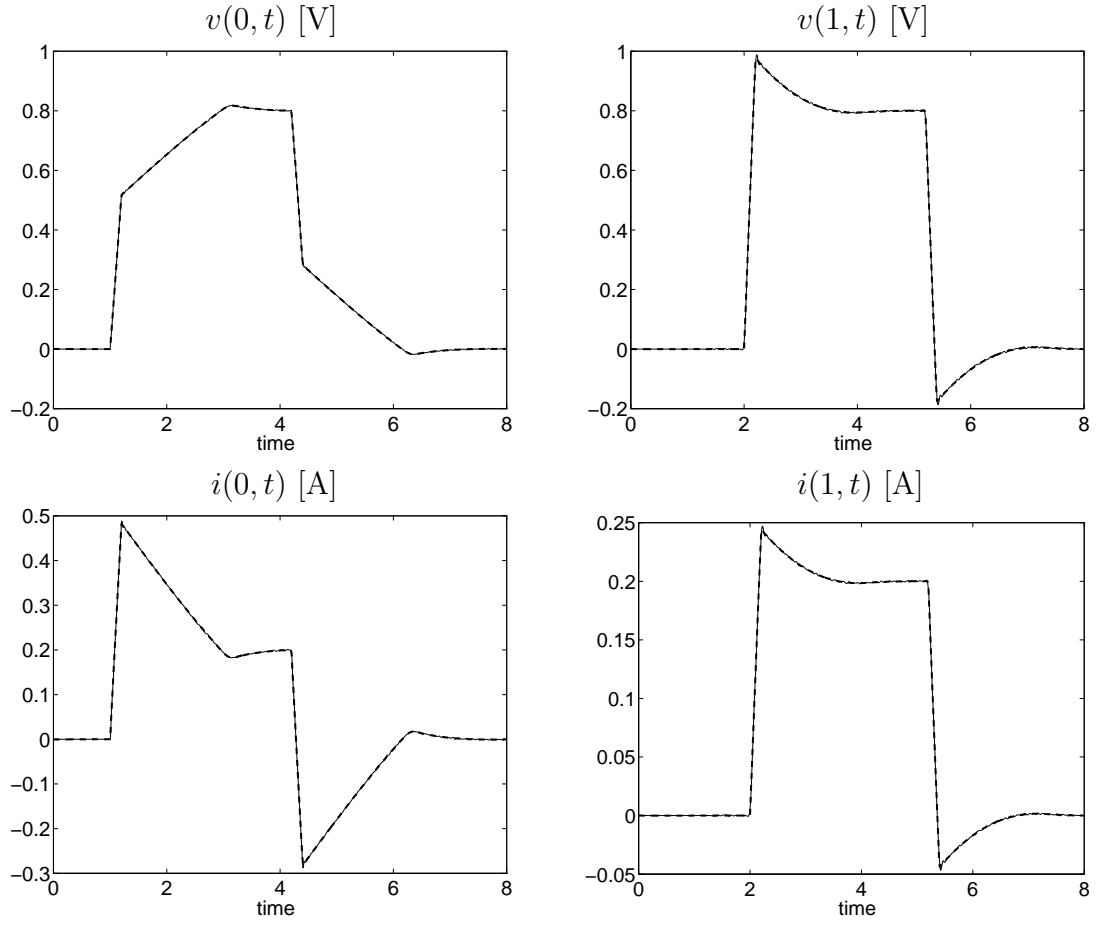
**Figure 1.5:** As in Figure 1.3, but with  $R_L = 16 \Omega$  (high impedance load).

A Thévenin voltage source is applied between the two side conductors, one of which is the reference, and the near-end crosstalk is determined on the middle conductor. The load matrices are

$$\mathbf{R}_S = \begin{bmatrix} 50 & 0 \\ 0 & 50 \end{bmatrix} \Omega, \quad \mathbf{R}_L = \begin{bmatrix} 50 & 0 \\ 0 & 50 \end{bmatrix} \Omega,$$

and the source waveform is a 1 MHz , 50% duty cycle trapezoidal pulse train with raise and fall times equal to 6.25 ns.

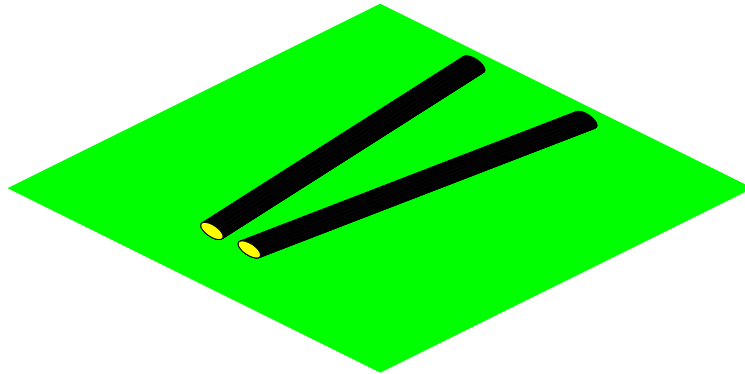
A reference solution is obtained in this case through inverse FFT from the exact solution in the frequency domain [7]. Figure 1.8 shows the magnitude of the voltage on the middle conductor obtained with our method (with  $N = 65$ ) compared with the reference solution (obtained using 2048 points in the evaluation of the inverse FFT). The two curves are almost undistinguishable.



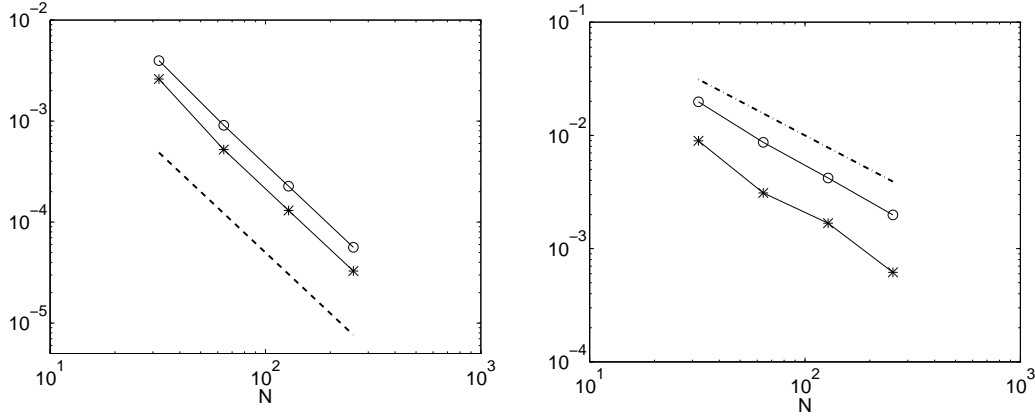
**Figure 1.6:** As in Figure 1.3, but with a trapezoidal pulse voltage source. The dimension of the approximation space is  $N = 65$ .

### 1.3.3 Nonparallel wires above a ground plane

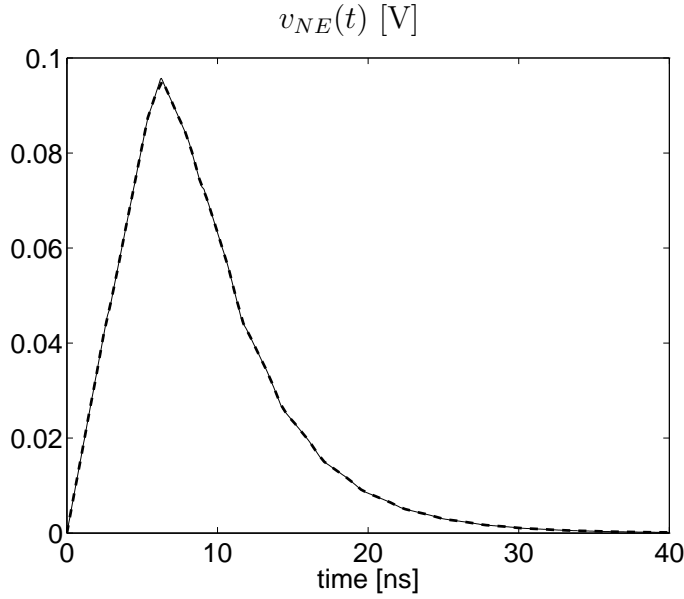
This section will examine the crosstalk on a nonuniform line made of two wires above a ground plane, sketched below.



The two wires are supposed to be parallel to the ground plane, but their distance



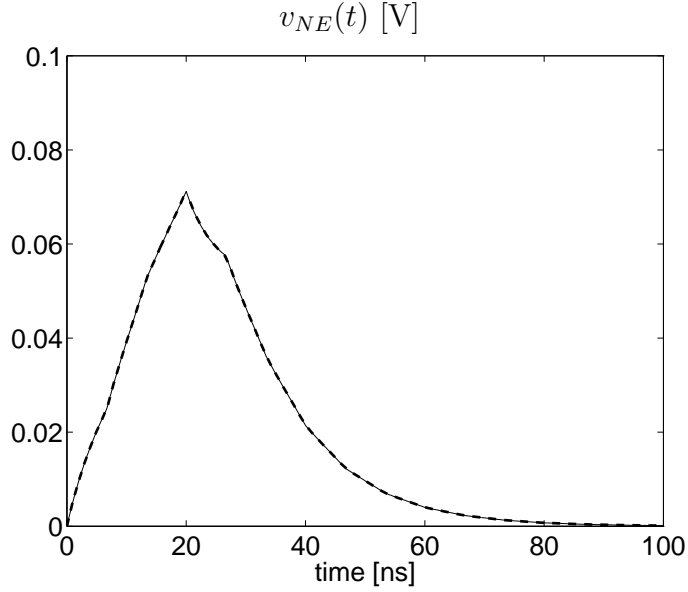
**Figure 1.7:** Approximation error on voltage  $E_v(N)$  (circles) and current  $E_i(N)$  (stars) as functions of  $N$  for the 1:4 exponential line with matched loads. The left and right panels refer to the gaussian and the trapezoidal pulse source, respectively. The dashed line corresponds to a slope  $N^{-2}$ , while the dash-dotted line corresponds to a slope  $N^{-1}$ .



**Figure 1.8:** Near end crosstalk for the PCB of Section 1.3.2. The continuous line represents the weak solution with  $N = 65$ , and the thick dashed line is the solution obtained through inverse FFT with 2048 points.

increases linearly along the length of the line. Lines of this type have been studied in [6, 2].

Let us consider two wires with radius  $r = 1$  mm placed at  $h = 3$  cm above a ground plane. Their separation is  $D_0 = 5$  mm at  $z = 0$  and  $D_1 = 15$  mm at  $z = 1$ . The medium is supposed to be free space. Due to these particular conditions,



**Figure 1.9:** Near end crosstalk for the structure described in Section 1.3.3. The continuous line represents the weak solution with  $N = 65$ . The thick dashed line is the solution obtained for the piecewise uniform line with 32 sections through inverse FFT (using 1024 points).

the expressions for the per-unit-length inductance and capacitance matrices can be obtained with the wide-separation approximation [7]. Taking the ground plane as the reference conductor, we have the approximate expressions

$$\begin{aligned} L_{11} &= L_{22} = \frac{\mu_0}{2\pi} \log \left( \frac{2h}{r} \right) \\ L_{12}(z) &= L_{21}(z) = \frac{\mu_0}{4\pi} \log \left( 1 + \frac{4h^2}{D^2(z)} \right), \end{aligned}$$

where  $\mu_0$  indicates the permeability of free space and the distance along the line is

$$D(z) = D_0 + \frac{z}{\mathcal{L}}(D_1 - D_0),$$

with  $\mathcal{L} = 1$  m. As the surrounding medium is homogeneous we can easily derive the per-unit-length capacitance matrix,

$$\mathbf{C}(z) = \varepsilon_0 \mu_0 \mathbf{L}^{-1}(z),$$

where  $\varepsilon_0$  is the permittivity of free space.

We will apply a voltage source consisting of a 1 MHz , 50% duty cycle trapezoidal pulse train with raise and fall times equal to 20 ns to one of the two wires at the edge  $z = 0$ , and calculate the near end crosstalk on the other wire.

The load matrices will be the same already used in Section 1.3.2, i.e. diagonal loads with  $50\ \Omega$  resistances.

As there is no closed form solution for lines of this type, we will have to use an approximate method to obtain a reference solution for this problem. The standard approach is to divide the line into  $N_z$  uniform subsections and to perform the analysis in the frequency domain [7, 2]. Each subsection is analyzed separately by deriving its chain matrix, which can be evaluated in closed form. The chain matrix of the overall structure is then obtained by multiplying the chain matrices of each subsection of the line, and the solution for the voltages and currents at the line ends is found by incorporating the terminal conditions. Finally, inverse FFT is applied to get the time domain waveform. This method converges to the exact solution when  $N_z$  increases. Some numerical tests on the convergence have been conducted to obtain the minimum number of subdivisions that insures a good approximation for the nonuniformity of the line. A number of  $N_z = 32$  subdivisions resulted beyond this limit and will be used in the following.

Figure 1.9 shows the results of the simulations with our method ( $N = 65$ ) and with the approximate piecewise uniform solution. The number of points for the evaluation of the inverse FFT in the latter was set to 1024. We notice that both methods give practically the same results.



# Chapter 2

## An introduction to multilevel decompositions

The TDSE method for the solution of the NMTL equations introduced in Chapter 1 is based on three sets of functions  $\{\zeta_n\}$ ,  $\{\phi_k\}$ , and  $\{\eta_m\}$ . However, only piecewise linear functions have been used to solve the examples presented in Sec. 1.3. On the other hand, the accuracy of any discretization method for PDEs is dependent on the particular choice of approximation spaces for the solution and the parameters, and consequently on their basis functions. The example in Section 1.3.1 showed that the decay of the approximation error with the number  $N$  of basis functions when piecewise linear functions are employed is at most  $O(N^{-2})$ . The rest of this work is devoted to the improvement of the TDSE method by using different approximation spaces, which allow a faster decay of the approximation errors with  $N$ , or equivalently a smaller error when  $N$  is fixed.

The key point in the choice of new basis sets is the determination of the characteristics of the solution to be represented. This solution must be represented with a small approximation error in an efficient way, i.e., with as few as possible basis functions. In the following we will show that the features of signals that are commonly found on transmission lines are well captured by the so-called *multilevel approximation spaces*. These spaces are obtained through multilevel decompositions of functional spaces like  $L^2(\mathbb{R})$ .

The purpose of this chapter is to introduce the multilevel decomposition of functional spaces. Section 2.1 will introduce qualitatively the properties of the simplest multilevel decomposition, the Haar system. An example will illustrate the representation of a function in terms of the canonical and hierarchical bases. We will derive empirically the two-scale relations, which are the milestone of any multilevel decomposition. In Section 2.2 the same aspect will be considered from an abstract point of view, by introducing projection operators that fully characterize multiresolution decompositions of a general space of functions  $V$ . The following two chapters will describe general biorthogonal multilevel decom-

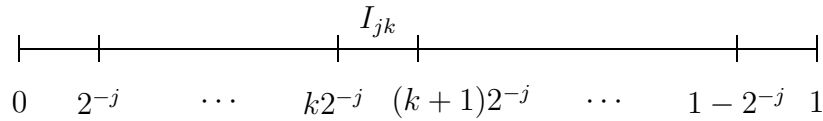


positions of  $L^2(\mathbb{R})$  (Chap. 3) and of  $L^2(\Omega)$ , where  $\Omega$  is a bounded domain like the unit interval  $(0, 1)$  (Chap. 4). This last construction is exactly what we need for the improvement of the TDSE method, as we will show in Chapter 6.

## 2.1 Multilevel representation of functions: a primer

In this section we consider the qualitative aspects of multilevel decompositions of functional spaces. Due to its simplicity, we will use the Haar decomposition of  $L^2(0, 1)$ . The aim is to illustrate through simple examples the basic properties of nested approximation spaces and the improvement that can be obtained by using hierarchical bases instead of canonical bases to represent a given function.

The Haar decomposition, dated back to 1910 [26], is based on piecewise constant approximations of a function  $f$  with domain in the unit interval. Let us suppose that the unit interval is subdivided in  $2^j$  intervals of length  $2^{-j}$  each. This corresponds to a collocation of separation points  $\{x_{jk} = k2^{-j}, k = 0, \dots, 2^j\}$ . We will work with approximation spaces  $V_j$  defined as the spaces of piecewise constant functions in any interval  $I_{jk} = [k2^{-j}, (k+1)2^{-j})$ ,  $k = 0, \dots, 2^j - 1$ .



We define now a function  $f_j \in V_j$ ,  $\forall j \geq 0$ , as the “closest” to the function  $f$  among all the functions in  $V_j$ . It is natural to define it as that particular function that minimizes the  $L^2$  norm,

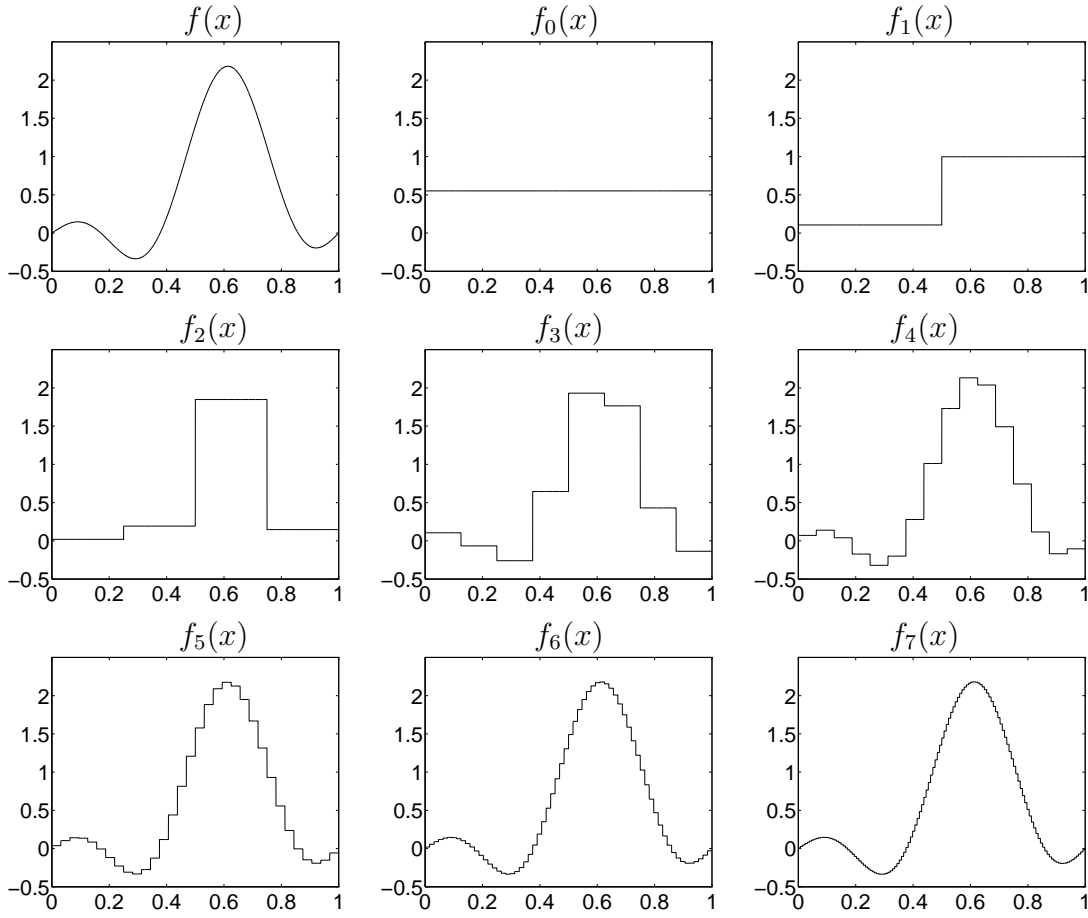
$$\|f - f_j\| = \min_{g_j \in V_j} \|f - g_j\|. \quad (2.1)$$

It is not difficult to prove that the approximation error tends to zero in  $L^2$  when the approximation level approaches to infinity,

$$\|f - f_j\| \rightarrow 0, \quad j \rightarrow +\infty. \quad (2.2)$$

Intuitively, even if the function  $f$  has fast variations, we can make the intervals  $I_{jk}$  small enough by increasing  $j$ , to reduce the approximation error below any fixed threshold  $\epsilon$ . Figure 2.1 shows a function  $f$  with its approximations  $f_j$  for some  $j$ . The figure clearly illustrates the fact that the approximation spaces are nested,

$$V_0 \subset \cdots \subset V_j \subset V_{j+1} \subset \cdots \subset L^2.$$



**Figure 2.1:** The top left panel shows a function  $f$  on the unit interval. The other panels show some piecewise constant approximations  $f_j$ .

We can interpret now the definition of the approximations  $f_j$  in an alternative way, through projection operators  $P_j : L^2 \rightarrow V_j$  such that

$$f_j = P_j f.$$

How do these operators act? As the functions  $f_j$  are constants on any subinterval  $I_{jk}$ , it is easy to see that the value that minimizes the norm in Eq. (2.1) is the average of  $f$  on  $I_{jk}$ ,

$$(P_j f)(x)|_{I_{jk}} = f_j(x)|_{I_{jk}} = \frac{1}{2^{-j}} \int_{I_{jk}} f(x) dx = c_{jk}. \quad (2.3)$$

The following interpretation of Eq. (2.3) is the basis for the definition of what we will call *scaling function spaces* in the next sections. Let us define the box function

$$\varphi(x) = \begin{cases} 1 & x \in [0, 1] \\ 0 & \text{otherwise,} \end{cases} \quad (2.4)$$

which is the indicator function of the unit interval. The  $L^2$  norm of  $\varphi$  is clearly equal to one. Let us now define a rescaled and dilated box function as

$$\varphi_{jk}(x) = 2^{j/2} \varphi(2^j x - k).$$

This function has still unitary  $L^2$  norm, and assumes the value  $2^{j/2}$  on any interval  $I_{jk}$ , being identically zero elsewhere. Any function  $f_j \in V_j$  can then be expressed as a finite superposition of these rescaled box functions, henceforth denoted as *scaling functions*, according to

$$f_j(x) = \sum_{k=0}^{2^j-1} f_{jk} \varphi_{jk}(x), \quad \forall x \in [0, 1). \quad (2.5)$$

The expansion coefficients (the *scaling function coefficients*) are such that

$$f_{jk} = 2^{-j/2} c_{jk} \quad (2.6)$$

on each subinterval  $I_{jk}$ . If we evaluate now the inner product between  $f$  and  $\varphi_{jk}$  we get

$$\langle f, \varphi_{jk} \rangle = \int_{I_{jk}} f(x) 2^{j/2} dx = 2^{-j/2} c_{jk} = f_{jk}.$$

Therefore, the series expansion can be expressed in a more abstract way as

$$f_j(x) = \sum_{k=0}^{2^j-1} \langle f, \varphi_{jk} \rangle \varphi_{jk}(x), \quad \forall x \in [0, 1), \quad (2.7)$$

where the functions  $\varphi_{jk}$  are orthonormal by construction,

$$\langle \varphi_{jk}, \varphi_{jl} \rangle = \delta_{kl}, \quad \forall j \geq 0.$$

In summary, we have characterized the projection operators  $P_j$  through an orthonormal basis  $\{\varphi_{jk}, k = 0, \dots, 2^j - 1\}$  of the approximation spaces  $V_j$ .

We address now the question on how we can obtain a coarse approximation  $f_j$  starting from the one at the immediately finer level  $j + 1$ . Recalling from Eq. (2.3) that the approximations are based on averages, and noting that any interval  $I_{jk}$  can be decomposed in the union of two intervals at a finer level according to

$$I_{jk} = I_{j+1,2k} \cup I_{j+1,2k+1},$$

we can easily see that

$$\begin{aligned} c_{jk} &= \frac{1}{2^{-j}} \int_{I_{jk}} f(x) dx \\ &= \frac{1}{2} \left[ \frac{1}{2^{-j-1}} \int_{I_{j+1,2k}} f(x) dx + \frac{1}{2^{-j-1}} \int_{I_{j+1,2k+1}} f(x) dx \right] \\ &= \frac{1}{2} (c_{j+1,2k} + c_{j+1,2k+1}). \end{aligned}$$

This means that a two-point average produces a coarser approximation of one level. Recalling the relation (2.6), we obtain a recurrence relation between the expansion coefficients at different levels,

$$f_{jk} = \frac{1}{\sqrt{2}} (f_{j+1,2k} + f_{j+1,2k+1}), \quad (2.8)$$

which can be also written as

$$f_{jk} = \sum_n h_n f_{j+1,2k+n}, \quad (2.9)$$

where

$$h_n = \begin{cases} 1/\sqrt{2} & \text{for } n = 0, 1 \\ 0 & \text{otherwise} \end{cases} \quad (2.10)$$

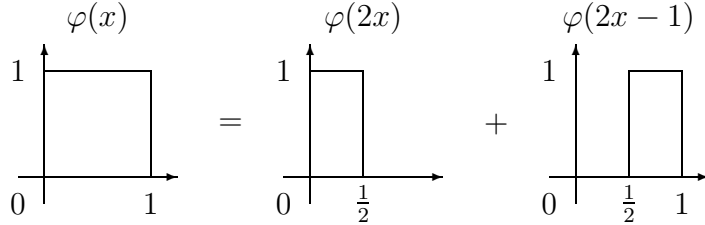
Eq. (2.9) can also be obtained directly from the basis functions  $\varphi_{jk}$ , by noting that they satisfy the so-called *two-scale relation*,

$$\varphi_{jk} = \sum_n h_n \varphi_{j+1,2k+n}, \quad (2.11)$$

which stems from

$$\varphi(x) = \varphi(2x) + \varphi(2x - 1).$$

This is illustrated graphically in the picture below.



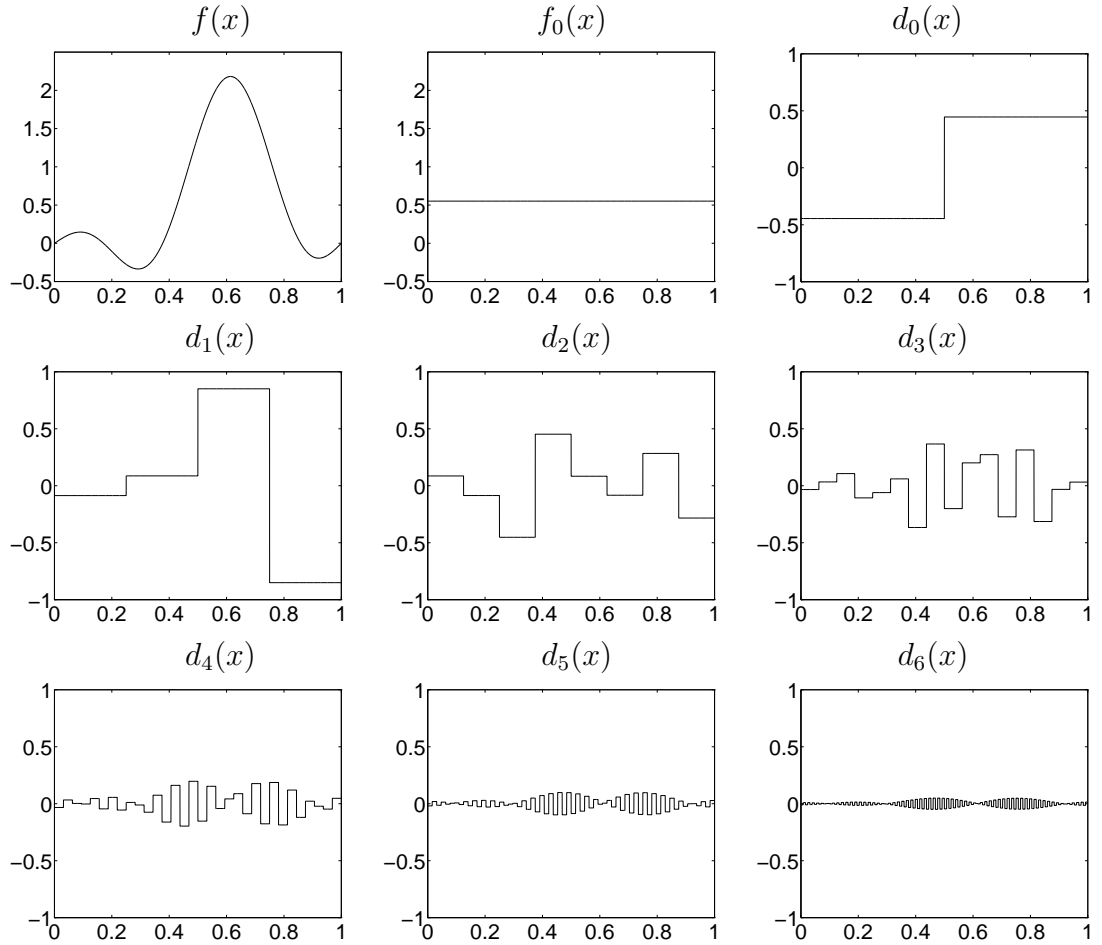
In the following sections Eq. (2.11) will be denoted as *refinement equation*, while the coefficients  $h_n$  will be denoted *filters*.

We want now to obtain the approximation  $f_{j+1}$  from the approximation at the coarser level  $f_j$ . An idea is to define what is the detail  $d_j$  we need to add to  $f_j$  in order to get a better approximation. This leads to another set of projection operators  $Q_j$ , obtained as

$$d_j = f_{j+1} - f_j = (P_{j+1} - P_j)f = Q_j f.$$

Clearly the function  $d_j$  belongs to the space  $V_{j+1}$ . We can also see that, having removed its component in  $V_j$ , the detail belongs to a complement space  $W_j$  such that

$$V_{j+1} = V_j \oplus W_j$$



**Figure 2.2:** Detail functions  $d_j = f_{j+1} - f_j$  for the example in Fig. 2.1. The function  $f$  is reported in the top left panel, and the coarsest approximation  $f_0$  is shown in the top middle panel.

The functions  $d_j$  for the example of Fig. 2.1 are reported in Fig. 2.2. Note that each function  $f_j$  can be recovered from an initial coarser approximation  $f_0$  by adding all the details at different levels

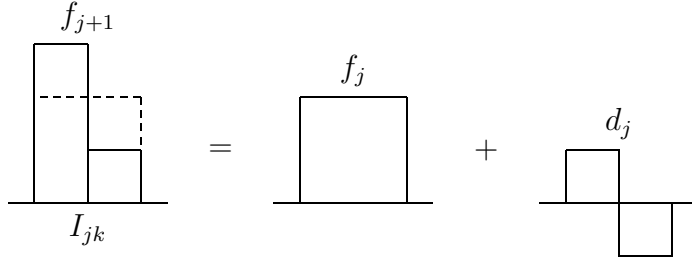
$$f_j = f_0 + \sum_{i=0}^{j-1} d_i.$$

Recalling the convergence relation (2.2), we can then write the *multilevel decomposition* of a function  $f$  as

$$f = f_0 + \sum_{i=0}^{\infty} d_i.$$

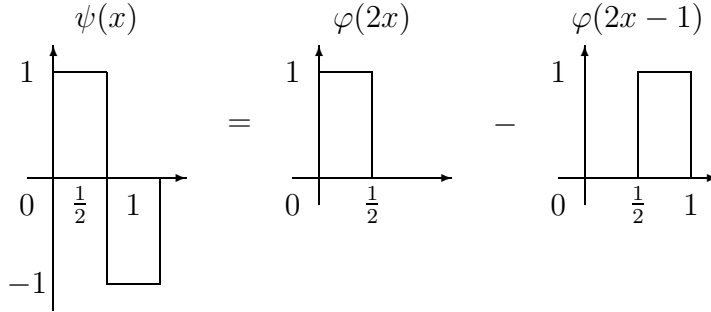
The above representation becomes useful when we are able to characterize the detail spaces  $W_j$  and the corresponding projection operators  $Q_j$  with a

suitable basis set. This can be easily accomplished if we recall that the approximations at any level  $j$  are obtained through averages. The detail function  $d_j$  on an interval  $I_{jk}$  can then be obtained from the approximation  $f_{j+1}$  by removing its average on  $I_{jk}$ . This average is exactly the approximation at level  $j$ , as shown in the picture below.



As a result, we see that the detail  $d_j$  can be represented on any subinterval  $I_{jk}$  as a zero-mean function that is constant in the two halves of the interval. It is then easy to express  $d_j$  as a superposition of zero-mean basis functions that we will call *wavelets*. These are obtained through dilations and translations of a single function  $\psi$ , the *mother wavelet*. In our case, the mother wavelet is defined as

$$\psi(x) = \varphi(2x) - \varphi(2x - 1). \quad (2.12)$$



The wavelets at any level  $j$  and location  $k$  are obtained as

$$\psi_{jk} = 2^{j/2} \psi(2^j x - k). \quad (2.13)$$

In addition, the wavelet at any  $j$  and  $k$  can then be expressed in terms of the scaling functions at level  $j + 1$  as

$$\psi_{jk} = \frac{1}{\sqrt{2}} (\varphi_{j+1,2k} - \varphi_{j+1,2k+1}),$$

which is analogous to the refinement equation (2.11) for the scaling function.

The detail function at level  $j$  will then be expressed as

$$d_j(x) = \sum_{k=0}^{2^j-1} w_{jk} \psi_{jk}(x), \quad (2.14)$$

where  $w_{jk}$  are the *wavelet coefficients*. It is easy to show that also the wavelets are orthonormal,

$$\langle \psi_{jk} \psi_{j'l} \rangle = \delta_{jj'} \delta_{kl}, \quad \forall j, j', k, l,$$

so that the series above can be expressed in a formal way as

$$d_j(x) = \sum_{k=0}^{2^j-1} \langle f_{j+1}, \psi_{jk} \rangle \psi_{jk}(x),$$

In addition, the wavelets are orthogonal to any scaling function at coarser levels,

$$\langle \psi_{jk}, \varphi_{j'l} \rangle = 0, \quad \forall j' \leq j.$$

In summary, we have constructed two different representations of the function  $f_{j+1}$ ,

$$f_{j+1}(x) = \sum_{k=0}^{2^{j+1}-1} f_{j+1,k} \varphi_{j+1,k}(x) \quad (2.15)$$

$$f_{j+1}(x) = \sum_{k=0}^{2^j-1} f_{j,k} \varphi_{j,k}(x) + \sum_{k=0}^{2^j-1} w_{j,k} \psi_{j,k}(x) \quad (2.16)$$

The first is the expansion in the *canonical* basis  $\{\varphi_{j+1,k}\}$ , while the second uses the *hierarchical* basis  $\{\{\varphi_{j,k}\} \cup \{\psi_{j,k}\}\}$ . The advantage in using the hierarchical basis is that the information already included in the approximation of the function at level  $j$  is reused at level  $j+1$ . The multilevel decomposition of the original  $L^2$  function  $f$  is then expressed in terms of the hierarchical basis functions as

$$f(x) = \sum_{k=0}^{2^{j_0}-1} f_{j_0,k} \varphi_{j_0,k} + \sum_{j \geq j_0} \sum_{k=0}^{2^j-1} w_{j,k} \psi_{j,k},$$

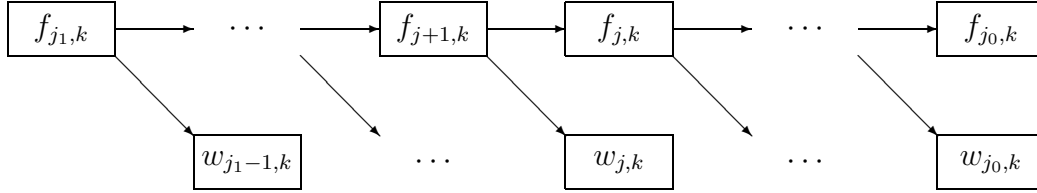
where  $j_0$  is an initial level.

The last step will be to show that also the coefficients  $w_{jk}$  can be evaluated in a recursive way through the approximation levels. If we combine the Equations (2.5), (2.8) and (2.14), we can derive the expression of the wavelet coefficients in terms of the scaling function coefficients at the next level,

$$w_{jk} = \frac{1}{\sqrt{2}} (f_{j+1,2k} - f_{j+1,2k+1}). \quad (2.17)$$

This equation, together with Eq. (2.8), constitutes the so-called *wavelet analysis*, i.e., the determination of wavelet coefficients starting from the approximation of the finest level  $j_1$  down to the coarsest level  $j_0$ . This operation can be visualized by the block diagram below, where the horizontal arrows use Eq. (2.8) and the

oblique arrows use Eq. (2.17).



It is important to note that the operations involved in the wavelet analysis are nothing else than convolutions with two filters (this will be made more precise in the following two chapters), together with a downsampling of a factor of two. Indeed, the number of scaling functions and wavelet coefficients at any step of the analysis is reduced by half. These operations can then be performed in  $O(N)$  operations, where  $N$  is the total number of initial data, and can be implemented as fast recursive algorithms. The wavelet analysis is then faster than the standard FFT, which requires  $O(N \log N)$  operations.

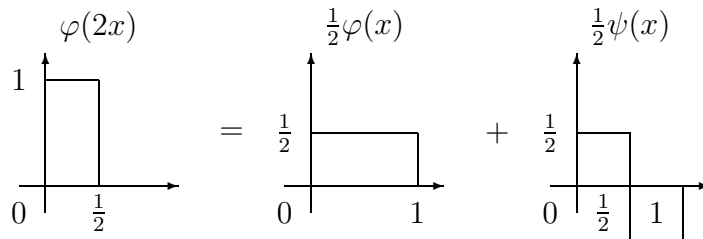
The inverse operation of the wavelet analysis is called *wavelet synthesis*, and corresponds to the summation of finer and finer details to the coarsest approximation at level  $j_0$  to get the approximation at the finest level  $j_1$ . It is easy to show that in the Haar case the coefficients  $f_{j+1,k}$  can be expressed in terms of the coefficients  $f_{j,k}$  and  $w_{j,k}$  as

$$f_{j+1,k} = \frac{1}{\sqrt{2}} \begin{cases} f_{j,l} + w_{j,l} & \text{if } k = 2l \\ f_{j,l} - w_{j,l} & \text{if } k = 2l + 1. \end{cases} \quad (2.18)$$

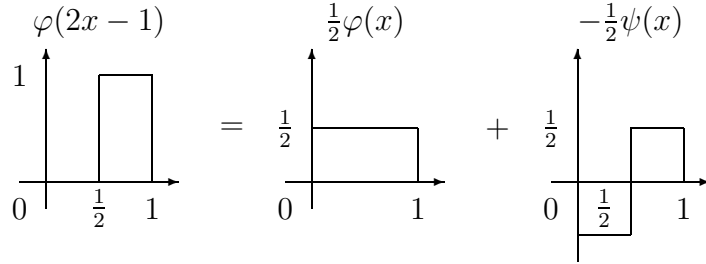
The proof of these relations is straightforward when we consider the expressions

$$\begin{aligned} \varphi(2x) &= \frac{1}{2}(\varphi(x) + \psi(x)), \\ \varphi(2x - 1) &= \frac{1}{2}(\varphi(x) - \psi(x)), \end{aligned}$$

which are graphically illustrated in the pictures below.







Note that the wavelet synthesis in Eq. (2.18) corresponds to merging at each level the scaling function and wavelet coefficients by placing their sum at even locations and their difference at odd locations in the finer level approximation. The diagram depicted above for the analysis is then valid also for the synthesis once the arrows are reversed.

In the next section we will generalize the foregoing example by constructing abstract multilevel decompositions. This generalization will be further detailed in Chapters 3 and 4, where multilevel decompositions with better approximation properties will be introduced. Basically, the generalization is due to a different choice for the filter  $h_n$ . All the properties of a given multilevel decomposition, like regularity and polynomials reproduction, will be determined only by the filter, provided it satisfies certain assumptions (see Sec. 3.1).

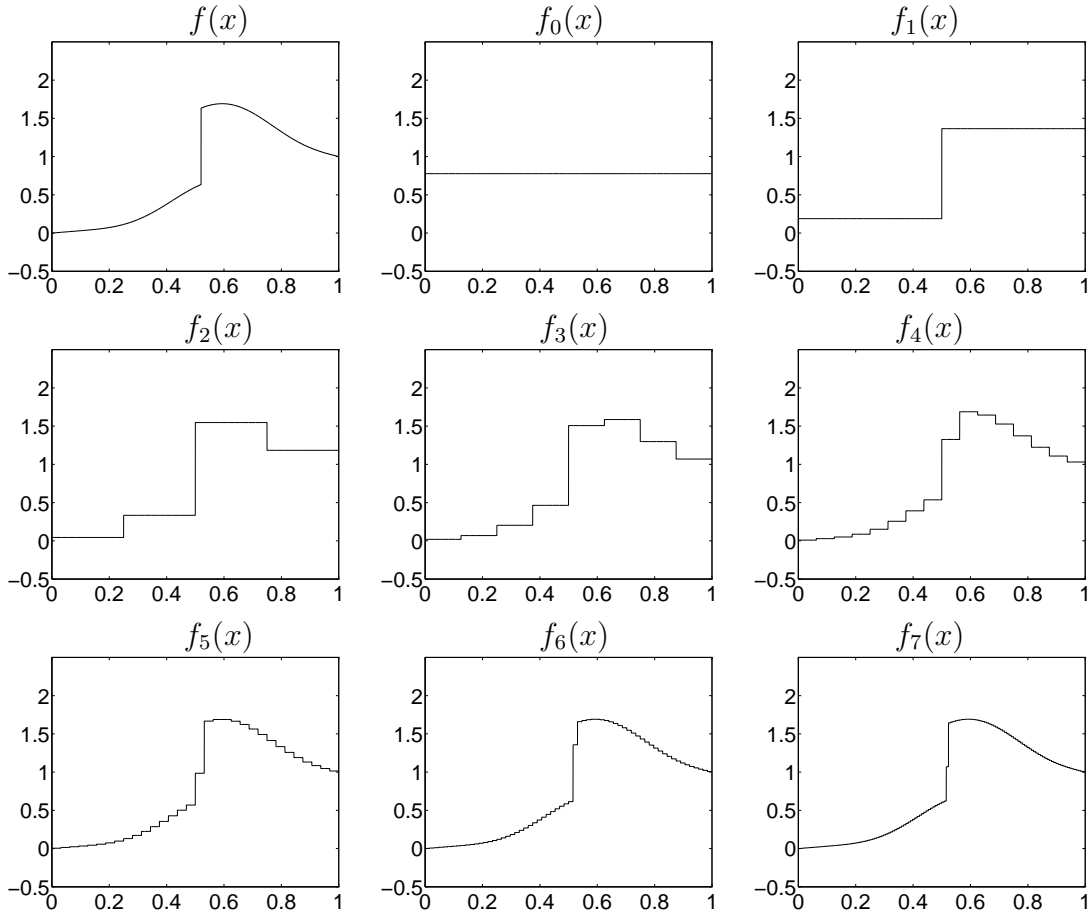
We conclude this section by showing with an example that the use of hierarchical bases with respect to canonical bases constitutes a very powerful technique for adaptive representation of functions and for data compression. This will be the feature used in the construction of numerical schemes for the solution of the NMTL equations (1.1)-(1.2), and is the main justification for the use of wavelets in this work.

Let us first begin to note that due to the  $L^2$  normalization of the basis functions, the wavelet coefficients  $w_{jk}$  can be related to the values assumed by the details  $d_j$  in the subintervals  $I_{jk}$  through multiplication by a factor of  $2^{j/2}$ . More precisely, if we indicate with  $|d_{jk}|$  the (constant) value assumed by the modulus of  $d_j$  in  $I_{jk}$ , we can see that

$$|w_{jk}| = 2^{-j/2}|d_{jk}|.$$

This expression is exactly the same as Eq. 2.6, which links scaling function coefficients and approximations  $f_j$  with the same level-dependent normalization constant. This shows that, apart from this normalization, which is the same for both  $c_{jk}$  and  $|d_{jk}|$ , we can interpret the plots in Figs. 2.1 and 2.2 as the expansion coefficients  $f_{jk}$  and  $w_{jk}$ .

In the mentioned example, dealing with a smooth function  $f$ , we can note that the absolute values of the wavelet coefficients have a faster decay than the scaling function coefficients when the level  $j$  increases. We will see in the following how this rate can be increased when more regular wavelets are used. This is due to Jackson-type inequalities, which are able to predict the rate of

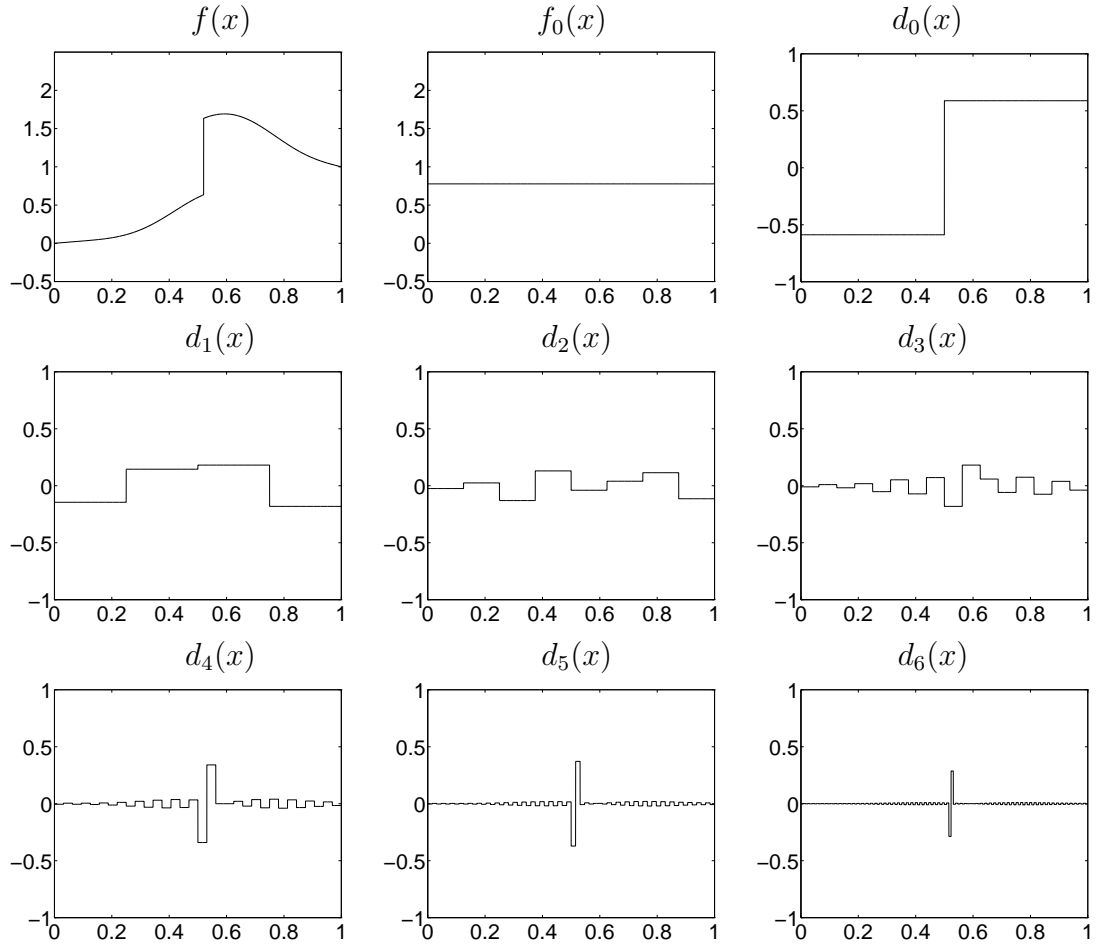


**Figure 2.3:** The top left panel shows a function  $f$  with a sharp discontinuity. The other panels show some piecewise constant approximations  $f_j$ .

decay of the wavelet coefficients as a function of the regularity of the function being analyzed.

We turn now to the representation of a function that presents a sharp discontinuity. Let us consider the function in the top left panel of Fig. 2.3. The approximations  $f_j$  and details  $d_j$  of this function at different levels are shown in Figs. 2.3 and 2.4, respectively.

These plots illustrate clearly that while all the scaling function coefficients must always be included to get a correct representation of  $f$ , only few wavelet coefficients need to be considered, especially for increasing  $j$ . The wavelet coefficients peak in the point where the discontinuity is located, and can be neglected elsewhere. Therefore, the total information needed to represent the function can be significantly reduced when using hierarchical representations by simply cutting away all coefficients  $w_{jk}$  below a certain threshold. This is obviously not possible using canonical representations. Examples will be provided in Sec. 4.5.2.



**Figure 2.4:** The top left panel shows the function  $f$  of Fig. 2.3. The top middle panel shows its coarsest approximation  $f_0$ , and the other panels show the details  $d_j$ .

## 2.2 Multilevel decompositions: the abstract setting

This section describes the abstract setting that underlies any construction of multilevel decompositions. This abstract setting can be applied to fairly general Banach or quasi-Banach spaces [16]. However, this work focuses on decompositions of Hilbert spaces, and the following sections will only deal with  $L^2$  spaces, or with Sobolev spaces  $H^s$  when regularity is needed. In any case, we will keep the working space unspecified in this section, and we will name it  $V$ . This will allow us to use the concepts described below to the multilevel decompositions of both  $L^2(\mathbb{R})$  and  $L^2([0, 1])$ .

We introduce a family of closed subspaces  $\{V_j\}_{j \in \mathcal{J}}$ , with  $\mathcal{J} = \mathbb{Z}$  or  $\mathbb{N}$  such that

$$V_j \subset V_{j+1}, \quad \forall j \in \mathcal{J}. \quad (2.19)$$

These spaces will be the approximation spaces in the multilevel decomposition, with the index  $j$  indicating the level of approximation. The approximation of a function  $v \in V$  can be defined through suitable linear projection operators  $P_j : V \rightarrow V_j$ , which must be bounded and satisfy the following two conditions

$$P_j v = v, \quad \forall v \in V_j, \quad (2.20)$$

$$P_j \circ P_{j+1} = P_j, \quad \forall j \in \mathcal{J}. \quad (2.21)$$

These conditions insure that  $P_j v$  is the best approximation of  $v$  in the space  $V_j$ , in the sense that the norm of the approximation error is minimized.

Let us define now another projection operator

$$Q_j v := P_{j+1} v - P_j v, \quad \forall v \in V. \quad (2.22)$$

This operator, which is linear and continuous, extracts from the function  $v$  the “details” to be added to its approximation at level  $j$  in order to get its approximation at level  $j + 1$ . Indeed, it can be shown that, indicating with  $W_j$  the closure of the range of this operator in  $V$ , we have

$$V_{j+1} = V_j \oplus W_j. \quad (2.23)$$

In other words, the approximation of a function at level  $j + 1$  can be obtained in two different ways. We can use the operator  $P_{j+1}$ , or we can refine the approximation at level  $j$  by using the correction obtained through the operator  $Q_j$ .

This procedure can be iterated: for any two indices  $j_0 < j_1$ , we can write

$$V_{j_1} = V_{j_0} \oplus \bigoplus_{j_0 \leq j < j_1} W_j. \quad (2.24)$$

In addition, if convergence holds according to

$$P_j v \rightarrow v \quad \text{for } j \rightarrow +\infty; \quad (2.25)$$

$$\begin{cases} P_j v \rightarrow 0 & \text{if } \mathcal{J} = \mathbb{Z} \\ P_0 v = 0 & \text{if } \mathcal{J} = \mathbb{N} \end{cases} \quad \text{for } j \rightarrow -\infty, \quad (2.26)$$

we get the multilevel decomposition of  $V$ ,

$$V = \bigoplus_{j \in \mathcal{J}} W_j, \quad \text{with} \quad v = \sum_{j \in \mathcal{J}} Q_j v, \quad \forall v \in V. \quad (2.27)$$

Let us now introduce basis sets  $\{\varphi_{jk} \mid k \in \check{\mathcal{K}}_j\}$  and  $\{\psi_{jk} \mid k \in \hat{\mathcal{K}}_j\}$  for  $V_j$  and  $W_j$  respectively, where  $\check{\mathcal{K}}_j, \hat{\mathcal{K}}_j$  are suitable sets of indices. From Eq. (2.23), we have two different basis sets for  $V_{j+1}$ . The first, called *canonical*, is  $\{\varphi_{j+1,k} \mid k \in$

$\check{\mathcal{K}}_{j+1}\}$ . The second, called *hierarchical*, is  $\{\varphi_{jk} \mid k \in \check{\mathcal{K}}_j\} \cup \{\psi_{jk} \mid k \in \hat{\mathcal{K}}_j\}$ . Each  $v \in V_{j+1}$  can then be represented in two ways,

$$v = \sum_{k \in \check{\mathcal{K}}_{j+1}} \check{v}_{j+1,k} \varphi_{j+1,k} = \sum_{k \in \check{\mathcal{K}}_j} \check{v}_{jk} \varphi_{jk} + \sum_{k \in \hat{\mathcal{K}}_j} \hat{v}_{jk} \psi_{jk}. \quad (2.28)$$

The evaluation of the expansion coefficients into one basis from the expansion coefficients into the other basis is the so-called two-scale analysis/synthesis, and corresponds to the transformation

$$\{\check{v}_{j+1,k}\}_{k \in \check{\mathcal{K}}_{j+1}} \leftrightarrow \{\{\check{v}_{jk}\}_{k \in \check{\mathcal{K}}_j}, \{\hat{v}_{jk}\}_{k \in \hat{\mathcal{K}}_j}\}. \quad (2.29)$$

According to Eq. (2.23), there are coefficients  $h_{km}^{(j)}$  ( $k \in \check{\mathcal{K}}_{j+1}$ ,  $m \in \check{\mathcal{K}}_j$ ) and  $g_{km}^{(j)}$  ( $k \in \check{\mathcal{K}}_{j+1}$ ,  $m \in \hat{\mathcal{K}}_j$ ) such that

$$\varphi_{jm} = \sum_{k \in \check{\mathcal{K}}_{j+1}} h_{km}^{(j)} \varphi_{j+1,k}, \quad \forall m \in \check{\mathcal{K}}_j, \quad (2.30)$$

$$\psi_{jm} = \sum_{k \in \check{\mathcal{K}}_{j+1}} g_{km}^{(j)} \varphi_{j+1,k}, \quad \forall m \in \hat{\mathcal{K}}_j. \quad (2.31)$$

On the other hand, there are also coefficients  $\chi_{mk}^{(j)}$  ( $m \in \check{\mathcal{K}}_j$ ,  $k \in \check{\mathcal{K}}_{j+1}$ ) and  $\gamma_{mk}^{(j)}$  ( $m \in \hat{\mathcal{K}}_j$ ,  $k \in \check{\mathcal{K}}_{j+1}$ ) such that

$$\varphi_{j+1,k} = \sum_{m \in \check{\mathcal{K}}_j} \chi_{mk}^{(j)} \varphi_{jm} + \sum_{m \in \hat{\mathcal{K}}_j} \gamma_{mk}^{(j)} \psi_{jm}, \quad \forall k \in \check{\mathcal{K}}_{j+1}. \quad (2.32)$$

The same coefficients can then be used to evaluate the transformation (2.29),

$$\check{v}_{jm} = \sum_{k \in \check{\mathcal{K}}_{j+1}} \chi_{mk}^{(j)} \check{v}_{j+1,k}, \quad \forall m \in \check{\mathcal{K}}_j, \quad (2.33)$$

$$\hat{v}_{jm} = \sum_{k \in \check{\mathcal{K}}_{j+1}} \gamma_{mk}^{(j)} \check{v}_{j+1,k}, \quad \forall m \in \hat{\mathcal{K}}_j, \quad (2.34)$$

$$\check{v}_{j+1,k} = \sum_{m \in \check{\mathcal{K}}_j} h_{km}^{(j)} \check{v}_{jm} + \sum_{m \in \hat{\mathcal{K}}_j} g_{km}^{(j)} \hat{v}_{jm}, \quad \forall k \in \check{\mathcal{K}}_{j+1}, \quad (2.35)$$

Note that this transformation implies that the approximation at level  $j$  can be obtained by coarsening (or *filtering*) the approximation at level  $j+1$ . On the other hand, the approximation at level  $j$  can be refined using the same approximation plus some details. Finally, the multilevel decomposition of any function  $v \in V$  can be expressed in the hierarchical basis as

$$v = \sum_{j \in \mathcal{J}} \sum_{k \in \hat{\mathcal{K}}_j} \hat{v}_{jk} \psi_{jk}. \quad (2.36)$$

The following two chapters will determine the conditions under which the abstract setting can be applied to insure convergence and stability of the canonical and hierarchical bases, and will give explicit rules for the determination of the coefficients of the above expressions.

# Chapter 3

## Biorthogonal decomposition of $L^2(\mathbb{R})$

This chapter introduces the biorthogonal multilevel decomposition of the space  $L^2(\mathbb{R})$ . This is essential for the construction of a multilevel decomposition of spaces of functions defined on bounded domains, like the solutions to the NMTL equations. This will be covered extensively in Chapter 4.

The construction of the scaling function and wavelet spaces is derived here from a restricted set of hypotheses. We will see that the properties of the decomposition are hidden in two sequences of real numbers which are denoted *filters* in the wavelet literature. If these filters satisfy certain conditions, then the multilevel decomposition, the approximation spaces, the canonical and the hierarchical bases are uniquely determined. In the following, we will detail only the main results, together with practical examples, and we will omit the proofs. Additional details and proofs can be found e.g. in Refs. [20, 17].

### 3.1 The basic axioms

Let us consider two sequences of real numbers  $\{h_n\}$  and  $\{\tilde{h}_n\}$ , with  $n \in \mathbb{Z}$ . The four axioms from which the whole construction starts are the following:

**M1.** *The two functions*

$$m_0(\xi) = \frac{1}{\sqrt{2}} \sum_{n=-\infty}^{\infty} h_n e^{-in\xi}, \quad \tilde{m}_0(\xi) = \frac{1}{\sqrt{2}} \sum_{n=-\infty}^{\infty} \tilde{h}_n e^{-in\xi}, \quad (3.1)$$

*are  $2\pi$ -periodic and belong to  $\mathcal{C}^r$  with  $r \geq 1$ ;*

**M2.**  *$m_0, \tilde{m}_0$  satisfy*

$$m_0(\xi) \overline{\tilde{m}_0(\xi)} + m_0(\xi + \pi) \overline{\tilde{m}_0(\xi + \pi)} = 1, \quad \forall \xi \in \mathbb{R}, \quad (3.2)$$

and

$$m_0(0) = \widetilde{m}_0(0) = 1, \quad m_0(\pi) = \widetilde{m}_0(\pi) = 0. \quad (3.3)$$

**M3.**  $m_0, \widetilde{m}_0$  have in  $\xi = \pi$  zeros of order  $L - 1$  and  $\widetilde{L} - 1$  ( $\leq r$ ), respectively, i.e.

$$\begin{aligned} \frac{d^l m_0}{d\xi^l}(\pi) &= 0 & \text{for each } l \in \mathbb{N} \text{ with } 0 \leq l \leq L - 1 \\ \frac{d^l \widetilde{m}_0}{d\xi^l}(\pi) &= 0 & \text{for each } l \in \mathbb{N} \text{ with } 0 \leq l \leq \widetilde{L} - 1. \end{aligned} \quad (3.4)$$

Note that  $m_0$  e  $\widetilde{m}_0$  can be factorized as

$$m_0(\xi) = \left( \frac{1 + e^{-i\xi}}{2} \right)^L \mathcal{F}(\xi) \quad \text{e} \quad \widetilde{m}_0(\xi) = \left( \frac{1 + e^{-i\xi}}{2} \right)^{\widetilde{L}} \widetilde{\mathcal{F}}(\xi) \quad (3.5)$$

where  $\mathcal{F}, \widetilde{\mathcal{F}}$  are  $2\pi$ -periodic.

**M4.** there are two integers  $\ell, \widetilde{\ell} > 0$  such that, setting

$$\max_{\xi} |\mathcal{F}(\xi) \cdot \dots \cdot \mathcal{F}(2^{\ell-1}\xi)| = 2^{\ell\tau} \quad \max_{\xi} |\widetilde{\mathcal{F}}(\xi) \cdot \dots \cdot \widetilde{\mathcal{F}}(2^{\widetilde{\ell}-1}\xi)| = 2^{\widetilde{\ell}\tilde{\tau}} \quad (3.6)$$

we have  $0 \leq \tau < L - \frac{1}{2}$ ,  $0 \leq \tilde{\tau} < \widetilde{L} - \frac{1}{2}$ . In the following we will set

$$\sigma = L - \frac{1}{2} - \tau > 0 \quad \text{e} \quad \tilde{\sigma} = \widetilde{L} - \frac{1}{2} - \tilde{\tau} > 0. \quad (3.7)$$

The condition M1 simply states that the filters are the Fourier coefficients of the functions  $m_0$  and  $\widetilde{m}_0$ . The condition M2 translates into the biorthogonality of the multiresolution spaces that we are going to build. The condition M3 will determine the properties of the approximation spaces, like the local reproduction of algebraic polynomials. The condition M4 determines the regularity of the basis functions for the approximation spaces, and determines consequently the spaces of functions that can be characterized. The necessity for these conditions will be cleared in the following sections. It should be noted that the orthogonal case can be easily derived from this more general setting by choosing  $\widetilde{h}_m = h_m$ ,  $\forall m$ . In this case, all the quantities with a tilde  $\sim$  coincide with the same quantities without tilde.

Some properties of the filters can be immediately derived from M2. We have the following identities

$$\sum_n h_n \widetilde{h}_{n-2k} = \delta_{k0}, \quad \forall k \in \mathbb{Z}, \quad (3.8)$$

$$\sum_n h_n = \sum_n \widetilde{h}_n = \sqrt{2}, \quad (3.9)$$

$$\sum_n (-1)^n h_n = \sum_n (-1)^n \widetilde{h}_n = 0. \quad (3.10)$$

## 3.2 Scaling function spaces in $\mathbb{R}$

The functions  $m_0$  and  $\widetilde{m}_0$  defined in the foregoing section allow to introduce the scaling functions  $\varphi$  and  $\widetilde{\varphi}$  in the Fourier domain,

$$\widehat{\varphi}(\xi) = \frac{1}{\sqrt{2\pi}} \prod_{j=1}^{\infty} m_0(2^{-j}\xi), \quad \widehat{\widetilde{\varphi}}(\xi) = \frac{1}{\sqrt{2\pi}} \prod_{j=1}^{\infty} \widetilde{m}_0(2^{-j}\xi). \quad (3.11)$$

From these definitions we can immediately derive the basic relations

$$\left\{ \begin{array}{l} \widehat{\varphi}(2\xi) = \widehat{\varphi}(\xi)m_0(\xi) \\ \widehat{\varphi}(0) = \frac{1}{\sqrt{2\pi}} \end{array} \right. \quad \text{and} \quad \left\{ \begin{array}{l} \widehat{\widetilde{\varphi}}(2\xi) = \widehat{\widetilde{\varphi}}(\xi)\widetilde{m}_0(\xi) \\ \widehat{\widetilde{\varphi}}(0) = \frac{1}{\sqrt{2\pi}} \end{array} \right. \quad (3.12)$$

that, when transformed back to the natural domain, lead to the well known refinement equations for the scaling functions,

$$\left\{ \begin{array}{l} \varphi(x) = \sqrt{2} \sum_{n \in \mathbb{Z}} h_n \varphi(2x - n) \\ \widetilde{\varphi}(x) = \sqrt{2} \sum_{n \in \mathbb{Z}} \widetilde{h}_n \widetilde{\varphi}(2x - n) \end{array} \right. \quad (3.13)$$

together with the normalization conditions

$$\int_{\mathbb{R}} \varphi(x) dx = 1 \quad \text{e} \quad \int_{\mathbb{R}} \widetilde{\varphi}(x) dx = 1. \quad (3.14)$$

It can be shown that the Fourier transforms of the two scaling functions are of class  $C^r$  on  $\mathbb{R}$ , and decay at infinity according to

$$|\widehat{\varphi}(\xi)| \leq C(1 + |\xi|)^{-1/2-\sigma} \quad (3.15)$$

$$|\widehat{\widetilde{\varphi}}(\xi)| \leq C(1 + |\xi|)^{-1/2-\widetilde{\sigma}}, \quad (3.16)$$

where  $\sigma, \widetilde{\sigma}$  are the same as in M4. Therefore, recalling the characterization of the Sobolev spaces in Eq. (9), we can immediately link the axiom M4 to the regularity of the scaling functions through these two exponents. We obtain  $\varphi \in H^s, \forall s < \sigma$  and  $\widetilde{\varphi} \in H^s, \forall s < \widetilde{\sigma}$ . In addition, this fact insures that  $\varphi$  and  $\widetilde{\varphi}$  belong to the Hilbert space  $L^2(\mathbb{R})$ .

A biorthogonality relation can be derived from M2. More precisely, we have

$$\langle \varphi, \widetilde{\varphi}(\cdot - k) \rangle = \delta_{0k}. \quad (3.17)$$

This means that the scaling function  $\varphi$  is orthogonal in  $L^2(\mathbb{R})$  to the integer translates of the dual scaling function  $\widetilde{\varphi}$ , except in the case  $k = 0$ .

Let us now introduce the translated and dilated versions of the scaling functions, that we will indicate as

$$\varphi_{jk}(x) = 2^{j/2} \varphi(2^j x - k) \quad \text{and} \quad \widetilde{\varphi}_{jk}(x) = 2^{j/2} \widetilde{\varphi}(2^j x - k), \quad (3.18)$$



with  $k, j$  in  $\mathbb{Z}$ . From Eq. (3.17) we see that the biorthogonality holds also at any level  $j$ ,

$$\langle \varphi_{jk}, \tilde{\varphi}_{jl} \rangle = \delta_{kl}. \quad (3.19)$$

This implies that the sets  $\{\varphi_{jk} | k \in \mathbb{Z}\}$  and  $\{\tilde{\varphi}_{jk} | k \in \mathbb{Z}\}$  are constituted by linearly independent functions, and allows to introduce the approximation spaces

$$V_j = \text{span}_{L^2(\mathbb{R})} \{\varphi_{jk} | k \in \mathbb{Z}\}, \quad (3.20)$$

$$\tilde{V}_j = \text{span}_{L^2(\mathbb{R})} \{\tilde{\varphi}_{jk} | k \in \mathbb{Z}\}. \quad (3.21)$$

Note that it is possible to introduce the spaces  $V_j$  and  $\tilde{V}_j$  from  $V_0$  and  $\tilde{V}_0$  by using the isometry  $T_j : L^2(\mathbb{R}) \rightarrow L^2(\mathbb{R})$  defined as

$$T_j f(\cdot) = 2^{j/2} f(2^j \cdot), \quad \forall f \in L^2(\mathbb{R}). \quad (3.22)$$

We have then

$$\begin{aligned} V_j &= T_j V_0 \\ \tilde{V}_j &= T_j \tilde{V}_0. \end{aligned}$$

Moreover, it can be easily proved that the spaces  $V_j$  and  $\tilde{V}_j$  are generated by linear combinations of the basis functions  $\{\varphi_{jk}\}$  and  $\{\tilde{\varphi}_{jk}\}$ , with coefficients in  $\ell^2$ , and that these bases are uniformly 2-stable,

$$V_j = \left\{ \sum_k \alpha_k \varphi_{jk} \mid \{\alpha_k\} \in \ell^2 \right\} \quad (3.23)$$

$$C_1 \|\{\alpha_k\}\|_{\ell^2} \leq \left\| \sum_k \alpha_k \varphi_{jk} \right\|_{L^2(\mathbb{R})} \leq C_2 \|\{\alpha_k\}\|_{\ell^2}. \quad (3.24)$$

This allows to show also that

$$v(x) \in V_j \iff v(x - 2^{-j}k) \in V_j \quad \forall k \in \mathbb{Z} \quad (3.25)$$

$$v(x) \in V_j \iff v(2^{-j}x) \in V_0. \quad (3.26)$$

In summary, using the biorthogonality relation (3.19), we have the representation of any element of  $V_j$  and  $\tilde{V}_j$  as a superposition of basis functions, with coefficients obtained through inner products with the corresponding dual functions,

$$v = \sum_k \langle v, \tilde{\varphi}_{jk} \rangle \varphi_{jk}, \quad \forall v \in V_j \quad (3.27)$$

$$\tilde{v} = \sum_k \langle \tilde{v}, \varphi_{jk} \rangle \tilde{\varphi}_{jk}, \quad \forall \tilde{v} \in \tilde{V}_j. \quad (3.28)$$

In Section 2.2 we showed that the multilevel decompositions of spaces of functions are based on sequences of encapsulated spaces, which characterize

better and better approximations when the level  $j$  tends to infinity. The construction of scaling functions and approximation spaces of this section provides the framework in which a multilevel decomposition can be built. Indeed, the approximation spaces  $V_j$  and  $\tilde{V}_j$  satisfy the embedding relations

$$V_j \subset V_{j+1}, \quad \tilde{V}_j \subset \tilde{V}_{j+1}$$

for all  $j \in \mathbb{Z}$ . This is evident by restating the refinement equations (3.13) and interpreting the right sides as functions of  $V_1$ . At any level  $j$  we get

$$\varphi_{jk} = \sum_n h_n \varphi_{j+1,2k+n}, \quad \tilde{\varphi}_{jk} = \sum_n \tilde{h}_n \tilde{\varphi}_{j+1,2k+n}. \quad (3.29)$$

As introduced in the abstract setting, the approximation of any function  $v \in V = L^2(\mathbb{R})$  at level  $j$  can be obtained through suitable projection operators. It is natural to define these operators  $P_j : L^2 \rightarrow V_j$  and  $\tilde{P}_j : L^2 \rightarrow \tilde{V}_j$  as

$$P_j v = \sum_k \langle v, \tilde{\varphi}_{jk} \rangle \varphi_{jk} \quad (3.30)$$

$$\tilde{P}_j v = \sum_k \langle v, \varphi_{jk} \rangle \tilde{\varphi}_{jk}. \quad (3.31)$$

The resulting operators are linear, bounded, and satisfy the conditions (2.20) and (2.21). Moreover, the projections are biorthogonal in the sense that

$$\langle v - P_j v, \tilde{v} \rangle = 0, \quad \forall \tilde{v} \in \tilde{V}_j \quad (3.32)$$

$$\langle \tilde{v} - \tilde{P}_j \tilde{v}, v \rangle = 0, \quad \forall v \in V_j \quad (3.33)$$

The projection operators  $P_j$  and  $\tilde{P}_j$  are the adjoints of each other, because

$$\langle P_j v, \tilde{v} \rangle = \sum_k \langle v, \tilde{\varphi}_{jk} \rangle \langle \tilde{v}, \varphi_{jk} \rangle = \langle v, \tilde{P}_j \tilde{v} \rangle. \quad (3.34)$$

for all  $v, \tilde{v} \in L^2(\mathbb{R})$ . It can also be proved that the convergence for  $j \rightarrow -\infty$  and  $j \rightarrow +\infty$  holds for each  $v \in L^2(\mathbb{R})$ ,

$$\begin{cases} P_j v \rightarrow v \\ \tilde{P}_j v \rightarrow v \end{cases} \quad \text{for } j \rightarrow +\infty$$

$$\begin{cases} P_j v \rightarrow 0 \\ \tilde{P}_j v \rightarrow 0 \end{cases} \quad \text{for } j \rightarrow -\infty$$

Moreover, stronger conditions hold when the functions being analyzed are sufficiently regular. As we are dealing with Hilbert spaces in this work, we will mention here a particular form of a general result, the Jackson inequality, which holds for more general Banach spaces like Besov spaces [17]. If the scaling function  $\varphi$  belongs to a Sobolev space  $H^{s_0}$ , it can be shown that, for all  $s < \min(s_0, L)$ ,

$$\|v - P_j v\|_{L^2(\mathbb{R})} \lesssim 2^{-js} |v|_{H^s}. \quad (3.35)$$

This means that the rate of convergence for the approximation  $P_j v$  as  $j$  increases is of exponential type, and is controlled by the regularity index of the function being analyzed, as long as the scaling functions are sufficiently regular.

It should be noted that in the orthogonal case the sequence of subspaces  $V_j$  determines a Multiresolution Analysis (MRA) in the sense introduced by Meyer [28] and Mallat [27], i.e., the following five conditions are verified,

1.  $\cdots \subset V_{j-1} \subset V_j \subset V_{j+1} \subset \cdots$ ;
2.  $\bigcap_j V_j = \{0\}$ ,  $\bigcup_j V_j$  is dense in  $L^2(\mathbb{R})$ ;
3. for each  $v \in L^2(\mathbb{R})$  and  $j \in \mathbb{Z}$ ,  $v(x) \in V_j \iff v(2x) \in V_{j+1}$ ;
4. for each  $v \in L^2(\mathbb{R})$  and  $k \in \mathbb{Z}$ ,  $v(x) \in V_0 \iff v(x - k) \in V_0$ ;
5. there exists  $\varphi \in V_0$  such that  $\{\varphi(\cdot - k)\}_{k \in \mathbb{Z}}$  is an orthonormal basis of  $V_0$ .

### 3.3 Wavelet spaces in $\mathbb{R}$

This section will introduce the wavelet spaces, following the same guidelines of the abstract setting in Section 2.2. Recalling the definition of the detail operators  $Q_j$  (and similarly for the dual system  $\tilde{Q}_j$ ), in Eq. (2.22), we can define the detail subspaces  $W_j$  and  $\tilde{W}_j$  as

$$W_j := \text{Im } Q_j, \quad \text{and} \quad \tilde{W}_j := \text{Im } \tilde{Q}_j.$$

Using the additional biorthogonality properties listed in Section 3.2, we are led to the decomposition

$$V_{j+1} = V_j \oplus W_j \quad W_j \perp \tilde{V}_j \quad (3.36)$$

$$\tilde{V}_{j+1} = \tilde{V}_j \oplus \tilde{W}_j \quad \tilde{W}_j \perp V_j. \quad (3.37)$$

In other words, the detail spaces are such that

$$W_j = \{v \in V_{j+1} \mid \langle v, \tilde{v} \rangle = 0, \forall \tilde{v} \in \tilde{V}_j\} \quad (3.38)$$

$$\tilde{W}_j = \{\tilde{v} \in \tilde{V}_{j+1} \mid \langle \tilde{v}, v \rangle = 0, \forall v \in V_j\} \quad (3.39)$$

We will now determine a basis for  $W_j$  and  $\tilde{W}_j$ . Only the case for  $j = 0$  needs to be explicitly studied, because we can obtain the spaces at any level  $j$  through  $W_j = T_j W_0$  and  $\tilde{W}_j = T_j \tilde{W}_0$ , where  $T_j$  is the isometry defined in Eq. (3.22). We omit here the proof that leads to the definition of the *mother wavelet* and its dual in the Fourier domain,

$$\hat{\psi}(\xi) = -e^{-i\xi/2} \overline{\hat{m}_0(\xi/2 + \pi)} \hat{\varphi}(\xi/2) \quad (3.40)$$

$$\hat{\tilde{\psi}}(\xi) = -e^{-i\xi/2} \overline{\hat{m}_0(\xi/2 + \pi)} \hat{\varphi}(\xi/2). \quad (3.41)$$

Passing to the natural domain through Fourier inversion we obtain

$$\psi(x) = \sqrt{2} \sum_n g_n \varphi(2x - n) \quad (3.42)$$

$$\tilde{\psi}(x) = \sqrt{2} \sum_n \tilde{g}_n \tilde{\varphi}(2x - n), \quad (3.43)$$

where the filters  $g_n$  and  $\tilde{g}_n$  are defined as

$$g_n = (-1)^n \tilde{h}_{1-n} \quad \text{and} \quad \tilde{g}_n = (-1)^n h_{1-n}. \quad (3.44)$$

We have then obtained two 2-stable bases for the detail spaces, which can be characterized as

$$W_0 = \left\{ \sum_k \alpha_k \psi(\cdot - k) \mid \{\alpha_k\} \in \ell^2 \right\} \quad (3.45)$$

$$\tilde{W}_0 = \left\{ \sum_k \alpha_k \tilde{\psi}(\cdot - k) \mid \{\alpha_k\} \in \ell^2 \right\} \quad (3.46)$$

The biorthogonality of the approximation spaces  $V_0, \tilde{V}_0$  with the detail spaces  $\tilde{W}_0$  and  $W_0$ , respectively, of Eqs. (3.36) and (3.37) can be restated in terms of the basis functions as

$$\langle \psi(\cdot - k), \tilde{\varphi}(\cdot - l) \rangle = 0 \quad \forall k, l \in \mathbb{Z} \quad (3.47)$$

$$\langle \tilde{\psi}(\cdot - k), \varphi(\cdot - l) \rangle = 0 \quad \forall k, l \in \mathbb{Z}, \quad (3.48)$$

and in terms of the filters as

$$\begin{aligned} \sum_n g_n \tilde{h}_{n-2k} &= 0 \\ \sum_n \tilde{g}_n h_{n-2k} &= 0 \end{aligned} \quad \forall k \in \mathbb{Z}. \quad (3.49)$$

We also have the following relation between the filters  $h, \tilde{h}, g, \tilde{g}$ ,

$$\sum_m [h_{k-2m} \tilde{h}_{n-2m} + g_{k-2m} \tilde{g}_{n-2m}] = \delta_{kn}, \quad \forall k, n \in \mathbb{Z}. \quad (3.50)$$

In addition, the wavelet bases are biorthogonal, so that

$$\langle \psi(\cdot - k), \tilde{\psi} \rangle = \delta_{0k}, \quad \forall k \in \mathbb{Z}, \quad (3.51)$$

or equivalently

$$\sum_n g_n \tilde{g}_{n-2k} = \delta_{0k}, \quad \forall k \in \mathbb{Z}. \quad (3.52)$$

If we introduce, as usual, the wavelet functions at level  $j$ ,

$$\begin{aligned} \psi_{jk}(x) &= 2^{j/2} \psi(2^j x - k) \\ \tilde{\psi}_{jk}(x) &= 2^{j/2} \tilde{\psi}(2^j x - k), \end{aligned}$$

we obtain the following characterization for the detail spaces,

$$W_j = \left\{ \sum_k \alpha_k \psi_{jk} \mid \{\alpha_k\} \in \ell^2 \right\} \quad (3.53)$$

$$\widetilde{W}_j = \left\{ \sum_k \alpha_k \widetilde{\psi}_{jk} \mid \{\alpha_k\} \in \ell^2 \right\}. \quad (3.54)$$

The bases  $\{\psi_{jk} \mid k \in \mathbb{Z}\}$  and  $\{\widetilde{\psi}_{jk} \mid k \in \mathbb{Z}\}$  are 2-stable, and satisfy the more general biorthogonality relation

$$\langle \psi_{jk}, \widetilde{\psi}_{j'l} \rangle = \delta_{jj'} \delta_{kl}, \quad \forall j, j', k, k' \in \mathbb{Z}. \quad (3.55)$$

As a consequence, we are able to express the detail operators  $Q_j$  and  $\widetilde{Q}_j$  in terms of the basis functions,

$$Q_j v = \sum_k \langle v, \widetilde{\psi}_{jk} \rangle \psi_{jk} \quad (3.56)$$

$$\widetilde{Q}_j v = \sum_k \langle v, \psi_{jk} \rangle \widetilde{\psi}_{jk}. \quad (3.57)$$

The biorthogonal multilevel decomposition of  $L^2(\mathbb{R})$  is then simply expressed, for any  $v \in L^2(\mathbb{R})$ , as

$$v = \sum_{j,k} \langle v, \widetilde{\psi}_{jk} \rangle \psi_{jk} = \sum_{j,k} \langle v, \psi_{jk} \rangle \widetilde{\psi}_{jk}, \quad (3.58)$$

with the stability condition

$$\|v\|_{L^2(\mathbb{R})} \asymp \left( \sum_{j,k} |\langle v, \widetilde{\psi}_{jk} \rangle|^2 \right)^{1/2} \asymp \left( \sum_{j,k} |\langle v, \psi_{jk} \rangle|^2 \right)^{1/2}. \quad (3.59)$$

If we need to start the decomposition from a fixed level  $j_0$ , as for multilevel decompositions of bounded domains, we can use the approximation spaces  $V_{j_0}$  and  $\widetilde{V}_{j_0}$  and add the wavelets for  $j \geq j_0$ , obtaining  $\forall v \in L^2(\mathbb{R})$

$$\begin{aligned} v &= P_{j_0} v + \sum_{j \geq j_0} Q_j v = \sum_k \langle v, \widetilde{\varphi}_{j_0,k} \rangle \varphi_{j_0,k} + \sum_{j \geq j_0} \sum_k \langle v, \widetilde{\psi}_{jk} \rangle \psi_{jk} = \\ &= \widetilde{P}_{j_0} v + \sum_{j \geq j_0} \widetilde{Q}_j v = \sum_k \langle v, \varphi_{j_0,k} \rangle \widetilde{\varphi}_{j_0,k} + \sum_{j \geq j_0} \sum_k \langle v, \psi_{jk} \rangle \widetilde{\psi}_{jk}. \end{aligned} \quad (3.60)$$

The stability condition (3.59) becomes then

$$\|v\|_{L^2(\mathbb{R})}^2 \asymp \|P_{j_0} v\|^2 + \left( \sum_{j \geq j_0} \sum_k |\langle v, \widetilde{\psi}_{jk} \rangle|^2 \right) \asymp \|\widetilde{P}_{j_0} v\|^2 + \left( \sum_{j \geq j_0} \sum_k |\langle v, \psi_{jk} \rangle|^2 \right). \quad (3.61)$$

We can now give an explicit expression for the coefficients of the two-scale analysis/synthesis described formally in Section 2.2. Using the canonical and hierarchical basis for  $V_{j+1}$  we have

$$v = \sum_{k \in \mathbb{Z}} \check{v}_{j+1,k} \varphi_{j+1,k} = \sum_{k \in \mathbb{Z}} \check{v}_{jk} \varphi_{jk} + \sum_{k \in \mathbb{Z}} \hat{v}_{jk} \psi_{jk}, \quad (3.62)$$

with  $\check{v}_{j+1,k} = \langle v, \tilde{\varphi}_{j+1,k} \rangle$ ,  $\check{v}_{jk} = \langle v, \tilde{\varphi}_{j,k} \rangle$ ,  $\hat{v}_{jk} = \langle v, \tilde{\psi}_{j,k} \rangle$ . Equations (2.30) and (2.31) become

$$\varphi_{jm} = \sum_{k \in \mathbb{Z}} h_{k-2m} \varphi_{j+1,k}, \quad \forall m \in \mathbb{Z} \quad (3.63)$$

$$\psi_{jm} = \sum_{k \in \mathbb{Z}} g_{k-2m} \varphi_{j+1,k}, \quad \forall m \in \mathbb{Z}, \quad (3.64)$$

while Eq. (2.32) reads

$$\varphi_{j+1,k} = \sum_{m \in \mathbb{Z}} \tilde{h}_{k-2m} \varphi_{jm} + \sum_{m \in \mathbb{Z}} \tilde{g}_{k-2m} \psi_{jm}, \quad \forall k \in \mathbb{Z}. \quad (3.65)$$

Finally, the corresponding analysis and synthesis on the expansion coefficients into the approximation and detail spaces become

$$\check{v}_{jm} = \sum_{k \in \mathbb{Z}} \tilde{h}_{k-2m} \check{v}_{j+1,k}, \quad \forall m \in \mathbb{Z}, \quad (3.66)$$

$$\hat{v}_{jm} = \sum_{k \in \mathbb{Z}} \tilde{g}_{k-2m} \check{v}_{j+1,k}, \quad \forall m \in \mathbb{Z}, \quad (3.67)$$

$$\check{v}_{j+1,k} = \sum_{m \in \mathbb{Z}} h_{k-2m} \check{v}_{jm} + \sum_{m \in \mathbb{Z}} g_{k-2m} \hat{v}_{jm}, \quad \forall k \in \mathbb{Z}. \quad (3.68)$$

### 3.3.1 Vanishing moments and polynomials reproduction

The preceding sections constructed sequences of nested subspaces of  $L^2(\mathbb{R})$  such that the approximation error vanishes when a characteristic index  $j$  tends to infinity. However, the question of “how good” is an approximation at a given level  $j_0$  was not answered. In particular, we still do not know how many levels need to be included in the approximation to get a small error. The answer to this question is hidden in the property M3 of Section 3.1. It can be proved that this condition is equivalent to the property of the local reproduction of all polynomials of degree up to  $L - 1$  (and  $\tilde{L} - 1$  for the dual system). More precisely, if we denote with  $\mathcal{P}_{L-1}$  the set of polynomials of degree at most  $L - 1$ , we can prove that M3 is equivalent to the two conditions

- the functions  $\{\varphi(x - k)\}_{k \in \mathbb{Z}}$  generate  $\mathcal{P}_{L-1}$  on  $\mathbb{R}$ ;
- $\int_{\mathbb{R}} x^l \tilde{\psi}(x) dx = 0, \quad 0 \leq l \leq L - 1.$

The first condition insures that at a fixed  $x$  any polynomial  $p(x)$  can be obtained as  $p(x) = \sum_k c_k \varphi(x - k)$  for a suitable choice of coefficients  $c_k$ . Note that the second condition represents a direct link between approximation properties of the spaces  $V_j$  and the number of vanishing moments of the dual wavelet  $\tilde{\psi}$ . With obvious substitutions the same facts apply for the dual system. Of course, the integrals representing the moments of order  $l$  must converge. This requires a sufficient fast decay of the wavelet functions for  $|x| \rightarrow \infty$ . This is insured when the wavelets are compactly supported.

### 3.3.2 Compactly supported wavelets

The use of compactly supported wavelets is of great interest under both a theoretical and practical standpoint. On one hand they allow to characterize spaces of functions like  $L^p$  and Besov spaces, which are far more general than  $L^2$  or Sobolev spaces [16]. On the other hand, multiresolution decompositions based on compactly supported scaling functions and wavelets can be extended to bounded domains, as we will show in Chap. 4.

The compactness of the support for scaling functions and wavelets can be obtained by imposing a finite length for the filters  $h_n$  and  $\tilde{h}_n$ . Let us assume that the functions  $m_0$  and  $\tilde{m}_0$  are trigonometric polynomials

$$m_0(\xi) = \frac{1}{\sqrt{2}} \sum_{n=n_0}^{n_1} h_n e^{-in\xi}, \quad \tilde{m}_0(\xi) = \frac{1}{\sqrt{2}} \sum_{n=\tilde{n}_0}^{\tilde{n}_1} \tilde{h}_n e^{-in\xi}. \quad (3.69)$$

This means that the refinement equations can be expressed as

$$\varphi(x) = \sqrt{2} \sum_{n=n_0}^{n_1} h_n \varphi(2x - n) \quad (3.70)$$

$$\tilde{\varphi}(x) = \sqrt{2} \sum_{n=\tilde{n}_0}^{\tilde{n}_1} \tilde{h}_n \tilde{\varphi}(2x - n). \quad (3.71)$$

From these expression we see that

$$\text{supp } \varphi \subseteq [n_0, n_1], \quad \text{supp } \tilde{\varphi} \subseteq [\tilde{n}_0, \tilde{n}_1].$$

It can be shown that the support of the scaling functions is indeed exactly

$$\text{supp } \varphi = [n_0, n_1], \quad \text{supp } \tilde{\varphi} = [\tilde{n}_0, \tilde{n}_1], \quad (3.72)$$

and that their polynomial order is such that  $L \leq n_1 - n_0 - 1$ . Only upper bounds for the supports of the wavelets can be obtained without imposing further conditions on  $n_0$ ,  $n_1$ ,  $\tilde{n}_0$ , and  $\tilde{n}_1$ . We have

$$\text{supp } \psi \subseteq [\nu_0, \nu_1] \quad \text{supp } \tilde{\psi} \subseteq [\tilde{\nu}_0, \tilde{\nu}_1], \quad (3.73)$$

where

$$\begin{aligned}\nu_0 &= \frac{1}{2}[1 - (\tilde{n}_1 - n_0)], & \nu_1 &= \frac{1}{2}[1 + (n_1 - \tilde{n}_0)], \\ \tilde{\nu}_0 &= \frac{1}{2}[1 - (n_1 - \tilde{n}_0)], & \tilde{\nu}_1 &= \frac{1}{2}[1 + (\tilde{n}_1 - n_0)].\end{aligned}$$

### 3.4 Examples of wavelets

We give in this section a few examples of orthogonal and biorthogonal wavelets which satisfy the basic axioms M1-M4 of Section 3.1 and therefore lead to multilevel decompositions of  $L^2(\mathbb{R})$ . Many examples of different wavelets can be found in the literature, and a complete list would be difficult to provide. Therefore, only three types of wavelets will be described in the following. The Haar system, the Daubechies wavelets, and the biorthogonal splines wavelets.

#### 3.4.1 The Haar wavelets

The Haar system constitutes the simplest case of orthogonal decomposition of  $L^2(\mathbb{R})$  (see Section 2.1). Being an orthogonal system, the dual spaces and basis functions coincide with the primal ones, so only the filter  $h$  needs to be specified. Let us define

$$m_0(\xi) = \tilde{m}_0(\xi) = \frac{1 + e^{-i\xi}}{2}. \quad (3.74)$$

This leads to the filter

$$h_n = \begin{cases} \frac{1}{\sqrt{2}} & n = 0, 1, \\ 0 & \text{otherwise.} \end{cases} \quad (3.75)$$

It is straightforward to verify that the axioms M1-M4 are satisfied, with  $L = 1$  and a regularity exponent  $\sigma = 1/2$ . From Eq. (3.11), we get the expression for the scaling function in the Fourier domain,

$$\hat{\varphi}(\xi) = \frac{1}{\sqrt{2\pi}} \frac{1 - e^{-i\xi}}{i\xi}, \quad (3.76)$$

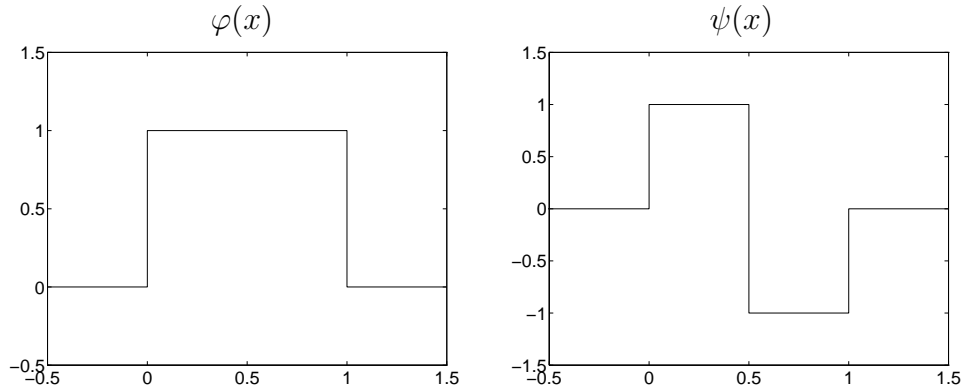
that when transformed back to the natural domain leads to

$$\varphi(x) = \begin{cases} 1 & \text{if } 0 \leq x < 1, \\ 0 & \text{otherwise.} \end{cases} \quad (3.77)$$

This is nothing else than the indicator function of the unit interval,  $\varphi = \chi_{[0,1]}$ . The corresponding wavelet can be determined from Eq. (3.40), obtaining

$$\hat{\psi}(\xi) = \frac{1}{\sqrt{2\pi}} \frac{1}{i\xi} (1 - 2e^{-i\xi/2} + e^{-i\xi}) \quad (3.78)$$





**Figure 3.1:** Haar scaling function (left panel) and wavelet (right panel).

in the Fourier domain, and

$$\psi(x) = \begin{cases} 1 & \text{if } 0 \leq x < 1/2, \\ -1 & \text{if } 1/2 \leq x < 1, \\ 0 & \text{otherwise} \end{cases} \quad (3.79)$$

Figure 3.1 shows the Haar scaling function and wavelet. As these functions are not regular, their approximation properties are not very good. It is then natural to generalize the Haar system to more regular wavelets, which can reproduce polynomials of higher degree. This can be done in many ways. If we want to keep the support compact for scaling functions and wavelets, we are forced to use the Daubechies systems for orthonormal decompositions. More flexibility, however, is provided by biorthogonal systems, like the B-splines, which are characterized by useful symmetry properties.

### 3.4.2 The Daubechies wavelets

The Daubechies wavelets form an orthogonal system. The basis functions have a compact support and are characterized by more regularity than the Haar functions. The regularity can be parameterized, and the support increases with the regularity. This explains the great popularity that these wavelets are having in the literature. The disadvantages of these systems are that scaling functions and wavelets are not known in closed form, but can only be generated through iterative algorithms. Moreover, they are not symmetric around the center of their support.

We do not give here the details of their construction (see Ref. [23]). Basically, the scaling functions are constructed from a trigonometric polynomial function  $m_0(\xi)$ , which is built so that the property M2 is satisfied. Figure 3.2 shows the scaling functions and wavelets for different values of  $L$  ranging from 2 to 6. The case  $L = 1$  is exactly the Haar system. Note that the regularity and the support increase with  $L$ . The regularity index  $\sigma$  can be evaluated numerically and is reported in Table 3.1 for the first values of  $L$ .

L	2	3	4	5	6
$\sigma$	0.84	1.14	1.41	1.68	1.91

**Table 3.1:** Regularity index for the Daubechies scaling functions

### 3.4.3 Biorthogonal spline wavelets

The biorthogonal spline wavelets generalize the Haar system to higher approximation order and regularity. As the Haar scaling functions can only reproduce polynomials of degree 0, i.e. constants, the B-spline scaling functions can reproduce locally polynomials of any fixed degree  $L - 1$ . In addition, these wavelets are symmetric with respect to the center of their support. This cannot be achieved with orthogonal systems like the Daubechies' wavelets. Indeed, to achieve symmetry the orthogonality conditions must be “relaxed” to the more general biorthogonality conditions. We will show in Chapter 4 that symmetry is a fundamental requirement for the construction of wavelets on the unit interval, because it allows to avoid different constructions at the edges of the domain.

We recall that the B-spline of order  $l$  is obtained through  $l$  convolutions of the box function  $\chi_{[0,1)}$ , through

$$\chi_{[0,1)}^{*l} = \underbrace{\chi_{[0,1)} * \cdots * \chi_{[0,1)}}_{l \text{ convolutions}} \quad (3.80)$$

It has compact support in  $[0, l + 1)$  and is expressed in any interval  $[m, m + 1)$  as a polynomial of degree  $l$ . As these functions are symmetric with respect to the center of their support, it is convenient to translate them to have the center in 0 for  $l$  odd and in  $1/2$  for  $l$  even. Let us set then

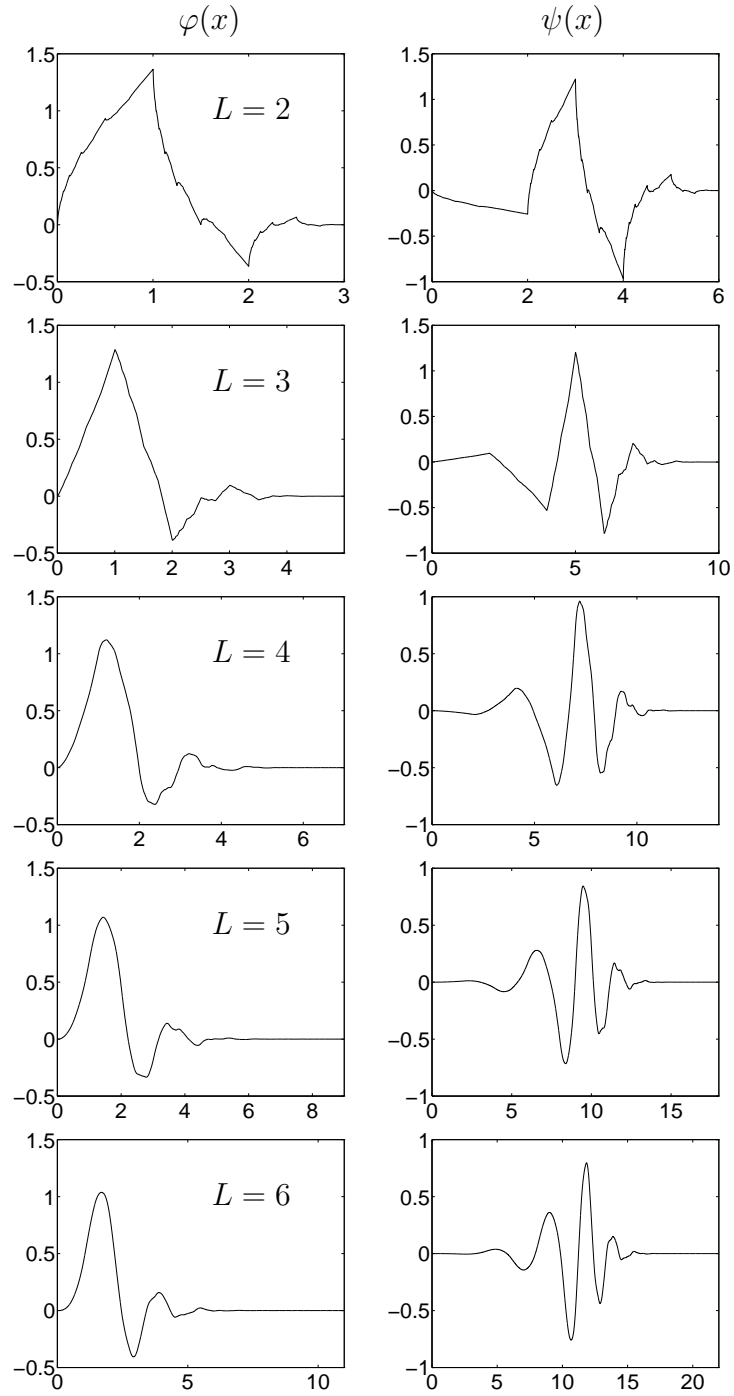
$$\Phi_l(x) = \chi_{[0,1)}^{*l} \left( x + \left\lfloor \frac{l+1}{2} \right\rfloor \right). \quad (3.81)$$

This will be the primal scaling function for the biorthogonal system.

We can easily see that the function  $m_0(\xi)$  of axiom M1 is expressed as

$$\begin{aligned} m_{0,l}(\xi) &= e^{i \lfloor \frac{l+1}{2} \rfloor \xi} \left( \frac{1 + e^{-i\xi}}{2} \right)^{l+1} \\ &= \begin{cases} (\cos \xi/2)^{l+1} & \text{for } l \text{ odd,} \\ e^{-i\frac{\xi}{2}} (\cos \xi/2)^{l+1} & \text{for } l \text{ even.} \end{cases} \end{aligned} \quad (3.82)$$

From this expression we see that, setting  $L = l + 1$ , we obtain a regularity index  $\sigma = L - 1/2$ . Therefore, both the polynomial reproduction and the regularity of the primal scaling functions increase with  $L$ .



**Figure 3.2:** Daubechies scaling functions (left column) and wavelets (right column) for different values of  $L$  ranging from 2 (top) to 6 (bottom).

As the integer translates of the primal scaling function are not orthogonal, we need to construct a dual scaling function so that its integer translates are biorthogonal to  $\Phi_l(x)$ . Again, we skip the details of the derivation (see Ref. [20]),

$L$	$\tilde{L}$	$\tilde{\sigma}$
2	2	0.18
2	4	0.87
2	6	1.33
2	8	1.85
3	5	0.33
3	7	0.85
4	8	0.32

**Table 3.2:** Regularity index  $\tilde{\sigma}$  of biorthogonal spline dual scaling functions, for some pairs  $L, \tilde{L}$ .

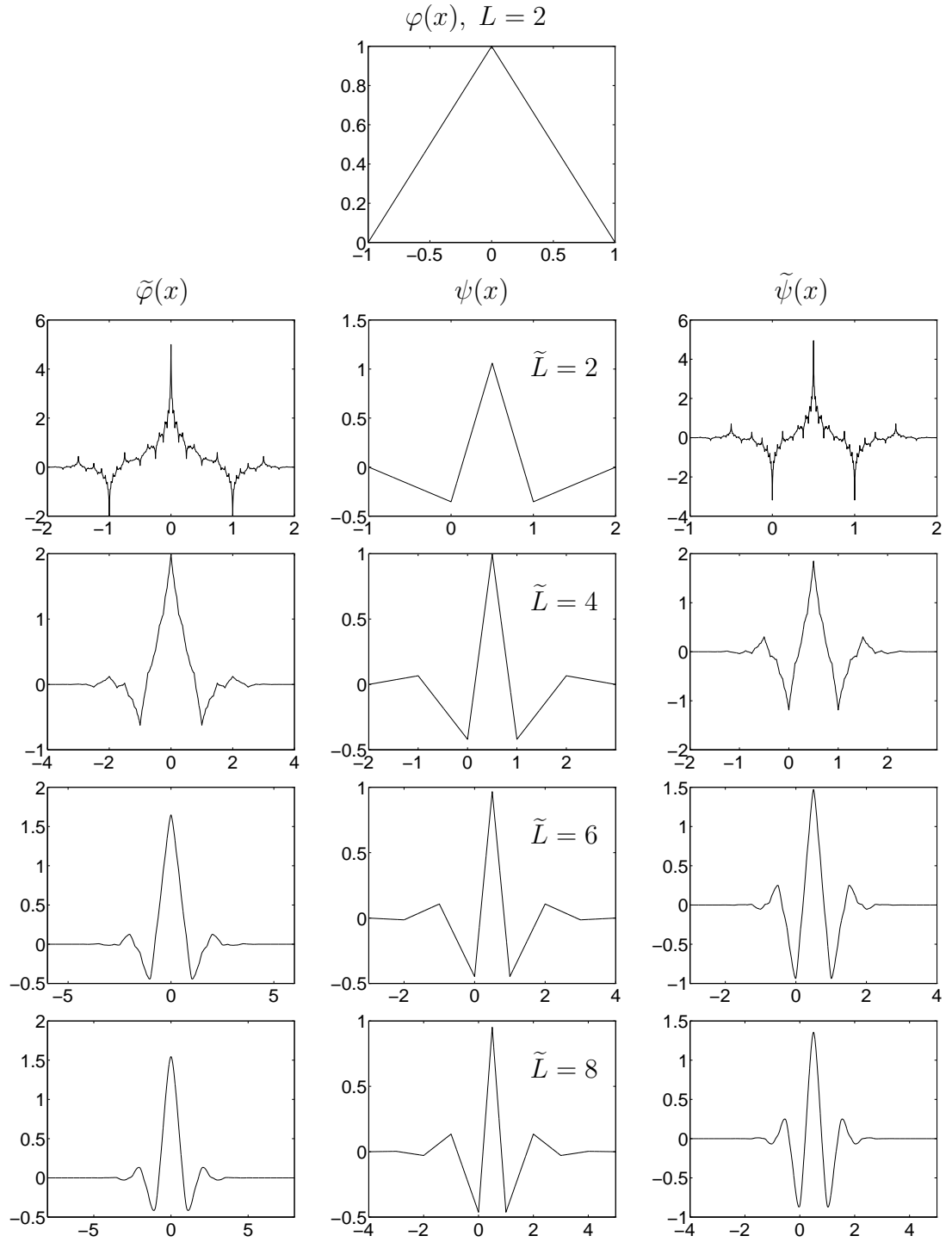
and we show some examples. The dual functions can be determined by finding solutions to M2. It is obvious that the solution is not unique, because the function  $\tilde{m}_0(\xi)$  can be chosen with an arbitrary number of zeros in  $\xi = \pi$ . However, once this number  $\tilde{L}$  has been fixed, it can be proved that there is a unique solution of minimal degree (i.e. minimal length of the dual filter). This solution is

$$\tilde{m}_0(\xi) = e^{-jr\xi/2} (\cos \xi/2)^{\tilde{L}} \left[ \sum_{n=0}^{k-1} \binom{k-1+n}{n} (\sin \xi/2)^{2n} \right],$$

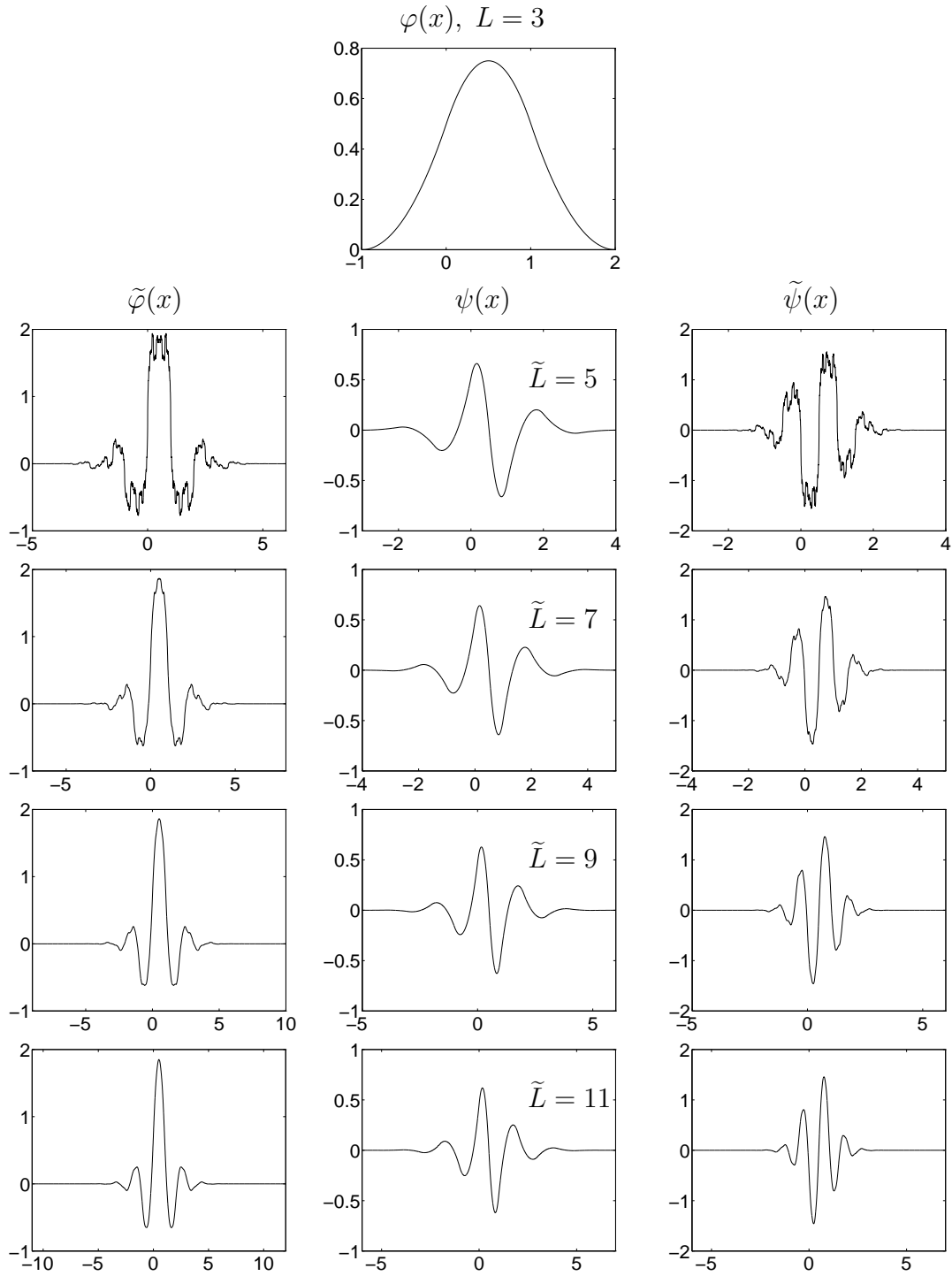
where  $2k = L + \tilde{L}$  and  $r = \text{rem}(L, 2)$ . This holds only when  $\tilde{L} \geq L$  and  $\tilde{L} + L$  is even. In summary, once  $L$  has been chosen, a corresponding  $\tilde{L} \geq L$  must be chosen, with the same parity as  $L$ . Then, the filter  $\tilde{h}$  is uniquely determined. The length of the filters  $h$  and  $\tilde{h}$  can be explicitly calculated in terms of  $L$  and  $\tilde{L}$ , obtaining

$$n_0 = -\left\lfloor \frac{L}{2} \right\rfloor, \quad n_1 = \left\lceil \frac{L}{2} \right\rceil, \quad \tilde{n}_0 = -\left\lfloor \frac{L}{2} \right\rfloor - \tilde{L} + 1, \quad \tilde{n}_1 = \left\lceil \frac{L}{2} \right\rceil + \tilde{L} - 1.$$

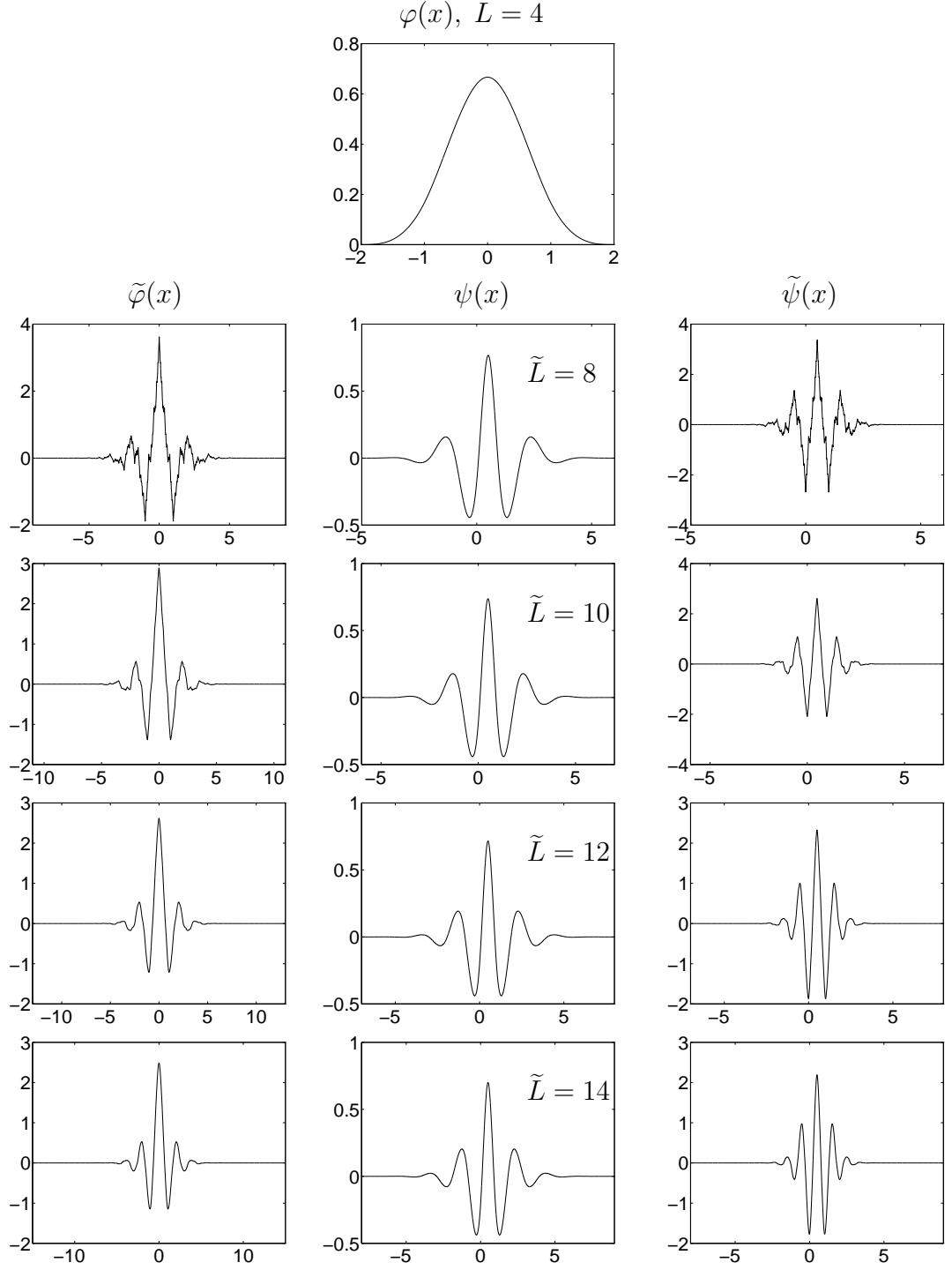
Table 3.2 shows the regularity index  $\tilde{\sigma}$  for the first pairs  $L, \tilde{L}$ , while Fig. 3.3 shows examples of primal and dual scaling functions and wavelets.



**Figure 3.3:** Biorthogonal splines for  $L = 2$  and  $\tilde{L} = 2, 4, 6, 8$ : primal scaling function (top panel), dual scaling function (left column), primal wavelet (middle column) and dual wavelet (right column).



**Figure 3.3:** (continued). Biorthogonal splines for  $L = 3$  and  $\tilde{L} = 5, 7, 9, 11$ .



**Figure 3.3:** (continued). Biorthogonal splines for  $L = 4$ ,  $\tilde{L} = 8, 10, 12, 14$ .

## Chapter 4

# Biorthogonal decomposition on bounded domains

This chapter focuses on the construction of biorthogonal multilevel decompositions on bounded domains. The need for this construction stems from the finite length of all NMTL structures. Clearly, if we want to use wavelet bases to numerically solve the NMTL equations, these bases must be defined on a bounded interval of the real line. It is obviously not possible to use directly the scaling function and the wavelet systems described in Chapter 3, because these bases are inherently translation invariant on  $\mathbb{R}$ .

The problem of an efficient and stable construction of wavelets on bounded domains is fundamental for the numerical solution of differential or integral equations stemming from many fields of application. For this reason, several papers are appearing in the literature about this subject. The reader is referred to the references [33]-[38]. In our derivation, we will follow the guidelines of the construction in Refs. [17, 25], although the main point for all the most recent constructions is the same, namely the preservation of the polynomial approximation properties at the edges of the domain.

The starting point in our construction is the biorthogonal decomposition of  $L^2(\mathbb{R})$  already described in Chapter 3. The aim is to preserve as much as possible the properties of the decomposition on the real line, and to modify the scaling functions and wavelet spaces only near the edges of the domain. Only a few basis functions of these spaces will need to be re-defined. The procedure illustrated in the following sections allows to preserve the approximation properties (i.e. the local generation of all the polynomials up to a certain degree) and the 2-stability.

As we will be dealing with bounded domains, it is convenient to start with finite length filters  $h = \{h_n\}_{n=n_0}^{n_1}$  and  $\tilde{h} = \{\tilde{h}_n\}_{n=\tilde{n}_0}^{\tilde{n}_1}$ , with  $n_0, \tilde{n}_0 \leq 0$  and  $n_1, \tilde{n}_1 \geq 0$ . This insures that both the scaling functions  $\varphi, \tilde{\varphi}$  and the wavelets  $\psi, \tilde{\psi}$  have compact support (see Section 3.3.2).

The following sections will use extensively the results of Chapter 3. We will



start with the simplest case, the half line  $[0, +\infty)$ , and we will construct approximation spaces  $V_j, \tilde{V}_j$ , wavelet spaces  $W_j, \tilde{W}_j$ , and corresponding 2-stable bases. Even if this domain is bounded only on one side, all the aspects in the treatment of the edges must be dealt with. We will see in Sec. 4.4 that the construction on the unit interval  $[0, 1]$  can be derived from the construction on  $[0, +\infty)$  through little modifications. Finally, Sec. 4.5 introduces the nonlinear approximations that will be used in next chapter to construct adaptive numerical schemes for the solution of the NMTL equations.

## 4.1 Scaling function spaces in $\mathbb{R}^+$

This section will show how a sequence of nested approximation spaces  $V_j(\mathbb{R}^+)$  and  $\tilde{V}_j(\mathbb{R}^+)$  can be constructed from the biorthogonal system already described in Chapter 3. To avoid ambiguity, a suffix  $\mathbb{R}$  will be appended to all the scaling functions defined on the real line. Also, this notation will be used for all the functions defined on  $\mathbb{R}$  and vanishing for  $x < 0$ . The domain of definition of these functions can also be interpreted to be  $\mathbb{R}^+$  with a slight abuse of notation. It will be assumed that all the functions without this suffix are intrinsically defined on  $\mathbb{R}^+$ . The primal decomposition spaces will be described in detail, while only the main results will be given for the dual spaces.

The philosophy underlying this construction is to localize the differences between the approximation spaces on  $\mathbb{R}$  and  $\mathbb{R}^+$  around the edge  $x = 0$ . This will allow to use most of the basis functions  $\varphi_{jk}^{\mathbb{R}}$  and  $\tilde{\varphi}_{jk}^{\mathbb{R}}$  to generate  $V_j(\mathbb{R}^+)$  and  $\tilde{V}_j(\mathbb{R}^+)$ . We will detail our construction for the scale  $j = 0$ , because the approximation spaces at any scale  $j$  can be generated from the spaces at scale 0 by using the isometry  $T_j$  already introduced in Eq. (3.22).

From the definitions in Eq. (3.72) we can determine the support of the dilated scaling functions

$$\text{supp } \varphi_{0k}^{\mathbb{R}} = [n_0 + k, n_1 + k].$$

Therefore the functions characterized by  $k \geq -n_0$  have the support included in  $[0, +\infty)$ . Our construction will preserve the scaling functions having support in  $[\delta, +\infty)$ , with  $\delta$  integer,  $\delta \geq 0$ . The corresponding value of  $k$  will be denoted by

$$k_0^* = \min\{k \in \mathbb{Z} : \text{supp } \varphi_{0k} \subset [\delta, +\infty)\} = -n_0 + \delta.$$

This will determine a corresponding subspace that will be left unchanged,

$$V^{(+)} = \text{span } \{\varphi_{0k}^{\mathbb{R}} : k \geq k_0^*\}. \quad (4.1)$$

This space clearly corresponds to a subspace of  $V_0(\mathbb{R})$ .

The scaling function space  $V_0(\mathbb{R}^+)$  will be obtained from  $V^{(+)}$  by adding some border functions in a finite number. The criterion to build these additional

functions will be the preservation of the approximation properties of the scaling function spaces on  $\mathbb{R}$ , i.e., the space  $\mathbb{P}_{L-1}$  of the polynomials of degree at most  $L-1$  will be locally reproduced. Let  $\{p_\alpha : \alpha = 0, \dots, L-1\}$  be a basis for  $\mathbb{P}_{L-1}$ . Without loss of generality we will use the basis of monomials, i.e.,  $p_\alpha = x^\alpha$ . Different choices will be discussed in the following sections. Each element of the basis can be represented for  $x \geq 0$  as

$$\begin{aligned} p_\alpha(x) &= \sum_{k \geq -n_1+1} c_{\alpha k} \varphi_{0k}^{\mathbb{R}}(x) \\ &= \sum_{k=-n_1+1}^{k_0^*-1} c_{\alpha k} \varphi_{0k}^{\mathbb{R}}(x) + \sum_{k \geq k_0^*} c_{\alpha k} \varphi_{0k}^{\mathbb{R}}(x), \end{aligned} \quad (4.2)$$

where

$$c_{\alpha k} = \langle p_\alpha, \tilde{\varphi}_{0k}^{\mathbb{R}} \rangle = \int_{\mathbb{R}} y^\alpha \tilde{\varphi}(y-k) dy, \quad \alpha = 0, \dots, L-1. \quad (4.3)$$

As the second sum in Eq. (4.2) belongs to  $V^{(+)}$ , the reproduction of polynomials on  $\mathbb{R}^+$  will be maintained if we choose as border functions the linear combinations

$$\theta_\alpha(x) = \sum_{k=-n_1+1}^{k_0^*-1} c_{\alpha k} \varphi_{0k}^{\mathbb{R}}(x), \quad x \geq 0, \quad \alpha = 0, \dots, L-1. \quad (4.4)$$

It can be proved that, as  $L \leq n_1 - n_0 - 1$ , the functions  $\theta_\alpha$ ,  $\alpha = 0, \dots, L-1$  are linearly independent, and also that the two sets of functions  $\theta_\alpha$ ,  $\alpha = 0, \dots, L-1$ , and  $\varphi_{0k}^{\mathbb{R}}$ ,  $k \geq k_0^*$ , are linearly independent. It is then natural to define

$$V_0(\mathbb{R}^+) = \text{span} \{ \theta_\alpha : \alpha = 0, \dots, L-1 \} \oplus V^{(+)}. \quad (4.5)$$

The same definition can be given for the dual space,

$$\tilde{V}_0(\mathbb{R}^+) = \text{span} \{ \tilde{\theta}_\beta : \beta = 0, \dots, \tilde{L}-1 \} \oplus \tilde{V}^{(+)},$$

where

$$\tilde{\theta}_\beta(x) = \sum_{k=-\tilde{n}_1+1}^{\tilde{k}_0^*-1} \tilde{c}_{\beta k} \tilde{\varphi}_{0k}^{\mathbb{R}}(x), \quad x \geq 0, \quad \beta = 0, \dots, \tilde{L}-1 \quad (4.6)$$

and

$$\tilde{c}_{\beta k} = \langle p_\beta, \varphi_{0k}^{\mathbb{R}} \rangle = \int_{\mathbb{R}} y^\beta \varphi(y-k) dy, \quad \beta = 0, \dots, \tilde{L}-1. \quad (4.7)$$

Let us set now

$$k^* = \max\{k_0^*, \tilde{k}_0^*\} = \max\{-n_0 + \delta, -\tilde{n}_0 + \tilde{\delta}\}. \quad (4.8)$$

From this point on, without loss of generality, we will set  $\tilde{L} \geq L$ , such as in the B-spline case. Under this assumption, we have  $k^* = \tilde{k}_0^* = -\tilde{n}_0 + \tilde{\delta}$ . The scaling

function spaces at level 0 will then have the form

$$\begin{aligned} V_0(\mathbb{R}^+) &= \text{span} \left\{ \{\theta_\alpha : \alpha = 0, \dots, L-1\} \cup \{\varphi_{0k}^R : k = \tilde{k}_0^*, \dots, k^* - 1\} \right. \\ &\quad \left. \cup \{\varphi_{0k}^R : k \geq k^*\} \right\} \\ \tilde{V}_0(\mathbb{R}^+) &= \text{span} \left\{ \{\tilde{\theta}_\beta : \beta = 0, \dots, \tilde{L}-1\} \cup \{\tilde{\varphi}_{0l}^R : l \geq k^*\} \right\}. \end{aligned} \quad (4.9)$$

From the biorthogonal decomposition on the real line we know that the last sets of functions in the expressions above are biorthogonal. These define the sets of scaling functions that are not modified by the construction,

$$\begin{aligned} V_0^I &= \text{span} \{\varphi_{0k}^R : k \geq k^*\} \\ \tilde{V}_0^I &= \text{span} \{\tilde{\varphi}_{0l}^R : l \geq k^*\}. \end{aligned}$$

The biorthogonality of the border functions has instead been destroyed by the construction procedure. It is necessary then to perform a change of basis for the border sets

$$\begin{aligned} V_0^B &= \text{span} \left\{ \{\theta_\alpha : \alpha = 0, \dots, L-1\} \cup \{\varphi_{0k}^R : k = \tilde{k}_0^*, \dots, k^* - 1\} \right\}, \\ \tilde{V}_0^B &= \text{span} \{\tilde{\theta}_\beta : \beta = 0, \dots, \tilde{L}-1\}. \end{aligned}$$

This is only possible when the dimensions of these two spaces match, with  $\dim V_0^B = \dim \tilde{V}_0^B = \tilde{L}$ . This leads to a relation between the constants  $\delta$  and  $\tilde{\delta}$ ,

$$\tilde{\delta} - \delta = \tilde{L} - L + \tilde{n}_0 - n_0. \quad (4.10)$$

It is convenient to set  $\tilde{\delta} = 0$ , so that the number of modified border functions is minimized. Also, it is required that  $\delta < L$ , in order to preserve the linear independence of the basis functions. Note the difference between the spaces labeled  $V_0^I$  and  $V_0^{(+)}$ . The latter includes also some of the primal scaling functions on the real line with support in  $\mathbb{R}^+$  that are “shifted” into the space  $V_0^B$ . This is only needed for the biorthogonalization procedure, that requires the dimension of  $V_0^B$  to be  $\tilde{L}$ .

Let us now relabel the primal and dual basis functions as

$$\theta_{0k} = \begin{cases} \theta_k & k = 0, \dots, L-1, \\ \varphi_{0, k_0^* + k - L}^R & k \geq L, \end{cases} \quad (4.11)$$

$$\tilde{\theta}_{0k} = \begin{cases} \tilde{\theta}_k & k = 0, \dots, \tilde{L}-1, \\ \tilde{\varphi}_{0, k^* + k - \tilde{L}}^R & k \geq \tilde{L}. \end{cases} \quad (4.12)$$

This notation leads to the simplified expressions

$$\begin{aligned} V_0^I &= \text{span} \{\theta_{0k} : k \geq \tilde{L}\}, & \tilde{V}_0^I &= \text{span} \{\tilde{\theta}_{0k} : k \geq \tilde{L}\} \\ V_0^B &= \text{span} \{\theta_{0k} : 0 \leq k < \tilde{L}\}, & \tilde{V}_0^B &= \text{span} \{\tilde{\theta}_{0k} : 0 \leq k < \tilde{L}\}. \end{aligned}$$

The final decomposition is obtained through biorthogonalization. From the results on the real line we have that  $V_0^I \perp \tilde{V}_0^I$ ,  $V_0^B \perp \tilde{V}_0^I$ , and  $V_0^I \perp \tilde{V}_0^B$ . Therefore we only need to biorthogonalize the border spaces. We need to determine new basis sets

$$\varphi_{0k'} = \sum_{k=0}^{\tilde{L}-1} d_{k'k} \theta_{0k}, \quad \tilde{\varphi}_{0l'} = \sum_{l=0}^{\tilde{L}-1} \tilde{d}_{l'l} \tilde{\theta}_{0l} \quad (4.13)$$

such that

$$\langle \varphi_{0k}, \tilde{\varphi}_{0l} \rangle = \delta_{kl}, \quad k, l = 0, \dots, \tilde{L} - 1.$$

This is equivalent to finding two real matrices  $D$  and  $\tilde{D}$  such that

$$DX\tilde{D}^T = I, \quad (4.14)$$

where  $X$  is the Gramian matrix

$$X_{kl} = \langle \theta_{0k}, \tilde{\theta}_{0l} \rangle, \quad \forall k, l = 0, \dots, \tilde{L} - 1. \quad (4.15)$$

Although there is no general result stating the invertibility of the matrix  $X$ , it can be proved that this matrix is non-singular at least in the B-spline case. The solution of Eq. (4.14) is obviously not unique. It is possible, for example, to preserve the primal functions by setting  $D = I$ , and solving for the dual system obtaining  $\tilde{D} = (X^{-1})^T$ .

In summary, whatever be the choice of the matrices  $D$  and  $\tilde{D}$ , we have constructed two different bases for  $V_0(\mathbb{R}^+)$  and  $\tilde{V}_0(\mathbb{R}^+)$ . The pair of basis sets indicated with  $\{\theta_{0k}, \tilde{\theta}_{0k}, k \geq 0\}$  are not biorthogonal in the first  $\tilde{L}$  functions, while the sets  $\{\varphi_{0k}, \tilde{\varphi}_{0k}, k \geq 0\}$ , obtained by setting  $\varphi_{0k} = \theta_{0k}$  and  $\tilde{\varphi}_{0k} = \tilde{\theta}_{0k}$  for  $k \geq \tilde{L}$ , are fully biorthogonal.

Section 4.1.1 shows how a modified refinement equation for the scaling functions on  $\mathbb{R}^+$  can be derived, while Section 4.1.2 details the calculation of the Gramian matrix  $X$ . A particular choice of the biorthogonalization matrices  $D$  and  $\tilde{D}$  preserving the boundary values of the scaling functions will be derived in Section 4.1.3.

### 4.1.1 The refinement equation

Let us consider the primal system  $\{\theta_{0\alpha}, \alpha \geq 0\}$  before the biorthogonalization. There are two cases  $\alpha \geq L$  and  $\alpha < L$ , which must be studied separately.

**Case  $\alpha \geq L$ .** In this case we can adapt the refinement equation (3.70) for the corresponding functions on the real line,

$$\theta_{0\alpha}(x) = \varphi_{0, k_0^* + \alpha - L}^R(x) = \sum_m h_{m-2(k_0^* + \alpha - L)} \varphi_{1m}^R. \quad (4.16)$$

This is possible because  $\varphi_{1m}^R \in V_1^{(+)} = T_1 V_0^{(+)}$  for all the values of  $m$  in the sum.

**Case  $\alpha < L$ .** As in this case  $\theta_{0\alpha} = \theta_\alpha$ , for  $x \geq 0$  we have, according to Eq. (4.2) and (4.3),

$$x^\alpha = \theta_\alpha + \sum_{k \geq k_0^*} c_{\alpha k} \varphi_{0k}^R. \quad (4.17)$$

Evaluating this expression in  $2x$  and multiplying for  $2^{1/2}$  we get

$$\begin{aligned} 2^{\alpha+1/2} x^\alpha &= 2^{1/2} \theta_\alpha(2x) + 2^{1/2} \sum_{k \geq k_0^*} c_{\alpha k} \varphi_{0k}^R(2x) \\ &= \theta_{1\alpha}(x) + \sum_{k \geq k_0^*} c_{\alpha k} \varphi_{1k}^R. \end{aligned}$$

Substituting again the expression (4.17) for  $x^\alpha$ , we find

$$\begin{aligned} \theta_{0\alpha} &= 2^{-\alpha-1/2} \left[ \theta_{1\alpha} + \sum_{k \geq k_0^*} c_{\alpha k} \varphi_{1k}^R \right] - \sum_{k \geq k_0^*} c_{\alpha k} \varphi_{0k}^R \\ &= 2^{-\alpha-1/2} \left[ \theta_{1\alpha} + \sum_{k \geq k_0^*} c_{\alpha k} \varphi_{1k}^R \right] - \sum_{k \geq L} c_{\alpha, k_0^*+k-L} \theta_{0k} \end{aligned}$$

Using the refinement equation (4.16) for  $\varphi_{0k}$ ,  $k \geq L$ , we have

$$\theta_{0\alpha} = 2^{-\alpha-1/2} \theta_{1\alpha} + \sum_{k \geq k_0^*} H_{\alpha k} \varphi_{1k}^R,$$

where

$$H_{\alpha k} = 2^{-\alpha-1/2} c_{\alpha k} - \sum_{l \geq k_0^*} c_{\alpha l} h_{k-2l}, \quad \alpha < L. \quad (4.18)$$

In summary, collecting the two cases, we can write

$$\theta_{0\alpha} = H_{\alpha\alpha} \theta_{1\alpha} + \sum_{k \geq k_0^*} H_{\alpha k} \varphi_{1k}^R, \quad (4.19)$$

where

$$H_{\alpha\alpha} = \begin{cases} 2^{-\alpha-1/2} & \alpha = 0, \dots, L-1, \\ 0 & \alpha \geq L, \end{cases} \quad (4.20)$$

$$H_{\alpha k} = \begin{cases} 2^{-\alpha-1/2} c_{\alpha k} - \sum_{l \geq k_0^*} c_{\alpha l} h_{k-2l}, & \alpha = 0, \dots, L-1, \\ h_{k-2(k_0^*+\alpha-L)} & \alpha \geq L. \end{cases} \quad (4.21)$$

This refinement equation shows that each border function at level 0 can be expressed in terms of the corresponding border function at level 1 plus a suitable linear combination of the internal functions. The final refinement equation for

the scaling functions in  $V_0(\mathbb{R}^+)$  can be obtained by applying the definition in Eq. (4.11) to the functions at level 1,

$$\begin{aligned}\theta_{0\alpha} &= H_{\alpha\alpha}\theta_{1\alpha} + \sum_{k \geq k_0^*} H_{\alpha k}\theta_{1,k+L-k_0^*} \\ &= H_{\alpha\alpha}\theta_{1\alpha} + \sum_{l \geq L} H_{\alpha,k_0^*+l-L}\theta_{1l}.\end{aligned}$$

This refinement equation can be written in a more compact form by introducing the (infinite) filtering matrix  $\mathcal{H}^\angle$  defined as

$$\theta_{0k} = \sum_{l \geq 0} \mathcal{H}_{kl}^\angle \theta_{1l}, \quad \forall k \geq 0. \quad (4.22)$$

The structure of this matrix can be easily visualized as

$$\mathcal{H}^\angle = \begin{array}{|c|c|c|c|c|c|c|} \hline & & & & & & \\ \hline & H_{\alpha\alpha} & & H_{\alpha,k_0^*+l-L} & & & \\ \hline L & \delta & \boxed{\phantom{0}} & \boxed{\phantom{0}} & h & \boxed{\phantom{0}} & \boxed{\phantom{0}} \\ \hline & & & \boxed{\phantom{0}} & \boxed{\phantom{0}} & h & \boxed{\phantom{0}} \\ \hline & & & & \boxed{\phantom{0}} & \boxed{\phantom{0}} & h \\ \hline & & & & & \boxed{\phantom{0}} & \boxed{\phantom{0}} \\ \hline \end{array}$$

The  $L \times L$  diagonal block in the upper left corner is due to the choice of monomials as the basis for  $\mathcal{P}_{L-1}$ . This block is not diagonal for other choices of basis sets. Note that the number of columns  $N'$  with nonzero entries in the first  $L$  rows can be exactly determined from Eq. (4.18) by imposing the limits  $n_0, n_1$  for the filter  $h$ . The result is  $N' = n_1 + L + k_0^* - 1$ . The rows starting from  $L$  consist simply of translations of the filter  $h$ . The offset of two adjacent rows is 2, while the first nonzero entry in row  $L$  corresponds to column  $L + \delta$  (the indexing of this matrix starts from 0).

The dual filter has a similar structure, with a  $\tilde{L} \times \tilde{L}$  diagonal block in the upper left corner, a number of nonzero entries in the first  $\tilde{L}$  rows equal to  $N' = n_1 + \tilde{L} + k^* - 1$ , and  $\tilde{L} + \tilde{\delta}$  zeros in row  $\tilde{L}$ . The matrix is visualized in the picture below.

$$\widetilde{\mathcal{H}}^\angle = \begin{array}{|c|c|c|c|c|c|c|} \hline & & & & & & \\ \hline & \widetilde{H}_{\alpha\alpha} & & \widetilde{H}_{\alpha,k^*+l-\tilde{L}} & & & \\ \hline \tilde{L} & \tilde{\delta} & \boxed{\phantom{0}} & \boxed{\phantom{0}} & \tilde{h} & \boxed{\phantom{0}} & \boxed{\phantom{0}} \\ \hline & & & \boxed{\phantom{0}} & \boxed{\phantom{0}} & \tilde{h} & \boxed{\phantom{0}} \\ \hline & & & & \boxed{\phantom{0}} & \boxed{\phantom{0}} & \tilde{h} \\ \hline & & & & & \boxed{\phantom{0}} & \boxed{\phantom{0}} \\ \hline \end{array}$$

Let us derive now the refinement equation for the biorthogonal basis set  $\varphi_{0\alpha}$ . As the biorthogonalization process involves only the first  $\tilde{L}$  functions, we must distinguish the two cases  $\alpha < \tilde{L}$  and  $\alpha \geq \tilde{L}$ .

- $\alpha < \tilde{L}$ . In this case we can use the matrix  $D$  to express  $\varphi_{0\alpha}$  as a linear combination of the non biorthogonal basis functions, according to Eq. (4.13). Then, we can apply the refinement equation in the non biorthogonal case, by using Eq. (4.22). Note that if we consider only the first  $\tilde{L}$  rows, the upper limit in the sum is finite. Let us call it  $N - 1$ . This limit can be easily determined, obtaining  $N = \tilde{L} + k^* + n_1 - 1$ . We have

$$\theta_{0s} = \sum_{l=0}^{N-1} \mathcal{H}_{sl}^{\angle} \theta_{1l}, \quad \forall s < \tilde{L}. \quad (4.23)$$

Finally, we express each function  $\varphi_{1l}$  as a combination of the biorthogonal functions,

$$\theta_{1l} = \begin{cases} \sum_{r=0}^{\tilde{L}-1} [D^{-1}]_{lr} \varphi_{1r} & \text{if } l < \tilde{L}, \\ \varphi_{1l} & \text{if } l \geq \tilde{L}. \end{cases}$$

This can be expressed in a compact form as

$$\theta_{1l} = \sum_{r=0}^{\tilde{L}-1} [D_a]_{lr} \varphi_{1r},$$

where the  $N \times N$  matrix  $D_a$  is defined as

$$D_a = \begin{bmatrix} D^{-1} & 0 \\ 0 & I \end{bmatrix}.$$

Putting all together we have, for  $\alpha < \tilde{L}$ ,

$$\begin{aligned} \varphi_{0\alpha} &= \sum_{r=0}^{N-1} \left\{ \sum_{s=0}^{\tilde{L}-1} \sum_{l=0}^{N-1} D_{ks} \mathcal{H}_{sl}^{\angle} [D_a]_{lr} \right\} \varphi_{1r} \\ &= \sum_{r=0}^{N-1} H_{\alpha r}^0 \varphi_{1r}, \end{aligned} \quad (4.24)$$

where the matrix  $H_0$  is simply the product of the three matrices  $D$ ,  $\mathcal{H}^{\angle}$  (the upper  $\tilde{L} \times N$  block) and  $D_a$ .

- $\alpha \geq \tilde{L}$ . In this case we do not need to change basis because we already have  $\varphi_{0\alpha} = \theta_{0\alpha}$ . Therefore, the refinement equation does not change from Eq. (4.16). After substitution of the correct limits in the sum, obtained from the support of the filter  $h$ , we have

$$\varphi_{0\alpha} = \sum_{m=\delta+2\alpha-L}^{n_1+k_0^*+2\alpha-L} h_{m-2\alpha+L-k_0^*} \varphi_{1m}. \quad (4.25)$$





We will show in the following that the terms in the summation above are identically zero, so that the evaluation of the matrix  $X$  reduces to the computation of the moments of the primal scaling functions.

Using Eq. (4.11), we see that

$$\langle \theta_{0\alpha}, \tilde{\varphi}_{0k}^R \rangle = \begin{cases} \langle \theta_\alpha, \tilde{\varphi}_{0k}^R \rangle & \alpha = 0, \dots, L-1 \\ \langle \varphi_{0, k_0^* + \alpha - L}^R, \tilde{\varphi}_{0k}^R \rangle & \alpha = L, \dots, \tilde{L}-1. \end{cases} \quad (4.28)$$

The second row in the preceding expression vanishes identically due to the biorthogonality of the scaling functions on the real line. Indeed, when  $\alpha < \tilde{L}-1$  we have  $k_0^* + \alpha - L < k^*$ , and the inner product evaluates to  $\langle \varphi_{0, k_0^* + \alpha - L}^R, \tilde{\varphi}_{0k}^R \rangle = \delta_{k_0^* + \alpha - L, k} = 0$  for  $k \geq k^*$ . It is easily shown that also the first row in Eq. (4.28) vanishes. Recalling that  $\text{supp} \{\tilde{\varphi}_{0k}^R, k \geq k^*\} \subseteq [\tilde{\delta}, +\infty)$ , we can notice that the restriction to positive values of  $x$  in the definition of  $\theta_\alpha$  (Eq. (4.4)) is not necessary for the evaluation of the inner product. Therefore, when  $\alpha < L$ , we have

$$\begin{aligned} \langle \theta_\alpha, \tilde{\varphi}_{0k}^R \rangle &= \left\langle \sum_{l=-n_1+1}^{k_0^*-1} c_{\alpha l} \varphi_{0l}^R \Big|_{[0, +\infty)}, \tilde{\varphi}_{0k}^R \right\rangle \\ &= \left\langle \sum_{l=-n_1+1}^{k_0^*-1} c_{\alpha l} \varphi_{0l}^R, \tilde{\varphi}_{0k}^R \right\rangle \\ &= \sum_{l=-n_1+1}^{k_0^*-1} c_{\alpha l} \delta_{kl} = 0 \end{aligned}$$

because  $k \geq k^*$ . We can conclude that the calculation of the Gramian matrix  $X$  reduces to

$$X_{\alpha\beta} = \langle \theta_{0\alpha}, \tilde{\theta}_{0\beta} \rangle = \langle \theta_{0\alpha}, x^\beta \rangle, \quad \alpha, \beta = 0, \dots, \tilde{L}-1.$$

Let us now apply the refinement equation (4.19) to  $\theta_{0\alpha}$  in the inner product. We get

$$\begin{aligned} \langle \theta_{0\alpha}, x^\beta \rangle &= H_{\alpha\alpha} \langle T_1 \theta_{0\alpha}, x^\beta \rangle + \sum_{k \geq k_0^*} H_{\alpha k} \langle \varphi_{1k}^R, x^\beta \rangle \\ &= H_{\alpha\alpha} 2^{-\beta-1/2} \langle \theta_{0\alpha}, x^\beta \rangle + \sum_{k \geq k_0^*} H_{\alpha k} 2^{-\beta-1/2} \langle \varphi_{0k}^R, x^\beta \rangle \end{aligned}$$

The last passage is easily obtained through change of variable in the integrals. Using now the definition of  $\tilde{c}_{\beta k}$  (see Eq. (4.7)), we can write

$$X_{\alpha\beta} = \langle \theta_{0\alpha}, x^\beta \rangle = \frac{2^{-\beta-1/2}}{1 - H_{\alpha\alpha} 2^{-\beta-1/2}} \sum_{k \geq k_0^*} H_{\alpha k} \tilde{c}_{\beta k}.$$

Note that the modified border filters  $H_{\alpha k}$  are known once the  $c_{\alpha k}$  have been evaluated. Therefore the elements of the matrix  $X$  can be expressed in terms of the coefficients  $c_{\alpha k}$  and  $\tilde{c}_{\beta k}$ , which are easy to evaluate because they are defined through functions on the real line. Indeed we have

$$\begin{aligned} c_{\alpha k} &= \int_{\mathbb{R}} x^{\alpha} \tilde{\varphi}_{0k}^R(x) dx = \int_{\mathbb{R}} (x+k)^{\alpha} \tilde{\varphi}(x) dx \\ &= \sum_{m=0}^{\alpha} \binom{\alpha}{m} k^{\alpha-m} \int_{\mathbb{R}} x^m \tilde{\varphi}(x) dx = \sum_{m=0}^{\alpha} \binom{\alpha}{m} k^{\alpha-m} c_{m0}. \end{aligned} \quad (4.29)$$

The value of  $c_{\alpha k}$  can therefore be evaluated for each  $k$  from the set of coefficients with  $k = 0$ . These can be evaluated by applying the refinement equation to the dual scaling function on  $\mathbb{R}$ ,

$$\begin{aligned} c_{\alpha 0} &= \int_{\mathbb{R}} x^{\alpha} \tilde{\varphi}(x) dx = \sum_{m=\tilde{n}_0}^{\tilde{n}_1} \tilde{h}_m 2^{1/2} \int_{\mathbb{R}} x^{\alpha} \tilde{\varphi}(2x-m) dx \\ &= \sum_{m=\tilde{n}_0}^{\tilde{n}_1} \tilde{h}_m 2^{-\alpha-1/2} \int_{\mathbb{R}} (x+m)^{\alpha} \tilde{\varphi}(x) dx \\ &= \sum_{m=\tilde{n}_0}^{\tilde{n}_1} \tilde{h}_m 2^{-\alpha-1/2} \sum_{s=0}^{\alpha} \binom{\alpha}{s} m^{\alpha-s} \int_{\mathbb{R}} x^s \tilde{\varphi}(x) dx \\ &= 2^{-\alpha-1/2} \sum_{s=0}^{\alpha} \binom{\alpha}{s} \left[ \sum_{m=\tilde{n}_0}^{\tilde{n}_1} \tilde{h}_m m^{\alpha-s} \right] c_{s0}. \end{aligned}$$

Recalling that  $\sum_m \tilde{h}_m = \sqrt{2}$ , we obtain

$$c_{\alpha 0} = \frac{2^{-\alpha-1/2}}{1-2^{-\alpha}} \sum_{s=0}^{\alpha-1} \binom{\alpha}{s} \left[ \sum_{m=\tilde{n}_0}^{\tilde{n}_1} \tilde{h}_m m^{\alpha-s} \right] c_{s0}, \quad \alpha > 0. \quad (4.30)$$

The evaluation of  $c_{00}$  is trivial because

$$c_{00} = \int_{\mathbb{R}} \tilde{\varphi}(x) dx = 1. \quad (4.31)$$

In conclusion, Eq. (4.30) can be applied recursively to calculate  $c_{\alpha 0}$ ,  $\forall \alpha > 0$ .

### 4.1.3 Boundary value preserving biorthogonalization

In the foregoing sections a set of modified border functions have been constructed for the primal and dual systems. These functions,  $\theta_{\alpha}$  and  $\tilde{\theta}_{\beta}$  respectively, were derived from the basis of monomials of the spaces of polynomials  $\mathcal{P}_{L-1}$  and  $\mathcal{P}_{\tilde{L}-1}$ . This basis is such that only the first monomial, i.e., the constant, is different from zero in the point  $x = 0$ . All the others vanish when evaluated in the origin. This property is obviously true also for the functions

$\theta_\alpha$  and  $\tilde{\theta}_\beta$ , of which only the two with  $\alpha = 0$  and  $\beta = 0$  do not vanish in zero. In general, the matrices  $D$  and  $\tilde{D}$  in Eq. (4.14) do not necessarily preserve this property for the biorthogonal basis functions  $\varphi_{0k}, \tilde{\varphi}_{0k}$ . This section shows how the system (4.14) can be solved with the constraint that only one function for the primal and only one for the dual system be nonvanishing in  $x = 0$ .

The procedure is based on two steps. First, all the functions except the first are biorthogonalized. Second, a border function is added for the primal and dual system so that they are biorthogonal to the rest of the basis functions. The two steps are now examined in further detail.

Let us partition the matrices  $X$ ,  $D$  and  $\tilde{D}$  by extracting the first row and column,

$$X = \begin{bmatrix} X_0 & X_r \\ X_c & X_1 \end{bmatrix} \quad D = \begin{bmatrix} D_0 & D_r \\ D_c & D_1 \end{bmatrix} \quad \tilde{D} = \begin{bmatrix} \tilde{D}_0 & \tilde{D}_r \\ \tilde{D}_c & \tilde{D}_1 \end{bmatrix}. \quad (4.32)$$

The suffix  $_0$  indicates scalars,  $_r$  stands for a row vector of length  $\tilde{L} - 1$  and  $_c$  stands for a column vector of length  $\tilde{L} - 1$ . The matrices with the suffix  $_1$  are square with dimension  $\tilde{L} - 1$ , and link quantities that involve all the basis functions except the first. The biorthogonalization of the sets  $\{\theta_{0k}, k = 1, \dots, \tilde{L} - 1\}$ ,  $\{\tilde{\theta}_{0k}, k = 1, \dots, \tilde{L} - 1\}$  corresponds to finding two real matrices  $D_1$  and  $\tilde{D}_1$  satisfying

$$D_1 X_1 \tilde{D}_1^T = I.$$

In the following we will suppose that these matrices have already been determined. The rest of the biorthogonalization does not depend on a particular choice for  $D_1$  and  $\tilde{D}_1$ . As a result, we get two sets of biorthogonal functions  $\{\varphi_{0k}, k = 1, \dots, \tilde{L} - 1\}$ ,  $\{\tilde{\varphi}_{0k}, k = 1, \dots, \tilde{L} - 1\}$  spanning the same space of the corresponding non biorthogonal functions.

We add now two new functions  $\theta_{00}$  and  $\tilde{\theta}_{00}$ , with the only constraint that they are independent from the others. This is obviously true for the basis of monomials, for which these new border functions are identically equal to 1. We want to define two new functions  $\varphi_{00}$  and  $\tilde{\varphi}_{00}$  that satisfy the following conditions,

$$\langle \varphi_{00}, \tilde{\varphi}_{0s} \rangle = 0 \quad \text{for } s > 0 \quad (4.33)$$

$$\langle \tilde{\varphi}_{00}, \varphi_{0s} \rangle = 0 \quad \text{for } s > 0 \quad (4.34)$$

$$\langle \varphi_{00}, \tilde{\varphi}_{00} \rangle = 1 \quad (4.35)$$

Let us write the new border functions as

$$\varphi_{00} = \alpha_0 \theta_{00} + \sum_{s=1}^{\tilde{L}-1} \alpha_s \varphi_{0s}, \quad \tilde{\varphi}_{00} = \beta_0 \tilde{\theta}_{00} + \sum_{s=1}^{\tilde{L}-1} \beta_s \tilde{\varphi}_{0s}. \quad (4.36)$$

Substituting these expressions in (4.33) and (4.34) we obtain,

$$\begin{cases} \alpha_s = -\alpha_0 \langle \theta_{00}, \tilde{\varphi}_{0s} \rangle, \\ \beta_s = -\beta_0 \langle \tilde{\theta}_{00}, \varphi_{0s} \rangle, \end{cases} \quad \forall s = 1, \dots, \tilde{L} - 1. \quad (4.37)$$

Substituting now Eq. (4.36) in Eq. (4.35) and using Eq. (4.37) we get, after a straightforward calculation,

$$\alpha_0 \beta_0 K = 1, \quad (4.38)$$

where

$$K = \langle \varphi_{00}, \tilde{\varphi}_{00} \rangle - \sum_{s=1}^{\tilde{L}-1} \langle \theta_{00}, \tilde{\varphi}_{0s} \rangle \langle \tilde{\theta}_{00}, \varphi_{0s} \rangle. \quad (4.39)$$

If we define now the column vectors

$$\underline{\alpha} = \begin{bmatrix} \alpha_1 \\ \vdots \\ \alpha_{\tilde{L}-1} \end{bmatrix}, \quad \underline{\beta} = \begin{bmatrix} \beta_1 \\ \vdots \\ \beta_{\tilde{L}-1} \end{bmatrix},$$

we can express the solution in a matrix form,

$$\begin{aligned} \underline{\alpha} &= -\alpha_0 \widetilde{D}_1 X_r^T \\ \underline{\beta} &= -\beta_0 D_1 X_c \\ K &= X_0 - X_r \widetilde{D}_1^T D_1 X_c. \end{aligned}$$

Note that there is still one degree of freedom, i.e. the ratio between  $\alpha_0$  and  $\beta_0$ . Once this is fixed, both  $\alpha_0$  and  $\beta_0$  are uniquely determined by Eq. (4.38), and so are the other coefficients.

Finally, the remaining blocks of the matrices  $D$  and  $\widetilde{D}$  can be determined. Using their definition in Eq. (4.13) we obtain

$$\begin{cases} D_r = \underline{\alpha}^T D_1 \\ D_c = 0 \\ D_0 = \alpha_0 \end{cases} \quad \begin{cases} \widetilde{D}_r = \underline{\beta}^T \widetilde{D}_1 \\ \widetilde{D}_c = 0 \\ \widetilde{D}_0 = \beta_0. \end{cases} \quad (4.40)$$

#### 4.1.4 Other polynomial bases

The construction of the border scaling functions is based on a particular choice of basis set for the polynomials to be reproduced. The construction in the foregoing sections was performed using the basis of monomials  $p_\alpha$  due to its simplicity. On the other hand, it is well known that the basis of monomials often leads to ill-conditioned problems and loss of accuracy in the evaluation of related quantities. This fact occurs also in the present construction of the border scaling functions. In particular the condition number of the biorthogonalization matrix  $X$  grows very large when  $L$  and  $\tilde{L}$  increase. Therefore, the gain in accuracy due to better approximation spaces is “spoiled” by the ill-conditioned biorthogonalization.

These problems could be partially avoided by using different bases for the polynomial spaces  $\mathcal{P}_{L-1}$  and  $\mathcal{P}_{\tilde{L}-1}$ . In this section we list the expressions for the basic quantities introduced in the construction of the scaling function

spaces when a different polynomial basis is chosen. Their derivation requires a few straightforward passages, and is not detailed here. In addition, only the primal quantities are listed, because the dual case can be recovered through obvious substitutions.

Let us assume a choice of basis functions for  $\mathcal{P}_{L-1}$ , different from the monomials  $p_\alpha = x^\alpha$ , which we will denote by  $\rho_\alpha(x)$ , with  $\alpha = 0, \dots, L-1$ . This new basis set can be expressed in terms of the monomials through a change of basis matrix  $Z$ ,

$$\begin{aligned}\rho_\alpha(x) &= \sum_{r=0}^{L-1} Z_{\alpha r} p_r(x) \\ p_\alpha(x) &= \sum_{r=0}^{L-1} [Z^{-1}]_{\alpha r} \rho_r(x)\end{aligned}\tag{4.41}$$

The matrix  $Z$  allows to determine the following parameters and functions in the new polynomial basis (denoted in the following with the superscript <sup>(new)</sup>) in terms of the same quantities based on the basis of monomials

- expansion coefficients of the polynomial basis functions into the scaling functions on the real line,

$$c_{\alpha,k}^{(\text{new})} = \langle \rho_\alpha, \tilde{\varphi}_{0k}^R \rangle = \sum_{r=0}^{L-1} Z_{\alpha r} c_{r,k},$$

- border scaling functions,

$$\theta_\alpha^{(\text{new})}(x) = \sum_{r=0}^{L-1} Z_{\alpha r} \theta_r(x),$$

- refinement equation for border scaling functions

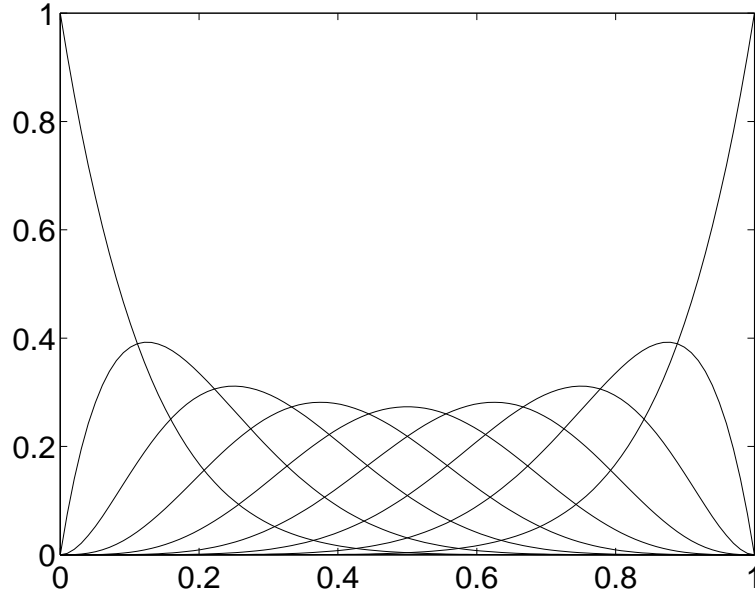
$$\mathbf{H}^{(\text{new})} = Z\mathbf{H}\hat{Z}^{-1},$$

where  $\mathbf{H}$  is the non-biorthogonal border filter matrix made of the two blocks in Eq. (4.20) and (4.21), and  $\hat{Z}^{-1}$  is a matrix with  $Z^{-1}$  in the upper-left block and filled with ones in the main diagonal to reach the correct dimension required by the matrix product,

- biorthogonalization matrix

$$X^{(\text{new})} = \hat{Z}X\tilde{Z}^T,$$

where  $\hat{Z}$  is a matrix with  $Z$  in the upper-left block and filled with ones in the main diagonal to reach the correct dimension required by the matrix product;



**Figure 4.1:** Bernstein polynomials for  $n = 8$  and  $b = 1$ .

Among the many polynomial bases there is one in particular that is known to be better behaved than the monomials, the Bernstein or Bezier polynomials [13, 14, 36, 37]. The general expression of the Bernstein polynomials of degree  $n$  is

$$\mathcal{B}_{r,b}^n(x) = b^{-n} \binom{n}{r} x^r (b-x)^{n-r}, \quad r = 0, \dots, n, \quad (4.42)$$

where  $b$  is a parameter that controls the definition interval  $[0, b]$  of the polynomials. Figure 4.1.4 shows the Bernstein polynomials of degree  $n = 8$  defined for  $b = 1$ .

The Bernstein polynomials are also convenient because only one is nonvanishing at  $x = 0$ , thus allowing for boundary adaption of the border scaling functions, and because the change of basis matrix  $Z$  and its inverse are known analytically,

$$Z_{r,s} = \begin{cases} (-1)^{r-s} \binom{n}{s} \binom{s}{r} b^{-s} & s \geq r \\ 0 & s < r \end{cases},$$

$$[Z^{-1}]_{r,s} = \begin{cases} \binom{s}{r} \left[ \binom{n}{r} \right]^{-1} b^r & s \geq r \\ 0 & s < r \end{cases}.$$

It should be noted that both  $Z$  and  $Z^{-1}$  are upper triangular. The usefulness of the Bernstein polynomials will be shown in the example of Section 4.4.4.

### 4.1.5 Approximation properties

We will list here the main approximation properties of the spaces  $V_j(\mathbb{R}^+)$  constructed in the previous sections. We recall that the unitary operator  $T_j$ , defined in Eq. (3.22), can be applied to the space  $V_0(\mathbb{R}^+)$  to obtain the corresponding space at level  $j$ . Throughout this section we will use the functions  $\{\varphi_{0k}, k \geq 0\}$  as the biorthogonalized basis set for  $V_0(\mathbb{R}^+)$ .

First of all, it can be demonstrated that the functions  $\{\varphi_{0k}, k \geq 0\}$  constitute a 2-stable basis for  $V_0(\mathbb{R}^+)$ . This result is extremely important because it states the equivalence of the norm in  $V_0(\mathbb{R}^+)$  with the  $\ell^2(\mathbb{N})$  norm of the expansion coefficients into this basis. More precisely, we have

$$V_0(\mathbb{R}^+) = \left\{ v = \sum_{k \geq 0} \alpha_k \varphi_{0k} : \{\alpha_k\} \in \ell^2(\mathbb{N}) \right\},$$

with

$$\|v\|_{L^2(\mathbb{R}^+)} \asymp \|\{\alpha_k\}\|_{\ell^2(\mathbb{N})}, \quad \forall v \in V_0(\mathbb{R}^+). \quad (4.43)$$

As for the construction on the real line, we can define projection operators  $P_j : L^2(\mathbb{R}^+) \rightarrow V_j(\mathbb{R}^+)$  that, given any function in  $L^2(\mathbb{R}^+)$ , can extract its approximation at level  $j$ . Again, only the operator  $P_0$  needs to be explicitly defined, because the operators at any other level  $j$  can be obtained through the use of the isometry  $T_j$  according to

$$P_j = T_j \circ P_0 \circ T_j^{-1}.$$

It is natural to define

$$P_0 v = \sum_{k \geq 0} \check{v}_{0k} \varphi_{0k}, \quad \text{with} \quad \check{v}_{0k} = \int_{\mathbb{R}^+} v(x) \tilde{\varphi}_{0k}(x) dx. \quad (4.44)$$

It can be easily proven that  $P_0 v \in V_0(\mathbb{R}^+)$  and that the operator is continuous. Moreover, we have

$$\begin{aligned} P_j v &= v, & \forall v \in V_j(\mathbb{R}^+) \\ P_j \circ P_{j+1} &= P_j, & \forall j \geq 0. \end{aligned}$$

With obvious substitutions, also the operators  $\tilde{P}_j$  acting on the dual spaces can be constructed.

The last important approximation property that we list here is the Jackson inequality. As we are dealing with Hilbert spaces in this work, we will use a particular case. The reader should note that this inequality (as well as the Bernstein inequality, which we will not mention here) is true for far more general Banach spaces, like Besov spaces. The Jackson inequality can be stated in a formal way as follows. If the scaling function  $\varphi$  belongs to a Sobolev space  $H^{s_0}(\mathbb{R}^+)$ , it can be shown that, for all  $0 < s < \min(s_0, L)$ ,

$$\|v - P_j v\|_{L^2(\mathbb{R}^+)} \lesssim 2^{-js} |v|_{H^s}, \quad \forall v \in H^s, \forall j \in \mathbb{N}. \quad (4.45)$$

This means that if we choose a scaling function that is sufficiently regular, we are sure to get a good approximation for all classes of less regular functions. As the level  $j$  increases, the norm of the approximation error tends exponentially to zero at a rate controlled by the regularity index  $s$  of the function being analyzed.

## 4.2 Wavelet spaces in $\mathbb{R}^+$

This section introduces the wavelet spaces on the half line, using the same procedure adopted on  $\mathbb{R}$ . Basically, we want to determine a complement space  $W_0(\mathbb{R}^+)$  such that  $V_1(\mathbb{R}^+) = V_0(\mathbb{R}^+) \oplus W_0(\mathbb{R}^+)$ . The wavelet space  $W_0(\mathbb{R}^+)$  will include the details needed to refine the approximation of a given function. Note that the sum must be a direct sum to insure uniqueness, but does not need to be orthogonal. Indeed, we will look for biorthogonal wavelet spaces, such that

$$V_1(\mathbb{R}^+) = V_0(\mathbb{R}^+) \oplus W_0(\mathbb{R}^+), \quad (4.46)$$

$$\tilde{V}_1(\mathbb{R}^+) = \tilde{V}_0(\mathbb{R}^+) \oplus \tilde{W}_0(\mathbb{R}^+), \quad (4.47)$$

$$W_0(\mathbb{R}^+) \perp \tilde{V}_0(\mathbb{R}^+), \quad (4.48)$$

$$\tilde{W}_0(\mathbb{R}^+) \perp V_0(\mathbb{R}^+). \quad (4.49)$$

The orthogonal setting is obviously a particular case.

The starting point is the scaling functions basis set  $\{\theta_{0k}, k \geq 0\}$  constructed in Section 4.1. Note that this is the basis before the biorthogonalization process, so the first  $\tilde{L}$  functions are, in general, not biorthogonal. By applying the operator  $T_1$  to this set we get a basis for  $V_1(\mathbb{R}^+)$ . We need to express each basis function of  $V_1(\mathbb{R}^+)$  in terms of the basis functions of  $V_0(\mathbb{R}^+)$ , plus some other functions. These functions will generate the wavelet space  $W_0(\mathbb{R}^+)$ . The space  $W_j(\mathbb{R}^+)$  will then be obtained applying again  $T_j$ , as we did to generate the scaling function spaces.

Let us recall two identities that hold on the real line, and that will be needed in the following,

$$\varphi_{1k}^{\mathbb{R}} = \sum_m \tilde{h}_{k-2m} \varphi_{0m}^{\mathbb{R}} + \sum_m \tilde{g}_{k-2m} \psi_{0m}^{\mathbb{R}}, \quad (4.50)$$

$$\psi_{0m}^{\mathbb{R}} = \sum_l g_{l-2m} \varphi_{1l}^{\mathbb{R}}, \quad (4.51)$$

with  $g_n = (-1)^n \tilde{h}_{1-n}$  and  $\tilde{g}_m = (-1)^m h_{1-m}$ . As for the scaling function spaces, also for the wavelet spaces the construction will preserve the internal basis functions, already defined on  $\mathbb{R}$ , and will define the smallest set of modified border functions to achieve completeness.

Let us consider first the internal wavelet functions. From Eq. (4.1) we know that the set  $T_1(V^{(+)}) = \{\varphi_{1k}^{\mathbb{R}} : k \geq k_0^*\}$  was left unchanged in the construction



of the scaling function spaces. Therefore, the wavelet functions that can be generated from this set through Eq. (4.51) will also be left unchanged. Recalling that  $\tilde{h}_n$  is not vanishing only when  $\tilde{n}_0 \leq n \leq \tilde{n}_1$ , we can evaluate the lower bound in the sum of Eq. (4.51) to be  $l = 1 + 2m - \tilde{n}_1$ . This lower bound must also be  $l \geq k_0^*$ . This condition allows to determine the smallest integer  $m = m_0^*$  such that the wavelet functions on the real line can be used without modifications. This integer evaluates to

$$m_0^* := \left\lceil \frac{k_0^* + \tilde{n}_1 - 1}{2} \right\rceil \quad (4.52)$$

and the corresponding internal wavelet space will be

$$W_0^I = \text{span} \{ \xi_{0m} = \psi_{0m}^R : m \geq m_0^* \}. \quad (4.53)$$

We need to define additional functions that, together with the basis functions of  $W_0^I$  and  $V_0(\mathbb{R}^+)$ , will generate all  $V_1(\mathbb{R}^+)$ .

Let us consider now which functions of  $V_1(\mathbb{R}^+)$  we are able to generate with the internal scaling function and wavelet spaces only. These spaces are generated by  $\{\varphi_{0m}^R, m \geq k_0^*\}$  and  $\{\psi_{0m}^R, m \geq m_0^*\}$ , respectively. Substituting these functions in Eq. (4.50) and enforcing the bounds on  $m$ , we get a lower bound on  $k$ . A straightforward calculation shows that this bound is

$$\bar{k} = 2k_0^* + \tilde{n}_1 - 1. \quad (4.54)$$

Therefore, using the internal spaces, we can only construct the internal functions  $\{\varphi_{1k}^R, k \geq \bar{k}\}$ .

However, we do not need to enforce the reconstruction for all the functions  $k < \bar{k}$ , because of the special form of the refinement equation for the non-biorthogonal scaling functions. Indeed, Eq. (4.19) can be rewritten as

$$\theta_{1\alpha} = H_{\alpha\alpha}^{-1} \left[ \theta_{0\alpha} - \sum_{k \geq k_0^*} H_{\alpha k} \varphi_{1k}^R \right],$$

showing that the border functions  $\{\theta_{1k}, k < k_0^*\}$  can be expressed in terms of the corresponding functions at level 0 plus the functions at level 1 with  $k \geq k_0^*$ . Therefore, only the gap  $k_0^* \leq k < \bar{k}$  needs to be filled. The border scaling functions will then be generated automatically through the refinement equation.

It can be proved that the dimension of the space  $W_0^B$ , which collects the border wavelet functions, is exactly  $m_0^*$ . When this space is added to  $W_0^I$ , also the remaining  $\bar{k} - k_0^*$  functions of  $V_1(\mathbb{R}^+)$  will be generated. However, from the definition of  $m_0^*$ , we can see that  $m_0^* = \lceil \frac{\bar{k} - k_0^*}{2} \rceil$ . The cardinality of the wavelet border functions set is then approximately half of the cardinality of the set of

functions that are to be generated. This simply means that these functions are not linearly independent, because part of them can still be expressed as linear combinations of functions in  $V_0(\mathbb{R}^+)$  plus some “internal wavelets”. We omit here the details of the proof, which leads to the definition of the border wavelet functions as

$$\xi_{0,m_0^*-k} := \varphi_{1,\bar{k}-2k+1}^{\mathbb{R}} - P_0 \varphi_{1,\bar{k}-2k+1}^{\mathbb{R}}, \quad \forall k = 1, \dots, m_0^*. \quad (4.55)$$

The corresponding space will then be

$$W_0^B = \{\xi_{0m} : m = 0, \dots, m_0^* - 1\}, \quad (4.56)$$

and the complete wavelet space on the half line will be

$$\begin{aligned} W_0(\mathbb{R}^+) &= \text{span} \{\xi_{0m} : m = 0, \dots, m_0^* - 1\} \oplus \text{span} \{\xi_{0m} = \psi_{0m}^{\mathbb{R}} : m \geq m_0^*\} \\ &= W_0^B \oplus W_0^I. \end{aligned} \quad (4.57)$$

The same scheme leads to the construction of the dual wavelet space,

$$\widetilde{W}_0(\mathbb{R}^+) = \text{span} \{\tilde{\xi}_{0m} : m \geq 0\} = \widetilde{W}_0^B \oplus \widetilde{W}_0^I,$$

where

$$\begin{aligned} \widetilde{W}_0^I &= \{\tilde{\xi}_{0m} = \tilde{\psi}_{0m}^{\mathbb{R}} : m \geq \widetilde{m}_0^*\}, \quad \widetilde{m}_0^* := \left\lceil \frac{k^* + n_1 - 1}{2} \right\rceil \\ \widetilde{W}_0^B &= \{\tilde{\xi}_{0,\widetilde{m}_0^*-k} := \tilde{\varphi}_{1,\bar{l}-2k+1}^{\mathbb{R}} - \tilde{P}_0 \tilde{\varphi}_{1,\bar{l}-2k+1}^{\mathbb{R}} : \forall k = 1, \dots, \widetilde{m}_0^*\} \\ \bar{l} &= 2k^* + n_1 - 1 \end{aligned}$$

The last step consists of the biorthogonalization of the border wavelet functions, in order to satisfy Eqs. (4.48) and (4.49). The biorthogonalization of wavelets is slightly more complex than the corresponding biorthogonalization of the scaling functions. This is due to the fact that both for the primal and dual systems some internal wavelets must be included in the sets to be modified in order to achieve biorthogonality. It is not difficult to prove that the total number of primal and dual wavelets to be modified is

$$m^* = \max\{m_1^*, \widetilde{m}_1^*\},$$

where

$$m_1^* = \left\lceil \frac{\bar{l} + \tilde{n}_1 - 1}{2} \right\rceil, \quad \widetilde{m}_1^* = \left\lceil \frac{\bar{k} + n_1 - 1}{2} \right\rceil.$$

Without loss of generality, and consistently with the assumption  $\tilde{L} \geq L$ , we will set  $\widetilde{m}_1^* \leq m_1^*$ . This is true when the primal scaling functions are B-splines. In

conclusion, the first  $m^*$  primal and dual wavelets need to be redefined through two basis changes,

$$\begin{aligned}\psi_{0m'} &= \sum_{m=0}^{m^*-1} e_{m'm} \xi_{0m} \\ \tilde{\psi}_{0n'} &= \sum_{n=0}^{m^*-1} \tilde{e}_{n'n} \tilde{\xi}_{0n}\end{aligned}\tag{4.58}$$

such that

$$\langle \psi_{0m}, \tilde{\psi}_{0n} \rangle = \delta_{mn} \quad m, n = 0, \dots, m^* - 1$$

This is equivalent to finding two real matrices  $E$  and  $\tilde{E}$  such that

$$EY\tilde{E}^T = I,\tag{4.59}$$

where  $Y$  is the Gramian matrix

$$Y_{mn} = \langle \xi_{0m}, \tilde{\xi}_{0n} \rangle, \quad \forall m, n = 0, \dots, m^* - 1.\tag{4.60}$$

As for the scaling functions, we have constructed two different basis sets for  $W_0(\mathbb{R}^+)$  and  $\tilde{W}_0(\mathbb{R}^+)$ . The first, indicated with  $\{\xi_{0m}, \tilde{\xi}_{0m}, \forall m \geq 0\}$  is not biorthogonal in the first  $m^*$  functions. The second, labelled  $\{\psi_{0m}, \tilde{\psi}_{0m}, \forall m \geq 0\}$  with  $\psi_{0m} = \xi_{0m}$  and  $\tilde{\psi}_{0m} = \tilde{\xi}_{0m}$  for  $m \geq m^*$ , is fully biorthogonal. Indeed, we can easily show that

$$\langle \psi_{jm}, \tilde{\psi}_{j'n} \rangle = \delta_{jj'} \delta_{mn}, \quad \forall j, j', m, n \geq 0.$$

Finally, it can be shown that the basis  $\{\psi_{jm} : m \geq 0\}$  for  $W_j(\mathbb{R}^+)$  (and similarly for the dual system) is uniformly 2-stable for each  $j \geq 0$ .

The following sections give further details about the construction of the wavelet spaces. Section 4.2.1 describes the construction of the border wavelets through the characterization of the projection operators  $P_0$  and  $\tilde{P}_0$ . Section 4.2.2 details the derivation of the filters of the refinement equation for the wavelets, and Section 4.2.3 determines the expression for the Gramian matrix  $Y$ .

### 4.2.1 Projection operators

If we recall the definition of the border wavelets in Eq. (4.55), we need to characterize the projection operators  $P_0$  by evaluating the quantities

$$P_0 \varphi_{1k}^R = \sum_{m \geq 0} \langle \varphi_{1k}^R, \tilde{\varphi}_{0m} \rangle \varphi_{0m}.$$

There are two different cases, depending on whether the index  $k$  is smaller or larger than  $k^*$ . We recall that we only need to project those functions with  $k \geq k_0^*$ . Let us examine these cases separately.

- $k_0^* \leq k < k^*$ . In this case we have  $\varphi_{1k}^R = \theta_{1,L+k-k_0^*}$ , which belongs to the space  $V_1^+ \setminus V_1^I = V_1^+ \cap V_1^B$  of the internal functions that have been included in the border functions for the biorthogonalization. Recalling the definition of the matrix  $D$  in Eq. (4.13), we can write

$$\theta_{1,L+k-k_0^*} = \sum_{n=0}^{\tilde{L}-1} [D^{-1}]_{L+k-k_0^*,n} \varphi_{1n}.$$

We have then

$$\begin{aligned} \langle \varphi_{1k}^R, \tilde{\varphi}_{0m} \rangle &= \sum_{n=0}^{\tilde{L}-1} [D^{-1}]_{L+k-k_0^*,n} \langle \varphi_{1n}, \tilde{\varphi}_{0m} \rangle \\ &= \sum_{n=0}^{\tilde{L}-1} [D^{-1}]_{L+k-k_0^*,n} \sum_{l \geq 0} \tilde{\mathcal{H}}_{ml}^\perp \langle \varphi_{1n}, \tilde{\varphi}_{1l} \rangle \\ &= \sum_{n=0}^{\tilde{L}-1} [D^{-1}]_{L+k-k_0^*,n} \tilde{\mathcal{H}}_{mn}^\perp; \end{aligned}$$

- $k \geq k^*$ . In this case we have  $\varphi_{1k}^R = \theta_{1,\tilde{L}+k-k^*} = \varphi_{1,\tilde{L}+k-k^*}$ . We can use directly the refinement equation on  $\tilde{\varphi}_{0m}$ , obtaining

$$\begin{aligned} \langle \varphi_{1k}^R, \tilde{\varphi}_{0m} \rangle &= \sum_{l \geq 0} \tilde{\mathcal{H}}_{ml}^\perp \langle \varphi_{1,\tilde{L}+k-k^*}, \tilde{\varphi}_{1l} \rangle \\ &= \tilde{\mathcal{H}}_{m,\tilde{L}+k-k^*}^\perp. \end{aligned}$$

Putting the two cases together we get

$$P_0 \varphi_{1k}^R = \sum_{m \geq 0} P_{km} \varphi_{0m}, \quad (4.61)$$

where

$$P_{km} = \begin{cases} \sum_{n=0}^{\tilde{L}-1} [D^{-1}]_{L+k-k_0^*,n} \tilde{\mathcal{H}}_{mn}^\perp & \text{if } k_0^* \leq k < k^* \\ \tilde{\mathcal{H}}_{m,\tilde{L}+k-k^*}^\perp & \text{if } k \geq k^*. \end{cases} \quad (4.62)$$

The projection matrix  $P_{km}$  is then formed by a column of the dual biorthogonal scaling function filter  $\tilde{\mathcal{H}}^\perp$ , eventually multiplied by the matrix  $D^{-1}$ . Recalling the structure of this filter from Section 4.1.1, we can derive the length of the row  $k$ ,

$$P_{km} = 0 \quad \forall m > \tilde{L} + \max \left\{ 0, \left\lceil \frac{k - k^* - \tilde{\delta} + 1}{2} \right\rceil \right\}.$$

We can now express the wavelets in terms of known functions, as

$$\begin{aligned} \xi_{0m} &= \varphi_{1,2m+\bar{k}-2m_0^*+1}^R - P_0 \varphi_{1,2m+\bar{k}-2m_0^*+1}^R \\ &= \varphi_{1,2m+\bar{k}-2m_0^*+1}^R - \sum_{n \geq 0} P_{2m+\bar{k}-2m_0^*+1,n} \varphi_{0n}. \end{aligned} \quad (4.63)$$

In the same way, we can characterize the projectors on the dual spaces. For the dual system, however, due to the initial choice  $\tilde{L} \geq L$ , we do not have to care for the special case  $k < k^*$ , because we only need to project the functions for  $k \geq k^*$ . We have

$$\tilde{P}_{km} = \langle \tilde{\varphi}_{1k}^R, \varphi_{0m} \rangle = \mathcal{H}_{m, \tilde{L}+k-k^*}^\perp, \quad k \geq k^*,$$

with nonzero entries determined by

$$\tilde{P}_{km} = 0 \quad \forall m > \tilde{L} + \max \left\{ 0, \left\lceil \frac{k - k^* - \tilde{L} + L - \delta + 1}{2} \right\rceil \right\},$$

while the corresponding dual wavelets read

$$\tilde{\xi}_{0m} = \tilde{\varphi}_{1, 2m + \tilde{L} - 2\tilde{m}_0^* + 1}^R - \sum_{n \geq 0} \tilde{P}_{2m + \tilde{L} - 2\tilde{m}_0^* + 1, n} \tilde{\varphi}_{0n}. \quad (4.64)$$

### 4.2.2 Wavelet filters

This section is devoted to the derivation of the filters for the primal and dual wavelets, for both the non biorthogonal and the biorthogonal system. The derivation of the non biorthogonal filters is straightforward from the decomposition in Eq. (4.63), which we recall here setting  $k_m = 2m + \bar{k} - 2m_0^* + 1$ ,

$$\xi_{0m} = \varphi_{1, k_m}^R - \sum_{n \geq 0} P_{k_m, n} \varphi_{0n}.$$

We can apply the refinement equation (4.26) to express the sum in terms of the scaling functions at level 1. Also, recalling the expression for  $\tilde{\varphi}_{1, k_m}^R$  already determined in Section 4.2.1, we can immediately write

$$\xi_{0m} = \sum_{l \geq 0} \mathcal{G}_{ml}^\angle \varphi_{1l}, \quad (4.65)$$

where the non biorthogonal filter has the expression

$$\mathcal{G}_{ml}^\angle = \begin{cases} [D^{-1}]_{L+k_m-k_0^*, l} - \sum_{n \geq 0} P_{k_m, n} \mathcal{H}_{nl}^\perp & \text{if } k_0^* \leq k_m < k^*, \\ \delta_{k_m, l} - \sum_{n \geq 0} P_{k_m, n} \mathcal{H}_{nl}^\perp & \text{if } k_m \geq k^*. \end{cases} \quad (4.66)$$

The length of the row  $m$  of this matrix is determined by the product between  $P$  and  $\mathcal{H}^\perp$ , and is given by the last column in  $\mathcal{H}^\perp$  with at least a nonzero entry in the first  $N_m$  rows, where  $N_m$  is the length of row  $m$  in  $P$ . The overall structure of this matrix is obtained by adding the ladder of the internal wavelets filters, as we did for the scaling functions. The number of leading zeros in row  $m_0^*$ , which corresponds to the first internal wavelet, evaluates to  $n_z = \nu_0 + \tilde{L} + 2m_0^* - k^*$ . The structure of  $\mathcal{G}^\angle$  is depicted below.



Putting the two cases together we obtain the biorthogonal refinement equation

$$\psi_{0m} = \sum_{l \geq 0} \mathcal{G}_{ml}^\perp \varphi_{1l}, \quad (4.69)$$

where the wavelet filter is

$$\mathcal{G}_{ml}^\perp = \begin{cases} \sum_{s=0}^{m^*-1} e_{ms} \mathcal{G}_{sl}^\perp & \text{if } m < m^* \\ \mathcal{G}_{ml}^\perp & \text{if } m \geq m^*, \end{cases} \quad (4.70)$$

obtained from the non biorthogonal filter by left multiplying its first  $m^*$  rows by the matrix  $E$  and leaving the others unchanged. The biorthogonal primal filter has the same structure as the non-biorthogonal one, except the number of rows in the upper-left block, which is now equal to  $m^*$ , and consequently the number of vanishing entries at the beginning of row  $m^*$ , which evaluates to  $n'_z = \nu_0 + \tilde{L} + 2m^* - k^*$ . Its global structure is shown below.

$$\mathcal{G}^\perp = \begin{array}{c} \begin{array}{|c|} \hline G_b^0 \\ \hline \end{array} \\ \begin{array}{c} n'_z \quad \begin{array}{|c|c|c|c|} \hline g \\ \hline \end{array} \\ \begin{array}{|c|c|c|c|c|c|} \hline g \\ \hline \end{array} \\ \begin{array}{|c|c|c|c|c|c|c|c|} \hline g \\ \hline \end{array} \end{array} \end{array}$$

The same expressions hold, with obvious substitutions, for the dual filter, depicted below.

$$\tilde{\mathcal{G}}^\perp = \begin{array}{c} \begin{array}{|c|} \hline \tilde{G}_b^0 \\ \hline \end{array} \\ \begin{array}{c} \tilde{n}'_z \quad \begin{array}{|c|c|c|c|} \hline \tilde{g} \\ \hline \end{array} \\ \begin{array}{|c|c|c|c|c|c|} \hline \tilde{g} \\ \hline \end{array} \\ \begin{array}{|c|c|c|c|c|c|c|c|} \hline \tilde{g} \\ \hline \end{array} \end{array} \end{array}$$

### 4.2.3 The Gramian matrix $Y$

The calculation of the Gramian matrix  $Y$  is much simpler than for the corresponding matrix  $X$  used to biorthogonalize the scaling functions. Indeed, we can use the expression of the non-biorthogonal filters  $\mathcal{G}^\perp$  and  $\tilde{\mathcal{G}}^\perp$  to reduce the calculation of the elements of  $Y$  to inner products of scaling functions that are biorthogonal. More precisely, we have

$$Y_{mn} = \langle \xi_{0m}, \tilde{\xi}_{0n} \rangle = \sum_{h \geq 0} \sum_{l \geq 0} \mathcal{G}_{mh}^\perp \tilde{\mathcal{G}}_{nl}^\perp \langle \varphi_{1h}, \tilde{\varphi}_{1l} \rangle = \sum_{l \geq 0} \mathcal{G}_{ml}^\perp \tilde{\mathcal{G}}_{nl}^\perp. \quad (4.71)$$

The matrix  $Y$  is then simply obtained by multiplying the primal and (transposed) dual non orthogonal wavelet filters (the first  $m^*$  rows), according to

$$Y = \mathcal{G}^\angle [\tilde{\mathcal{G}}^\angle]^T.$$

#### 4.2.4 Boundary adaption of wavelets

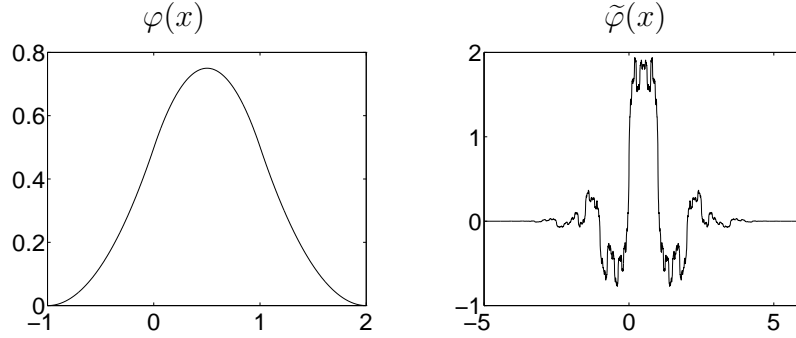
As for the scaling functions, it may be convenient in the applications to have as few as possible nonvanishing wavelets at the edge  $x = 0$ . As a matter of fact, when the scaling functions are boundary adapted, i.e., only one is nonzero at  $x = 0$ , the non-biorthogonal wavelets can be constructed so that only one primal and one dual wavelet is nonzero at  $x = 0$ . The non-biorthogonal filters  $\mathcal{G}^\angle$  and  $\tilde{\mathcal{G}}^\angle$  express the wavelets in terms of the biorthogonal scaling functions, which we suppose to be boundary adapted. If a row has a vanishing entry in the first column, the corresponding wavelets will not depend on the first scaling function but only on the other scaling functions, which all vanish at  $x = 0$ . Consequently, if we perform one loop of gaussian elimination in order to force the first column (starting from the second row) of  $\mathcal{G}^\angle$  and  $\tilde{\mathcal{G}}^\angle$  to be made of zeros, we get automatically boundary adaption for wavelets with a nonzero value at  $x = 0$ . Finally, if the biorthogonalization is performed with care according to the guidelines of Section 4.1.3, also the biorthogonal wavelets will be boundary adapted.

### 4.3 An example

This section will apply the results of the preceding sections to the biorthogonal spline decomposition of the half line for the case  $L = 3$ ,  $\tilde{L} = 5$ . We will show the border scaling functions, the wavelets, and the corresponding filters for the primal and dual systems, before and after the biorthogonalization. Before proceeding with the example, we summarize the main computational steps of the construction.

1. Pick a biorthogonal multilevel decomposition on the real line. This will determine the scaling functions  $\varphi$  and  $\tilde{\varphi}$ , which are uniquely defined by the filters  $h$  and  $\tilde{h}$ .
2. Choose appropriate values of the arbitrary constants  $\delta$  and  $\tilde{\delta}$ , such that Eq. (4.10) is satisfied.
3. Evaluate  $c_{\alpha 0} \forall \alpha$  (resp.  $\tilde{c}_{\beta 0} \forall \beta$ ) through Eq. (4.30).
4. Evaluate the coefficients  $c_{\alpha k}$  (resp.  $\tilde{c}_{\beta k}$ ) through Eq. (4.29).
5. Evaluate the modified filters for the border scaling functions  $H_{\alpha \alpha}$  and  $H_{\alpha k}$  (resp.  $\tilde{H}_{\beta \beta}$  and  $\tilde{H}_{\beta k}$ ) through Eqs. (4.20) and (4.21).





**Figure 4.2:** Primal and dual scaling functions for the biorthogonal spline decomposition of  $\mathcal{R}$  in the case  $L = 3$ ,  $\tilde{L} = 5$

6. Evaluate the Gramian matrix  $X_{\alpha\beta}$  through Eq. (4.1.2).
7. Apply the biorthogonalization procedure to the border scaling functions by solving Eq. (4.14)
8. Determine the filters for the biorthogonal bases with Eq. (4.26).
9. Evaluate the projection matrices  $P$  and  $\tilde{P}$  through Eq. (4.62).
10. Build the filters for the non-biorthogonal border wavelets using Eq. (4.66).
11. Calculate the Gramian matrix  $Y$  from Eq. (4.71).
12. Apply the biorthogonalization procedure to the border wavelets by solving Eq. (4.59)
13. Determine the biorthogonal filters for the border wavelets using Eq. (4.70).

#### 4.3.1 The case $L = 3$ , $\tilde{L} = 5$

In this section we will examine in detail the case  $L = 3$ ,  $\tilde{L} = 5$ . We know that this is a valid choice of the parameters from Section 3.4.3. The characteristic constants of the corresponding multiresolution scheme on  $\mathcal{R}$  are

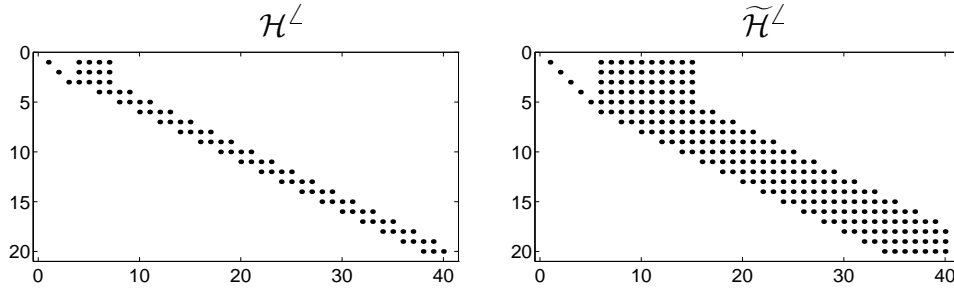
$$\begin{array}{ll} L = 3 & \tilde{L} = 5 \\ n_0 = -1 & \tilde{n}_0 = -5 \\ n_1 = 2 & \tilde{n}_1 = 6. \end{array} \quad (4.72)$$

We show for convenience plots of the primal and dual scaling functions (Fig. 4.2). These scaling functions are symmetric with respect to  $x = 1/2$  and can reproduce locally all polynomials of degree 2 and 4, respectively.

To start the construction of the boundary scaling functions we need to pick a pair of positive integers  $\delta$  and  $\tilde{\delta}$  that determine the left edge of the support

$k$	$\tilde{c}_{0k}$	$\tilde{c}_{1k}$	$\tilde{c}_{2k}$	$\tilde{c}_{3k}$	$\tilde{c}_{4k}$
0	1	0.5	0.5	0.5	0.6
1	1	1.5	2.5	4.5	8.6
2	1	2.5	6.5	17.5	48.6
3	1	3.5	12.5	45.5	168.6
4	1	4.5	20.5	94.5	440.6
5	1	5.5	30.5	170.5	960.6
6	1	6.5	42.5	279.5	1848.6
7	1	7.5	56.5	427.5	3248.6
8	1	8.5	72.5	620.5	5328.6
9	1	9.5	90.5	864.5	8280.6
10	1	10.5	110.5	1165.5	12320.6

**Table 4.1:** Coefficients  $\tilde{c}_{\beta k}$  for  $\beta = 1, \dots, \tilde{L} - 1$ .



**Figure 4.3:** Structure of the filter matrices for the non biorthogonal primal and dual scaling functions.

for the first primal and dual scaling functions that will not be modified (in the non-biorthogonal system). As  $\tilde{L} > L$  we can choose  $\tilde{\delta} = 0$ . The corresponding value of  $\delta$  is then automatically determined by Eq. (4.10), as well as  $k_0^*$  and  $k^*$ ,

$$\begin{aligned} \delta &= 2 & \tilde{\delta} &= 0 \\ k_0^* &= 3 & k^* &= 5. \end{aligned} \tag{4.73}$$

For illustration, we list the coefficients  $\tilde{c}_{\beta k}$  needed in the construction in table 4.1. Due to the choice of the monomials as the basis for the polynomials to be reproduced at the border, the values of these coefficients grows with  $\beta$  and  $k$ .

The structure of the non-biorthogonal scaling function filter matrices  $\mathcal{H}^L$  and  $\tilde{\mathcal{H}}^L$  is visualized in Fig. 4.3. We recall that these are infinite matrices (on one side only), so the number of rows and columns shown in the pictures has been determined to visualize appropriately their global structure. Note the  $L \times L$  and  $\tilde{L} \times \tilde{L}$  upper diagonal blocks, which are due to the choice of the basis

3.0000	4.5833	9.5000	22.5167	57.8000
4.3333	8.7500	20.2333	50.7500	134.6952
8.0000	18.0000	44.0000	114.0000	308.8000
1.0000	3.5000	12.5000	45.5000	168.6000
1.0000	4.5000	20.5000	94.5000	440.6000

**Table 4.2:** The Gramian matrix  $X$ .

of monomials. The number of nonzero entries in the first  $L$  and  $\tilde{L}$  rows for the primal and dual matrices, respectively, is

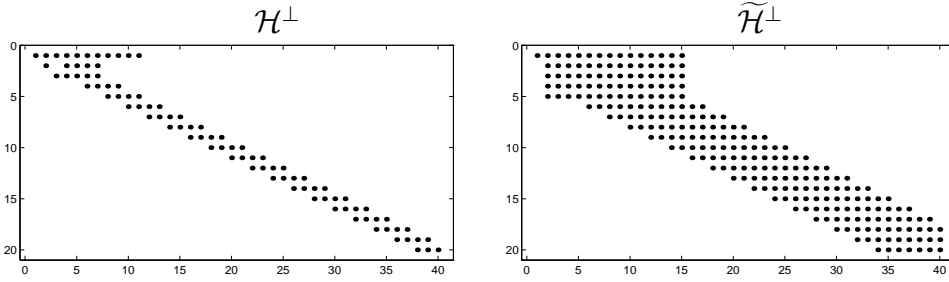
$$\begin{aligned} n_1 + L + k_0^* - 1 &= 7 \\ \tilde{n}_1 + \tilde{L} + k^* - 1 &= 15 \end{aligned}$$

We turn now to the biorthogonalization of the scaling function spaces. The Gramian matrix  $X$  is shown in Table 4.2. Its condition number is about  $2 \times 10^4$ . Therefore, a significant loss of precision is expected in the biorthogonalization of the border scaling functions. As a result, the biorthogonality relations will not be satisfied at the machine precision, because the least significant digits will be lost.

The biorthogonalization system has been solved by setting the entries in the main diagonal of the matrix  $D$  to ones. The structure of the resulting biorthogonal filters  $\mathcal{H}^\perp$  and  $\tilde{\mathcal{H}}^\perp$  for the primal and dual scaling functions is depicted in Fig. 4.4. Note that due to the particular choice of  $D$  in the biorthogonalization, the upper diagonal block has disappeared from  $\tilde{\mathcal{H}}^\perp$ , but not completely from  $\mathcal{H}^\perp$ . The first column has only one nonzero entry, because the boundary-value preserving biorthogonalization was obtained by separating the first border function from the others (see Section 4.1.3). The number of nonzero entries in the first  $\tilde{L}$  rows for the primal and dual filters is, respectively,

$$\begin{aligned} \tilde{L} + k^* + n_1 - 1 &= 11 \\ \tilde{L} + k^* + \tilde{n}_1 - 1 &= 15. \end{aligned}$$

We show now plots of primal and dual scaling functions before and after the biorthogonalization. Figure 4.5 reports the primal scaling functions, while Figure 4.6 reports the duals. Note that all the primal scaling functions except the first do not change with the biorthogonalization. The biorthogonal dual scaling functions, instead, are completely different from their non biorthogonal counterpart, because the matrix  $\tilde{D}$  is nearly full. Note that the support of the first  $L$  primal scaling functions before the biorthogonalization is  $[0, n_1 + k_0^* - 1]$ , while for the first  $\tilde{L}$  duals we have  $[0, \tilde{n}_1 + k^* - 1]$ . After the biorthogonalization, the support of the duals does not change, while for the primal scaling functions



**Figure 4.4:** Structure of the filter matrices for the biorthogonal primal and dual scaling functions.

we get  $[0, n_1 + k^* - 1]$  due to the inclusion of some internal functions (two in this case) to reach a number of  $\tilde{L}$ .

Next, we consider the generation of wavelets and their filters. The constants involved in the construction of the non-biorthogonal wavelets are

$$\begin{aligned} m_0^* &= 4 & \tilde{m}_0^* &= 3 \\ \bar{k} &= 11 & \bar{l} &= 11, \end{aligned}$$

while the total number of primal and dual biorthogonal border wavelets must be

$$m^* = 8.$$

This means that four primal wavelets and five dual wavelets need to be added for the biorthogonalization.

The first step is the construction of the projection matrices  $P$  and  $\tilde{P}$ . We do not need to determine all the entries in these matrices, but only the rows needed for the construction of the border wavelets. These rows correspond in the present case to the indices  $k = 4, 6, 8, 10$  for primals and  $k = 4, 6, 8$  for duals.

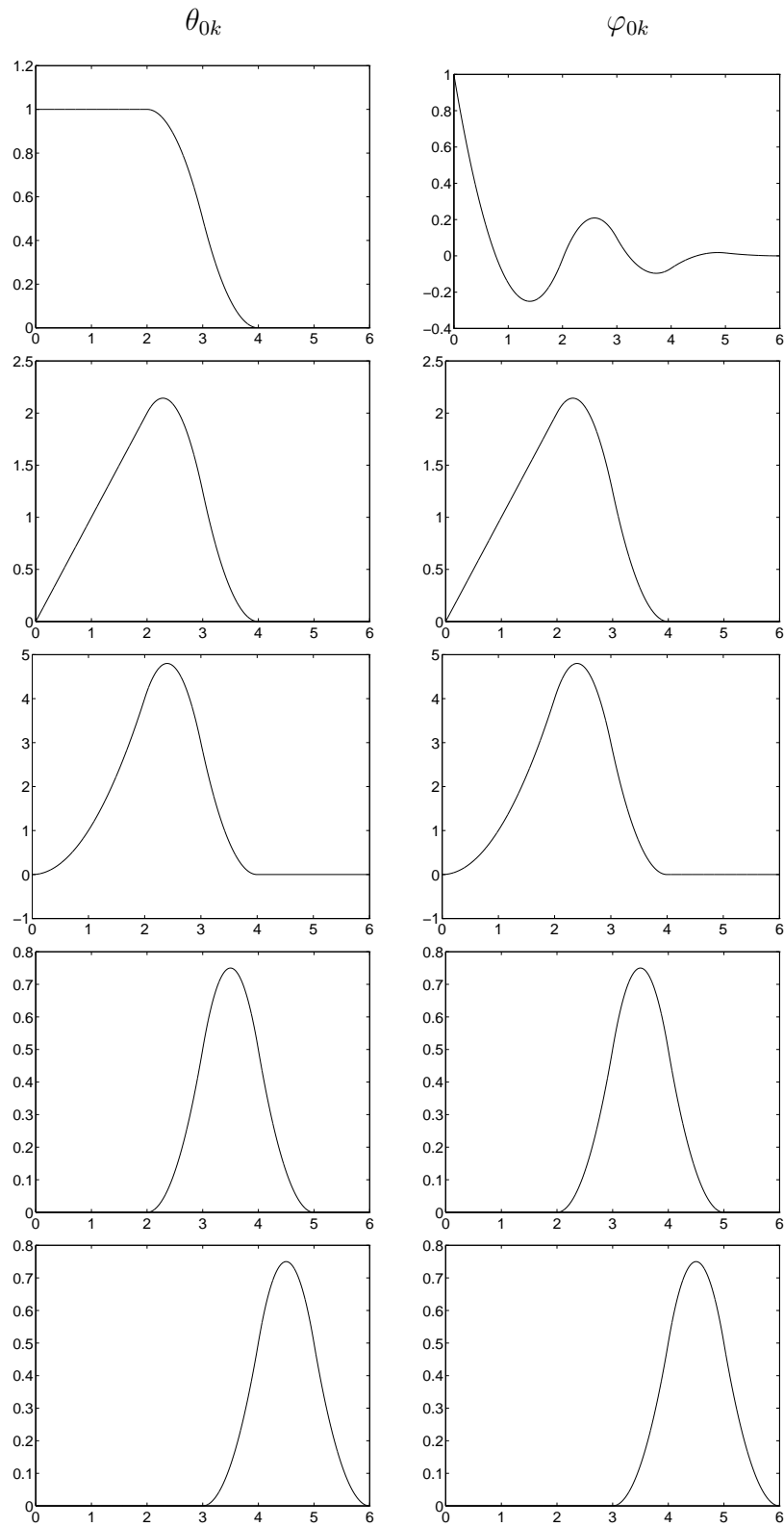
The non-biorthogonal wavelet filters  $\mathcal{G}^\perp$  and  $\tilde{\mathcal{G}}^\perp$  are shown in Fig 4.7. The structure is similar to the one of the scaling function filters. Note that there is only one nonvanishing entry in the first column of the filters. These are indeed the filters after the boundary adaption of the border wavelets, which is possible because also the scaling functions are boundary adapted.

The biorthogonal wavelet filters  $\mathcal{G}^\perp$  and  $\tilde{\mathcal{G}}^\perp$  are shown in Fig. 4.8. The biorthogonalization was performed as for the scaling functions, i.e., setting the diagonal entries in the primal change of basis matrix  $E$  to ones. This choice preserves the primal wavelets except the first and modifies all the dual wavelets. It should be mentioned that the condition number of the biorthogonalization matrix  $Y$  is about  $8.4 \times 10^6$ . This means that also for the wavelets the biorthogonalization relations will be satisfied with reduced accuracy.

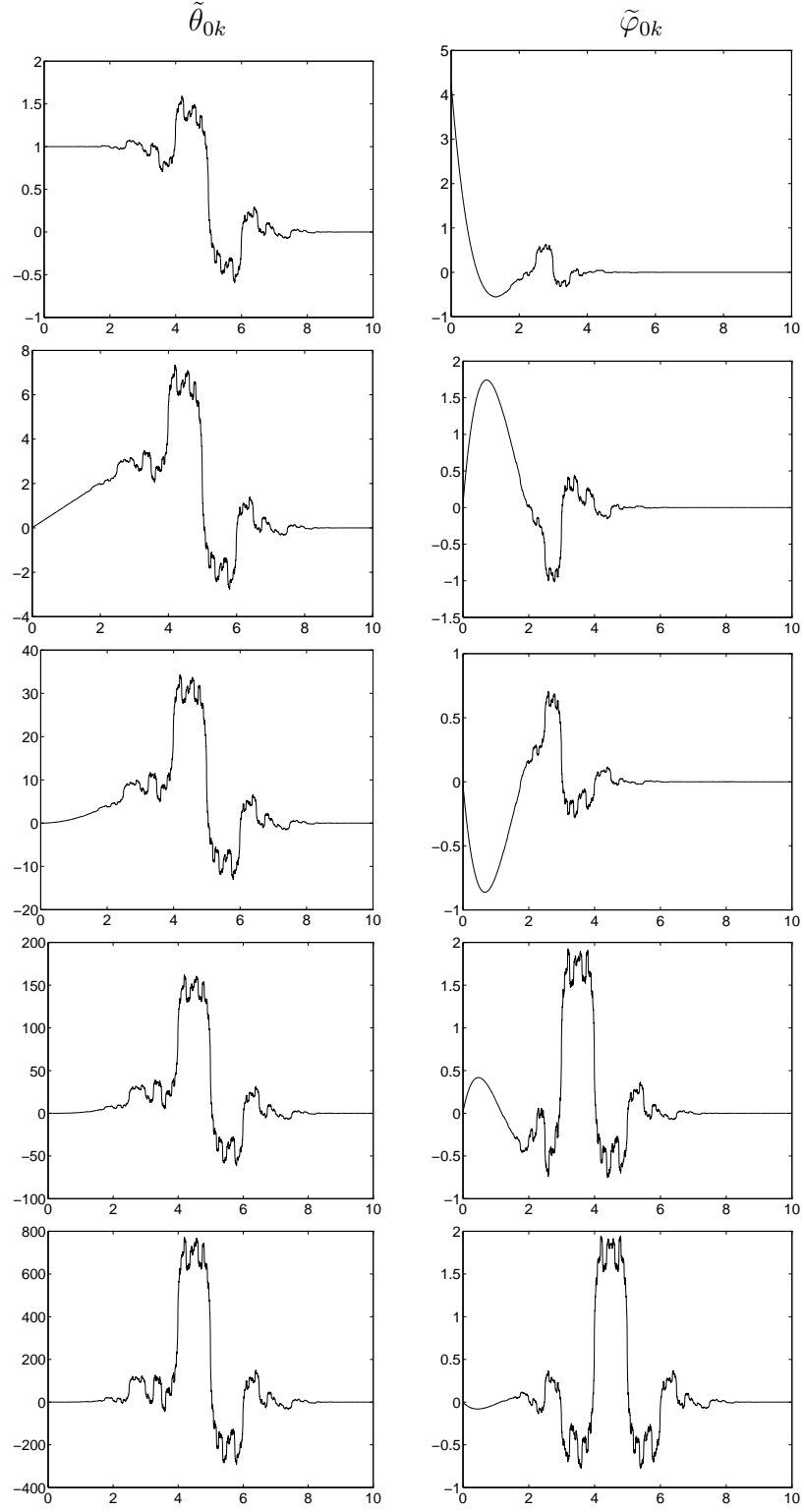
Finally, we show the plots of the border wavelets. Figure 4.9 shows the primal and dual non-biorthogonal wavelets, Fig. 4.10 shows the biorthogonal

primal wavelets, and Fig. 4.11 shows the biorthogonal dual wavelets. Also in this case the biorthogonal primal wavelets are identical to the non-biorthogonal ones except the first due to the choice of matrix  $E$ . Due to the bad condition number of the matrix  $Y$  and the initial choice of the monomials, the resulting dual biorthogonal wavelets have a bad behavior, with highly varying values between one function and the other.

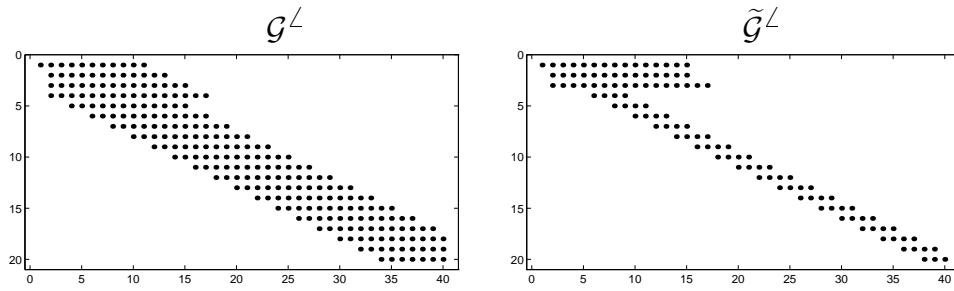
At the end of this chapter, in Section 4.4.4, we will show that the Bernstein polynomials and a careful biorthogonalization can be employed to generate better behaved scaling function and wavelet systems.



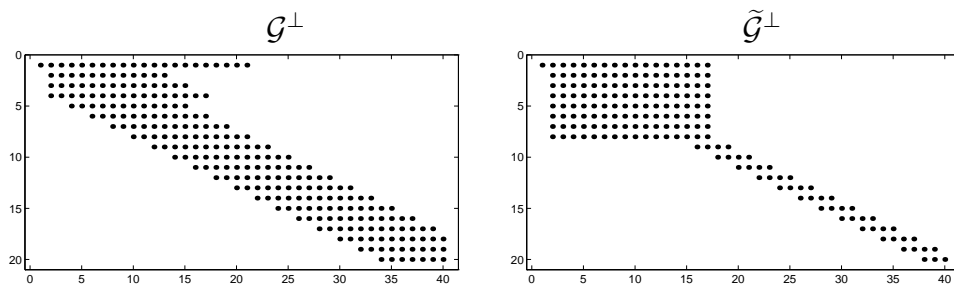
**Figure 4.5:** Primal border scaling functions  $\theta_{0k}$  (left column) and  $\varphi_{0k}$  (right column). The index  $k$  ranges from 0 (top) to  $\tilde{L} - 1$  (bottom).



**Figure 4.6:** Dual border scaling functions  $\tilde{\theta}_{0k}$  (left column) and  $\tilde{\varphi}_{0k}$  (right column). The index  $k$  ranges from 0 (top) to  $\tilde{L} - 1$  (bottom).

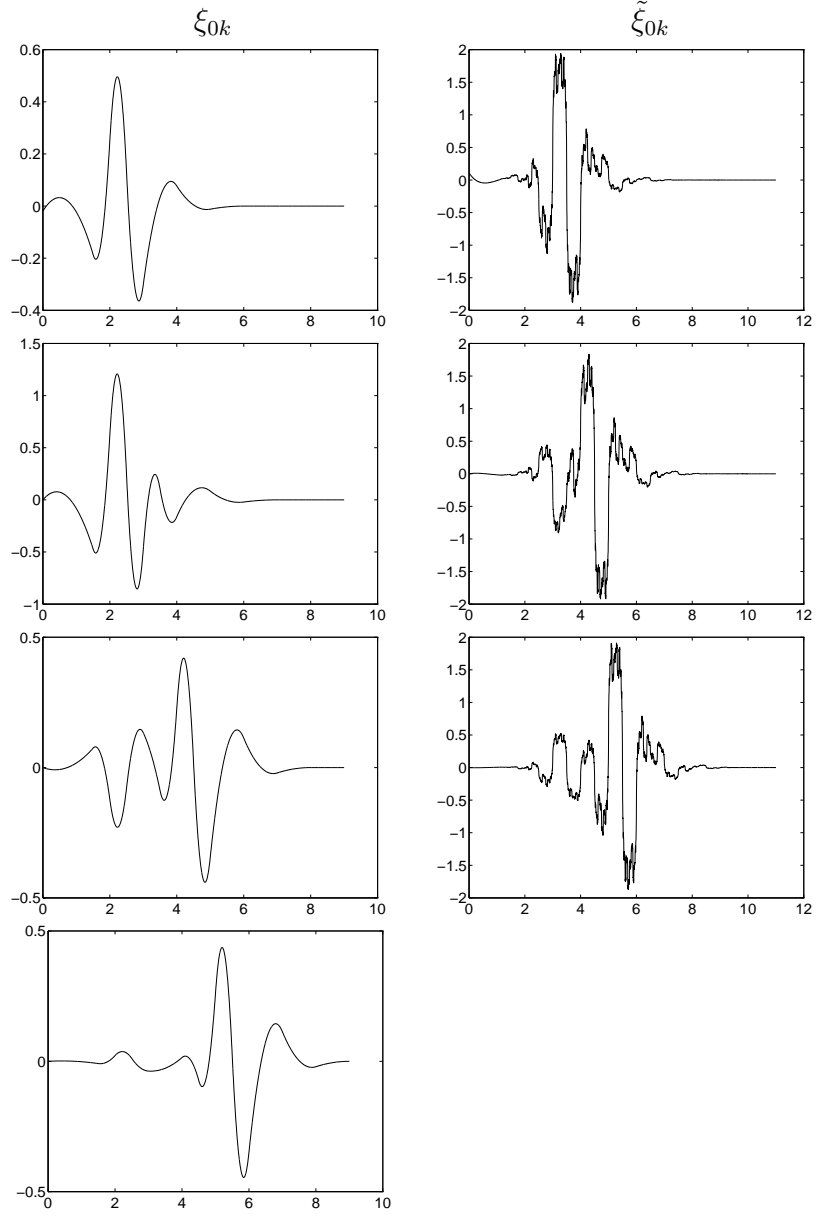


**Figure 4.7:** Structure of the filter matrices for the non-biorthogonal primal and dual wavelets.

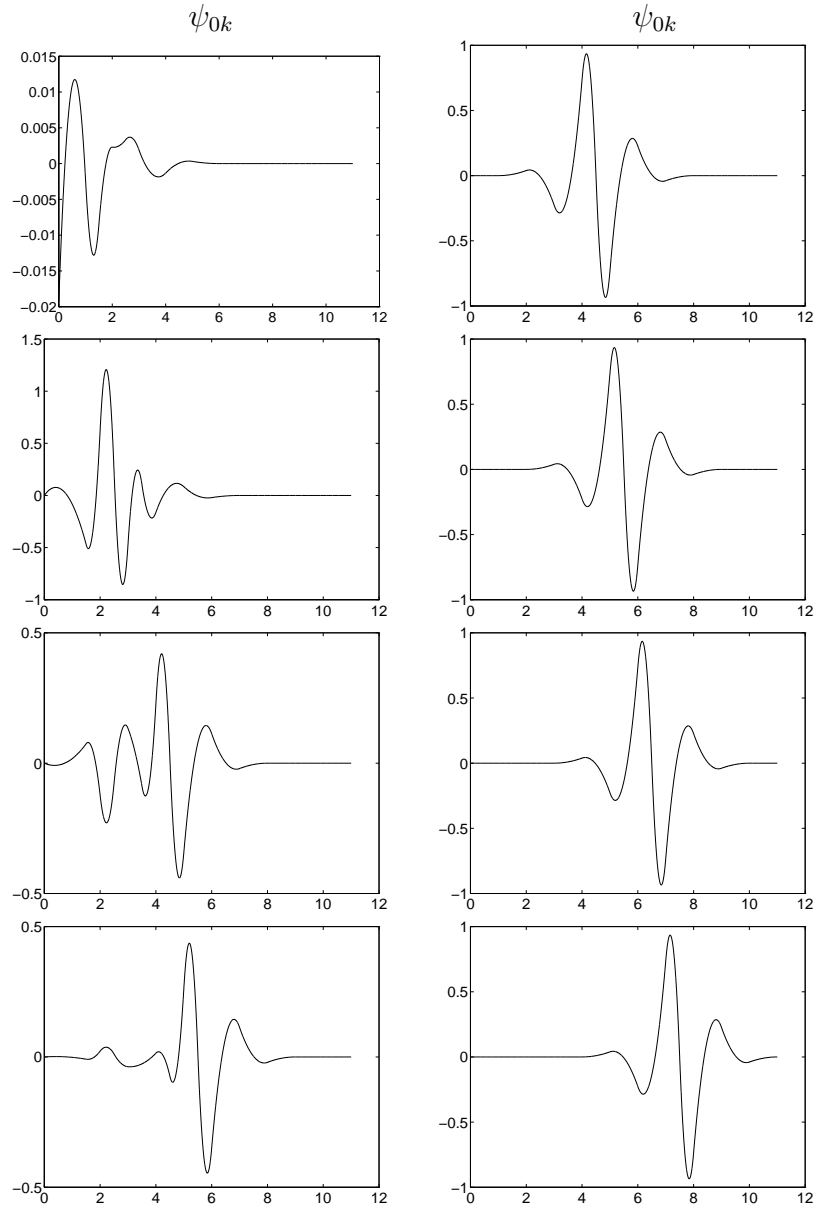


**Figure 4.8:** Structure of the filter matrices for the biorthogonal primal and dual wavelets.

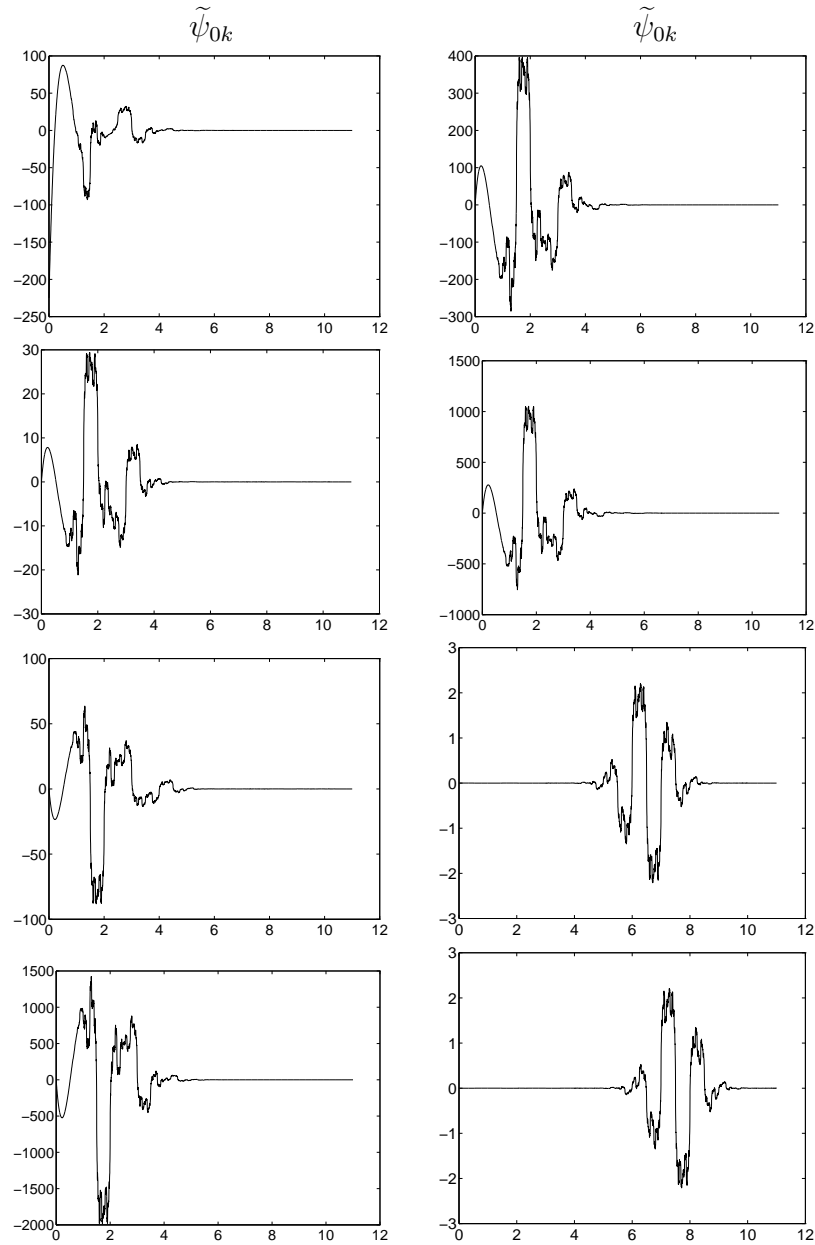




**Figure 4.9:** Primal (left column) and dual (right column) non-biorthogonal border wavelets. The index  $k$  ranges from 0 (top) to  $m_0^* - 1$  (bottom) for primal wavelets and from 0 (top) to  $\tilde{m}_0^* - 1$  (bottom) for dual wavelets.



**Figure 4.10:** Primal biorthogonal border wavelets. The index  $k$  ranges from 0 (top) to 3 (bottom) in the left column and from 4 (top) to 7 (bottom) in the right column.



**Figure 4.11:** Dual biorthogonal border wavelets. The index  $k$  ranges from 0 (top) to 3 (bottom) in the left column and from 4 (top) to 7 (bottom) in the right column.

## 4.4 The unit interval

This section describes how the multilevel decomposition of the half line derived in the preceding sections can be adapted to build a decomposition of the unit interval  $[0, 1]$ . The main point is to merge two parallel decompositions of the positive and negative half line, after a translation of the latter by 1. In addition, we want to decouple the effects of the two boundary points, so that each of them can be studied separately. This requires to start with an initial level  $j_0 > 0$ , in order to have the supports of the scaling functions small enough with respect to the length of the whole domain  $[0, 1]$ . More precisely, we will set  $j_0$  so that the scaling function spaces on the unit interval will be constructed as

$$\begin{aligned} V_j(0, 1) = & \text{span} \{ \varphi_{jl}^{(0)} : l \in \mathcal{I}_L \} \oplus \\ & \text{span} \{ \varphi_{jk} : k \in \mathcal{I}_I \} \oplus \\ & \text{span} \{ \varphi_{jr}^{(1)} : r \in \mathcal{I}_R \}, \forall j \geq j_0, \end{aligned}$$

with  $\mathcal{I}_I \neq \emptyset$  and where the boundary functions  $\varphi_{jl}^{(0)}$  and  $\varphi_{jr}^{(1)}$  are constructed independently. Here, and from now on, the suffix  $^{(0)}$  or  $^{(1)}$  refers to the boundary point 0 or 1, respectively.

### 4.4.1 Scaling function spaces

The first step is the construction of the scaling function spaces at level 0 for the primal and dual systems on the half line  $\mathbb{R}^-$ . We will parallel our construction of Section 4.1, adapting it to the case  $x \leq 0$ . Therefore, we will not show all the details of the construction.

We need to specify two bases for the spaces of polynomials  $\mathcal{P}_{L-1}$  and  $\mathcal{P}_{\tilde{L}-1}$ . These can be the same of the ones used for  $\mathbb{R}^+$  or different. In the following, we will use a particular basis to obtain symmetric scaling functions at the two edges. After fixing these bases, we can define modified border functions by imposing the reproduction of the polynomials. The right edge of the support for the first scaling functions that will not be changed in the construction is tuned by two nonnegative integers  $\delta_1$  and  $\tilde{\delta}_1$ . Therefore, the internal scaling function spaces will be formed by functions with support in  $(-\infty, -\delta_1]$  and  $(-\infty, -\tilde{\delta}_1]$  for the primal and dual system, respectively. The newly defined border functions will then have the expression

$$\theta_\alpha^{(0-)}(x) = \sum_{k=1-\delta_1-n_1}^{-n_0-1} c_{\alpha k}^{(1)} \varphi_{0k}^{\mathcal{R}}(x), \quad x \leq 0, \quad \forall \alpha = 0, \dots, L-1,$$

and similarly for the duals. The scaling function spaces can be expressed as the linear combination of internal and border functions, just like in Eq. (4.9),

$$V_0(\mathbb{R}^-) = \text{span} \{ \theta_\alpha^{0-} : \alpha = 0, \dots, L-1 \} \oplus$$

$$\begin{aligned}
& \text{span } \{\varphi_{0k}^R : k = 1 - \tilde{\delta}_1 - \tilde{n}_1, \dots, -n_1 - \delta_1\} \oplus \\
& \text{span } \{\varphi_{0k}^R : k \leq -\tilde{\delta}_1 - \tilde{n}_1\} \\
\tilde{V}_0(\mathbb{R}^-) = & \text{span } \{\tilde{\theta}_\beta^{0-} : \beta = 0, \dots, \tilde{L} - 1\} \oplus \\
& \text{span } \{\tilde{\varphi}_{0k}^R : k \leq -\tilde{\delta}_1 - \tilde{n}_1\}.
\end{aligned}$$

Matching the dimensions of the border spaces for the subsequent biorthogonalization, we get a relation between  $\delta_1$  and  $\tilde{\delta}_1$  similar to Eq. (4.10),

$$\tilde{\delta}_1 - \delta_1 = \tilde{L} - L - (\tilde{n}_1 - n_1). \quad (4.74)$$

The spaces at any level  $j \geq 0$  can be obtained through the operator  $T_j$  applied to  $V_0(\mathbb{R}^-)$  and  $\tilde{V}_0(\mathbb{R}^-)$ .

We need now to translate the construction to have its origin in the point  $x = 1$  and domain in  $(-\infty, 1]$ . We recall that the translation of the spaces at level  $j$  requires exactly  $2^j$  steps to shift the point  $x = 0$  in  $x = 1$ . It is then straightforward to get the expressions for the new border functions at  $x = 1$ ,

$$\theta_{j\alpha}^{(1)}(x) = \sum_{k=2^j+1-\delta_1-n_1}^{2^j-n_0-1} c_{\alpha,k-2^j}^{(1)} \varphi_{jk}^R(x), \quad x \leq 1,$$

and similarly for the duals. The scaling function spaces will be

$$\begin{aligned}
V_0(-\infty, 1) &= \text{span } \{\theta_{j\alpha}^{(1)} : \alpha = 0, \dots, L - 1\} \oplus \\
& \text{span } \{\varphi_{jk}^R : k = 2^j + 1 - \tilde{\delta}_1 - \tilde{n}_1, \dots, 2^j - n_1 - \delta_1\} \oplus \\
& \text{span } \{\varphi_{jk}^R : k \leq 2^j - \tilde{\delta}_1 - \tilde{n}_1\} \\
\tilde{V}_0(-\infty, 1) &= \text{span } \{\tilde{\theta}_{j\beta}^{(1)} : \beta = 0, \dots, \tilde{L} - 1\} \oplus \\
& \text{span } \{\tilde{\varphi}_{0k}^R : k \leq 2^j - \tilde{\delta}_1 - \tilde{n}_1\}.
\end{aligned}$$

We impose now that the two sets of border functions in  $x = 0$  (see Eq. (4.9)) and  $x = 1$  do not overlap. As we are still working under the hypothesis  $\tilde{L} \geq L$ , we only need to use the dual system. We obtain the relation

$$-\tilde{n}_0 + \tilde{\delta}_0 \leq 2^j - \tilde{\delta}_1 - \tilde{n}_1,$$

that can be used to evaluate the minimal starting level  $j_0$ ,

$$j_0 = \lceil \log_2(\tilde{n}_1 - \tilde{n}_0 + \tilde{\delta}_0 + \tilde{\delta}_1) \rceil. \quad (4.75)$$

In summary, we have constructed the scaling function spaces

$$\begin{aligned}
V_j(0, 1) &= \text{span } \{\theta_{jk}^{(0)} : k = 0, \dots, L - 1\} \oplus \\
& \text{span } \{\varphi_{jk}^R : k = -n_0 + \delta_0, \dots, 2^j - \delta_1 - n_1\} \oplus \\
& \text{span } \{\theta_{jk}^{(1)} : k = 0, \dots, L - 1\},
\end{aligned}$$

and similarly for the dual system. These expressions are valid when  $j \geq j_0$ . Note that the dimensions of the primal and dual scaling function spaces are the same,

$$\dim V_j(0, 1) = 2^j + 2L + 1 - \delta_0 - \delta_1 - n_1 + n_0, \quad \forall j \geq j_0.$$

Since  $V_{j+1}(0, 1) = V_j(0, 1) \oplus W_j(0, 1)$ , this implies that the dimension of the wavelet spaces will be

$$\dim W_j(0, 1) = 2^{j+1} - 2^j = 2^j$$

for both primal and dual systems.

We describe now how to choose the polynomial basis sets in order to obtain symmetric scaling functions at the two edges of the interval. The key point is obviously to start with symmetric polynomials and with symmetric scaling functions on  $\mathbb{R}$ . Therefore, we will particularize our derivation to the biorthogonal B-spline system, and we will choose a basis set  $p_{1,\alpha}$  for  $\mathbb{P}_{L-1}$  in  $\mathbb{R}^-$  such that

$$p_{1,\alpha}(y) = p_{0,\alpha}(-y) \quad \alpha = 0, \dots, L,$$

where  $p_{0,\alpha}$  are the basis elements used in  $0^+$ . This leads immediately to a relation between the modified scaling functions in 0 and 1,

$$\theta_{j\alpha}^{(1)}(x) = \theta_{j\alpha}^{(0)}(1 - x).$$

To obtain now symmetric filters we set  $\delta_1 = \delta_0$  and  $\tilde{\delta}_1 = \tilde{\delta}_0$ . This is possible because Eqs. (4.10) and (4.74) can still be satisfied. Therefore, we will not distinguish between these integers at the two edges, and we will indicate them with  $\delta$  and  $\tilde{\delta}$ , respectively.

In the symmetric B-spline case the coefficients  $c_{\alpha k}^{(1)}$  can be easily related to  $c_{\alpha k}^{(0)}$ . Indeed we have

$$\begin{aligned} c_{\alpha k}^{(1)} &= \int_{\mathbb{R}} p_{1,\alpha}(y) \tilde{\varphi}_{0k}(y) dy = \int_{\mathbb{R}} p_{0,\alpha}(-y) \tilde{\varphi}(y - k) dy \\ &= \int_{\mathbb{R}} p_{0,\alpha}(y) \tilde{\varphi}(y + k + r) dy = c_{\alpha, -k-r}^{(0)}, \end{aligned}$$

where  $r = \text{rem}(L, 2)$ . From the B-spline properties shown in Section 3.4.3 we know that  $r$  and the limits of the filters on  $\mathbb{R}$  are related through

$$n_1 = r - n_0.$$

This leads to the following expression of the border functions in terms of the scaling functions of the multiresolution on  $\mathbb{R}$ ,

$$\theta_{j\alpha}^{(0)}(x) = \sum_{k=1}^{n_1 - n_0 + \delta - 1} C_{\alpha, k} \varphi_{j, k - n_1}^{\mathbb{R}}(x) \Big|_{[0, 1]}, \quad (4.76)$$

$$\theta_{j\alpha}^{(1)}(x) = \sum_{k=1}^{n_1 - n_0 + \delta - 1} C_{\alpha, k} \varphi_{j, 2^j - k - n_0}^{\mathbb{R}}(x) \Big|_{[0, 1]}, \quad (4.77)$$

where the matrix  $C_{\alpha k}$  corresponds to the coefficients evaluated at the left edge  $x = 0^+$ ,

$$C_{\alpha k} = \left\{ c_{\alpha, k-n_1}^{(0)}, \alpha = 0, \dots, L, k = 1, \dots, n_1 - n_0 - 1 + \delta \right\}. \quad (4.78)$$

Note that in Eq. (4.77) the scaling functions on  $\mathbb{R}$  are superimposed with the same matrix  $C$ , but with a decreasing translation index. Equivalently their support shifts towards left when  $k$  increases. If we reorder now the sequence of basis functions of  $V_j(0, 1)$  as

$$\begin{aligned} V_j(0, 1) &= \text{span} \{ \theta_{jk}^{(0)} : k = 0, \dots, L-1 \} \oplus \\ &\quad \text{span} \{ \varphi_{jk}^{\mathbb{R}} : k = -n_0 + \delta_0 \dots, 2^j - \delta_1 - n_1 \} \oplus \\ &\quad \text{span} \{ \theta_{j, L-k}^{(1)} : k = 0, \dots, L-1 \} = \\ &= \text{span} \{ \theta_{jl}, l = 0, \dots, \dim V_j(0, 1) - 1 \}, \end{aligned}$$

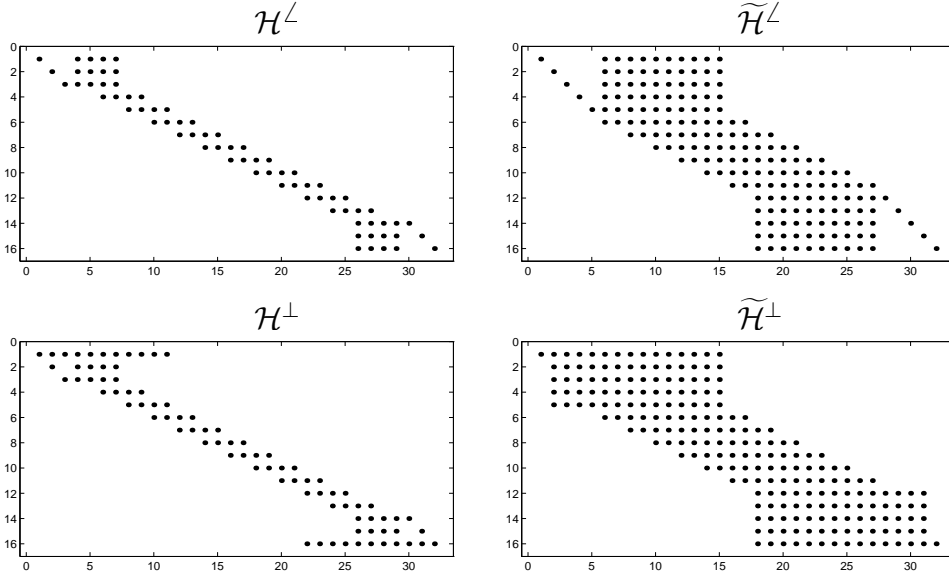
we can express the refinement equation with the filter matrix  $\mathcal{H}^\angle$  in a simple form,

$$\theta_{jl} = \sum_m \mathcal{H}_{lm}^\angle \theta_{j+1, m}.$$

The structure of  $\mathcal{H}^\angle$  is shown in the picture below,

$$\mathcal{H}^\angle = \begin{array}{|c|c|} \hline \mathcal{H}_0^\angle & \\ \hline & \mathcal{H}_c^\angle \\ \hline & \mathcal{H}_1^\angle \\ \hline \end{array}$$

where  $\mathcal{H}_1^\angle$  is obtained from  $\mathcal{H}_0^\angle$  by reversing the order of rows and columns and the center block  $\mathcal{H}_c^\angle$  is formed by a ladder of filters  $h$  as already shown in Section 4.1.1 for the half line filters. It should be noted that the bottom right block does not need to be explicitly computed. In addition, as the filters  $h$  are also symmetric, the overall matrix  $\mathcal{H}^\angle$  is invariant under an inversion of the order of its rows and columns. The filters for the dual scaling functions and their biorthogonal versions for primals and duals are derived in the same way as above, therefore the details will not be presented here. The structure of the filter matrices remains unchanged. In conclusion, once the filter matrices for the half line have been obtained, the filter matrices for the scaling functions on  $[0, 1]$



**Figure 4.12:** Structure of the filter matrices for the scaling functions on the unit interval at level  $j_0 = 4$ .

are easily derived by imposing that their dimension matches the dimension of  $V_j$  (rows) and  $V_{j+1}$  (columns), and by forcing the symmetry conditions described above.

As an example, we show in Fig. 4.12 the structure of the filter matrices in the case of the B-spline multiresolution with  $L = 3$  and  $\tilde{L} = 5$  already treated in Section 4.3.1 for the half line. In the figure the level  $j$  is set equal to the minimum allowed level  $j_0$ . This leads to the dimensions of the scaling function spaces  $\dim V_j(0, 1) = 16$  and  $\dim V_{j+1}(0, 1) = 32$ . These are the number of rows and columns, respectively, in the filter matrices.

We turn now to the derivation of a direct link between the biorthogonal scaling functions on the interval and the biorthogonal scaling functions on  $\mathbb{R}$ . This will be used in Section 5.2 for the evaluation of integrals of refinable functions on the interval.

Let us consider the primal system of biorthogonal scaling functions on  $[0, 1]$  at refinement level  $j$ , with  $j \geq j_0$ . We consider them indexed as  $\varphi_{jk}$ , with  $k$  ranging from 0 to  $\dim V_j(0, 1) - 1$ . There are three different cases.

- $k < \tilde{L}$ . In this case the biorthogonalization matrix  $D$  can be used to express the functions in terms of the non biorthogonal scaling functions, through

$$\varphi_{jk} = \sum_{n=0}^{L-1} d_{kn} \theta_{jn} + \sum_{n=L}^{\tilde{L}-1} d_{kn} \theta_{jn}$$



$$\begin{aligned}
&= \sum_{n=0}^{L-1} d_{kn} \sum_{l=1}^{n_1-n_0+\delta-1} C_{nl} \varphi_{j,l-n_1}^R \Big|_{[0,1]} + \sum_{n=L}^{\tilde{L}-1} d_{kn} \varphi_{j,k_0^*+n-L}^R \\
&= \sum_{l=1}^{n_1-n_0+\delta-1} \left\{ \sum_{n=0}^{L-1} d_{kn} C_{nl} \right\} \varphi_{j,l-n_1}^R \chi_{[0,1]} + \sum_{n=L}^{\tilde{L}-1} d_{kn} \varphi_{j,k_0^*+n-L}^R,
\end{aligned}$$

where we used the definition of the matrix  $C$  in Eq. (4.78) and the indicator function of the unit interval  $\chi_{[0,1]}$ .

- $\tilde{L} \leq k < \dim V_j(0,1) - \tilde{L}$ . In this case the scaling functions are internal and correspond to

$$\varphi_{jk} = \theta_{jk} = \varphi_{j,k^*+k-\tilde{L}}^R$$

- $\dim V_j(0,1) - \tilde{L} \leq k \leq \dim V_j(0,1) - 1$ . This case can be analyzed from the first case through direct use of symmetry.

If we split the matrix  $D$  as

$$D = \left( D^l \mid D^r \right),$$

with  $D^l$  including the first  $L$  columns and  $D^r$  the remaining  $\tilde{L} - L$  columns, we can define a matrix  $\mathcal{M}^0$  as

$$\mathcal{M}^0 = \left( D^l C \mid D^r \right).$$

This matrix can be used to express the functions  $\varphi_{jk}$  in terms of the functions  $\varphi_{jk}^R$ ,

$$\begin{bmatrix} \underline{\varphi}_j^0 \\ \underline{\varphi}_j^I \\ \underline{\varphi}_j^1 \end{bmatrix} = \begin{bmatrix} \mathcal{M}^0 & 0 & 0 \\ 0 & I & 0 \\ 0 & 0 & \mathcal{M}^1 \end{bmatrix} \begin{bmatrix} \underline{\varphi}_j^{R,0} \\ \underline{\varphi}_j^{R,I} \\ \underline{\varphi}_j^{R,1} \end{bmatrix} \chi_{[0,1]}, \quad (4.79)$$

where the block  $\mathcal{M}^1$  is derived from the block  $\mathcal{M}^0$  by reversing the order of rows and columns and the arrays on the left and right side are defined formally as

$$\begin{aligned}
\underline{\varphi}_j^0 &= [\varphi_{j,0}, \dots, \varphi_{j,\tilde{L}-1}]^T \\
\underline{\varphi}_j^I &= [\varphi_{j,\tilde{L}}, \dots, \varphi_{j,\dim V_j - \tilde{L} - 1}]^T \\
\underline{\varphi}_j^1 &= [\varphi_{j,\dim V_j - \tilde{L}}, \dots, \varphi_{j,\dim V_j - 1}]^T \\
\underline{\varphi}_j^{R,0} &= [\varphi_{j,1-n_1}^R, \dots, \varphi_{j,k^*-1}^R]^T \\
\underline{\varphi}_j^{R,I} &= [\varphi_{j,k^*}^R, \dots, \varphi_{j,2^j - \tilde{\delta} - \tilde{n}_1}^R]^T \\
\underline{\varphi}_j^{R,1} &= [\varphi_{j,2^j+1-\tilde{\delta}-\tilde{n}_1}^R, \dots, \varphi_{j,2^j-1-n_0}^R]^T.
\end{aligned}$$

The size of the identity matrix in Eq. (4.79) is  $2^j - 2\tilde{\delta} - \tilde{n}_1 + \tilde{n}_0 + 1$ , and the number of columns in the top-left and bottom-right blocks is  $n_1 - \tilde{n}_0 + \tilde{\delta} - 1$ .

Similar results hold for the dual system, where the scaling functions  $\tilde{\varphi}_{jk}^{\mathcal{R}}$  involved range from  $k = 1 - \tilde{n}_1$  to  $k = 2^j - 1 - \tilde{n}_0$ . The size of the identity matrix in the dual case is the same, while the number of columns in the top-left and bottom-right blocks becomes  $\tilde{n}_1 - \tilde{n}_0 + \tilde{\delta} - 1$ .

#### 4.4.2 Wavelet spaces

Section 4.4.1 showed that the construction of the scaling functions on the unit interval is readily obtained from the construction on the half line, provided that the filters are symmetric and the refinement level  $j$  is larger than a minimum refinement level  $j_0$ . This last requirement cannot be relaxed, because the whole construction on the unit interval is based upon a decoupling of the effects of the left and right borders. Instead, the first requirement is not strictly necessary. The scaling functions and wavelets can also be constructed starting from asymmetric functions, like the Daubechies' ones. We will not discuss here this more general construction, because the symmetric setting is sufficient for our applications. In addition, as we plan to use scaling functions and wavelets on the unit interval to solve differential problems which treat the two boundaries in the same way, it does seem appropriate to use symmetric basis functions. Given these assumptions, the construction of wavelets on the unit interval can be readily obtained from the construction on the half line in the same way as we did for the scaling functions. We will not give all the details of the construction here, but only the main results which lead to the definition of the  $2^j$  wavelets and to the primal and dual wavelet filters. We will describe the derivation of the biorthogonal wavelets. The non-biorthogonal wavelets are derived in the same way through obvious substitutions.

We recall from Section 4.2 that the number of modified border wavelets at the edge  $x = 0$  is  $m^*$ . Therefore, the number of border wavelets at the edge  $x = 1$  will also be  $m^*$ . As the dimension of the wavelet space  $W_j(0, 1)$  must be  $2^j$ , there will be  $2^j - 2m^*$  internal wavelets that remain unchanged from the construction on  $\mathcal{R}$ . This is true both for the primal space and for the dual space. Therefore, we can construct the wavelet space on the unit interval as

$$\begin{aligned} W_j(0, 1) &= \{\psi_{jm}^{(0)} : m = 0, \dots, m^* - 1\} \oplus \\ &\quad \{\psi_{jm}^{\mathcal{R}} : m = m^*, \dots, 2^j - m^* - 1\} \oplus \\ &\quad \{\psi_{jm}^{(1)} : m = 0, \dots, m_0^* - 1\} \\ &= \{\psi_{jm} : m = 0, \dots, 2^j - 1\} \end{aligned} \quad (4.80)$$

where the border wavelets  $\psi_{jm}^{(1)}$  at the edge  $x = 1$  are expressed in terms of the border wavelets at  $x = 0$  through reflection and translation by 1,

$$\psi_{jm}^{(1)}(x) = \psi_{jm}^{(0)}(1 - x).$$

This definition leads to a refinement equation for the primal biorthogonal wavelets as

$$\psi_{jm} = \sum_l \mathcal{G}_{ml}^\perp \varphi_{j+1,l}, \quad (4.81)$$

where  $\{\varphi_{j+1,l}, l = 0, \dots, \dim V_{j+1}(0, 1)\}$  represent the biorthogonal scaling functions on the unit interval at level  $j+1$ . The matrix  $\mathcal{G}^\perp$  has a structure depicted below,

$$\mathcal{G}^\perp = \begin{array}{|c|c|} \hline \mathcal{G}_0^\perp & \\ \hline & \mathcal{G}_c^\perp \\ \hline & \mathcal{G}_1^\perp \\ \hline \end{array}$$

where the bottom right block  $\mathcal{G}_1^\perp$  is obtained from the top left block  $\mathcal{G}_0^\perp$  by reversing the order of rows and columns, and the center block  $\mathcal{G}_c^\perp$  is made of a ladder of wavelet filters on  $\mathbb{R}$ .

This construction is possible when the two systems of border wavelets do not interact with each other. This leads to the definition of a minimum refinement level  $j_0^w$ , which can be different from the minimum level  $j_0$  required by the scaling function construction. The level  $j_0^w$  can be determined by imposing that the support of the rightmost border wavelet at the edge  $x = 0$  is strictly included in the interval  $[0, 1/2]$ . For the B-spline wavelets this leads to the expression

$$j_0^w = \left\lceil 1 + \log_2 \left( 2\tilde{L} + \frac{3}{2}L + \text{rem}(L, 2) - 3 \right) \right\rceil. \quad (4.82)$$

The actual minimum refinement level will be

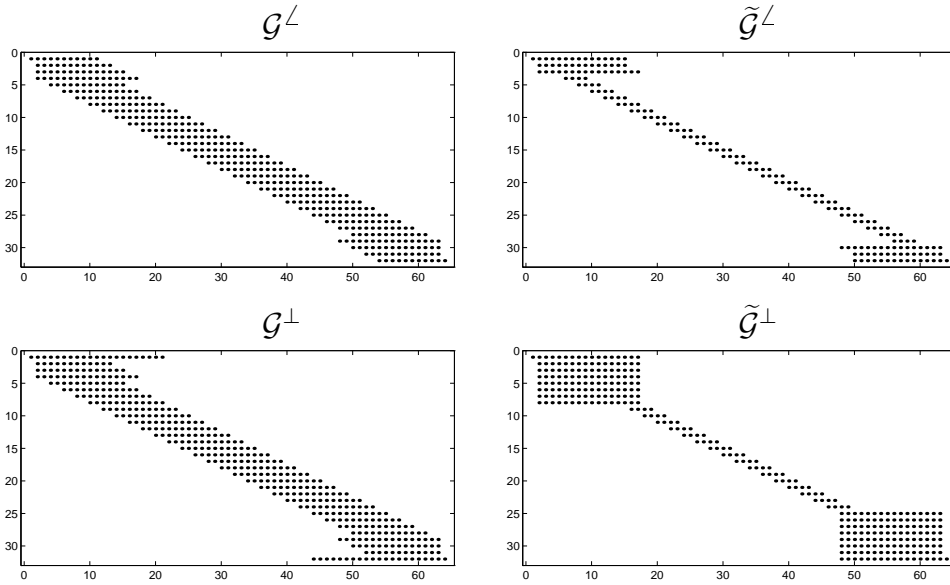
$$J_0 = \max\{j_0, j_0^w\} \quad (4.83)$$

Table 4.4.2 shows the value of the minimum refinement levels  $j_0$  and  $j_0^w$  for various pairs  $(L, \tilde{L})$ .

As an example, we report in Fig. 4.13 the structure of the non-biorthogonal and biorthogonal primal and dual wavelet filter matrices for the example  $L = 3$ ,  $\tilde{L} = 5$  treated in Section 4.3.1. The refinement level is set in the figure to the minimum allowed  $j_0^w = 5$ .

$L$	$\tilde{L}$	$j_0$	$j_0^w$
2	2	2	3
2	4	3	4
2	6	4	5
2	8	4	5
3	5	4	5
3	7	4	6
3	9	5	6
3	11	5	6
4	8	5	6
4	10	5	6
4	12	5	6
4	14	5	6

**Table 4.3:** Minimum refinement levels for the B-spline multiresolution on the unit interval for various pairs  $(L, \tilde{L})$ .



**Figure 4.13:** Structure of the filter matrices for the wavelets on the unit interval at level  $j = 5$ .

#### 4.4.3 Wavelet analysis and synthesis

A biorthogonal multiresolution on  $\mathbb{R}$  is fully characterized by the primal and dual scaling function and wavelet filters. Similarly, the multiresolution on the unit interval that has been constructed is fully characterized by the modified scaling functions and wavelets filter matrices. These matrices are listed below without the  $^\perp$  superscript because we will refer hereafter to the biorthogonal

filters. Instead, we emphasize with the superscript  $(j)$  the dependence of these matrices on the refinement level  $j$ .

- Primal scaling function filter matrix  $\mathcal{H}^{(j)}$ . The dimensions are  $\dim V_j(0, 1) \times \dim V_{j+1}(0, 1)$ . The matrix  $\mathcal{H}^{(j)}$  defines the refinement equation for the primal scaling functions

$$\varphi_{jk} = \sum_m \mathcal{H}_{km}^{(j)} \varphi_{j+1,m}.$$

- Dual scaling function filter matrix  $\widetilde{\mathcal{H}}^{(j)}$ . The dimensions are  $\dim \widetilde{V}_j(0, 1) \times \dim \widetilde{V}_{j+1}(0, 1)$ . The matrix  $\widetilde{\mathcal{H}}^{(j)}$  defines the refinement equation for the dual scaling functions

$$\widetilde{\varphi}_{jk} = \sum_m \widetilde{\mathcal{H}}_{km}^{(j)} \widetilde{\varphi}_{j+1,m}.$$

- Primal wavelet filter matrix  $\mathcal{G}^{(j)}$ . The dimensions are  $2^j \times \dim V_{j+1}(0, 1)$ . The matrix  $\mathcal{G}^{(j)}$  defines the refinement equation for the primal wavelets

$$\psi_{jk} = \sum_m \mathcal{G}_{km}^{(j)} \varphi_{j+1,m}.$$

- Dual wavelet filter matrix  $\widetilde{\mathcal{G}}^{(j)}$ . The dimensions are  $2^j \times \dim \widetilde{V}_{j+1}(0, 1)$ . The matrix  $\widetilde{\mathcal{G}}^{(j)}$  defines the refinement equation for the dual wavelets

$$\widetilde{\psi}_{jk} = \sum_m \widetilde{\mathcal{G}}_{km}^{(j)} \widetilde{\varphi}_{j+1,m}.$$

The biorthogonality of the scaling functions and wavelets bases can be restated in terms of these four matrices. In particular, we have the following identities, which are valid for any refinement level  $j \geq J_0$

- Biorthogonality of primal and dual scaling functions

$$\mathcal{H} \widetilde{\mathcal{H}}^T = I. \quad (4.84)$$

The equivalent relation on  $\mathcal{R}$  is expressed in Eqs. (3.8) and (3.9).

- Biorthogonality of primal and dual wavelets

$$\mathcal{G} \widetilde{\mathcal{G}}^T = I. \quad (4.85)$$

The equivalent relation on  $\mathcal{R}$  is Eq. (3.52).

- Orthogonality between dual scaling functions and primal wavelets

$$\mathcal{G} \widetilde{\mathcal{H}}^T = 0. \quad (4.86)$$

The equivalent relation on  $\mathcal{R}$  is Eq. (3.49)(first row).

- Orthogonality between primal scaling functions and dual wavelets

$$\tilde{\mathcal{G}} \mathcal{H}^T = 0. \quad (4.87)$$

The equivalent relation on  $\mathbb{R}$  is Eq. (3.49)(second row).

Finally, the reconstruction property stated formally in Eq. (2.32) and for the multiresolution on  $\mathbb{R}$  in Eq. (3.50) becomes

$$\tilde{\mathcal{H}}^T \mathcal{H} + \tilde{\mathcal{G}}^T \mathcal{G} = I. \quad (4.88)$$

Let us consider now an arbitrary function  $v \in L^2(\mathbb{R})$  and suppose that the expansion coefficients  $\check{v}_{Jk}$  into the primal scaling functions at refinement level  $J$  have been computed in some way,

$$P_J v(x) = \sum_k \check{v}_{Jk} \varphi_{Jk}.$$

The computation of these coefficients will be the subject of Section 5.1.4. We particularize now the expressions for the wavelet analysis and synthesis stated formally in Eqs. (2.33)-(2.35). The analysis expressions, valid  $\forall j = J_0, \dots, J-1$ , are

$$\check{v}_{jk} = \sum_m \tilde{\mathcal{H}}_{km}^{(j)} \check{v}_{j+1,m} \quad (4.89)$$

$$\hat{v}_{jk} = \sum_m \tilde{\mathcal{G}}_{km}^{(j)} \check{v}_{j+1,m}, \quad (4.90)$$

while the synthesis relation, valid  $\forall j = J_0, \dots, J-1$ , is

$$\check{v}_{j+1,m} = \sum_k \mathcal{H}_{km}^{(j)} \check{v}_{jk} + \sum_k \mathcal{G}_{km}^{(j)} \hat{v}_{jk}. \quad (4.91)$$

These relations are readily derived from the refinement equations for scaling functions and wavelets (see the items in the list above) and from the reconstruction identity

$$\varphi_{j+1,m} = \sum_k \tilde{\mathcal{H}}_{km}^{(j)} \varphi_{jk} + \sum_k \tilde{\mathcal{G}}_{km}^{(j)} \psi_{jk}$$

through use of the biorthogonality of primal and dual systems. Equations (4.89)-(4.91) can be viewed as simple matrix-vector products at any fixed refinement level  $j$ , although the practical implementation should take advantage of the particular structure of the filter matrices. An optimized code can perform the wavelet analysis and synthesis in  $O(N)$  operations, where  $N = \dim V_{j+1}(0, 1)$  is the starting number of coefficients.

For future reference, we will describe the full wavelet analysis and synthesis processes through all possible levels  $j = J_0, \dots, J-1$  with abstract operators, defined below. We introduce the following notations

- $\check{\mathbf{v}}$ : array with the  $\dim V_J$  scaling function coefficients of a general function  $v_J = P_J v$ ,  $v \in L^2$ .
- $\hat{\mathbf{v}}$ : array with the  $\dim V_{J_0}$  scaling function coefficients at level  $J_0$  followed by the wavelet coefficients at increasing levels  $j = J_0, \dots, J-1$ .
- $\hat{\mathbf{v}}_j$ : array with the  $\dim V_j$  scaling function coefficients at level  $j$  followed by the wavelet coefficients at increasing levels  $j, \dots, J-1$ . Clearly  $\hat{\mathbf{v}}_{J_0} = \hat{\mathbf{v}}$  and  $\hat{\mathbf{v}}_J = \check{\mathbf{v}}$ .

We can restate Eqs. (4.89)-(4.91) with these notations as

$$\hat{\mathbf{v}}_j = \widetilde{\mathcal{W}}_j \hat{\mathbf{v}}_{j+1}, \quad \hat{\mathbf{v}}_{j+1} = \mathcal{W}_j^T \hat{\mathbf{v}}_j$$

for  $j = J_0, \dots, J-1$ . The operators  $\widetilde{\mathcal{W}}_j$  and  $\mathcal{W}_j$  are depicted in the figure below, and obviously satisfy the identities

$$\widetilde{\mathcal{W}}_j \mathcal{W}_j^T = \mathcal{W}_j^T \widetilde{\mathcal{W}}_j = I.$$

$$\mathcal{W}_j = \begin{array}{c|c} \begin{array}{c} \mathcal{H}^{(j)} \\ \hline \mathcal{G}^{(j)} \end{array} & 0 \\ \hline 0 & I \end{array} \quad \widetilde{\mathcal{W}}_j = \begin{array}{c|c} \begin{array}{c} \widetilde{\mathcal{H}}^{(j)} \\ \hline \widetilde{\mathcal{G}}^{(j)} \end{array} & 0 \\ \hline 0 & I \end{array}$$

The full chain of wavelet analysis can finally be expressed through

$$\hat{\mathbf{v}} = \widetilde{\mathcal{W}} \check{\mathbf{v}}, \quad \check{\mathbf{v}} = \mathcal{W}^T \hat{\mathbf{v}}$$

where the operators  $\widetilde{\mathcal{W}}$  and  $\mathcal{W}$  are defined as

$$\mathcal{W} = \mathcal{W}_{J_0} \cdots \mathcal{W}_{J-1}, \quad \widetilde{\mathcal{W}} = \widetilde{\mathcal{W}}_{J_0} \cdots \widetilde{\mathcal{W}}_{J-1}. \quad (4.92)$$

It should be noted that, due to the highly sparse structure of these operators, the application of  $\widetilde{\mathcal{W}}$  and  $\mathcal{W}^T$  can be performed in  $O(\dim V_J)$  operations.

#### 4.4.4 An example (continued)

The foregoing section showed that the biorthogonal multiresolution on the unit interval is characterized by four filter matrices that satisfy Eqs (4.84)-(4.88). However, under a practical standpoint, these relations cannot be satisfied exactly. This is due to roundoff errors in the implementation of the algorithms in a computer code. In this section we show which are the major sources of loss of accuracy in the determination of the scaling function and wavelet filters,

and therefore in all applications using this construction of wavelets on the unit interval. We will use as a test case the example presented in Section 4.3.1, i.e. the B-spline biorthogonal multiresolution for  $L = 3$ ,  $\tilde{L} = 5$ .

The best achievable accuracy in numerical computations is the machine precision, which varies between platforms. The order of magnitude in double precision floating point arithmetic is generally  $\varepsilon \sim 10^{-16}$ . It is well known that some ill-conditioned problems lead to the amplification of small perturbations, thus reducing the overall accuracy in the solution. Even a simple linear system can only be solved at a reduced accuracy when the system matrix has a large condition number. The construction of scaling functions and wavelets on the unit interval is based upon the solution of two linear systems, namely Eq. (4.14) and Eq. (4.59). These are the biorthogonalization systems for scaling functions and wavelets. It is then important that the system matrices  $X$  and  $Y$  are well-conditioned. Otherwise, the biorthogonality and reconstruction identities in Eqs. (4.84)-(4.88) will be satisfied at a reduced accuracy.

We want to define significant quantities that can be immediately related to the losses of accuracy in the wavelet computations. We will consider these measures, where  $\max |\cdot|$  indicates the maximum magnitude among all the elements of the matrix in the argument.

- Maximum error in the biorthogonality between primal and dual scaling functions,

$$\varepsilon_1 = \max |\mathcal{H} \tilde{\mathcal{H}}^T - I|$$

- Maximum error in the biorthogonality between primal and dual wavelets,

$$\varepsilon_2 = \max |\mathcal{G} \tilde{\mathcal{G}}^T - I|$$

- Maximum error in the orthogonality between primal wavelets and dual scaling functions,

$$\varepsilon_3 = \max |\mathcal{G} \tilde{\mathcal{H}}^T|$$

- Maximum error in the orthogonality between primal scaling functions and dual wavelets,

$$\varepsilon_4 = \max |\tilde{\mathcal{G}} \mathcal{H}^T|$$

- Maximum error in the reconstruction,

$$\varepsilon_5 = \max |\tilde{\mathcal{H}}^T \mathcal{H} + \tilde{\mathcal{G}}^T \mathcal{G} - I|$$

In addition, we compute the condition number of the matrices  $X$  and  $Y$ .

Let us consider the aforementioned example. Table 4.4 shows the quantities listed above in the case  $L = 3$ ,  $\tilde{L} = 5$ . In column (a) the basis of monomials has been used for the polynomials, and the biorthogonalization has been performed



	(a)	(b)	(c)	(d)
$\varepsilon_1$	$2.7 \times 10^{-12}$	$5.7 \times 10^{-14}$	$5.7 \times 10^{-14}$	$5.7 \times 10^{-14}$
$\varepsilon_2$	$5.1 \times 10^{-9}$	$1.2 \times 10^{-8}$	$3.1 \times 10^{-9}$	$6.0 \times 10^{-10}$
$\varepsilon_3$	$1.6 \times 10^{-12}$	$3.2 \times 10^{-13}$	$1.3 \times 10^{-8}$	$1.4 \times 10^{-11}$
$\varepsilon_4$	$1.4 \times 10^{-9}$	$7.1 \times 10^{-10}$	$2.8 \times 10^{-14}$	$6.2 \times 10^{-11}$
$\varepsilon_5$	$1.4 \times 10^{-9}$	$9.8 \times 10^{-11}$	$1.2 \times 10^{-10}$	$1.6 \times 10^{-11}$
$\text{cond}(X)$	$2.1 \times 10^4$	$2.9 \times 10^1$	$2.9 \times 10^1$	$2.9 \times 10^1$
$\text{cond}(Y)$	$8.4 \times 10^6$	$8.4 \times 10^6$	$8.4 \times 10^6$	$8.4 \times 10^6$

**Table 4.4:** Errors in the case  $L = 3$ ,  $\tilde{L} = 5$ .

for both scaling functions and wavelets by setting the main diagonal of the primal change of basis matrices  $D$  and  $E$  to ones. It should be noted that the condition number of both  $X$  and  $Y$  is quite large, leading to large errors for all the biorthogonality relations. In column (b) the basis of Bernstein polynomials with  $b = 2$ ,  $\tilde{b} = 4$  has been used, with the same biorthogonalization procedure. The condition number of the matrix  $X$  is significantly reduced. The corresponding error on the biorthogonality between primal and dual scaling functions is small ( $\varepsilon_1$ ). However, the condition number of  $Y$  has not changed, so there is a large error in the biorthogonality between primal and dual wavelets ( $\varepsilon_2$ ). As the primal wavelets are unchanged (except the first for boundary adaption), the error involving the primal wavelet filter and not the dual wavelet filter ( $\varepsilon_3$ ) is also small. It can be concluded that the largest errors occur in the biorthogonalization of the border wavelets. This can be checked in column (c), where the same Bernstein polynomials have been used, but the biorthogonalization of wavelets was obtained by leaving the duals unchanged (except the first). The role of primals and duals in the magnitude of errors is exchanged with the results in column (b).

The weight of errors can be balanced between primal and dual wavelets given a badly conditioned matrix  $Y$ . In order to do so, the biorthogonalization of wavelets can be achieved through SVD decomposition. More precisely, the matrix  $Y$  is decomposed into two orthogonal matrices  $U$  and  $V$  and a diagonal matrix  $S$ ,

$$Y = U S V^T.$$

Then, the matrices  $E$  and  $\tilde{E}$  are computed according to

$$E = S^{-r} U^T, \quad \tilde{E} = S^{1-r} V^T,$$

where  $r$  is an additional parameter that can be used to attribute part of the singular values of  $Y$  to the primal wavelets and part to the dual wavelets. If  $r = 0.5$  we have a perfect balance between primal and dual wavelets. Column (d) of Table 4.4 was obtained with this biorthogonalization procedure (together with boundary adaption). Note that the errors are now more balanced between

all the measures considered in the first 5 rows of the table. In particular, the maximum error ( $6.0 \times 10^{-10}$ ) is smaller than the corresponding maximum errors in the other cases (a), (b), and (c) by at least one order of magnitude.

In conclusion, a careful biorthogonalization procedure together with a good choice of the polynomial basis can reduce the loss of accuracy in the computation of wavelets. We showed that the errors due to the biorthogonalization of border wavelets dominate on errors due to other sources. The main responsible is the biorthogonalization matrix  $Y$ , which has a large condition number. At present, it does not seem possible to reduce further the condition number of  $Y$  with this construction of wavelets. Therefore, the overall accuracy at which computations can be performed using these wavelets on the interval cannot be higher than the limits described in this section. However, when only scaling functions are employed in the computations, very high accuracies can be achieved by using the Bernstein polynomials instead of the monomials. When very high accuracies are required, it is advisory to use extended precision arithmetic for the computation of the filters.

## 4.5 Wavelet approximations

This section focuses on the approximation properties of the wavelet expansions on the unit interval. The first issue that will be discussed in Section 4.5.1 is the behavior of the approximation error of a given function  $f$  with different refinement levels. The second issue (Section 4.5.2) will be the concept of nonlinear approximation based on thresholding of the wavelet coefficients. This capability of wavelets is of paramount importance, because it allows very high compression rates in the representation of functions, it can limit memory usage and save computation time in numerical applications. This is the main reason why there is so much interest for wavelets in the literature. The reader is referred to [54, 53] and references therein for reviews of the theory and applications.

### 4.5.1 Linear approximations

The concept of linear approximation has already been discussed throughout the foregoing sections. The key is the characterization of functional spaces based on the projection operators  $P_j$  and  $\tilde{P}_j$  associated to biorthogonal multiresolution analyses on bounded or unbounded domains. Bearing in mind our application, we will focus here on the biorthogonal B-spline multiresolution on the unit interval constructed in Chapter 4.

The behavior of the approximation error of a given function  $f$  at a refinement level  $j$  is explicitly predicted by the Jackson inequality of Eq. (4.45), which we recall here for convenience. If the scaling function  $\varphi$  belongs to a Sobolev space

$H^{s_0}$ , it can be shown that, for all  $s < \min(s_0, L)$ ,

$$\|v - P_j v\|_{L^2} \lesssim 2^{-js} |v|_{H^s}, \quad \forall v \in H^s, \forall j \in \mathbb{N}.$$

This means that if the function  $f$  under investigation is sufficiently regular, the decay of the approximation error is controlled by the regularity of the scaling function  $\varphi$ .

We illustrate this fact on a simple example. Let the function  $f$  be a gaussian pulse centered at  $x = 1/2$ ,

$$f(x) = \exp \left\{ -\frac{(x - 1/2)^2}{0.002} \right\}. \quad (4.93)$$

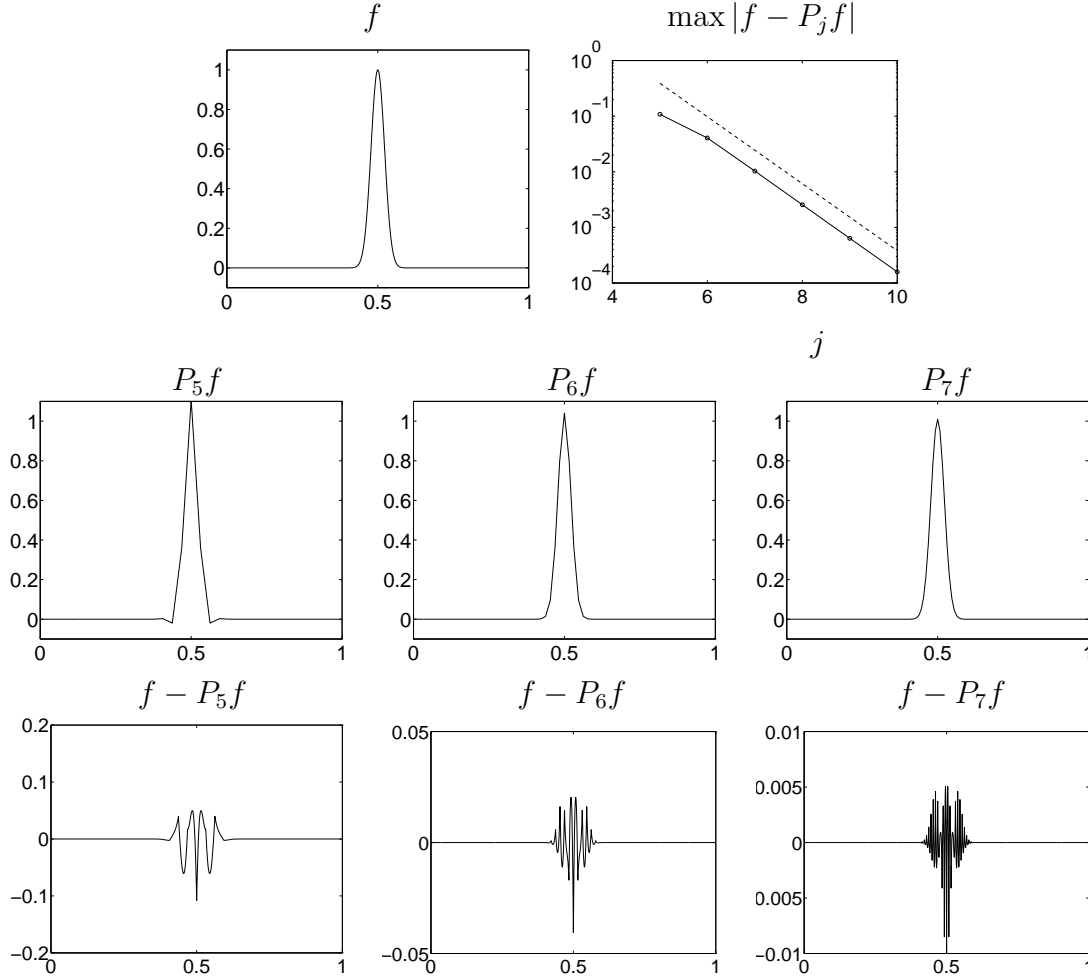
This function, plotted in Fig. 4.14 (top left panel) is  $\mathcal{C}^\infty$  on  $[0, 1]$ . Let us consider for example the biorthogonal B-spline multiresolution with  $L = 2$  and  $\tilde{L} = 4$ . The primal scaling functions are piecewise linear, and all polynomials of degree at most 1 can be locally reproduced. From the Jackson inequality we expect that the  $L^2$  approximation error should decay at least as  $2^{-2j}$  for increasing  $j$ . The top right panel shows the  $L^\infty$  approximation error (not predicted by Jackson, which only deals with the space  $L^2$ ) as a function of the refinement level  $j$ . Even if the  $L^\infty$  norm is stronger than the  $L^2$  norm, the approximation error decays with slope  $-2$ , as expected. The second row in the figure shows different approximations of the function  $f$ , and the third row depicts the corresponding approximation errors as functions of  $x$ .

A further investigation on the behavior of the approximation error was performed by repeating the same analysis with a different multiresolution, namely the biorthogonal B-spline system with  $L = 3$  and  $\tilde{L} = 5$ . In this case the primal scaling function is piecewise quadratic, so all polynomials of degree at most 2 can be locally reproduced. The behavior of the approximation error for increasing  $j$  should decay at least as  $2^{-3j}$ . This is confirmed by the top right panel of Fig. 4.15.

In conclusion, for regular functions  $f$ , the behavior of the approximation error can be explicitly bounded by the Jackson inequality. In particular, both the multiresolution system to be adopted and the maximum needed refinement level  $j$  can be chosen *a priori* once the functions to be represented are known. With the approximations described in this section, the choice of a multiresolution based on a scaling function that is more regular than the functions to be represented is useless. This is the limit of the so-called *linear approximations*. In the next section we show that highly optimized representations can be obtained even when the function  $f$  has low regularity.

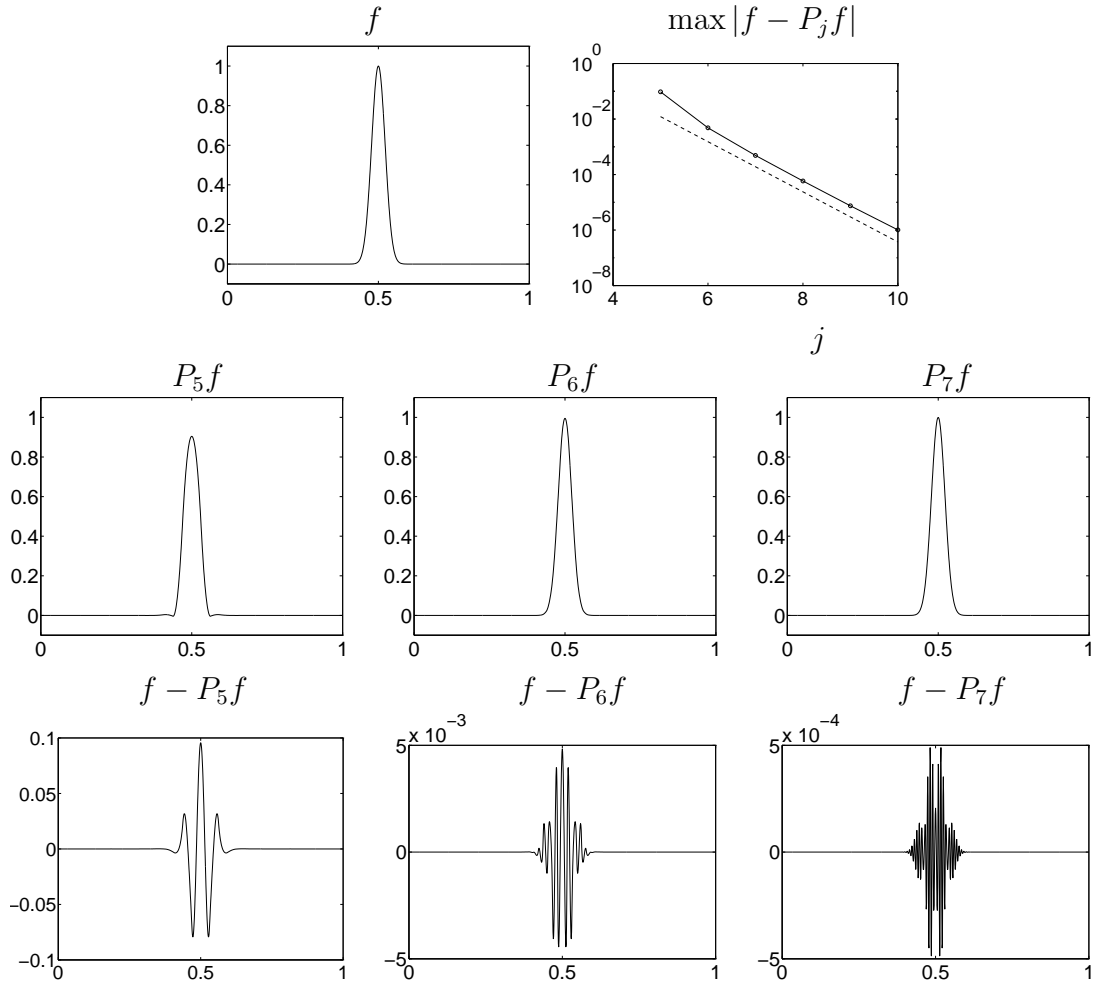
### 4.5.2 Nonlinear approximations

Let us consider again the Jackson inequality under a different perspective. If we are analyzing a function  $f$  that is less regular than the scaling function  $\varphi$ ,



**Figure 4.14:** Approximations of a gaussian pulse (top left panel) at various levels  $j$  in the case  $L = 2$ ,  $\tilde{L} = 4$ . The second row shows the projections onto the scaling function spaces  $P_j f$ , and the third row shows the approximation errors as functions of  $x$  with respect to the exact function  $f$ . The top right panel illustrates the decay of the approximation errors as the level  $j$  increases. The dashed line corresponds to a slope  $-2$ .

then the approximation error decays at a rate controlled by the regularity of the function under investigation. This rate can be unacceptably low, leading to a large number of refinement levels to be included in the representation. If the function  $f$  is poorly regular everywhere in its domain, nothing better can be done. Instead, if the function  $f$  presents few localized singularities and is sufficiently regular elsewhere, the *local* approximation error predicted by the Jackson inequality will have a fast decay in the regions of high regularity, and a slow decay in correspondence of the singular points. This leads to large wavelet coefficients only near the singularities of  $f$ , due to the sharp space localization



**Figure 4.15:** Approximations of a gaussian pulse (top left panel) at various levels  $j$  in the case  $L = 3$ ,  $\tilde{L} = 5$ . The second row shows the projections onto the scaling function spaces  $P_j f$ , and the third row shows the approximation errors as functions of  $x$  with respect to the exact function  $f$ . The top right panel illustrates the decay of the approximation errors as the level  $j$  increases. The dashed line corresponds to a slope  $-3$ .

of each single wavelet. The example of Fig. 2.4 for the Haar system is a simple illustration of this important fact.

Let us fix a coarse level  $J_0$ . The difference between  $f$  and its projection  $P_{J_0} f$  in the left hand side of the Jackson inequality can be expanded into a series of all the detail functions with levels larger than  $J_0$ . More precisely, in the wavelet basis we have

$$f - P_{J_0} f = \sum_{j \geq J_0} \sum_k \hat{f}_{j,k} \psi_{jk}.$$

Taking the norm of this expression, we obtain the approximation error in terms

of the magnitudes of the wavelet coefficients,

$$||f - P_{J_0}f||^2 \asymp \sum_{j \geq J_0} \sum_k |\hat{f}_{j,k}|^2$$

from which we can note that if the representation of  $f$  is “sparse” in the wavelet basis, i.e., few wavelet coefficients are large, only those coefficients will contribute significantly to the approximation error.

This argument naturally leads to the concept of *nonlinear approximation*. We introduce the finest refinement level  $J_{\max}$  at which we want to approximate a given function  $f$ . The approximation  $P_{J_{\max}}f$  can be expanded into scaling functions at the minimum refinement level  $J_0$  and wavelets up to a maximum level  $J_{\max} - 1$ ,

$$P_{J_{\max}}f(x) = \sum_k \check{f}_{J_0,k} \varphi_{J_0,k} + \sum_{j=J_0}^{J_{\max}-1} \sum_k \hat{f}_{j,k} \psi_{j,k}.$$

The nonlinear (adaptive) approximation is based on the thresholding of the wavelet coefficients. In this superposition we keep only those terms corresponding to wavelet coefficients above a given threshold  $\varepsilon$ . We define a threshold-dependent set of wavelet indices

$$\Lambda_\varepsilon = \{(j, k) : |\hat{f}_{j,k}| > \varepsilon\}. \quad (4.94)$$

The set of retained coefficients will be

$$\mathcal{S}_\varepsilon = \{\check{f}_{J_0,k}, \forall k\} \cup \{\hat{f}_{j,k} : (j, k) \in \Lambda_\varepsilon\}. \quad (4.95)$$

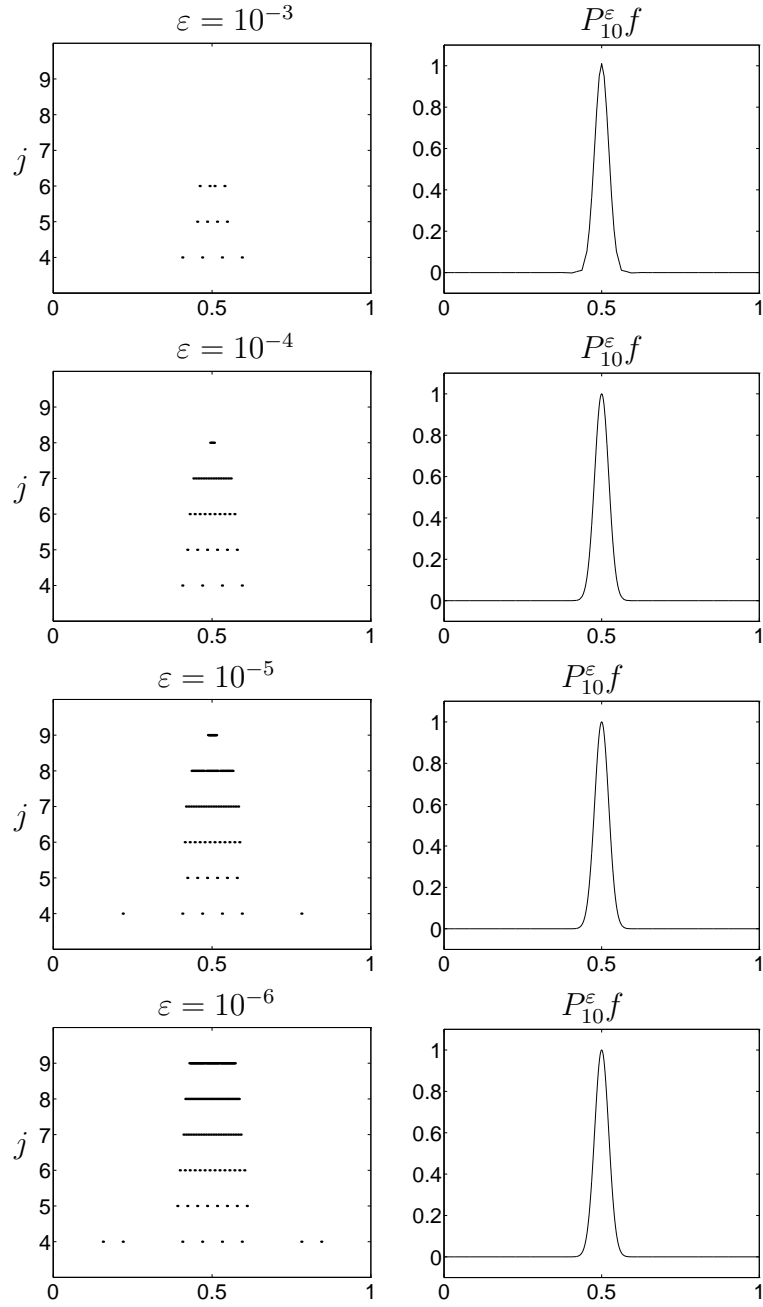
Note that we still include in this set all the scaling function coefficients at the minimum refinement level  $J_0$  for a correct representation of the “low frequency” portion of the function  $f$ . The nonlinear approximation of  $f$  will then be expressed by a projection operator that depends itself on  $f$ ,

$$P_{J_{\max}}^\varepsilon f(x) = \sum_k \check{f}_{J_0,k} \varphi_{J_0,k} + \sum_{\Lambda_\varepsilon} \hat{f}_{j,k} \psi_{j,k}. \quad (4.96)$$

We can also define a *sparsity index*  $S_I^{J_{\max}}(\varepsilon)$ , defined as the percentage of retained coefficients with respect to the total number of coefficients for a given threshold  $\varepsilon$  and maximum refinement level  $J_{\max}$ ,

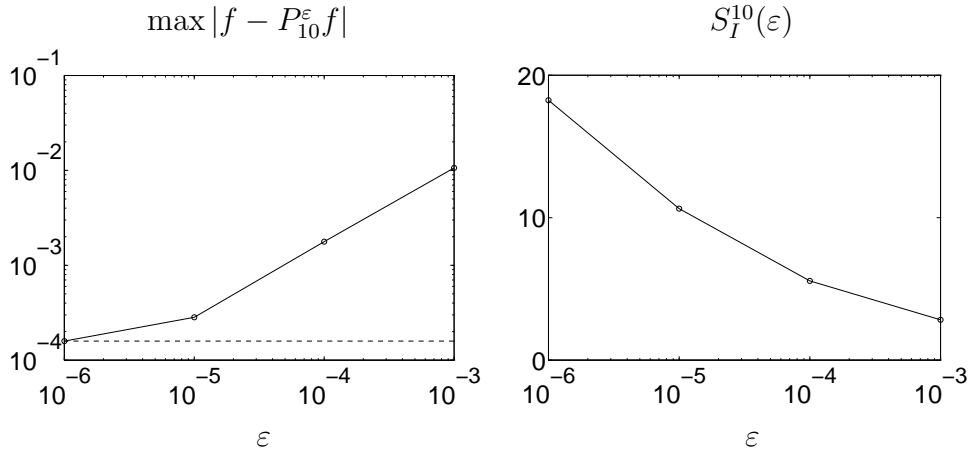
$$S_I^{J_{\max}}(\varepsilon) = 100 \cdot \frac{\text{card} \mathcal{S}_\varepsilon}{\dim V_{J_{\max}}} \quad (4.97)$$

Moreover, as each coefficient  $\hat{f}_{j,k}$  pertains to a wavelet function localized around the dyadic point  $x_{jk} = (2k+1)2^{-j-1}$ , the magnitude of the wavelet coefficients can be used as a “measure” of the regularity of the function  $f$  at a given



**Figure 4.16:** Adaptive approximations ( $L = 2$ ,  $\tilde{L} = 4$ ) of a gaussian pulse (Fig. 4.14, top left panel) obtained through suppression of wavelet coefficients below a given threshold  $\varepsilon$ . The left panels show the location of the wavelet coefficients with magnitude larger than  $\varepsilon$ , separated through refinement levels  $j$ . The right panels show the adaptive approximation obtained with the coefficients depicted in the left panels.

location  $x$  and a given scale  $j$ . Figure 4.16 illustrates these concepts using the gaussian pulse of Eq. 4.93. The left column shows the locations  $x_{jk}$  of

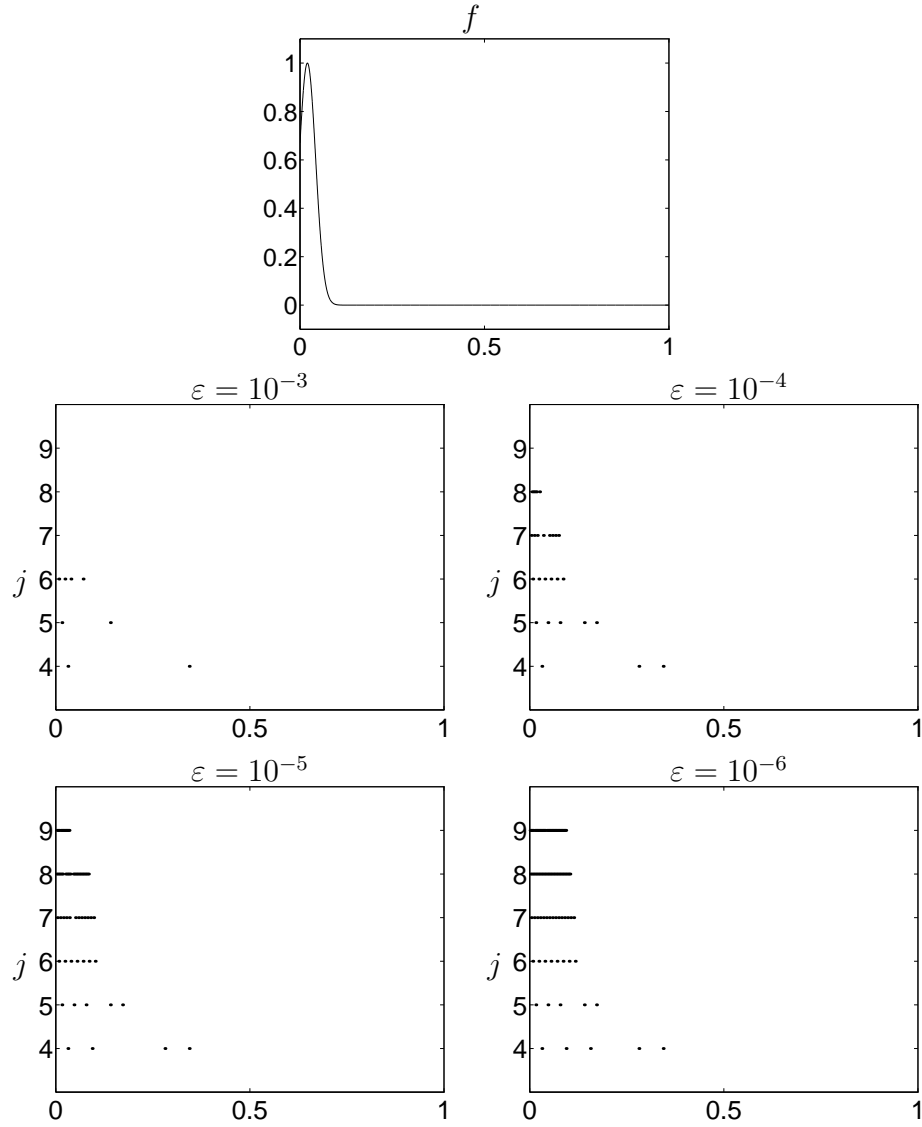


**Figure 4.17:** Adaptive approximations ( $L = 2$ ,  $\tilde{L} = 4$ ) of a Gaussian pulse (Fig. 4.14, top left panel) obtained through suppression of wavelet coefficients below a given threshold  $\varepsilon$ . The left panel shows the nonlinear approximation error (at refinement level  $j = 10$ ) as a function of the threshold  $\varepsilon$ . The dashed line indicates the error in the linear approximation retaining all wavelet coefficients. The right panel shows the sparsity index in the adaptive approximation (i.e. the percentage of retained coefficients) as a function of the threshold.

the wavelet coefficients with magnitude larger than  $\varepsilon$ . For convenience they have been separated through scales  $j$  on the  $y$ -axis. The right panel shows the corresponding approximation  $P_j^\varepsilon$ . Obviously the approximation gets better as the threshold decreases and the number of retained coefficients increases. However, the significant coefficients are concentrated where the “singularity” of  $f$  occurs. Figure 4.17 shows the behavior of the nonlinear approximation error and of the sparsity index as functions of the threshold. It is clear that quite sparse representations can be obtained with a small loss of details in the representation of  $f$ .

Part of this work has been dedicated to the construction of wavelet systems on the unit interval. As the modified border wavelets are not simple translations of a dilated mother wavelet, it should be checked that the space-scale localization is preserved also at the edges of the domain. Figure 4.18 shows that this localization is indeed preserved. The top panel shows the same Gaussian pulse of Fig. 4.16 centered now at  $x_c = 0.02$ . The bottom panels show the locations of wavelet coefficients (always for the B-splines with  $L = 2$ ,  $\tilde{L} = 4$ ), above a fixed threshold  $\varepsilon$ . The approximation errors and the sparsity index, shown in Fig. 4.19, top row, behave in the same way as in the case with  $x_c = 0.5$ . As the function is sufficiently regular everywhere, it is possible to gain accuracy and sparsity in the representation by increasing the regularity of the adopted multiresolution. This fact is illustrated in the bottom row of Fig. 4.19, where the approximation errors and the sparsity index obtained in the B-spline case with

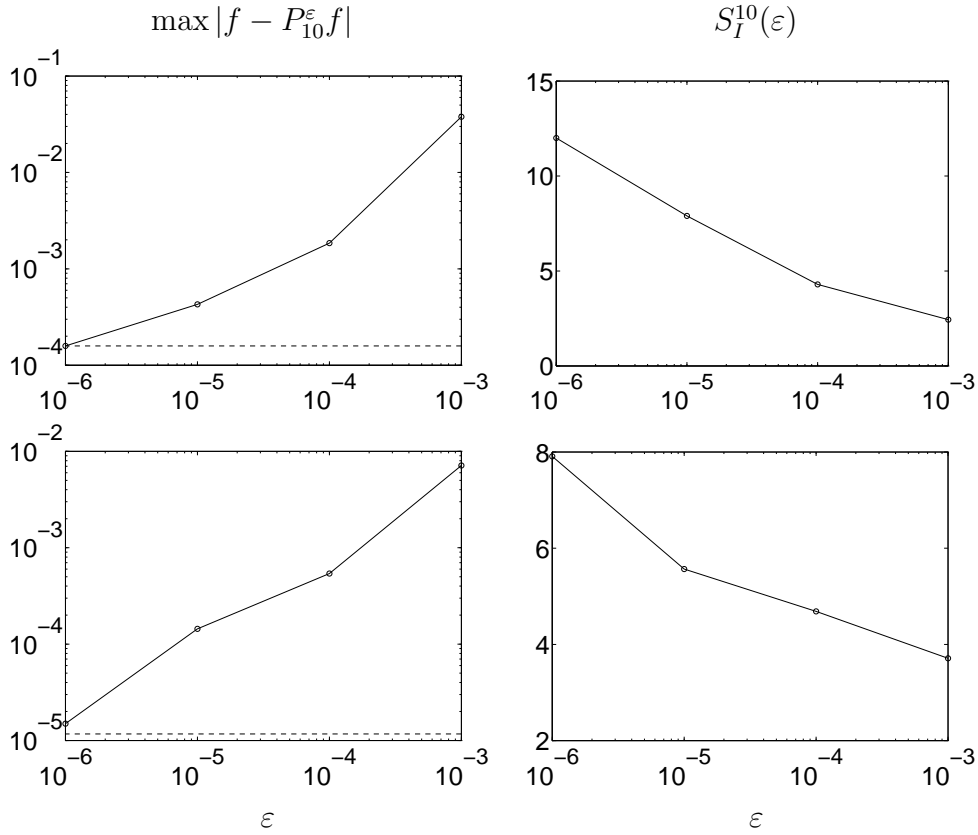




**Figure 4.18:** Gaussian pulse centered near the edge of the interval  $[0, 1]$  (top panel) and locations of its wavelet coefficients above different thresholds  $\varepsilon$ .

$L = 3$  and  $\tilde{L} = 5$  are plotted versus the same thresholds  $\varepsilon$  as in the top row. Both the errors and the sparsity indices are smaller than in the case  $L = 2, \tilde{L} = 4$ . It can be concluded that highly sparse representations can be achieved throughout the domain  $[0, 1]$  using the wavelet systems constructed in this chapter.

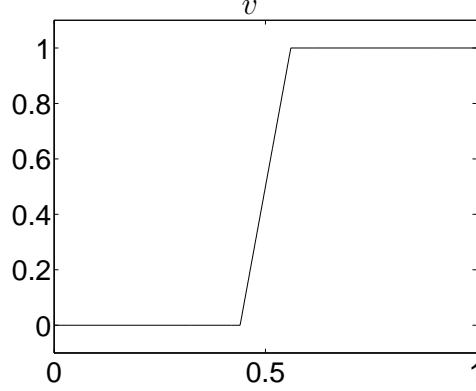
The following example further illustrates the high degree of sparsity in the representation of signals with isolated singularities. Figure 4.20 shows a step function  $v$  with finite rise time. This function has two localized singularities in the first derivative. Signals of this kind are very important in our applications, because the waveform of digital signals can often be modelled with functions like  $v$ . Performing the nonlinear approximations with different thresholds  $\varepsilon$  we



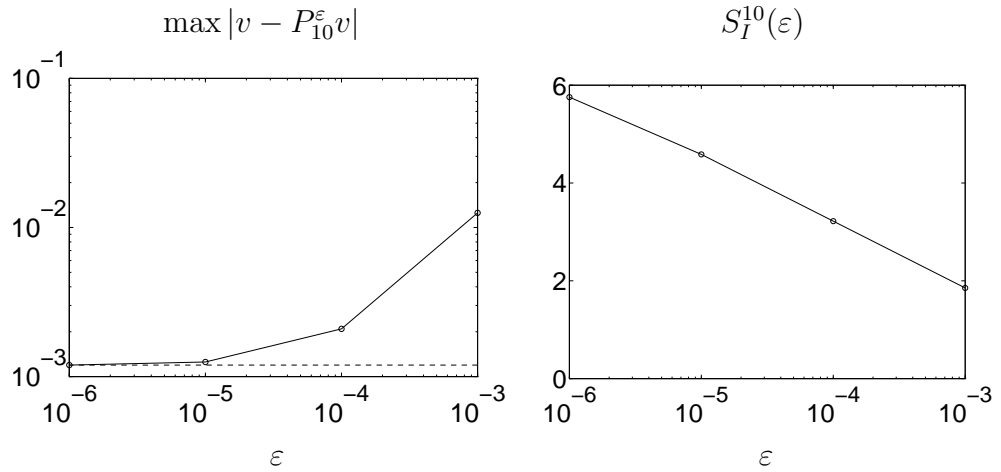
**Figure 4.19:** Nonlinear approximation of the gaussian pulse in Fig. 4.18 in the B-spline case  $L = 2, \tilde{L} = 4$  (top row) and  $L = 3, \tilde{L} = 5$  (bottom row). The left panels show the nonlinear approximation error (at refinement level  $j = 10$ ) as a function of the threshold  $\varepsilon$ . The dashed line indicates the error in the linear approximation retaining all wavelet coefficients. The right panels show the sparsity index in the nonlinear approximation (i.e. the percentage of retained coefficients) as a function of the threshold.

obtained the approximation error depicted in the left panel of Fig. 4.21, and the corresponding sparsity index in the right panel. It should be noted that retaining only less than 6% of the total number of coefficients in the wavelet expansion the approximation error is the smallest possible at the refinement level  $J_{\max}$ . Figure 4.22 shows the locations of the wavelet coefficients above a fixed threshold  $\varepsilon$  in the left panels, and the corresponding nonlinear approximation in the right panels.

The nonlinear approximation can also be regarded under a slightly different perspective. Instead of fixing the threshold  $\varepsilon$ , we can fix the total number  $N$  of coefficients to retain. Consequently, we need to choose *which* are the coefficients to be included. The answer is the set of the largest scaling function and wavelet coefficients. As each coefficient contributes independently to the  $L^2$  norm of

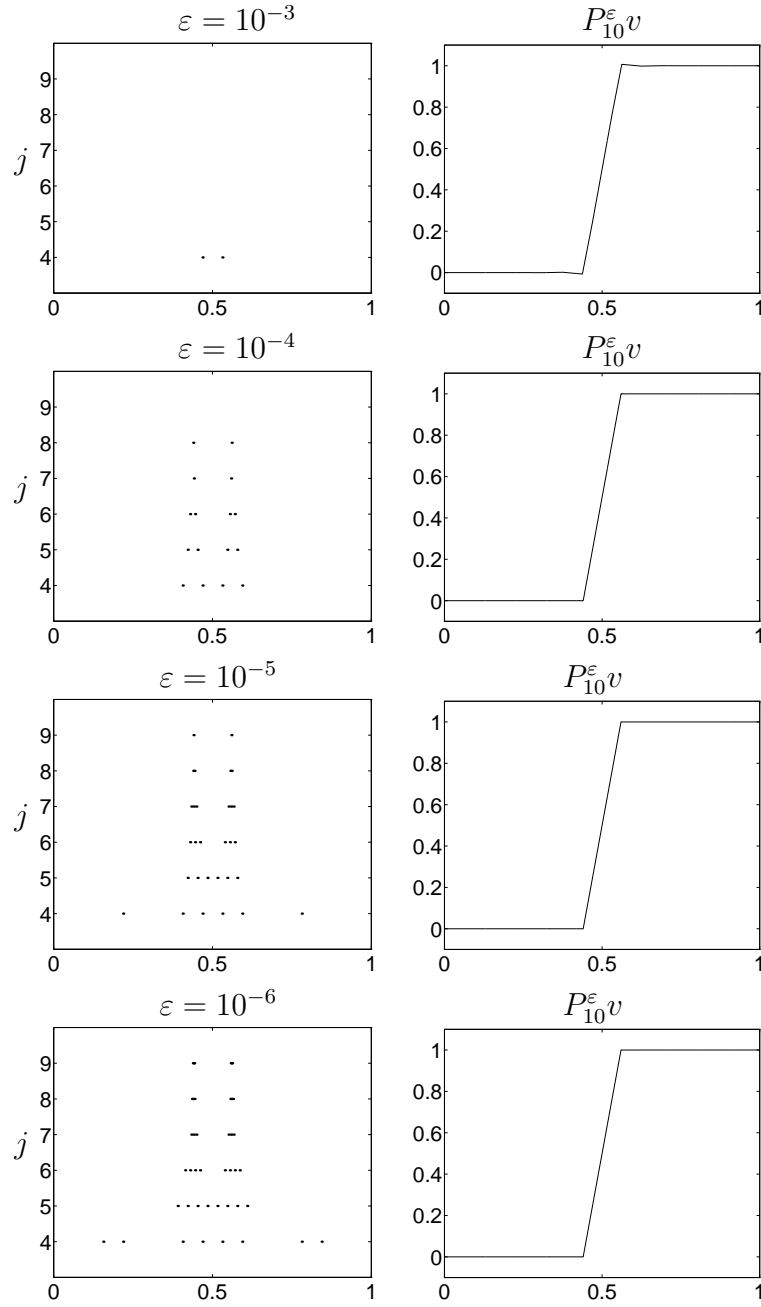


**Figure 4.20:** Step function with finite rise time



**Figure 4.21:** The left panel shows the nonlinear approximation error (at refinement level  $j = 10$ ) as a function of the threshold  $\varepsilon$  for the step function in Fig. 4.20. The dashed line indicates the error in the linear approximation retaining all wavelet coefficients. The right panel shows the sparsity index in the nonlinear approximation (i.e. the percentage of retained coefficients) as a function of the threshold.

the function  $f$ , we are sure that the  $L^2$  norm of the approximation error will be the smallest possible when the smallest coefficients are disregarded.



**Figure 4.22:** Adaptive approximations ( $L = 2$ ,  $\tilde{L} = 4$ ) of the step function in Fig. 4.20, top panel, obtained through suppression of wavelet coefficients below a given threshold  $\varepsilon$ . The left panels show the location of the wavelet coefficients with magnitude larger than  $\varepsilon$ , separated through refinement levels  $j$ . The right panels show the adaptive approximation obtained with the coefficients depicted in the left panels.



# Chapter 5

## Integrals of refinable functions

The formulation of the TDSE method in Chap. 1 derived a linear system of ODEs equivalent to the spatial discretization of the NMTL equations. This linear system (1.28) is fully characterized once the integrals of products of trial and test functions or their derivatives are known. The computation of these integrals in the case of piecewise linear functions used in Chap. 1 was trivial.

This chapter is devoted to the evaluation of integrals of refinable functions. We call *refinable function* any function satisfying a refinement equation. We will see that derivatives or primitives of the scaling functions are refinable functions. The results presented in this chapter will be applied in Chap. 6 to the solution of the NMTL equations by using as trial and test functions the scaling functions and wavelets on the unit interval constructed in Chap. 4.

The following sections will show that the computation of the aforementioned integrals can be performed with the only knowledge of the refinement equations of the functions to be integrated. These integrals are indeed multivariate refinable functions, and their computation can be reduced to an equivalent eigenvector problem.

Section 5.1 will summarize the results of Ref. [21] about the calculation of integrals of refinable functions on unbounded domains. Section 5.2 will then apply the construction of scaling functions and wavelets on the unit interval given in Chapter 4, and determine how to evaluate integrals of functions on bounded domains.

### 5.1 Unbounded domains

Throughout the following we will indicate with  $\varphi^m : \mathbb{R}^q \rightarrow \mathbb{R}^q$  a *refinable function* with associated mask  $\{a_k^m, k \in \mathbb{Z}^q\}$ ,

$$\varphi^m(x) = \sum_{k \in \mathbb{Z}^q} a_k^m \varphi^m(2x - k), \quad \forall x \in \mathbb{R}^q. \quad (5.1)$$

The suffix  $m$  is an index that we will use to distinguish different refinable functions. Even if our application is intrinsically 1D, we will work with a space of dimension  $q$  to keep the derivation as general as possible. Recalling the definition of scaling functions in Chapter 3, we see that the mask is nothing else than the rescaled associated filter,

$$a_k^m = \sqrt{2}h_k^m, \quad \forall k \in \mathbb{Z}^q.$$

We will suppose in the following that the masks have finite support.

The objective will be the evaluation of integrals of the type

$$\int_{\mathbb{R}^q} \varphi^0(2^{j_0}x - \alpha_0) \prod_{m=1}^M q^m(D)\varphi^m(2^{j_m}x - \alpha_m)dx, \quad (5.2)$$

where  $q^m$  is a homogeneous polynomial of degree  $d^m$  and  $D$  is the differentiation operator. We can get rid of the differentiation operator by defining new functions

$$\Upsilon^m(x) = q^m(D)\varphi^m(x), \quad m = 1, \dots, M. \quad (5.3)$$

It is easy to see that the functions  $\Upsilon^m$  also satisfy a refinement equation with associated masks  $\{b_k^m\}$ ,

$$b_k^m = 2^{d^m} a_k^m, \quad \forall k \in \mathbb{Z}^q, m = 1, \dots, M. \quad (5.4)$$

If we set  $\Upsilon^0 = \varphi^0$  and  $b_k^0 = a_k^0 \forall k$ , we are reduced to the evaluation of integrals of the type

$$I(\alpha_0, \dots, \alpha_M; j_0, \dots, j_M) = \int_{\mathbb{R}^q} \prod_{m=0}^M \Upsilon^m(2^{j_m}x - \alpha_m)dx. \quad (5.5)$$

Without loss of generality we will assume that the refinement levels are sorted in ascending order, i.e.

$$j_0 \leq j_1 \leq \dots \leq j_{M-1} \leq j_M.$$

The process of the evaluation of the integrals (5.5) is based on the iterative application of the refinement equation to each of the functions in the product, with repeated changes of variable. At the end of the process we will have to evaluate a set of integrals of type  $I(0, \alpha_1, \dots, \alpha_M; 0, \dots, 0)$ , i.e., only containing translated refinable functions with all the dilation factors equal to one. This reduction of refinement levels is dealt with in Section 5.1.1. Section 5.1.2 will show that the computation of integrals with all the functions at the same refinement level can be reduced to an eigenvector problem.

### 5.1.1 Refinement levels reduction

We begin with the consideration that the change of variable  $x \rightarrow 2^r x$  leads to the identity

$$I(\alpha_0, \dots, \alpha_M; j_0, \dots, j_M) = 2^{qr} I(\alpha_0, \dots, \alpha_M; j_0 + r, \dots, j_M + r).$$

Setting then  $r = -j_0$  in previous expression we obtain

$$I(\alpha_0, \dots, \alpha_M; j_0, \dots, j_M) = 2^{-j_0 q} I(\alpha_0, \dots, \alpha_M; 0, j_1 - j_0, \dots, j_M - j_0). \quad (5.6)$$

Therefore, we can consider directly the case with  $j_0 = 0$ .

Next step is to apply the refinement equation (5.1) to the function with  $m = 1$ . We get

$$\begin{aligned} I(\alpha_0, \dots, \alpha_M; 0, j_1, \dots, j_M) &= \\ &= \int_{\mathbb{R}^q} \left( \sum_{\alpha \in \mathbb{Z}^q} b_\alpha^0 \Upsilon^0(2x - 2\alpha_0 - \alpha) \right) \prod_{m=1}^M \Upsilon^m(2^{j_m} x - \alpha_m) dx \\ &= 2^{-q} \sum_{\alpha \in \mathbb{Z}^q} b_\alpha^0 \int_{\mathbb{R}^q} \Upsilon^0(x - 2\alpha_0 - \alpha) \prod_{m=1}^M \Upsilon^m(2^{j_m-1} x - \alpha_m) dx \\ &= 2^{-q} \sum_{\alpha \in \mathbb{Z}^q} b_\alpha^0 I(\alpha + 2\alpha_0, \alpha_1, \dots, \alpha_M; 0, j_1 - 1, \dots, j_M - 1). \end{aligned}$$

The same procedure can be repeated until the refinement level of the second function reaches 0, obtaining

$$I(\alpha_0, \dots, \alpha_M; 0, j_1, \dots, j_M) = 2^{-j_1 q} \sum_{\gamma_1} \sum_{\gamma_2} \dots \sum_{\gamma_{j_1}} b_{\gamma_1}^0 b_{\gamma_2}^0 \dots b_{\gamma_{j_1}}^0. \quad (5.7)$$

$$I(\gamma_{j_1} + 2\gamma_{j_1-1} + \dots + 2^{j_1-1} \gamma_1 + 2^{j_1} \alpha_0, \alpha_1, \dots, \alpha_M; 0, 0, j_2 - j_1, \dots, j_M - j_1).$$

We have then reduced the integral (5.5) to a sum of integrals with two functions at refinement level 0. Repeating this procedure we can reach the situation where all the functions are at the same refinement level. Let us then consider the general case, when  $p$  functions have already reached this level and the remaining ones are  $M - p + 1$ . We can split the integral as

$$\begin{aligned} I(\alpha_0, \dots, \alpha_M; 0, \dots, 0, j_p, \dots, j_M) &= \\ &= \int_{\mathbb{R}^q} \left( \prod_{m=0}^{p-1} \Upsilon^m(x - \alpha_m) \right) \left( \prod_{m=p}^M \Upsilon^m(2^{j_m} x - \alpha_m) \right) dx, \end{aligned} \quad (5.8)$$

where  $p = 1, \dots, M$ . Indeed, when  $p = 0$  we can apply Eq. (5.6), and when  $p = M + 1$  all the functions are at the level 0, and we can define a new function

$$H(\alpha_1 - \alpha_0, \dots, \alpha_M - \alpha_0) = I(\alpha_0, \dots, \alpha_M; 0, 0, \dots, 0) = \int_{\mathbb{R}^q} \prod_{m=0}^M \Upsilon^m(x - \alpha_m) dx. \quad (5.9)$$



This will be the exit condition from a recursive scheme that we are going to describe in the remaining part of this section. The evaluation of integrals of the type (5.9) will be the subject of Section 5.1.2.

Let us now apply the refinement equation to the first  $p$  functions in the integral (5.8),

$$\begin{aligned}
I(\alpha_0, \dots, \alpha_M; 0, \dots, 0, j_p, \dots, j_M) &= \\
&= \int_{\mathbb{R}^q} \left[ \prod_{m=0}^{p-1} \left( \sum_{n_m} b_{n_m}^m \Upsilon^m(2x - 2\alpha_m - n_m) \right) \right] \left( \prod_{m=p}^M \Upsilon^m(2^{j_m} x - \alpha_m) \right) dx \\
&= 2^{-q} \int_{\mathbb{R}^q} \left[ \prod_{m=0}^{p-1} \left( \sum_{n_m} b_{n_m}^m \Upsilon^m(x - 2\alpha_m - n_m) \right) \right] \left( \prod_{m=p}^M \Upsilon^m(2^{j_m-1} x - \alpha_m) \right) dx \\
&= 2^{-q} \sum_{n_0} \sum_{n_1} \dots \sum_{n_{p-1}} \left( \prod_{m=0}^{p-1} b_{n_m}^m \right) \cdot \\
&\quad \int_{\mathbb{R}^q} \left( \prod_{m=0}^{p-1} \Upsilon^m(x - 2\alpha_m - n_m) \right) \left( \prod_{m=p}^M \Upsilon^m(2^{j_m-1} x - \alpha_m) \right) dx.
\end{aligned}$$

This equation can be interpreted as a recursive formula, that can express the integral (5.8) as a superposition of integrals of the same kind, but with smaller refinement levels for all the functions with  $m \geq p$ . This formula reads

$$\begin{aligned}
I(\alpha_0, \dots, \alpha_M; 0, \dots, 0, j_p, \dots, j_M) &= 2^{-q} \sum_{n_0} \sum_{n_1} \dots \sum_{n_{p-1}} \left( \prod_{m=0}^{p-1} b_{n_m}^m \right) \cdot \quad (5.10) \\
&\quad I(n_0 + 2\alpha_0, \dots, n_{p-1} + 2\alpha_{p-1}, \alpha_p, \dots, \alpha_M; 0, \dots, 0, j_p - 1, \dots, j_M - 1)
\end{aligned}$$

In conclusion, the evaluation of the integral (5.5) in the general case can be reduced to a superposition of integrals of the type (5.9). This superposition is determined by a recursive algorithm whose steps are itemized below.

- Apply Eq. (5.6) to get rid of the first refinement level  $j_0$ ;
- Apply the recurrence relation (5.10) for  $p = 1, \dots, M$ ;
- Evaluate the integrals (5.9) (see next section).

### 5.1.2 The eigenvector equation

This section will show how the integral already introduced in Eq. (5.9),

$$H(x_1, \dots, x_M) = \int_{\mathbb{R}^q} \Upsilon^0(x) \prod_{m=1}^M \Upsilon^m(x - x_m) dx$$

can be evaluated at lattice points  $x_i = \alpha_i \in \mathbb{Z}^q$  without any use of quadrature formulas. The refinement equation of the functions in the product will be

used to determine an eigenvector problem equivalent to the computation of the integral. Note that the number of variables is here  $s = Mq$ .

Let us substitute the refinement equation for each  $\Upsilon^m$  in the integral. After few straightforward passages we get

$$\begin{aligned} H(x_1, \dots, x_M) &= \int_{\mathbb{R}^q} \left( \sum_{\gamma_0} b_{\gamma_0}^0 \Upsilon^0(2x - \gamma_0) \right) \prod_{m=1}^M \left( \sum_{\gamma_m} b_{\gamma_m}^m \Upsilon^m(2x - 2x_m - \gamma_m) \right) dx \\ &= \sum_{\gamma_0 \gamma_1 \dots \gamma_M} \prod_{m=0}^M b_{\gamma_m}^m \int_{\mathbb{R}^q} \Upsilon^0(2x - \gamma_0) \prod_{m=1}^M \Upsilon^m(2x - 2x_m - \gamma_m) dx. \end{aligned}$$

Changing now variable according to  $2x - \gamma_0 = \xi$ , and setting  $\mu_m = \gamma_0 - \gamma_m$ ,  $m = 1, \dots, M$  we get the expression

$$H(x_1, \dots, x_M) = 2^{-q} \sum_{\mu_1 \dots \mu_M} \sum_{\gamma_0} \left\{ b_{\gamma_0}^0 \prod_{m=1}^M b_{\gamma_0 - \mu_m}^m \right\} H(2x_1 - \mu_1, \dots, 2x_M - \mu_M).$$

If we evaluate this expression at lattice points  $\alpha_m \in \mathbb{Z}^q$ , we obtain

$$H(\alpha_1, \dots, \alpha_M) = \sum_{\mu_1 \dots \mu_M} h_{2\alpha_1 - \mu_1, \dots, 2\alpha_M - \mu_M} H(\mu_1, \dots, \mu_M), \quad (5.11)$$

where

$$h_{\mu_1, \dots, \mu_M} = h_{\mu} = 2^{-q} \sum_{\gamma} b_{\gamma}^0 \prod_{m=1}^M b_{\gamma - \mu_m}^m. \quad (5.12)$$

Equation (5.11), which can be rewritten in compact form

$$H(x) = \sum_{\mu \in \mathbb{Z}^{qM}} h_{\mu} H(2x - \mu), \quad x \in \mathbb{Z}^{qM} \quad (5.13)$$

represents an eigenvalue equation of which the integral  $H$  is an eigenvector. Note also that from Eq. (5.13) we immediately see that the integral  $H$  itself is a refinable function with filter  $h_{\mu} \in \mathbb{Z}^{qM}$ .

Let us now recall that the functions  $\Upsilon^m$  are the derivatives of the refinable functions  $\varphi^m$ , with their masks related by Eq. (5.4). We can introduce an integral in the same fashion as  $H$  but with the derivatives removed,

$$F(x_1, \dots, x_M) = \int_{\mathbb{R}^q} \varphi^0(x) \prod_{m=1}^M \varphi^m(x - x_m) dx. \quad (5.14)$$

Following the same procedure above we can easily prove that also  $F$  satisfies a refinement equation

$$F(x) = \sum_{\mu \in \mathbb{Z}^{qM}} c_{\mu} F(2x - \mu),$$

where the filter  $c_\mu$  is defined similarly to  $h_\mu$ ,

$$c_{\mu_1, \dots, \mu_M} = c_\mu = 2^{-q} \sum_{\gamma} a_{\gamma}^0 \prod_{m=1}^M a_{\gamma - \mu_m}^m. \quad (5.15)$$

The two functions  $F$  and  $H$  can be easily related if we differentiate  $F$  through the polynomials  $q^m(D)$  with respect to the arguments  $x^m$ . The result is

$$H(x) = (Q(D)F)(x),$$

where  $Q$  is a homogeneous polynomial on  $\mathbb{R}^{qM}$  expressed by

$$Q(y_1, \dots, y_M) = \prod_{m=1}^M q^m(-y_m).$$

Note that  $\deg[Q]$  is the total number of derivatives in the original integral (5.2). Using now the definition of the masks  $b^m$  in Eq. 5.4 we can see that

$$h_\mu = 2^{\deg[Q]} c_\mu. \quad (5.16)$$

We obtain then the two eigenvector equations

$$F(\alpha) = \sum_{\mu \in \mathbb{Z}^{qM}} c_{2\alpha - \mu} F(\mu), \quad (5.17)$$

$$2^{-\deg[Q]} H(\alpha) = \sum_{\mu \in \mathbb{Z}^{qM}} c_{2\alpha - \mu} H(\mu), \quad (5.18)$$

from which we see that once the matrix  $c_{2\alpha - \mu}$  is computed, both the integrals  $F$  and  $H$  (with any number of derivatives) are expressed by eigenvectors associated to the eigenvalues 1 and  $2^{-\deg[Q]}$ , respectively.

Two problems need still to be addressed. One is the unicity of the eigenvalues. The other is the suitable normalization to be considered for the numerical evaluation of the corresponding eigenvectors. As far as unicity is concerned, it is proved in Ref. [21] that when each  $\varphi^m$  is  $\mu_m$  times continuously differentiable the eigenvector in Eq. (5.18) is unique. The main theorem, which is based on asymptotic expansions of certain Stationary Subdivision Operators [18, 21], also proves that

$$H(\alpha) = (-1)^{|\mu|} W_\alpha,$$

where

$$\begin{aligned} \sum_{\beta \in \mathbb{Z}^{qM}} c_{2\alpha - \beta} W_\beta &= 2^{-|\mu|} W_\alpha, \quad \alpha \in \mathbb{Z}^{qM} \\ \sum_{\alpha \in \mathbb{Z}^{qM}} (-\alpha)^\nu W_\alpha &= \mu! \delta_{\nu\mu}, \quad |\nu| \leq |\mu|, \quad \nu, \mu \in \mathbb{Z}_+^{qM} \end{aligned}$$

This solves the problem of evaluating the integral  $H$ , because it indicates explicitly the correct normalization to be used for the eigenvector. We recall that the construction of the matrix whose eigenvectors must be evaluated is not difficult because its entries are sums of products of the original masks  $a_k^m$ .

### 5.1.3 An example: the scalar case

This section particularizes the results of previous sections to the scalar univariate case. We want to compute integrals of the type

$$I(\alpha_0, \alpha_1; j_0, j_1) = \int_{\mathbb{R}} \Upsilon^0(2_0^j x - \alpha_0) \Upsilon^1(2_1^j x - \alpha_1) dx$$

where the functions in the product are the derivatives of refinable functions

$$\varphi^m(x) = \sum_n a_n^m \varphi^m(2x - n)$$

according to

$$\Upsilon^m(x) = \frac{d^{\ell_m}}{dx^{\ell_m}} \varphi^m(x).$$

Note that the integral can only be evaluated when the order of differentiation of the first function is  $\ell_0 = 0$ . The application of the recursive procedure illustrated in Section 5.1.1 leads to the following three equations,

- $I(\alpha_0, \alpha_1; j_0, j_1) = 2^{-j_0} I(\alpha_0, \alpha_1; 0, j_1 - j_0),$
- $I(\alpha_0, \alpha_1; 0, j_1) = \frac{1}{2} \sum_n a_n^0 I(n + 2\alpha_0, \alpha_1; 0, j_1 - 1),$  for  $j_1 > 0,$
- $I(\alpha_0, \alpha_1; 0, 0) = H(\alpha_1 - \alpha_0).$

The function  $H$  is defined as

$$H(y) = \int_{\mathbb{R}} \Upsilon^0(x) \Upsilon^1(x - y) dx$$

and satisfies the refinement equation

$$2^{-\ell_1} H(y) = \sum_n c_n H(2y - n),$$

where

$$c_n = \frac{1}{2} \sum_{\nu} a_{\nu}^0 a_{\nu-\mu}^1.$$

We show the structure of the system matrix in the case when the mask is  $\{c_n, n = 0, \dots, N + 1\}$ . The refinement equation can also be written, when  $H$  is evaluated at lattice points, as

$$2^{-\ell_1} H(n) = \sum_{m=1}^N c_{2m-n} H(m).$$

From this expression we can visualize the structure of the eigenvector problem as

$$2^{-\ell_1} \begin{bmatrix} H(1) \\ H(2) \\ \vdots \\ H(N) \end{bmatrix} = \begin{bmatrix} c_1 & c_0 & 0 & \cdot & \cdot & 0 \\ c_3 & c_2 & c_1 & c_0 & \cdot & \cdot \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 0 & \cdot & \cdot & 0 & c_{N+1} & c_N \end{bmatrix} \cdot \begin{bmatrix} H(1) \\ H(2) \\ \vdots \\ H(N) \end{bmatrix}. \quad (5.19)$$

The eigenvector will be unique when  $\varphi^1$  is continuously differentiable at least  $\ell^1$  times, and the correct normalization will be given by

$$\sum_{m=1}^N m^\ell H(m) = \ell!$$

#### 5.1.4 Inner products

This section is devoted to the calculation of inner products between an arbitrary function and a refinable function. This case is not included in the discussion of the foregoing sections, because all the functions in the product to be integrated in Eq. (5.2) must be refinable. The problem of the evaluation of inner products stems from the need of finding the expansion coefficients of an arbitrary function into a biorthogonal basis of refinable functions. Given such a biorthogonal system on  $\mathbb{R}$ , we can write, for any function  $f \in V_j$ , where  $V_j$  is the primal scaling function space,

$$f = \sum_k \langle f, \tilde{\varphi}_{jk} \rangle \varphi_{jk}.$$

Conversely, when  $f \in \tilde{V}_j$ , the dual scaling function space, we have

$$f = \sum_k \langle f, \varphi_{jk} \rangle \tilde{\varphi}_{jk}.$$

The expansion coefficients are then evaluated in both cases through an integral of the function  $f$  times a refinable scaling function. In case  $f$  does not belong to the approximation spaces  $V_j$  or  $\tilde{V}_j$ , the evaluation of the coefficients can be performed in the same way. The result, however, will be the set of coefficients of the approximation of  $f$  at level  $j$ ,

$$\begin{aligned} P_j f &= \sum_k \langle f, \tilde{\varphi}_{jk} \rangle \varphi_{jk} \\ \tilde{P}_j f &= \sum_k \langle f, \varphi_{jk} \rangle \tilde{\varphi}_{jk}. \end{aligned}$$

We will consider first the integral

$$I_f = \int_{\mathbb{R}} f(x) \varphi(x) dx, \quad (5.20)$$

where  $f \in \mathcal{C}^0(\mathbb{R})$  is a continuous function on  $\mathbb{R}$  and  $\varphi$  is a refinable function with (finitely supported) mask  $\{a_n, n = n_0, \dots, n_1\}$ . The continuity of  $f$  is essential for the following derivation. Applying the refinement equation to  $\varphi$  in the integral we get

$$I_f = \frac{1}{2} \sum_{m=n_0}^{n_1} a_m \int f\left(\frac{y+m}{2}\right) \varphi(y) dy,$$

where the change of variable  $y = 2x - m$  has been used. If we iterate  $p$  times the refinement equation and the corresponding change of variable, we get

$$I_f = \frac{1}{2^p} \sum_{m_1=n_0}^{n_1} \cdots \sum_{m_p=n_0}^{n_1} a_{m_1} \cdots \alpha_{m_p} \int_{\mathbb{R}} f\left(\frac{x}{2^p} + \frac{m_p}{2^p} + \cdots + \frac{m_1}{2}\right) \varphi(x) dx.$$

When  $p$  gets large, the function in the integral viewed as a function of  $x$  approaches a constant value when evaluated in the support of  $\varphi$ . This can be showed in a rigorous way by using the continuity of  $f$ . As the integral of the scaling function  $\varphi$  is equal to one,  $\varphi$  behaves like a Dirac function centered at some point  $x_c$ . We recall that for the biorthogonal B-spline scaling functions we have  $x_c = \text{rem}(L, 2)$ . In summary, the function

$$\frac{1}{2^p} \sum_{m_1=n_0}^{n_1} \cdots \sum_{m_p=n_0}^{n_1} a_{m_1} \cdots \alpha_{m_p} f\left(\frac{x_c}{2^p} + \frac{m_p}{2^p} + \cdots + \frac{m_1}{2}\right) \quad (5.21)$$

converges uniformly to the constant  $I_f$  when  $p \rightarrow \infty$ . The level of approximation in the evaluation of the integral is controlled by  $p$ .

Let us assume now  $p$  fixed. The evaluation of Eq. (5.21) requires to sample the original function  $f(y)$  at some points on the  $y$  axis. Intuitively, looking at the main integral (5.20), it is clear that  $f$  needs only to be sampled where the refinable function  $\varphi$  is nonzero. The actual minimum and maximum values of  $y$  to be used can be obtained by considering the lower and upper limits in the sum for all indices  $m_i$ ,  $i = 1, \dots, p$ . We obtain

$$y_{min} = \frac{x_c}{2^p} + n_0 \sum_{i=1}^p \frac{1}{2^i}, \quad y_{max} = \frac{x_c}{2^p} + n_1 \sum_{i=1}^p \frac{1}{2^i}.$$

When  $p \rightarrow \infty$  we can sum the geometrical series obtaining

$$y_{min} = n_0, \quad y_{max} = n_1,$$

which are exactly the limits of the support of  $\varphi$ . As the set of points  $y_i$  are equally spaced with a sampling interval of  $2^{-p}$ , the required values of the original function are  $f(y_i)$ , where

$$y_i = n_0 + \frac{i - n_0 + x_c}{2^p}, \quad i = 0, \dots, 2^p(n_1 - n_0) - n_1 + n_0.$$

A straightforward generalization holds when the integral to be evaluated contains the refinable function at level  $j$  and location  $k$ . Indeed, a trivial change of variable leads to

$$I_{f,j,k} = \int_{\mathbb{R}} f(x) \varphi_{jk}(x) dx = \int_{\mathbb{R}} \hat{f}(x) \varphi(x) dx,$$

where

$$\hat{f}(x) = 2^{-j/2} f\left(\frac{x+k}{2^j}\right).$$

Therefore, the integral  $I_{f,jk}$  can be handled in the same way as above. The sequence of points  $y_i$  that are needed for its evaluation at the approximation level  $p$  will be

$$y_i = \frac{1}{2^j} \left( k + n_0 + \frac{i - n_0 + x_c}{2^p} \right), \quad i = 0, \dots, 2^p(n_1 - n_0) - n_1 + n_0,$$

which are strictly included in the support of  $\varphi_{jk}$ . Setting now

$$\{f_i\} = \{2^{-j/2} f(y_i)\},$$

it can be easily shown that the evaluation of the multiple sums in Eq. (5.21) can be achieved through the iterated application of the operator

$$I_{f,jk} \simeq (\downarrow *)^p \{f_i\},$$

where the star indicates the convolution with the filter

$$\left\{ \frac{a_n}{2}, n = n_0, \dots, n_1 \right\}$$

and the arrow indicates a downsampling of a factor of two. In other words, a sliding window of length  $n_1 - n_0 + 1$  is applied to the sequence  $f_i$  skipping one point every two, and the whole procedure is repeated  $p$  times.

## 5.2 Bounded domains

This section is devoted to the evaluation of integrals containing products of derivatives of scaling functions on bounded domains. Only the monodimensional case will be considered. The construction of biorthogonal multiresolutions on the unit interval derived in Section 4.4 will be used here. This will allow to perform the calculations with a simple extension of the evaluation scheme described in Section 5.1 for the integrals on the real line.

We consider the general form of the integrals to be evaluated,

$$I_{[0,1]} = \int_0^1 \varphi^0(2^{j_0}x - \alpha_0) \prod_{m=1}^M q^m(D) \varphi_{j_m, \alpha_m}^m(x) dx, \quad (5.22)$$

where the functions  $\varphi_{j_m, \alpha_m}^m$  are scaling functions belonging to  $V_{j_m}(0, 1)$ . As the support of all these functions is included in the unit interval, the integral can be also restated by substituting the indicator function of  $[0, 1]$  to the first function

$$\varphi^0(x) = \chi_{[0,1]}(x), \quad j_0 = \alpha_0 = 0,$$

obtaining

$$I_{[0,1]} = \int_{\mathbb{R}} \chi_{[0,1]}(x) \prod_{m=1}^M q^m(D) \varphi_{j_m, \alpha_m}^m(x) dx. \quad (5.23)$$

We recall from Eq. (3.77) that the indicator function of the unit interval is the scaling function of the Haar decomposition. Therefore, it is a refinable function with a mask having only two nonzero entries,

$$a_0^0 = a_1^0 = 1, \quad a_k^0 = 0 \text{ otherwise.}$$

Moreover, we can also express each scaling function  $\varphi^m$ ,  $m = 1, \dots, M$  as a superposition of refinable functions on the real line. This was already accomplished in Section 4.4.1. We restate here the main result, derived from Eq. 4.79, with the notations used in this chapter,

$$\varphi_{j_m, \alpha_m}^m = \chi_{[0,1]} \sum_{l_m} \mathcal{M}_{\alpha_m, l_m}^m \varphi_{j_m, l_m}^{m, \mathbb{R}}. \quad (5.24)$$

The limits in the sum are described explicitly in Section 4.4.1. The key point is that each scaling function at level  $j_m$  on the unit interval is expressed as a linear combination of (few) scaling functions of the underlying construction on  $\mathbb{R}$ , at the same refinement level, eventually restricted to the domain  $[0, 1]$ . As the indicator function obviously satisfies

$$\chi_{[0,1]} = \chi_{[0,1]} \cdot \chi_{[0,1]},$$

the expression of the integrals (5.22) becomes

$$I_{[0,1]} = \sum_{l_1} \cdots \sum_{l_m} \left\{ \prod_{m=1}^M \mathcal{M}_{\alpha_m, l_m}^m \right\} \int_{\mathbb{R}} \chi_{[0,1]}(x) \prod_{m=1}^M q^m(D) \varphi_{j_m, l_m}^{m, \mathbb{R}} dx. \quad (5.25)$$

Each of the integrals in the sum is now applied to the product of derivatives of refinable functions on  $\mathbb{R}$ , and can be evaluated through the scheme given in Section 5.1.

### 5.2.1 Inner products

We consider here the evaluation of the inner product of an arbitrary continuous function  $f \in \mathcal{C}^0([0, 1])$  with (biorthogonal) scaling functions on the unit interval. The aim is then the computation of

$$f_{jk} = \int_0^1 f(x) \tilde{\varphi}_{jk}(x) dx.$$

These can also be interpreted as the expansion coefficients of  $P_j f$  into the basis of the primal scaling function space  $V_j(0, 1)$ ,

$$P_j f(x) = \sum_k f_{jk} \varphi_{jk}(x),$$

with  $j \geq j_0$ , defined in Section 4.4.1. The expansion coefficients into the dual basis can be obtained by simply exchanging the role of primal and dual scaling functions.



As the translation invariance is lost on the unit interval (and on bounded domains in general), the procedure developed in Section 5.1.4 cannot be applied in the present case. Let us write then the function  $f$  as

$$f(x) = g(x) + r(x),$$

where  $g(0) = g(1) = 0$  and

$$r(x) = f(0) + [f(1) - f(0)]x.$$

Due to the linearity of this decomposition, we can obtain the expansion coefficients  $f_{jk}$  as the superposition

$$f_{jk} = g_{jk} + r_{jk},$$

where

$$P_j g(x) = \sum_k g_{jk} \varphi_{jk}(x), \quad r(x) = \sum_k r_{jk} \varphi_{jk}(x).$$

Note that  $P_j r = r$  because the scaling function spaces that we use in this work include the polynomials of degree at least one.

Let us consider the  $g_{jk}$  first. We recall that the dual scaling functions on the unit interval can be expressed as linear combinations of the dual scaling functions of the underlying biorthogonal multiresolution on  $\mathbb{R}$ . Therefore we can write (see Section 4.4.1)

$$g_{jk} = \int_0^1 g(x) \tilde{\varphi}_{jk}(x) dx = \int_{\mathbb{R}} \bar{g}(x) \tilde{\varphi}_{jk}(x) dx = \sum_l \tilde{\mathcal{M}}_{k,l} \int_{\mathbb{R}} \bar{g}(x) \tilde{\varphi}_{jl}^{\mathbb{R}}(x) dx,$$

where

$$\bar{g}(x) = \begin{cases} g(x), & x \in [0, 1] \\ 0 & \text{otherwise.} \end{cases}$$

As the function  $\bar{g}$  is continuous on  $\mathbb{R}$ , the results obtained in Section 5.1.4 can be immediately applied for the evaluation of  $g_{jk}$ .

We turn now to the computation of  $r_{jk}$ . These are readily obtained once the expansion coefficients of the monomials  $p_\alpha(x) = x^\alpha$  with  $\alpha \in \{0, 1\}$  are known, i.e.

$$p_\alpha(x) = \sum_k p_{\alpha k} \varphi_{jk}(x).$$

Referring to the scaling function spaces constructed in Section 4.4.1, at a fixed level  $j$  three cases can be distinguished,

1.  $0 \leq k < \tilde{L}$ ,
2.  $\tilde{L} \leq k < \dim V_j - \tilde{L} - 1$ ,
3.  $\dim V_j - \tilde{L} - 1 \leq k < \dim V_j$

The first case refers to the left edge, the second to the internal scaling functions, and the third to the right edge. The third case can be immediately obtained from the first by reflection and symmetry, and will not be detailed here.

Let us consider the internal functions first (case 2 in the list above). As the internal scaling functions are the translated biorthogonal scaling functions on the real line,  $\tilde{\varphi}_{jk} = \tilde{\varphi}_{j,k^*+k-\tilde{L}}^R$ , we have

$$\begin{aligned} p_{\alpha k} &= \langle p_{\alpha}, \tilde{\varphi}_{j,k^*+k-\tilde{L}}^R \rangle = \int_{\mathbb{R}} x^{\alpha} 2^{j/2} \tilde{\varphi}(2^j x - [k^* + k - \tilde{L}]) dx = \\ &= 2^{-j(\alpha+1/2)} \int_{\mathbb{R}} y^{\alpha} \tilde{\varphi}(y - [k^* + k - \tilde{L}]) dy = 2^{-j(\alpha+1/2)} c_{\alpha, k^*+k-\tilde{L}}. \end{aligned}$$

The coefficients  $c_{\alpha l}$  are defined in Eq. (4.3).

We consider now the left edge, i.e., case 1 in the list above. As the construction of the scaling function spaces is based on the separation of the left and right edge, we start with the study of the expansion of the monomials on the half line  $[0, \infty)$  at level  $j = 0$ . From the definition of the primal border functions in Eq. (4.4) we can write for  $x \geq 0$

$$\begin{aligned} p_{\alpha}(x) &= \theta_{\alpha}(x) + \sum_{k=k_0^*}^{k^*-1} c_{\alpha k} \varphi_{0k}^R(x) + \sum_{k \geq k^*} c_{\alpha k} \varphi_{0k}^R(x) = \\ &= \sum_{l=0}^{\tilde{L}-1} q_{\alpha l} \theta_{0l}(x) + \sum_{k \geq k^*} c_{\alpha k} \varphi_{0k}^R(x), \end{aligned}$$

where

$$q_{\alpha l} = \begin{cases} \delta_{\alpha l}, & 0 \leq l < L, \\ c_{\alpha, l-L+k_0^*}, & L \leq l < \tilde{L}, \end{cases}$$

and where  $\varphi_{0l}$  represent the primal scaling functions on the half line before the biorthogonalization. These are expressed in terms of the biorthogonal scaling functions through the inverse of the change of basis matrix  $D$ ,

$$\theta_{0l}(x) = \sum_{s=0}^{\tilde{L}-1} [D^{-1}]_{ls} \varphi_{0s}(x).$$

Therefore, we obtain

$$p_{\alpha k} = \sum_{s=0}^{\tilde{L}-1} p_{\alpha k}^{(0)} \varphi_{0s}(x) + \sum_{k \geq k^*} c_{\alpha k} \varphi_{0k}^R(x),$$

where

$$p_{\alpha k}^{(0)} = \sum_{l=0}^{\tilde{L}-1} q_{\alpha l} [D^{-1}]_{lk}.$$

Finally, the expansion coefficients at level  $j$  are found by applying the operator  $T_j$  to the equation

$$p_\alpha(x) = \sum_k p_{\alpha k}^{(0)} \varphi_{0k}(x),$$

obtaining

$$p_{\alpha k} = 2^{-j(\alpha+1/2)} p_{\alpha k}^{(0)} = 2^{-j(\alpha+1/2)} \sum_{l=0}^{\tilde{L}-1} q_{\alpha l} [D^{-1}]_{lk}.$$

Putting all the cases together, we get the final expression for the expansion coefficients  $r_{jk}$ ,

$$r_{jk} = f(0) p_{0k} + [f(1) - f(0)] p_{1k}.$$

### 5.2.2 Expansion of discontinuous functions

The evaluation of the expansion coefficients of an arbitrary continuous function defined on the unit interval into the scaling function bases has been solved in the foregoing section. This section focuses on the evaluation of the same expansion coefficients when the function presents some isolated discontinuities. It should be noted that the continuity is essential for the application of the method described in Section 5.2.1.

Any discontinuous function  $f$  can be expressed as a superposition of a continuous part  $f_c \in \mathcal{C}^0$  plus a discontinuous part  $\mu$ , which collects all the jumps of  $f$ ,

$$\mu(x) = \sum_{n=1}^{N_d} \alpha_n u(x - x_n),$$

where the total number of jumps is  $N_d$ , their height is expressed by the coefficients  $\alpha_n$ , and  $u(\cdot)$  is the Heaviside step function. Due to the linearity of the projection operator  $P_j$ , the computation of the expansion coefficients of  $f$  is solved once the expansion coefficients of the singular part  $\mu$  are known. These coefficients are expressed as inner products with the dual scaling functions,

$$\mu_{jk} = \int_0^1 \mu(x) \tilde{\varphi}_{jk}(x) dx, \quad \forall k = 0, \dots, \dim V_j(0, 1) - 1. \quad (5.26)$$

The same procedure applies with obvious substitutions for the computation of the expansion coefficients into the dual scaling function basis, therefore only the evaluation of the primal coefficients will be detailed here. By substituting the definition of the singular part  $\mu$  in Eq. (5.26), we get

$$\mu_{jk} = \sum_{n=1}^{N_d} \alpha_n \int_0^1 u(x - x_n) \tilde{\varphi}_{jk}(x) dx.$$

If we integrate each term in the sum by parts, we obtain

$$\mu_{jk} = \sum_{n=1}^{N_d} \{ \tilde{\mathcal{N}}_{jk}(1) - \tilde{\mathcal{N}}_{jk}(x_n) \},$$

where  $\widetilde{\mathcal{N}}_{jk}$  indicates the primitive of the scaling function  $\widetilde{\varphi}_{jk}$ ,

$$\widetilde{\mathcal{N}}_{jk}(x) = \int_0^x \widetilde{\varphi}_{jk}(s) ds.$$

We recall that the scaling functions on the unit interval are a linear combination of the scaling functions on  $\mathbb{R}$  restricted to the unit interval through some coefficients that we collect in a matrix  $\widetilde{\mathcal{M}}$  (see Eqs. (4.79) and (5.24)). Using this property we can also express the primitives of the scaling functions on the unit interval as a linear combination of the primitives of the scaling functions on  $\mathbb{R}$ . With few straightforward passages we get the expression

$$\mu_{jk} = \sum_{n=1}^{N_d} \alpha_n \sum_l \widetilde{\mathcal{M}}_{kl} \left[ \widetilde{\Phi}_{jl}^{\mathbb{R}}(1) - \widetilde{\Phi}_{jl}^{\mathbb{R}}(x_n) \right], \quad (5.27)$$

where  $\widetilde{\Phi}_{jl}^{\mathbb{R}}$  represents the primitive of the dual scaling function on  $\mathbb{R}$ ,

$$\widetilde{\Phi}_{jl}^{\mathbb{R}}(x) = \int_{-\infty}^x \widetilde{\varphi}_{jl}^{\mathbb{R}}(s) ds. \quad (5.28)$$

Combining Eqs. (5.28) and (5.27) we obtain the final expression for the expansion coefficients,

$$\mu_{jk} = 2^{-j/2} \sum_{n=1}^{N_d} \alpha_n \sum_l \widetilde{\mathcal{M}}_{kl} \left[ \widetilde{\Phi}(2^j - l) - \widetilde{\Phi}(2^j x_n - l) \right],$$

where  $\widetilde{\Phi}$  is the primitive of the dual scaling function on  $\mathbb{R}$  at the zeroth refinement level,

$$\widetilde{\Phi}(x) = \int_{-\infty}^x \widetilde{\varphi}(s) ds. \quad (5.29)$$

The last step is the evaluation of the primitive of the scaling function on  $\mathbb{R}$  at an arbitrary point. We assume as usual that the scaling function is compactly supported in  $[\widetilde{n}_0, \widetilde{n}_1]$ . This implies that

$$\begin{aligned} \widetilde{\Phi}(x) &= 0 & \forall x \leq \widetilde{n}_0, \\ \widetilde{\Phi}(x) &= 1 & \forall x \geq \widetilde{n}_1. \end{aligned}$$

It is also easy to prove that the primitive satisfies the refinement equation

$$\widetilde{\Phi}(x) = \sum_{n=\widetilde{n}_0}^{\widetilde{n}_1} \frac{\widetilde{a}_n}{2} \widetilde{\Phi}(2x - n), \quad (5.30)$$

where  $\widetilde{a}_n$  are the mask coefficients of the scaling function on  $\mathbb{R}$ .

Let us evaluate Eq. (5.30) at an integer  $k \in \{\tilde{n}_0 + 1, \dots, \tilde{n}_1 - 1\}$ . This is the set of integer points at which the primitive assumes nontrivial values. A simple change of variable leads to

$$\tilde{\Phi}(k) = \sum_{m=m_i}^{2k-\tilde{n}_0} \frac{1}{2} \tilde{a}_{2k-m} \tilde{\Phi}(m) = \sum_{m=m_i}^{m_f} \frac{1}{2} \tilde{a}_{2k-m} \tilde{\Phi}(m) + \sum_{m=m_f+1}^{2k-\tilde{n}_0} \frac{1}{2} \tilde{a}_{2k-m}, \quad (5.31)$$

where

$$m_i = \max\{2k - \tilde{n}_1, \tilde{n}_0 + 1\}, \quad m_f = \min\{2k - \tilde{n}_0, \tilde{n}_1 - 1\}.$$

The above equation can be restated in matrix form as

$$\left[ I - \frac{1}{2}A \right] \underline{\tilde{\Phi}} = \underline{\alpha}, \quad (5.32)$$

where the matrix  $A$  has the same structure as in Eq. (5.19),  $\underline{\tilde{\Phi}}$  is the array of the primitive values at nontrivial integers, and  $\underline{\alpha}$  collects the second sum in the right-hand side of Eq. (5.31). The system (5.32) is invertible because the maximum eigenvalue of the matrix  $A$  is unitary. Once the primitive values at integers are known, the value at any dyadic point  $x_{jk} = k2^{-j}$  can be computed by applying recursively the refinement equation (5.30).

# Chapter 6

## TDSE with scaling functions and wavelets

In this chapter we will apply the biorthogonal scaling function and wavelet bases on the unit interval constructed in Chapter 4 to the transient solution of Nonuniform Multiconductor Transmission Lines through the TDSE method described in Chapter 1.

Two different approaches will be followed, according to the nature of the voltage source waveforms exciting the line. If these waveforms are regular, we will show in Section 6.1 that very high accuracies can be obtained by increasing the regularity and the polynomial order of the trial functions. This holds because also the voltage and current along the line are highly regular, and can be represented efficiently when the approximation spaces match their regularity. On the other hand, when there are singularities in the forcing waveforms, increasing the regularity of the approximation spaces has no effects on the maximum approximation error. A large number of trial functions is needed to represent the solution with a good accuracy. However, we will show in Section 6.2 that when the singularities are isolated, the nonlinear approximation based on wavelet thresholding leads to adapted representations of the solution. The number of trial functions giving a significant contribution to the representation of the solution are concentrated near its singularities. Consequently, the overall representation of the solution results highly sparse.

### 6.1 TDSE with scaling functions

The TDSE method was introduced in Chapter 1 without reference to any specific choice of trial and test functions. It is easy to show that the formulation of the TDSE method of Sec. 1.1 can be used without modifications if the trial and test functions are the biorthogonal scaling functions on the unit interval constructed in Chapter 4. Indeed, the only requirements on the basis functions are expressed in Eq. (1.19) and (1.20). These conditions insure that only one basis

function is non-vanishing at  $z = 0$ , and similarly at  $z = 1$ . In addition, the supports of these two border functions must be disjoint. Recalling the properties of the scaling function systems on the unit interval constructed in Chapter 4, we see that these two conditions are automatically satisfied.

We recall here for convenience the implicit system of ODE's of Eq. (1.36), stemming from the aforementioned formulation of the TDSE method,

$$\widehat{\Psi} \frac{d}{dt} \hat{\mathbf{x}}(t) + \widehat{\Phi} \hat{\mathbf{x}}(t) = \widehat{\Delta}_S \mathbf{V}_S(t) + \widehat{\Delta}_{SD} \frac{d}{dt} \mathbf{V}_S(t) + \widehat{\Delta}_L \mathbf{V}_L(t) + \widehat{\Delta}_{LD} \frac{d}{dt} \mathbf{V}_L(t). \quad (6.1)$$

Each unknown in the array  $\hat{\mathbf{x}}(t)$  represents one expansion coefficient of voltage or current at time  $t$  into the trial functions system (see Eq. (1.35)). The approximation space for the solution is the scaling function space  $V_j(0, 1)$  at a refinement level  $j \geq j_0$ , where  $j_0$  is the minimum allowed level expressed by Eq. (4.75). Following the notations of Chapter 1 we set

$$\begin{aligned} \zeta_n &= \varphi_{j,n-1}, & n = 1, \dots, N_\zeta = \dim V_j \\ \phi_n &= \varphi_{j,n-1}, & n = 1, \dots, N_\phi = \dim V_j \\ \eta_n &= \varphi_{j,n-1}, & n = 1, \dots, N_\eta = \dim V_j. \end{aligned}$$

Note that all the three basis sets coincide with the primal scaling function systems, while the dual scaling functions are not used. Indeed, the primal scaling functions offer the best trade-off between regularity and length of the support. On one hand, the dual scaling functions are always less regular than the primal scaling functions, at least in the biorthogonal B-spline systems with practical values of  $L$  and  $\tilde{L}$ . On the other hand, they are characterized by longer supports and longer filters. Therefore, the primal scaling functions offer a better representation of the solution when employed as trial functions and, when employed as test functions, produce a larger number of vanishing entries, or equivalently a smaller bandwidth, in the system matrices of Eq. (6.1) with respect to the dual scaling functions.

The system matrices  $\widehat{\Psi}$  and  $\widehat{\Phi}$  are computed through the following operations.

- Expansion of the per-unit-length matrices  $\mathbf{L}(z)$ ,  $\mathbf{C}(z)$ ,  $\mathbf{R}(z)$ , and  $\mathbf{G}(z)$  into the basis functions  $\phi_n$  (see Eq. (1.11)). Due to the biorthogonality of the scaling function systems the expansion coefficients can be computed through inner products with the dual scaling functions. Section 5.2.1 shows how these inner products can be computed in a fast and accurate way.
- Computation of inner products of the basis functions. In particular, the two sets of inner products

$$A_{mn} = \left\langle \frac{d}{dz} \zeta_n, \eta_m \right\rangle$$

(see Eq. (1.16)) and

$$B_{mn}^{(k)} = \langle \zeta_n \phi_k, \eta_m \rangle$$

(see Eq. (1.18)) are needed. These are easily expressed as integrals of refinable functions, as shown in Section 5.2. The computation of these inner products can then be performed without any use of quadrature formulas through the algorithms developed in Chapter 5.

Figure 6.1 shows examples of the system matrices obtained by varying the values of  $L$  and  $\tilde{L}$  in the construction of the scaling function systems (i.e. the accuracy of the discretization scheme). The refinement level is set to  $j = 6$  for all the pairs  $(L, \tilde{L})$  in the figure. Note that the structure of these matrices is unchanged from the simple piecewise linear approximation case detailed in Chapter 1. In particular, both the matrices  $\hat{\Psi}$  and  $\hat{\Phi}$  have a banded structure. This allows to use fast factorization algorithms for banded matrices and leads to the computation of the quantities  $\hat{\Psi}^{-1}\hat{\Phi}\hat{\mathbf{x}}$  in  $O(N_\zeta)$  operations. The bandwidth of these matrices increases with increasing regularity (i.e., with the parameter  $L$ ). The modifications of the matrices near the edges are due to the special construction of the border scaling functions developed in Chapter 4.

The solution of the system of ODE's in Eq. (6.1) can be obtained through a suitable integration scheme in time. As in the case of piecewise linear approximations, we will use a 5<sup>th</sup> – 6<sup>th</sup> order Runge-Kutta scheme [10]. Once the expansion coefficients  $\hat{\mathbf{x}}(t)$  are computed by the time-stepping routine, the voltage and current solutions on each conductor can be computed at any  $z$  location along the line and any time  $t$  in a postprocessing stage, through the superpositions in Eq. (1.9) and (1.10). Also this postprocessing involves  $O(N_\zeta)$  arithmetic operations.

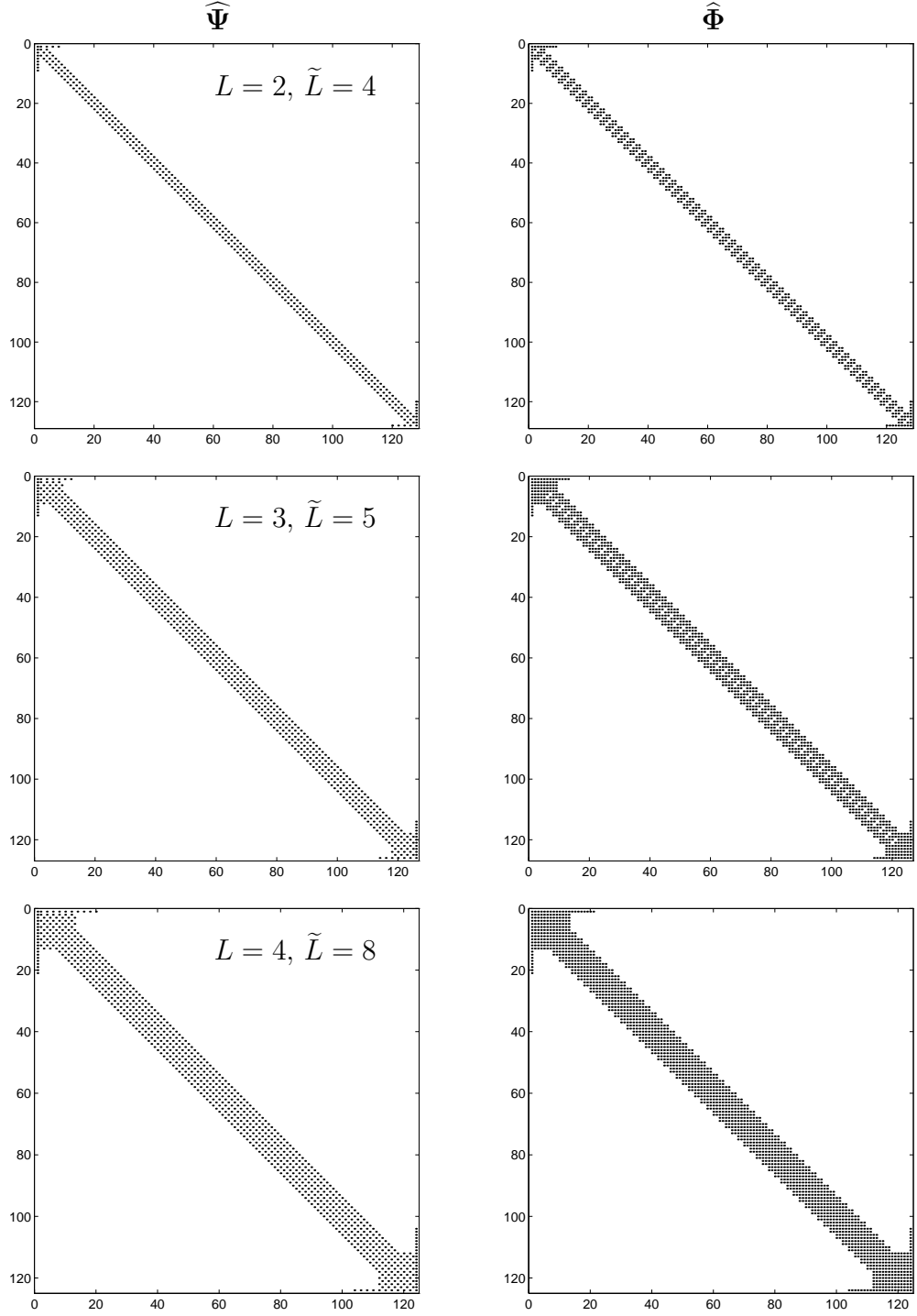
### 6.1.1 The exponential line

The improved approximation properties of the TDSE method with scaling functions are illustrated on the following example. The matched 1:4 scalar exponential line described and solved with piecewise linear trial and test functions in Sec. 1.3.1 is solved here with scaling function bases of various order. In particular, we will use the primal scaling functions corresponding to the pairs

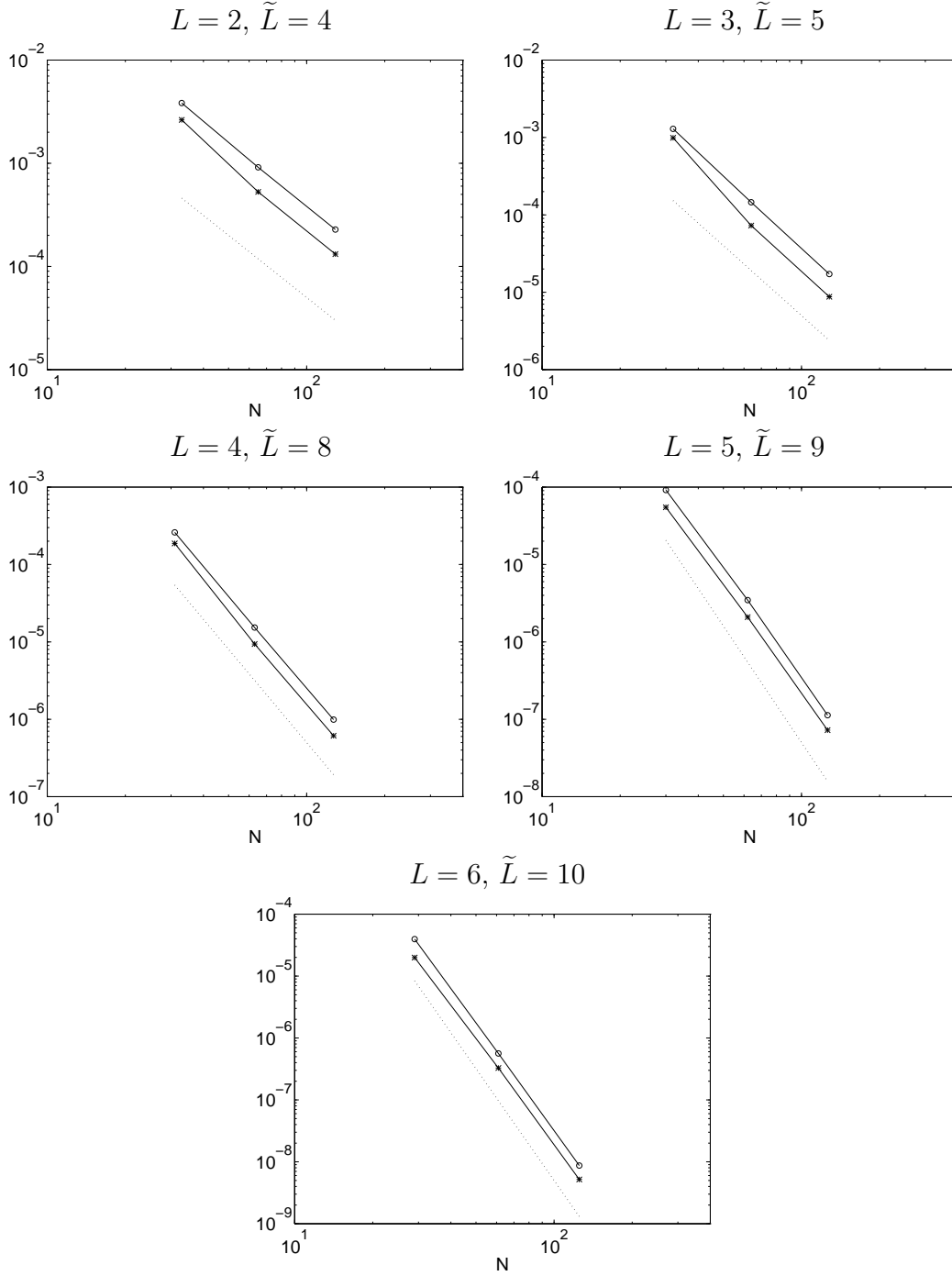
$$(L, \tilde{L}) \in \{(2, 4), (3, 5), (4, 8), (5, 9), (6, 10)\}$$

and will investigate the behavior of the approximation error of the solution obtained with the TDSE method with respect to the reference solution obtained by inverse FFT of the analytical frequency domain solution. Special care has been taken to insure that the periodization effects due to the FFT introduce insignificant deviations in the reference solution with respect to the exact solution. The approximation errors on voltage ( $E_v$ ) and current ( $E_i$ ) are defined in Eq. (1.38).

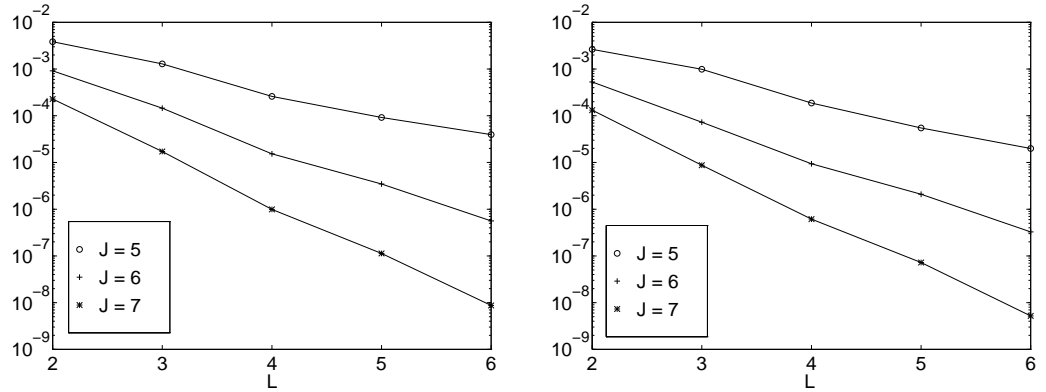




**Figure 6.1:** Structure of the system matrices of Eq. (6.1) using primal scaling functions from the biorthogonal B-spline systems on the unit interval with order  $(L, \tilde{L})$  as trial and test functions. Each dot represents a nonzero  $P \times P$  block.



**Figure 6.2:** Maximum absolute error on voltage (dots) and current (stars) for the matched 1:4 exponential line excited by a unitary gaussian pulse. Each panel corresponds to a different biorthogonal B-spline scaling function system of order  $(L, \tilde{L})$ . The slope of the dotted lines in the each panel corresponds to  $N^{-L}$ , where  $N$  is the total number of basis functions.



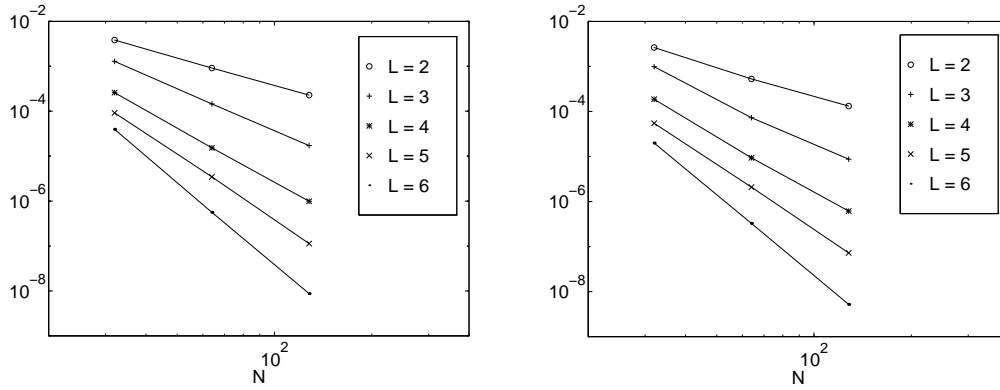
**Figure 6.3:** Matched 1:4 exponential line with gaussian excitation. Maximum absolute error on voltage (left) and current (right) as a function of the refinement level  $J$  and approximation order  $L$  of the basis functions. The total number of basis functions is approximately  $N = 2^J$ .

The source voltage generator is the same gaussian pulse as in Eq. (1.37). This waveform is  $\mathcal{C}^\infty$ . Therefore, the approximation error predicted by the Jackson inequality is only due to the choice of the approximation spaces and not to the functions to be approximated. Figure 6.2 shows the behavior of the voltage and current approximation errors as functions of the total number of trial functions  $N = \dim V_j$ . Each panel reports the errors for one of the pairs  $(L, \tilde{L})$  listed above. This results show that the decay of the approximation error follows a power law which is consistent with the polynomial order of the trial functions. More precisely,

$$E_v(N) \sim N^{-L}, \quad E_i(N) \sim N^{-L}.$$

Figures 6.3 and 6.4 depict the same errors on voltage and current as functions of the polynomial order  $L$  of the approximating spaces and the refinement levels  $j$  (we recall that  $N$  is approximately equal to  $2^j$ ). These plots show that a good approximation can be achieved either by using few functions with a high regularity or many functions with less regularity. The best results have of course been obtained with the largest number ( $N = 125$ ) of trial functions with the highest regularity (5<sup>th</sup> degree polynomials).

It should be noted that a high accuracy can be obtained by using high order scaling functions only when the source waveform is highly regular. When  $v_s(t)$  presents some localized singularities, like in the examples of Sec. 4.5.2, the decay of the approximation errors is dominated by the regularity of the solution, making the choice of higher order approximating functions useless. In this case, better results can be obtained through a nonlinear approximation by using hierarchical bases instead of canonical bases. This is dealt with in the forthcoming section.



**Figure 6.4:** Matched 1:4 exponential line with gaussian excitation. Maximum absolute error on voltage (left) and current (right) as a function of the approximation order  $L$  of the  $N$  basis functions.

## 6.2 TDSE with wavelets

This section shows that the nonlinear approximations with thresholding of the wavelet coefficients described in Section 4.5.2 can be employed to get adapted representations of the solution to the NMTL equations. This is particularly useful when the waveforms of the voltage sources at the line terminations present isolated singularities. For example, this is the case of trapezoidal pulse trains, which present discontinuities in their first derivative.

The key point is the characterization of the approximation spaces for the solution to be obtained with the TDSE method. Given a maximum refinement level  $J_{\max}$ , we consider only those functions in the space  $V_{J_{\max}}$  with significant wavelet coefficients. This requires to introduce a threshold  $\varepsilon$  that defines whether a wavelet coefficient is to be retained or not. This threshold, as shown in Sec. 4.5.2, determines also the overall accuracy at which the approximation can be obtained. Once the threshold is fixed, the approximation of the “true” solution is obtained by applying the nonlinear projection operator  $P_{J_{\max}}^{\varepsilon}$  introduced in Eq. (4.96). This operator selects the wavelet coefficients belonging to the set  $\mathcal{S}_{\varepsilon}$ , defined in Eq. (4.95), and disregards the others. The nonlinear approximation of the solution can then be described through a sparse representation, with a sparsity index  $S_I^{J_{\max}}(\varepsilon)$ , defined in Eq. (4.97), indicating the percentage of retained coefficients.

The procedure described in the above paragraphs leads naturally to operate directly with the wavelet bases instead of the scaling function bases. We recall from Sec. 4.4.3 that the wavelet analysis and synthesis processes allow to switch between scaling function and wavelet bases through iterative application of the filters  $\mathcal{H}$  and  $\mathcal{G}$ . These operations can be performed in  $O(N_{\zeta})$  operations. If we apply this basis change to the system (6.1), we obtain a new system where the unknowns are the expansion coefficients of voltage and current into the

scaling functions at level  $J_0$  (see Eq. (4.83)) and wavelets at the increasing levels  $j = J_0, \dots, J_{\max} - 1$ . Unfortunately, this new formulation destroys the banded structure of the system matrices. We illustrate this fact on a simple example, by showing how the structure of a linear operator acting in the canonical (scaling function) basis varies when the representation is changed into a hierarchical (wavelet) basis.

Let us consider an operator  $T : V_j \rightarrow V_j$  mapping a function  $v$  to a function  $w$ ,

$$w = Tv.$$

If we denote with  $\check{v}$ ,  $\check{w}$  the column arrays with scaling function coefficients of the functions  $v$ ,  $w$ , respectively, we have that the operator can be represented with a  $\dim V_j \times \dim V_j$  matrix  $\mathbf{T} = T_{mn}$ ,

$$\check{w} = \mathbf{T}\check{v}.$$

We indicate now with  $\hat{v}$  and  $\hat{w}$  the set of coefficients of  $v$  and  $w$  in the wavelet basis, according to the notations used in Sec. 4.4.3. We have the following formal identities,

$$\hat{v} = \widetilde{\mathcal{W}} \check{v}, \quad \check{v} = \mathcal{W}^T \hat{v},$$

where the operator  $\widetilde{\mathcal{W}}$  performs the wavelet analysis and the operator  $\mathcal{W}^T$  the synthesis. These two operators are defined in Eq. (4.92). The operator  $T$  in the wavelet basis will be represented as

$$\hat{w} = \widehat{\mathbf{T}} \hat{v},$$

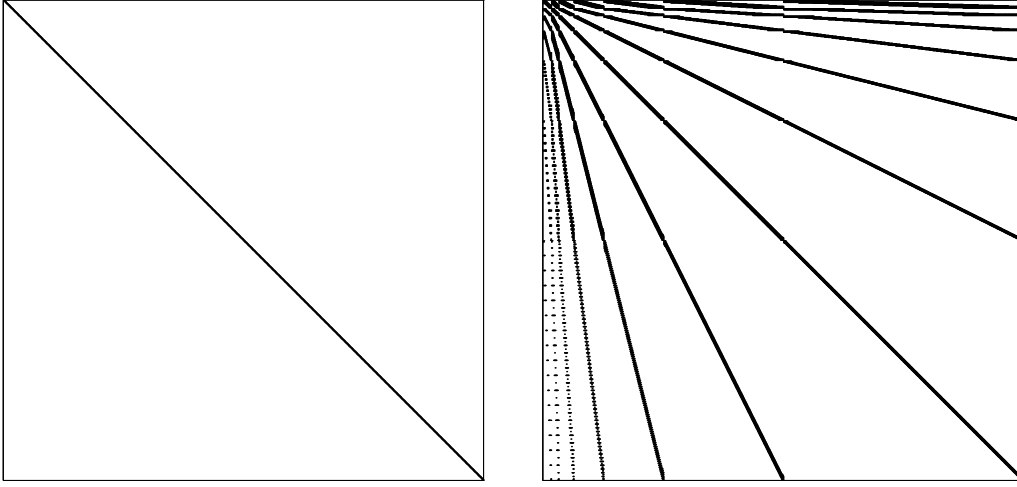
where

$$\widehat{\mathbf{T}} = \widetilde{\mathcal{W}} \mathbf{T} \mathcal{W}^T.$$

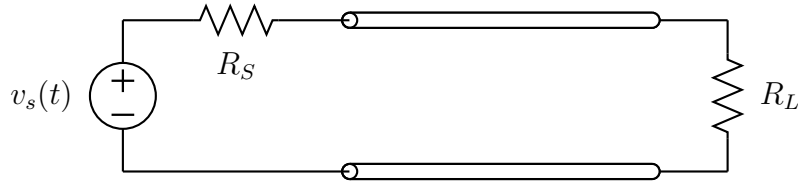
Let us consider a simple tridiagonal matrix  $T_{mn}$  with all the nonvanishing entries equal to 1. The structure of this operator is depicted in the left panel of Fig. 6.5. The right panel shows the structure of this operator in the wavelet basis. The operator is sparse in both cases, but the structure of its representation in the wavelet basis is much more complex.

The foregoing example showed that the structure of operators in wavelet bases is sparse and highly structured. However, this structure does not allow an efficient sparse factorization and, consequently, a fast inversion. Similarly, the structure of the system matrix  $\widehat{\Psi}$  when using a wavelet representation does not allow a sparse factorization and inversion as in the case of the scaling function representation. This results in a much larger computational effort and makes this representation highly inefficient for the solution of the NMTL equations with the TDSE method.

There are two possible alternatives to deal with this problem. The first is to derive an alternative formulation of the TDSE method, which is explicit in the time derivatives of the unknowns. This formulation will be the subject of



**Figure 6.5:** Operator  $T$  in the scaling function representation (left) and wavelet representation (right).



**Figure 6.6:** Scalar lossless uniform line

Sec. 6.3, where a detailed derivation is shown. In the remaining part of this section we will focus instead on a different approach, which is based on performing all the computations in the scaling function basis. Therefore, the system (6.1) is actually solved. However, at each time iteration (or every  $n_{\text{iter}}$  iterations) the solution is processed through wavelet analysis, thresholding, and synthesis. This process does not affect the overall computational effort, which remains  $O(N_\zeta)$ , but the structure of the solution is enforced to be sparse. Indeed, even if the time-stepping routine operates in the scaling function basis, the solution is non-linearly filtered at each time step and is therefore adaptively represented. We will illustrate the inclusion of nonlinear wavelet filtering in the TDSE method through simple examples, like the scalar lossless uniform line (Sec. 6.2.1) and the scalar exponential line (Sec. 6.2.2).

### 6.2.1 The uniform line

Let us consider the scalar lossless uniform line of Fig. 6.6. The (normalized) per-unit-length parameters are  $L = 1$  H and  $C = 1$  F. The line length is  $\mathcal{L} = 1$  m, consequently the one-way delay time is  $T = 1$  s. The line terminations are

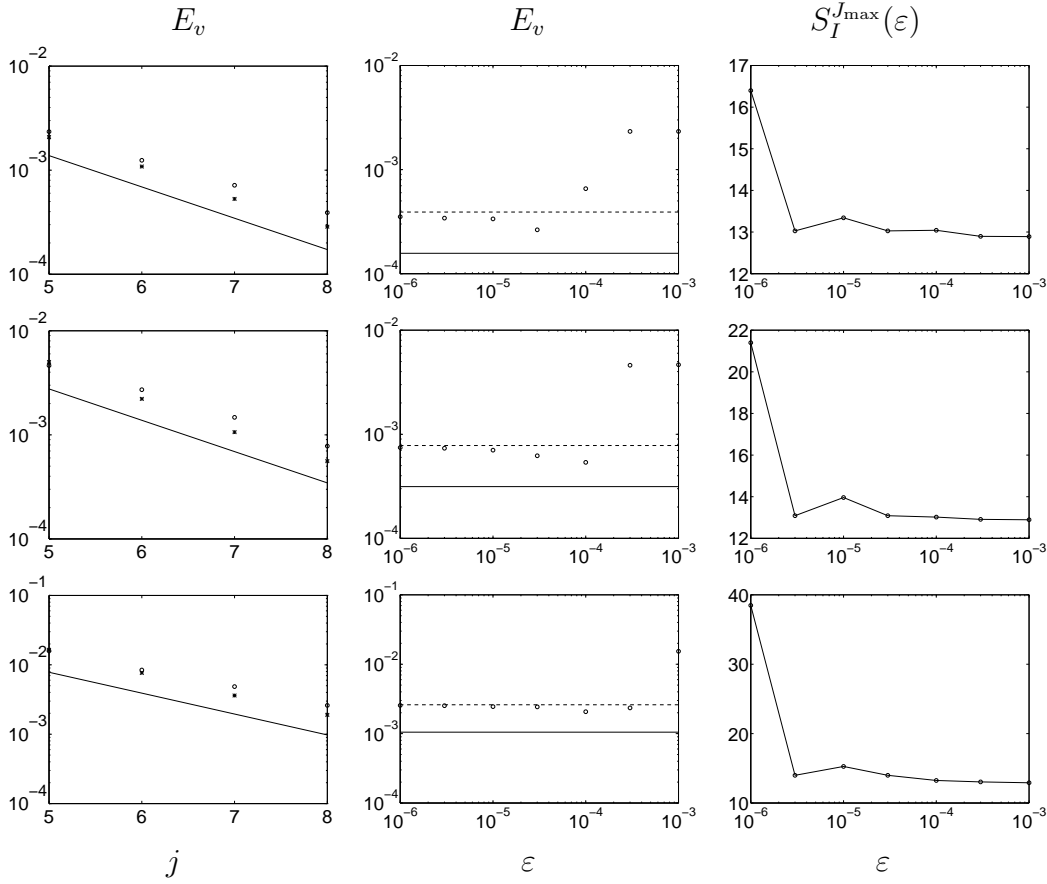
matched to the characteristic impedance of the line,  $R_s = 1 \, \Omega$ ,  $R_L = 1 \, \Omega$ . The voltage source is a 1 V trapezoidal pulse train with period  $T_0 = 8$  s, 50% duty cycle, and rise time  $t_r$  equal to the fall time  $t_f$ . We will examine three different values of the rise time, namely  $t_r = 0.3T$ ,  $t_r = T$ , and  $t_r = 2T$ . For each of these three cases, we will study the approximation error as a function of the wavelet threshold  $\varepsilon$  used in the nonlinear approximation.

We consider first the solution with no wavelet thresholding. This corresponds to a linear approximation of the solution in the sense of Sec. 4.5.1. The left panels of Fig. 6.7 show the behavior of the voltage approximation error  $E_v$  as a function of the refinement level  $J$  used in the simulations for two different pairs  $(L, \tilde{L}) = (2, 4)$  and  $(3, 5)$ . These figures show that the decay of the errors has a slope  $-1$  in both cases. This slope is dictated by the regularity of the solution and not by the regularity of the trial functions. Therefore, the use of higher order trial functions is useless for this type of forcing waveforms. For this reason, we will use the pair  $(L, \tilde{L}) = (2, 4)$  throughout the following.

Next, we consider the dependence of the approximation error on the threshold  $\varepsilon$ . As described in the foregoing section, at each time step we perform wavelet analysis to compute the wavelet coefficients, then we suppress all the wavelet coefficients with magnitude smaller than  $\varepsilon$ , and finally we reconstruct the solution through wavelet synthesis. This allows to compute both the overall sparsity index of the representation and the spatial location of the significant wavelet coefficients as a function of time. The middle panels of Fig. 6.7 show the approximation error  $E_v$  as a function of the threshold in the case  $J_{\max} = 8$ ,  $j_0 = 4$ . We can see that the errors are practically unchanged from the errors obtained by retaining all the wavelet coefficients (dashed lines). The only deviation can be noticed in the first two rows for large values of the threshold. In these cases the nonlinear approximation error becomes dominant with respect to the discretization error at the finest level. We can also notice that the minimum allowable error for the representation of the analytical solution (continuous line) is approximately one-half of the error obtained with the TDSE method.

The main advantage of the nonlinear approximation process is evident from the right panels of Fig. 6.7, where the sparsity index of the adapted representation is plotted versus the threshold  $\varepsilon$ . We see that, with the exception of very small values of the threshold, for which many wavelet coefficients become non-negligible, quite high degrees of sparsity can be achieved. Less than 15% of the total number of coefficients needed to linearly represent the solution are used in the nonlinear representation. This can be obtained at no loss of accuracy.

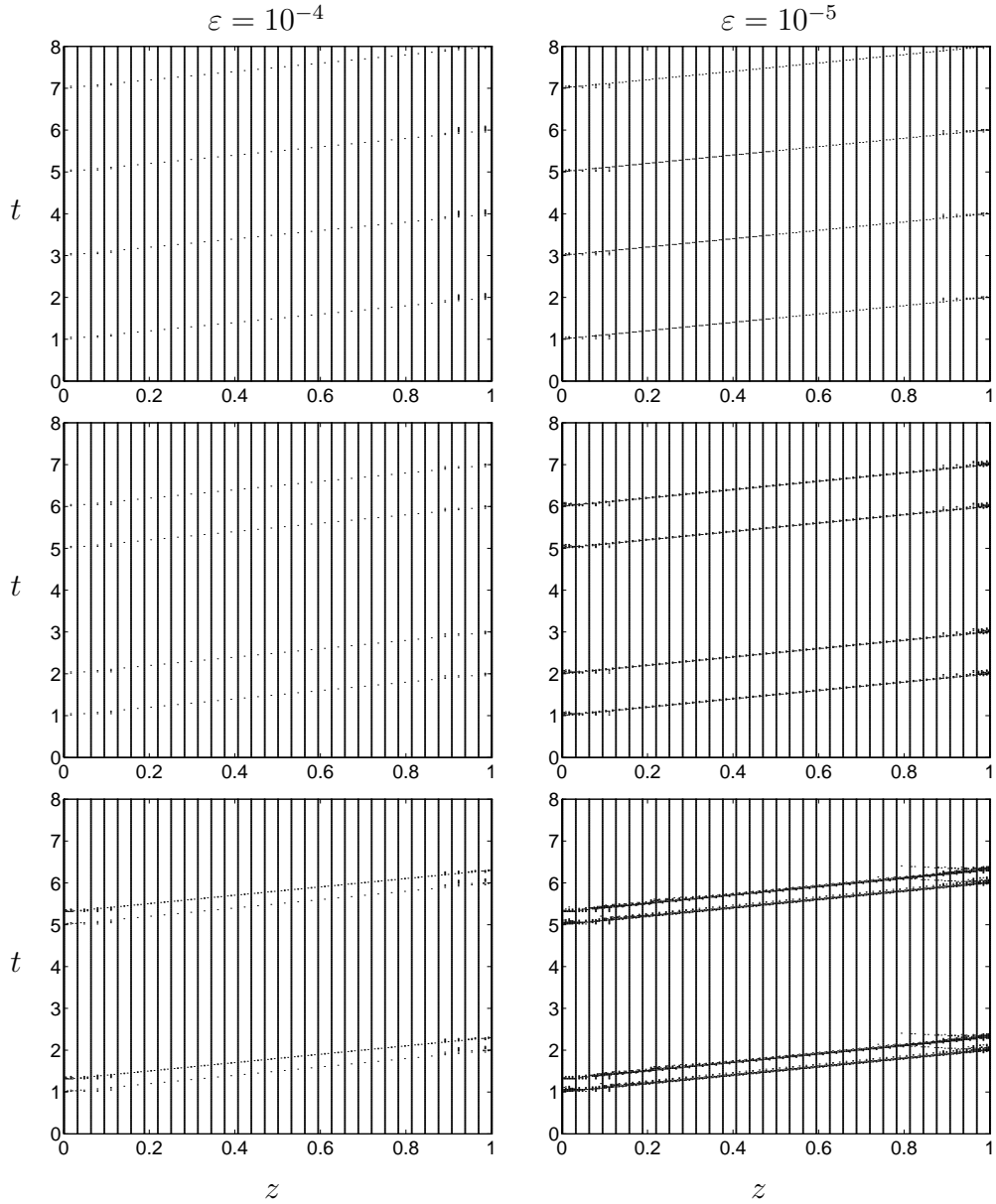
Figure 6.8 shows the location of the significant scaling function and wavelet coefficients in the  $(z, t)$  plane for two different values of the threshold  $\varepsilon$ . From these plots we notice that the wavelet representation is adapted to the travelling singularities along the line. In fact, the wavelet coefficients follow precisely the characteristics of the hyperbolic transmission line equation, which are lines in



**Figure 6.7:** Uniform matched lossless scalar line with a trapezoidal pulse train voltage source. The one-way delay time is  $T = 1$  s. The rise time is  $t_r = 2T$  in the first row,  $t_r = T$  in the second row, and  $t_r = 0.3T$  in the third row. The left panels show the linear approximation errors in the voltage obtained with the TDSE method with respect to the exact solution for different refinement levels  $j$  in the case  $L = 2$ ,  $\tilde{L} = 4$  (circles) and  $L = 3$ ,  $\tilde{L} = 5$  (stars). The slope of the continuous line is  $-1$ . The middle panels show the nonlinear voltage approximation errors obtained with different thresholds for the wavelet coefficients, and using  $J_{\max} = 8$ ,  $L = 2$ ,  $\tilde{L} = 4$ . The dashed line is the linear approximation error for the solution without wavelet thresholding, and the continuous line is the maximum linear approximation error of the exact solution. The right panels show the sparsity index in the representation of the solution as a function of the threshold  $\varepsilon$ .

the  $(t, z)$  plane with slope equal to the inverse of the propagation speed. It is evident from the plots that the number of wavelet coefficients increase when the threshold is decreased, thus leading to a larger sparsity index. Also, the number of wavelet coefficients increases when the rise time is decreased at a fixed threshold  $\varepsilon$ . This is due to the fact that the “strength” of the singularities, which can be formally defined as the jump in the first derivative, is inversely

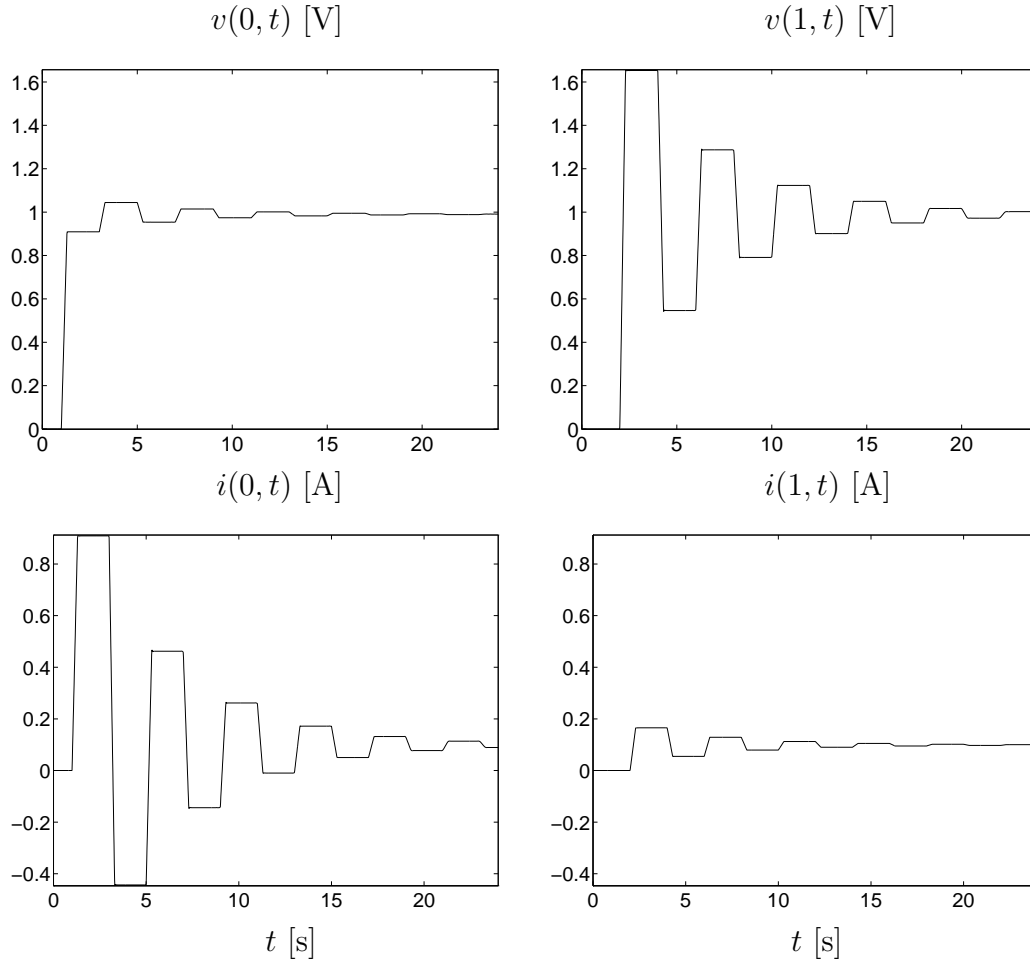




**Figure 6.8:** Uniform matched lossless scalar line with a trapezoidal pulse train voltage source. The one-way delay time is  $T = 1s$ . The rise time is  $t_r = 2T$  in the first row,  $t_r = T$  in the second row, and  $t_r = 0.3T$  in the third row. The locations of the significant wavelet coefficients in the nonlinear approximation of the voltage solution (with  $J_{\max} = 8$ ,  $L = 2$ ,  $\tilde{L} = 4$ ) are plotted in the  $(z, t)$  plane. The threshold for wavelet coefficients is  $\varepsilon = 10^{-4}$  in the left panels and  $\varepsilon = 10^{-5}$  in the right panels.

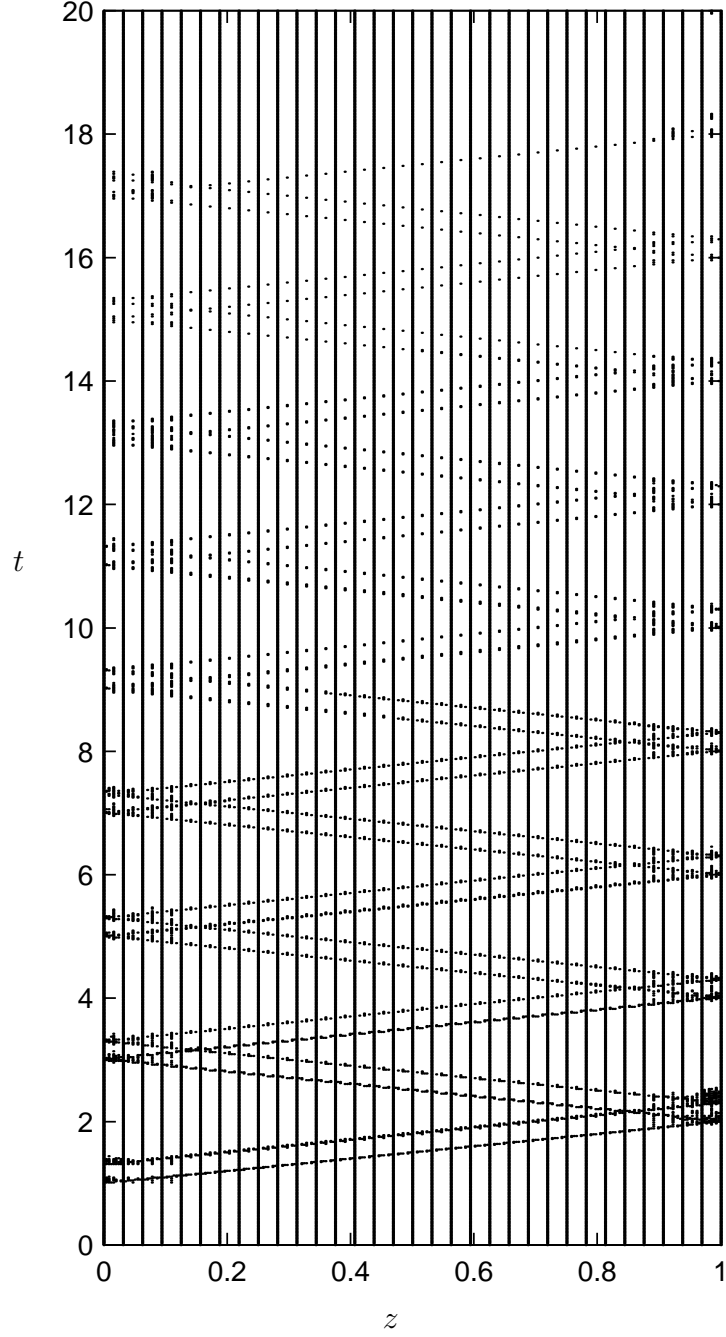
proportional to the rise time because the amplitude of the pulses is fixed.

The adaptivity of this method is further illustrated on another example. The same transmission line as above is terminated with unmatched loads, namely

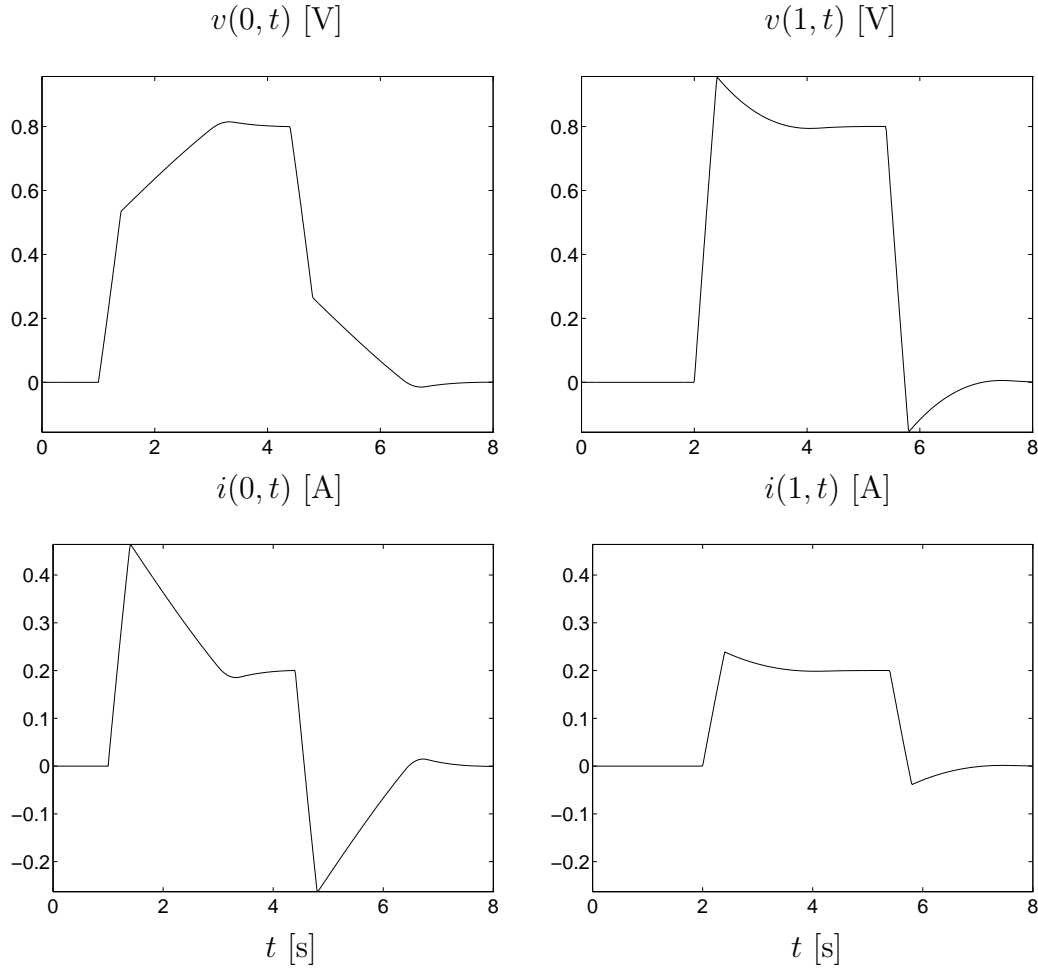


**Figure 6.9:** TDSE solution for the unmatched lossless scalar uniform line with  $J_{\max} = 8$ ,  $L = 2$ ,  $\tilde{L} = 4$ .

$R_S = 0.1$  (small driver impedance) and  $R_L = 10$  (high receiver impedance). The source waveform is a 1 V step with rise time  $t_r = 0.3T$ . With these load conditions the voltage and current at the line ends show a typical oscillating behavior due to multiple reflections of the input pulse. This is illustrated by Fig. 6.9, depicting voltage and current at the loads obtained with the TDSE method with nonlinear thresholding of the wavelet coefficients using  $\varepsilon = 10^{-5}$ . Figure 6.10 shows the location of the retained wavelet coefficients, which follow the bouncing singularities along the characteristic lines. As the travelling reflected wave reduces its amplitude at each reflection, the strength of the singularities dims down, leading to a smaller number of wavelet coefficients needed to represent the solution.



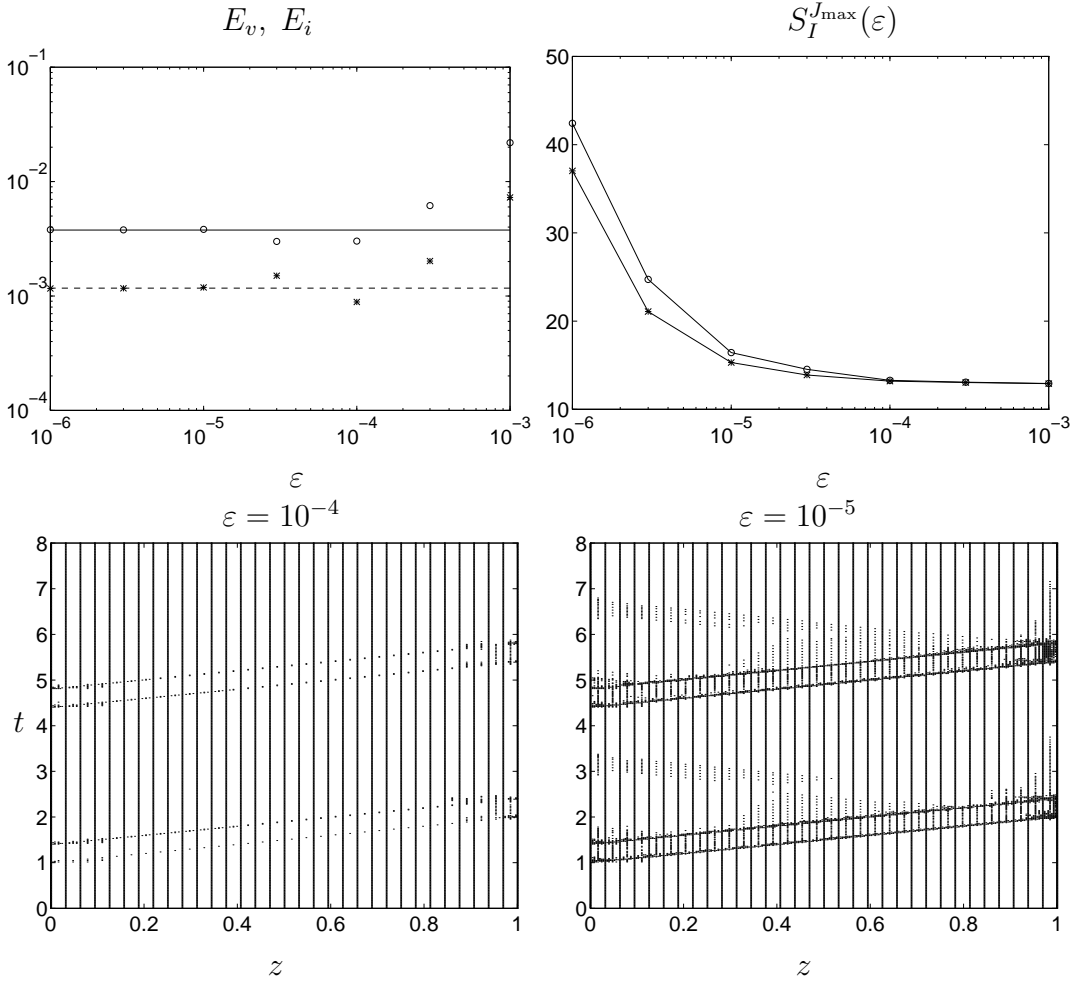
**Figure 6.10:** TDSE solution for the unmatched lossless scalar uniform line with  $J_{\max} = 8$ ,  $L = 2$ ,  $\tilde{L} = 4$ : location in the  $(z, t)$  plane of the voltage wavelet coefficients larger than  $\varepsilon = 10^{-5}$ .



**Figure 6.11:** TDSE solution for the matched 1:4 exponential line excited by a trapezoidal train voltage source, obtained with  $J_{\max} = 8$ ,  $L = 2$ ,  $\tilde{L} = 4$ .

### 6.2.2 The exponential line

This section solves the scalar exponential line already treated in Sections 1.3.1 and 6.1.1 with a trapezoidal pulse excitation (rise time  $t_r = 0.4$  s, duration  $\tau = 3.4$  s). We consider the matched line in order to compare the results with the reference solution obtained through inverse FFT. Figure 6.11 reports the results of the TDSE simulations with no wavelet thresholding. The top panels of Fig. 6.12 report the approximation errors on voltage and current obtained with wavelet thresholding for different values of the threshold  $\varepsilon$  together with the corresponding sparsity indices. Also in this case, as for the uniform line, highly sparse representations can be obtained at no loss of accuracy. Indeed, the errors obtained with thresholding are always comparable to the corresponding errors in the solution without thresholding, indicated in the top left panel of the figure with solid (voltage) and dashed (current) lines. The bottom panels of



**Figure 6.12:** Matched 1:4 exponential line with trapezoidal voltage source. Top left panel: maximum absolute error on voltage (circles) and current (stars) obtained with the TDSE method ( $J_{\max} = 8$ ,  $L = 2$ ,  $\tilde{L} = 4$ ) as a function of the threshold  $\varepsilon$  for the wavelet coefficients. The solid and dashed lines indicate the approximation errors obtained with no wavelet thresholding for voltage and current, respectively. Top right panel: sparsity index of the nonlinear approximation of voltage and current obtained with the same wavelet thresholds used in the left panel. Bottom panels: location of the voltage wavelet coefficients larger than  $\varepsilon = 10^{-4}$  (left) and  $\varepsilon = 10^{-5}$  (right).

Fig. 6.12 depict the location of significant wavelet coefficients in the  $(z, t)$  plane for two different values of the threshold  $\varepsilon$ . As for the uniform line, these plots give a quick interpretation of the solution in terms of travelling waves, which in this case are not simple translations of a single pulse as the time increases. It should be noted that the location of the significant coefficients is determined from the solution itself while it is generated by the time-stepping routine.

## 6.3 Adaptive TDSE

This section is devoted to the derivation of an alternative formulation of the TDSE method with respect to the derivation of Sec. 1.1. The advantages of the following formulation are fourfold. First, the resulting system of ODE's will be explicit in the time derivatives of the unknowns. Second, there will be no derivatives of the source voltage waveforms in the system of ODE's to be solved. Third, the use of biorthogonal trial and test functions allows to use wavelet bases and therefore to compute the solution directly in its adapted wavelet representation. Fourth, a preprocessing renormalization of the NMTL equations leads to solve non-dimensional equations. This will allow to perform wavelet thresholding without concerning about the physical nature of the unknowns (voltage or current). Section 6.3.1 details the formulation of the time-explicit TDSE method, and Sec. 6.3.2 shows its application with wavelet bases. Finally, Sec. 6.3.3 presents some applications of the method.

### 6.3.1 The general formulation

We begin with the derivation of the non-dimensional form of the NMTL equations. We introduce the following reference quantities,

- reference voltage  $V_0$ ,
- reference current  $I_0$ ,
- reference impedance  $R_0 = V_0/I_0$ ,
- reference time delay  $T$ ,
- reference length  $\mathcal{L}$  (the length of the line),
- reference speed  $v_0 = \mathcal{L}/T$ .

Even if these reference quantities are arbitrary, the choice of their numerical value must be guided by prior knowledge of the characteristics of the transmission line under investigation. For example,  $T$  should be chosen as close as possible to the one-way propagation delay time, while  $R_0$  should be approximately equal to the impedance level of the line. These normalization constants are used to define the following non-dimensional variables

- normalized space coordinate  $x = z/\mathcal{L}$ ,  $x \in [0, 1]$ ,
- normalized time coordinate  $\tau = t/T$ ,
- normalized voltage vector  $\mathbf{v}(x, \tau) = \mathbf{V}(\mathcal{L}x, T\tau)/V_0$ ,
- normalized current vector  $\mathbf{i}(x, \tau) = \mathbf{I}(\mathcal{L}x, T\tau)/I_0$ .

A simple change of variable in the NMTL equations (1.1)-(1.2) leads to

$$\frac{\partial}{\partial x} \mathbf{i}(x, \tau) = -\bar{\mathbf{L}}(x) \frac{\partial}{\partial \tau} \mathbf{i}(x, \tau) - \bar{\mathbf{R}}(x) \mathbf{i}(x, \tau), \quad (6.2)$$

$$\frac{\partial}{\partial x} \mathbf{i}(x, \tau) = -\bar{\mathbf{C}}(x) \frac{\partial}{\partial \tau} \mathbf{v}(x, \tau) - \bar{\mathbf{G}}(x) \mathbf{v}(x, \tau), \quad (6.3)$$

where

$$\bar{\mathbf{L}}(x) = \frac{v_0}{R_0} \mathbf{L}(\mathcal{L}x)$$

$$\bar{\mathbf{C}}(x) = v_0 R_0 \mathbf{C}(\mathcal{L}x)$$

$$\bar{\mathbf{R}}(x) = \frac{\mathcal{L}}{R_0} \mathbf{R}(\mathcal{L}x)$$

$$\bar{\mathbf{G}}(x) = \mathcal{L} R_0 \mathbf{G}(\mathcal{L}x).$$

The same procedure can be applied to the Thévenin terminations of the line (1.3)-(1.4), obtaining their non-dimensional form

$$\mathbf{v}(0, \tau) = \mathbf{v}_S(\tau) - \mathbf{r}_S \mathbf{i}(0, \tau) \quad (6.4)$$

$$\mathbf{v}(1, \tau) = \mathbf{v}_L(\tau) + \mathbf{r}_L \mathbf{i}(1, \tau), \quad (6.5)$$

where

$$\mathbf{r}_S = \frac{\mathbf{R}_S}{R_0}, \quad \mathbf{r}_L = \frac{\mathbf{R}_L}{R_0}.$$

We discussed in Sec. 6.2 that it would be convenient to obtain a time-explicit system of ODE's in order to save computation time. This can only be achieved when the partial differential equations to be discretized are explicit in the time derivatives. For this reason, we derive two sets of equations equivalent to the NMTL equations by inverting the normalized inductance and capacitance matrices  $\bar{\mathbf{L}}(x)$  and  $\bar{\mathbf{C}}(x)$ . This is possible because these two matrices are symmetric and positive definite for all  $x$  [7]. We can introduce then the normalized *elastance* and *reciprocal inductance* matrices as the inverse of the normalized capacitance and inductance matrices,

$$\mathbf{S}(x) = \bar{\mathbf{C}}^{-1}(x), \quad \mathbf{\Gamma}(x) = \bar{\mathbf{L}}^{-1}(x), \quad \forall x. \quad (6.6)$$

The NMTL equations become

$$\frac{\partial}{\partial \tau} \mathbf{i}(x, \tau) = -\mathbf{\Gamma}(x) \frac{\partial}{\partial x} \mathbf{v}(x, \tau) - \mathbf{\Gamma} \mathbf{R}(x) \mathbf{i}(x, \tau), \quad (6.7)$$

$$\frac{\partial}{\partial \tau} \mathbf{v}(x, \tau) = -\mathbf{S}(x) \frac{\partial}{\partial x} \mathbf{i}(x, \tau) - \mathbf{S} \mathbf{G}(x) \mathbf{v}(x, \tau), \quad (6.8)$$

where

$$\mathbf{S} \mathbf{G}(x) = \bar{\mathbf{C}}^{-1}(x) \bar{\mathbf{G}}(x), \quad \mathbf{\Gamma} \mathbf{R}(x) = \bar{\mathbf{L}}^{-1}(x) \bar{\mathbf{R}}(x), \quad \forall x. \quad (6.9)$$

Following the same procedure used in Section 1.1, we expand the voltages and currents into the functions  $\zeta_n$ ,

$$\mathbf{v}(x, \tau) = \sum_{n=1}^{N_\zeta} \zeta_n(x) \mathbf{v}_n(\tau), \quad (6.10)$$

$$\mathbf{i}(x, \tau) = \sum_{n=1}^{N_\zeta} \zeta_n(x) \mathbf{i}_n(\tau), \quad (6.11)$$

and the four matrices  $\mathbf{\Gamma}(x)$ ,  $\mathbf{S}(x)$ ,  $\mathbf{SG}(x)$ ,  $\mathbf{\Gamma R}(x)$  into the functions  $\phi_k$ ,

$$\begin{aligned} \mathbf{\Gamma} &= \sum_{k=1}^{N_\phi} \phi_k(x) \mathbf{\Gamma}_k, \\ \mathbf{S} &= \sum_{k=1}^{N_\phi} \phi_k(x) \mathbf{S}_k, \\ \mathbf{\Gamma R} &= \sum_{k=1}^{N_\phi} \phi_k(x) \mathbf{\Gamma R}_k, \\ \mathbf{SG} &= \sum_{k=1}^{N_\phi} \phi_k(x) \mathbf{SG}_k. \end{aligned} \quad (6.12)$$

Projecting now equations (6.7)-(6.8) onto the test functions  $\eta_m$ , we get the equations

$$\sum_{n=1}^{N_\zeta} \mathbf{E}_{mn} \frac{d}{d\tau} \mathbf{i}_n(\tau) + \sum_{n=1}^{N_\zeta} \mathbf{\Gamma R}_{mn} \mathbf{i}_n(\tau) + \sum_{n=1}^{N_\zeta} \mathbf{\Gamma}_{mn} \mathbf{v}_n(\tau) = 0, \quad (6.13)$$

$$\sum_{n=1}^{N_\zeta} \mathbf{E}_{mn} \frac{d}{d\tau} \mathbf{v}_n(\tau) + \sum_{n=1}^{N_\zeta} \mathbf{SG}_{mn} \mathbf{v}_n(\tau) + \sum_{n=1}^{N_\zeta} \mathbf{S}_{mn} \mathbf{i}_n(\tau) = 0, \quad (6.14)$$

valid  $\forall m = 1, \dots, N_\zeta$ . The matrices used in the above expressions can be expressed as

$$\begin{aligned} \mathbf{\Gamma R}_{mn} &= \sum_{k=1}^{N_\phi} \mathbf{\Gamma R}_k B_{mn}^{(k)}, \\ \mathbf{\Gamma}_{mn} &= \sum_{k=1}^{N_\phi} \mathbf{\Gamma}_k F_{mn}^{(k)}, \\ \mathbf{SG}_{mn} &= \sum_{k=1}^{N_\phi} \mathbf{SG}_k B_{mn}^{(k)}, \\ \mathbf{S}_{mn} &= \sum_{k=1}^{N_\phi} \mathbf{S}_k F_{mn}^{(k)}, \end{aligned} \quad (6.15)$$



where

$$B_{mn}^{(k)} = \langle \zeta_n \phi_k, \eta_m \rangle \quad (6.16)$$

$$F_{mn}^{(k)} = \left\langle \frac{\partial}{\partial x} \zeta_n \phi_k, \eta_m \right\rangle. \quad (6.17)$$

These equations are explicit only when the coefficients of the time derivatives are such that

$$\mathbf{E}_{mn} = \langle \zeta_n, \eta_m \rangle \mathcal{I}_P = \delta_{mn} \mathcal{I}_P.$$

This is true only when the two sets  $\{\zeta_n\}$  and  $\{\eta_m\}$  are biorthogonal. In this case we obtain

$$\frac{d}{d\tau} \mathbf{i}_m(\tau) = - \sum_{n=1}^{N_\zeta} \mathbf{\Gamma} \mathbf{R}_{mn} \mathbf{i}_n(\tau) - \sum_{n=1}^{N_\zeta} \mathbf{\Gamma}_{mn} \mathbf{v}_n(\tau), \quad (6.18)$$

$$\frac{d}{d\tau} \mathbf{v}_m(\tau) = - \sum_{n=1}^{N_\zeta} \mathbf{S} \mathbf{G}_{mn} \mathbf{v}_n(\tau) - \sum_{n=1}^{N_\zeta} \mathbf{S}_{mn} \mathbf{i}_n(\tau). \quad (6.19)$$

We proceed now to apply the line terminations. We will follow a procedure slightly different from the one we used in Sec. 1.1. As the choice of basis functions needed to obtain a time-explicit discretization of the NMTL equations is restricted to biorthogonal pairs, we will relax some of the assumptions made in Sec. 1.1. In fact, as we are planning to use the biorthogonal wavelet pairs as trial and test function, we must take into account that more than one wavelet function is nonvanishing at the edges of the unit interval. The trial and test functions will have hereafter the following properties

$$\begin{aligned} \zeta_1(0) &\neq 0, & \zeta_{N_\zeta}(1) &\neq 0 \\ \zeta_1(1) &= 0, & \zeta_{N_\zeta}(0) &= 0 \\ \text{supp } \zeta_1 \cup \text{supp } \eta_{N_\zeta} &= \text{supp } \zeta_{N_\zeta} \cup \text{supp } \eta_0 = \emptyset. \end{aligned} \quad (6.20)$$

It is important that the contributions coming from the two edges do not interact directly with each other, as stated by the second and third condition. It should be noted that both the biorthogonal scaling function and wavelet bases constructed in Chapter 4 satisfy these conditions. Note also that for this derivation we identified the main nonvanishing border functions with the first and last trial function. However, this ordering does not necessarily need to be preserved in the actual implementation.

Let us take now the load equations (6.4)-(6.5) and substitute the expansion of the voltage and current at the two edges,

$$\begin{aligned} \sum_{n=1}^{N_\zeta-1} \zeta_n(0) \mathbf{v}_n &= \mathbf{v}_S - \mathbf{r}_S \sum_{n=1}^{N_\zeta-1} \zeta_n(0) \mathbf{i}_n \\ \sum_{n=2}^{N_\zeta} \zeta_n(1) \mathbf{v}_n &= \mathbf{v}_L + \mathbf{r}_L \sum_{n=2}^{N_\zeta} \zeta_n(1) \mathbf{i}_n. \end{aligned}$$

We define now two pairs of new variables at the two edges,

$$\begin{aligned} \mathbf{a}_1 &= \sum_{n=1}^{N_\zeta-1} \lambda_n^0 [\mathbf{v}_n + \mathbf{r}_S \mathbf{i}_n] \\ \mathbf{b}_1 &= \sum_{n=1}^{N_\zeta-1} \lambda_n^0 [\mathbf{v}_n - \mathbf{r}_S \mathbf{i}_n] \\ \mathbf{a}_{N_\zeta} &= \sum_{n=2}^{N_\zeta} \lambda_n^1 [\mathbf{v}_n + \mathbf{r}_L \mathbf{i}_n] \\ \mathbf{b}_{N_\zeta} &= \sum_{n=2}^{N_\zeta} \lambda_n^1 [\mathbf{v}_n - \mathbf{r}_L \mathbf{i}_n]. \end{aligned}$$

These can be interpreted as progressive and regressive normalized voltage waves with respect to a reference impedance equal to the load normalized resistance matrices. The coefficients  $\lambda_n^0$  and  $\lambda_n^1$  are defined in terms of the boundary values of the trial functions,

$$\lambda_n^0 = \frac{\zeta_n(0)}{\zeta_1(0)}, \quad \lambda_n^1 = \frac{\zeta_n(1)}{\zeta_{N_\zeta}(1)}. \quad (6.21)$$

If the matrices  $\mathbf{r}_S$  and  $\mathbf{r}_L$  are invertible, we can express the voltage and current coefficients of the edge trial functions with  $n = 1$  and  $n = N_\zeta$  in terms of the newly defined variables and the “internal” ( $m = 2, \dots, N_\zeta - 1$ ) voltage and current coefficients,

$$\begin{aligned} \mathbf{v}_1 &= \frac{1}{2}(\mathbf{a}_1 + \mathbf{b}_1) - \sum_{n=2}^{N_\zeta-1} \lambda_n^0 \mathbf{v}_n \\ \mathbf{i}_1 &= \frac{1}{2}\mathbf{r}_S^{-1}(\mathbf{a}_1 - \mathbf{b}_1) - \sum_{n=2}^{N_\zeta-1} \lambda_n^0 \mathbf{i}_n \\ \mathbf{v}_{N_\zeta} &= \frac{1}{2}(\mathbf{a}_{N_\zeta} + \mathbf{b}_{N_\zeta}) - \sum_{n=2}^{N_\zeta-1} \lambda_n^1 \mathbf{v}_n \\ \mathbf{i}_{N_\zeta} &= \frac{1}{2}\mathbf{r}_L^{-1}(\mathbf{a}_{N_\zeta} - \mathbf{b}_{N_\zeta}) - \sum_{n=2}^{N_\zeta-1} \lambda_n^1 \mathbf{i}_n. \end{aligned}$$

The main advantage in using the new variables  $\mathbf{a}$ ,  $\mathbf{b}$  is in the implementation of the load equations. In fact, these equations are simply expressed as Dirichlet type conditions on  $\mathbf{a}_1$  and  $\mathbf{b}_{N_\zeta}$ . More precisely,

$$\mathbf{a}_1(\tau) = \frac{\mathbf{v}_S(\tau)}{\zeta_1(0)}, \quad (6.22)$$

$$\mathbf{b}_{N_\zeta}(\tau) = \frac{\mathbf{v}_L(\tau)}{\zeta_{N_\zeta}(1)}. \quad (6.23)$$

Therefore, only two border equations involving the variables  $\mathbf{b}_1$  and  $\mathbf{a}_{N_\zeta}$  need to be derived. These are easily obtained by forming the linear combinations

$$\begin{aligned} & \sum_{m=1}^{N_\zeta-1} \lambda_m^0 [(6.19) - \mathbf{r}_S(6.18)] \\ & \sum_{m=2}^{N_\zeta} \lambda_m^1 [(6.19) + \mathbf{r}_L(6.18)]. \end{aligned}$$

Finally, the elimination of  $\mathbf{v}_1$ ,  $\mathbf{i}_1$ ,  $\mathbf{v}_{N_\zeta}$ , and  $\mathbf{i}_{N_\zeta}$  from all the resulting equations ( $2P$  border equations plus  $2P(N-1)$  “internal equations”) leads to the system of ODE’s

$$\frac{d}{d\tau} \bar{\mathbf{x}}(\tau) = \mathbf{\Theta} \bar{\mathbf{x}}(\tau) + \mathbf{\Omega}_S \mathbf{v}_S(\tau) + \mathbf{\Omega}_L \mathbf{v}_L(\tau), \quad (6.24)$$

where

$$\bar{\mathbf{x}} = [\mathbf{b}_1^T, \mathbf{i}_2^T, \dots, \mathbf{i}_{N_\zeta-1}^T, \mathbf{a}_{N_\zeta}^T, \mathbf{v}_2^T, \dots, \mathbf{v}_{N_\zeta-1}^T]^T. \quad (6.25)$$

It should be noted that no derivatives of the source voltages are involved. This is impossible to achieve when the ODE system is derived from the non-explicit NMTL equations like in Sec. 1.1. This allows to process any source waveform of practical interest as long as it can be evaluated at an arbitrary time  $t$ . The drawback is that the load matrices  $\mathbf{r}_S$  and  $\mathbf{r}_L$  must be invertible. No short-circuit or open-circuit terminations can be used with this formulation. A permutation matrix  $\mathbf{T}$  can be applied to the above system as in Section 1.2 to rearrange the unknowns and obtain a more compact system matrix,

$$\frac{d}{d\tau} \hat{\mathbf{x}}(\tau) = \hat{\mathbf{\Theta}} \hat{\mathbf{x}}(\tau) + \hat{\mathbf{\Omega}}_S \mathbf{v}_S(\tau) + \hat{\mathbf{\Omega}}_L \mathbf{v}_L(\tau), \quad (6.26)$$

where

$$\hat{\mathbf{x}} = [\mathbf{b}_1^T, \mathbf{i}_2^T, \mathbf{v}_2^T, \dots, \mathbf{i}_{N_\zeta-1}^T, \mathbf{v}_{N_\zeta-1}^T, \mathbf{a}_{N_\zeta}^T]^T \quad (6.27)$$

and

$$\hat{\mathbf{\Theta}} = \mathbf{T}^T \mathbf{\Theta} \mathbf{T}, \quad \hat{\mathbf{\Omega}}_S = \mathbf{T}^T \mathbf{\Omega}_S, \quad \hat{\mathbf{\Omega}}_L = \mathbf{T}^T \mathbf{\Omega}_L.$$

### 6.3.2 Time-explicit TDSE with wavelets

This section will particularize the formulation of Sec. 6.3.1 to the use of wavelet bases on the unit interval constructed in Chapter 4 as trial and test functions. Two different implementations will be described. The first uses the standard biorthogonal wavelet bases. The ODE system that can be derived with these basis functions, however, is not fully equivalent to the system obtained by using the scaling functions at the maximum refinement level. Consequently, we present also a second scheme, which uses a slightly modified set of wavelets. This scheme leads to a system of ODE’s that is fully equivalent to the system derived in Sec. 6.1.

### Standard wavelet bases

This section illustrates the adaptive TDSE method through use of wavelet bases on the unit interval. As the trial and test functions must be biorthogonal to preserve the time-explicit form of the discretized NMTL equations, we will expand the solution into primal wavelets and test the equations with dual wavelets. The total number of trial and test functions will be equal to the dimension of the scaling function space at the maximum refinement level  $J_{\max}$ . We will set

$$\begin{aligned} \{\zeta_n\}_{n=1}^{N_\zeta} &= \{\varphi_{J_0,k} : k = 0, \dots, \dim V_{J_0} - 1\} \cup \bigcup_{j=J_0}^{J_{\max}-1} \{\psi_{jk} : k = 0, \dots, 2^j - 1\} \\ \{\phi_k\}_{k=1}^{N_\phi} &= \{\varphi_{J_{\max},k} : k = 0, \dots, \dim V_{J_{\max}} - 1\} \\ \{\eta_m\}_{m=1}^{N_\eta} &= \{\tilde{\varphi}_{J_0,k} : k = 0, \dots, \dim V_{J_0} - 1\} \cup \bigcup_{j=J_0}^{J_{\max}-1} \{\tilde{\psi}_{jk} : k = 0, \dots, 2^j - 1\}. \end{aligned} \quad (6.28)$$

The expansion functions for the per-unit-length parameters are the scaling functions at the maximum refinement level  $J_{\max}$ , because there is no need for an adapted (sparse) representation. The ordering sequence of the trial and test functions will be important in the following. The scaling functions at the minimum level  $J_0$  are placed first, followed by the wavelets with increasing refinement levels. Also, the indices  $n$  and  $m$  become now double indices,

$$\{n\} \leftrightarrow \{(j, k)\}, \quad \{m\} \leftrightarrow \{(j', k')\},$$

with the first indicating the refinement level and the second the number of function at level  $j$  according to the definitions in Eq. (6.28). To insure uniqueness in this representation, we conventionally indicate with the “dummy” refinement level  $j = J_0 - 1$  the pairs  $(j, k)$  referring to the scaling functions at the minimum level  $J_0$ .

With the foregoing definitions it is possible to compute the boundary values of the trial functions, i.e., the coefficients  $\lambda_n^0$  and  $\lambda_n^1$  of Eq. (6.21). We will suppose that the boundary adapted biorthogonalization has been used for both scaling functions and wavelets, and that the basis of the polynomials  $\{\rho_\alpha(x), \alpha = 0, \dots, L - 1\}$  is such that

$$\rho_0(0) = 1, \quad \rho_\alpha(0) = 0 \quad \forall \alpha > 0.$$

This assumption is true for the two polynomial bases considered in this work, i.e., monomials and Bernstein polynomials. A straightforward substitution leads to the boundary values of the non-biorthogonal border scaling functions,

$$\theta_{00}(0) = 1, \quad \theta_{0k}(0) = 0 \quad \forall k > 0,$$

and similarly for the duals. With the boundary adapted biorthogonalization we get the corresponding boundary values of the biorthogonal scaling functions,

$$\varphi_{00}(0) = D_{00}, \quad \varphi_{0k}(0) = 0 \quad \forall k > 0,$$

where  $D$  is the change of basis matrix for the biorthogonalization of the primal scaling functions, and similarly for the duals, for which  $\widetilde{D}_{00}$  is used. The boundary values of the scaling functions at a generic refinement level  $j$  is readily obtained as

$$\varphi_{j0}(0) = 2^{j/2} D_{00}, \quad \varphi_{jk}(0) = 0 \quad \forall j, \forall k > 0.$$

The boundary values of the wavelets can be derived from the biorthogonal wavelet filters  $\mathcal{G}^\perp$  and  $\tilde{\mathcal{G}}^\perp$  for primals and duals respectively. At any refinement level  $j$  we have

$$\psi_{jk}(0) = 0, \quad \forall k > 0$$

due to boundary adaption, while the value at 0 of the first wavelet is

$$\psi_{j0}(0) = \sum_{l \geq 0} \mathcal{G}_{j+1,l}^\perp \varphi_{j+1,l}(0) = \mathcal{G}_{00}^\perp \varphi_{j+1,0} = \sqrt{2} 2^{j/2} \mathcal{G}_{00}^\perp D_{00}.$$

The dual case can be obtained with obvious substitutions. We can now write explicitly the coefficients  $\lambda_n^0$ ,

$$\begin{aligned} \lambda_{J_0-1,0}^0 &= 1, \\ \lambda_{j,0}^0 &= \mathcal{G}_{00}^\perp 2^{\frac{j-J_0+1}{2}} \quad \forall j = J_0, \dots, J_{\max}-1, \\ \lambda_{j,k}^0 &= 0 \quad \text{otherwise.} \end{aligned}$$

Finally, recalling the symmetry of the construction on the unit interval, we can easily prove that the boundary values at the right edge  $x = 1$  are equal to the values obtained at  $x = 0$ . Therefore, the coefficients  $\lambda_n^1$  are

$$\begin{aligned} \lambda_{J_0-1, \dim V_{J_0}-1}^1 &= 1, \\ \lambda_{j, 2^j-1}^1 &= \mathcal{G}_{00}^\perp 2^{\frac{j-J_0+1}{2}} \quad \forall j = J_0, \dots, J_{\max}-1, \\ \lambda_{j,k}^1 &= 0 \quad \text{otherwise.} \end{aligned}$$

We turn now to the computation of the system matrices  $\Theta$  or  $\widehat{\Theta}$  in the wavelet basis. As the derivation of Sec. 6.3.1 was tailored for the application of wavelets, no modifications are necessary, except a reordering of the trial and test functions. The principal nonvanishing functions at the edges are indeed the first and last scaling functions at level  $J_0$ , and these are not the first and last functions in the ordering of Eq. (6.28). This operation is trivial. The only crucial point is the computation of the matrices  $\mathbf{\Gamma}_{mn}$ ,  $\mathbf{S}_{mn}$ ,  $\mathbf{\Gamma R}_{mn}$ , and  $\mathbf{S G}_{mn}$  of Eq. (6.15), which involve inner products of wavelets through Eqs. (6.16) and (6.17). In particular, as the expansion functions of the per-unit-length matrices are the scaling functions, like in the original formulation of the TDSE method in Sec. 6.1, the same procedure can be applied to compute the expansion coefficients  $\mathbf{\Gamma}_k$ ,  $\mathbf{S}_k$ ,  $\mathbf{\Gamma R}_k$ , and  $\mathbf{S G}_k$ . Therefore, we only need to focus on the scalar inner products  $B_{mn}^{(k)}$  and  $F_{mn}^{(k)}$ .

Under an abstract point of view, the aforementioned inner products can be regarded as the projection of two operators onto the spaces spanned by the trial and test functions. More precisely, these operators act on a generic function  $f$  as follows

$$(B^{(k)}f)(x) = \phi_k(x)f(x), \quad (F^{(k)}f)(x) = \phi_k(x)\frac{d}{dx}f(x),$$

If we choose the biorthogonal scaling functions as trial and test functions we have that the matrices  $B_{mn}^{(k)}$  and  $F_{mn}^{(k)}$  represent the approximated operators

$$B_{mn}^{(k)} = P_{J_{\max}} B^{(k)} P_{J_{\max}}, \quad F_{mn}^{(k)} = P_{J_{\max}} F^{(k)} P_{J_{\max}},$$

where  $J_{\max}$  is the refinement level for the scaling function spaces and  $P_{J_{\max}}$  is the corresponding projection operator. It follows that the wavelet change of basis process already applied in Section 6.2 to a general operator  $T$  can be applied also in this case to obtain the wavelet representation of the operators  $B^{(k)}$  and  $F^{(k)}$  in the wavelet basis, obtaining

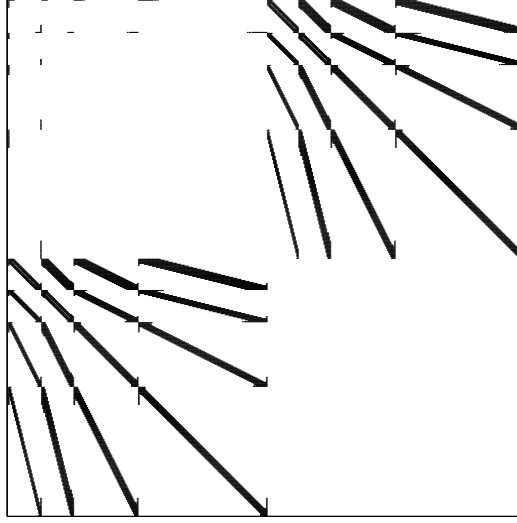
$$\hat{B}^{(k)} = \widetilde{\mathcal{W}} B^{(k)} \mathcal{W}^T, \quad \hat{F}^{(k)} = \widetilde{\mathcal{W}} F^{(k)} \mathcal{W}^T.$$

The matrices corresponding to these two operators can then be used to generate the system matrices of Eqs. (6.18)-(6.18), and consequently the system of ODE's (6.26) through the procedure indicated in Sec. 6.3.1. The typical structure of the system matrix is depicted in Fig. 6.13, where the unknowns are ordered with all the current expansion coefficients first followed by the voltage coefficients, and where no losses are included. The placement of the nonvanishing entries in this matrix shows the mutual interaction between different scales, indicated by non-zero elements outside a band around the main diagonal.

### Modified wavelet bases

The TDSE implementation with standard wavelets illustrated in the foregoing section is not fully equivalent to the scheme that uses the scaling functions at the maximum refinement level. This fact is due to a difference in the inclusion of the boundary conditions. We do not report here the formal proof, which involves a few straightforward calculations. Some numerical tests performed with the standard hierarchical wavelet bases showed indeed a significant loss of precision with respect to the canonical bases.

In order to have fully equivalent systems, the sets of trial and test functions should include the two nonvanishing scaling functions at the maximum refinement level  $J_{\max}$ . In addition, these should be the only nonvanishing functions at the borders. This is obviously not possible when using the standard hierarchical wavelet bases. However, we show in the following that it is possible to define a modified hierarchical system with these features.



**Figure 6.13:** Structure of the system matrix of Eq. (6.26) when hierarchical basis functions are used as trial and test functions.

Let us recall the multilevel decomposition of the scaling functions space at the maximum level  $J_{\max}$ ,

$$V_{J_{\max}} = V_{J_0} \oplus W_{J_0} \oplus \cdots \oplus W_{J_{\max}-1}. \quad (6.29)$$

We consider in the following only the left edge  $x = 0$ , because the same procedure can be applied to the right edge  $x = 1$  using the symmetry of the basis functions. Each space in Eq. (6.29) can be decomposed into a direct sum of a border space, indicated with the superscript  $^o$ , and an internal space, indicated with the superscript  $^i$ ,

$$\begin{aligned} V_{J_{\max}} &= V_{J_{\max}}^o \oplus V_{J_{\max}}^i \\ V_{J_0} &= V_{J_0}^o \oplus V_{J_0}^i \\ W_j &= W_j^o \oplus W_j^i, \quad j = J_0, \dots, J_{\max} - 1 \end{aligned}$$

The border spaces are generated by all the nonvanishing scaling functions and wavelets. If boundary adaption is used their dimension is exactly one. Conversely, the internal spaces are generated by scaling functions and wavelets that are all vanishing at the border. If we collect all the nonvanishing functions in a border space  $V_{J_{\max}}^b$ , we have that

$$V_{J_{\max}} = \left\{ V_{J_0}^i \oplus \bigoplus_{j=J_0}^{J_{\max}-1} W_j^i \right\} \oplus V_{J_{\max}}^b, \quad (6.30)$$

where

$$V_{J_{\max}}^b = V_{J_0}^o \oplus \bigoplus_{j=J_0}^{J_{\max}-1} W_j^o. \quad (6.31)$$

It can be proved that the border space at the maximum level  $V_{J_{\max}}^o$  is part of the border space  $V_{J_{\max}}^b$ . Therefore, we can define a modified detail space  $W^b$  as

$$V_{J_{\max}}^b = V_{J_{\max}}^o \oplus W^b. \quad (6.32)$$

The dimension of this detail space is  $J_{\max} - J_0$ . A set of basis functions for this space can be obtained by removing from the nonvanishing border wavelets their projection onto  $V_{J_{\max}}^o$ ,

$$\psi_{j,0}^{\text{mod}} = \psi_{j,0} - \langle \psi_{j,0}, \tilde{\varphi}_{J_{\max},0} \rangle \varphi_{J_{\max},0}, \quad j = J_0, \dots, J_{\max} - 1. \quad (6.33)$$

It can be shown that these functions are linearly independent. Moreover, also the modified scaling function at level  $J_0$ , obtained in the same way, can be expressed as a linear combination of these modified wavelets. It should be noted that these wavelets are not zero-mean functions and are not biorthogonal to the corresponding duals. However, they all vanish at the border. The gramian matrix and its inverse can be explicitly calculated, obtaining

$$\begin{aligned} \Psi_{j',j} &= \langle \tilde{\psi}_{j',0}^{\text{mod}}, \psi_{j,0}^{\text{mod}} \rangle = \delta_{jj'} - 2^{(j+j'-2J_{\max})/2}, \\ [\Psi^{-1}]_{j',j} &= \delta_{jj'} + 2^{(j+j'-2J_0)/2}, \end{aligned}$$

for  $j, j' = J_0, \dots, J_{\max} - 1$ .

Two different hierarchical representations of a function  $v \in V_{J_{\max}}$  have been constructed,

$$\begin{aligned} v &= \check{v}_{J_0,0} \varphi_{J_0,0} + \sum_{j=J_0}^{J_{\max}-1} \hat{v}_{j,0} \psi_{j,0} + \sum_{k \geq 1} \check{v}_{J_0,k} \varphi_{J_0,k} + \sum_{j=J_0}^{J_{\max}-1} \sum_{k \geq 1} \hat{v}_{j,k} \psi_{j,k}, \\ v &= \check{v}_{J_{\max},0} \varphi_{J_{\max},0} + \sum_{j=J_0}^{J_{\max}-1} \hat{v}_{j,0} \psi_{j,0}^{\text{mod}} + \sum_{k \geq 1} \check{v}_{J_0,k} \varphi_{J_0,k} + \sum_{j=J_0}^{J_{\max}-1} \sum_{k \geq 1} \hat{v}_{j,k} \psi_{j,k}. \end{aligned}$$

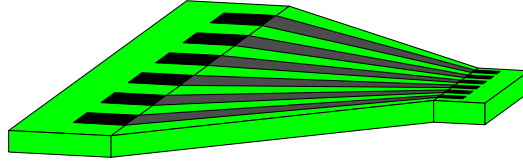
The first is the usual multilevel decomposition, and the second is a modified decomposition which uses the same internal functions, the nonvanishing scaling function at the maximum level  $J_{\max}$ , and the modified border wavelets of Eq. (6.33). A basis change between these two representations can be derived analytically by using the boundary adaption of wavelets, obtaining

$$\begin{aligned} \check{v}_{J_{\max},0} &= \mathcal{H}_{00}^{J_{\max}-J_0} \check{v}_{J_0,0} + \sum_{j=J_0}^{J_{\max}-1} \mathcal{G}_{00} \mathcal{H}_{00}^{J_{\max}-j-1} \hat{v}_{j,0}, \\ \hat{v}_{j,0}^{\text{mod}} &= \hat{v}_{j,0} - \tilde{\mathcal{G}}_{00} \tilde{\mathcal{H}}_{00}^{J_0-j-1} \check{v}_{J_0,0}, \quad j = J_0, \dots, J_{\max} - 1, \end{aligned}$$

and

$$\begin{aligned} \check{v}_{J_0,0} &= \tilde{\mathcal{H}}_{00}^{J_{\max}-J_0} \check{v}_{J_{\max},0} - \mathcal{G}_{00} \tilde{\mathcal{H}}_{00}^{J_{\max}-J_0} \sum_{j=J_0}^{J_{\max}-1} \mathcal{H}_{00}^{J_{\max}-j-1} \hat{v}_{j,0}^{\text{mod}}, \\ \hat{v}_{j,0} &= \tilde{\mathcal{G}}_{00} \tilde{\mathcal{H}}_{00}^{J_{\max}-j-1} \check{v}_{J_{\max},0} + \sum_{j'=J_0}^{J_{\max}-1} \Psi_{j,j'} \hat{v}_{j',0}^{\text{mod}}, \quad j = J_0, \dots, J_{\max} - 1. \end{aligned}$$





**Figure 6.14:** Geometry of a high-speed packaging interconnect.

This basis change can be applied to the system of Eqs. (6.18)-(6.19). The result can be expressed as

$$\begin{aligned} \frac{d}{d\tau} \mathbf{i}_m^{\text{mod}}(\tau) &= - \sum_{n=1}^{N_\zeta} \mathbf{\Gamma} \mathbf{R}_{mn}^{\text{mod}} \mathbf{i}_n^{\text{mod}}(\tau) - \sum_{n=1}^{N_\zeta} \mathbf{\Gamma}_{mn}^{\text{mod}} \mathbf{v}_n^{\text{mod}}(\tau), \\ \frac{d}{d\tau} \mathbf{v}_m^{\text{mod}}(\tau) &= - \sum_{n=1}^{N_\zeta} \mathbf{S} \mathbf{G}_{mn}^{\text{mod}} \mathbf{v}_n^{\text{mod}}(\tau) - \sum_{n=1}^{N_\zeta} \mathbf{S}_{mn}^{\text{mod}} \mathbf{i}_n^{\text{mod}}(\tau), \end{aligned}$$

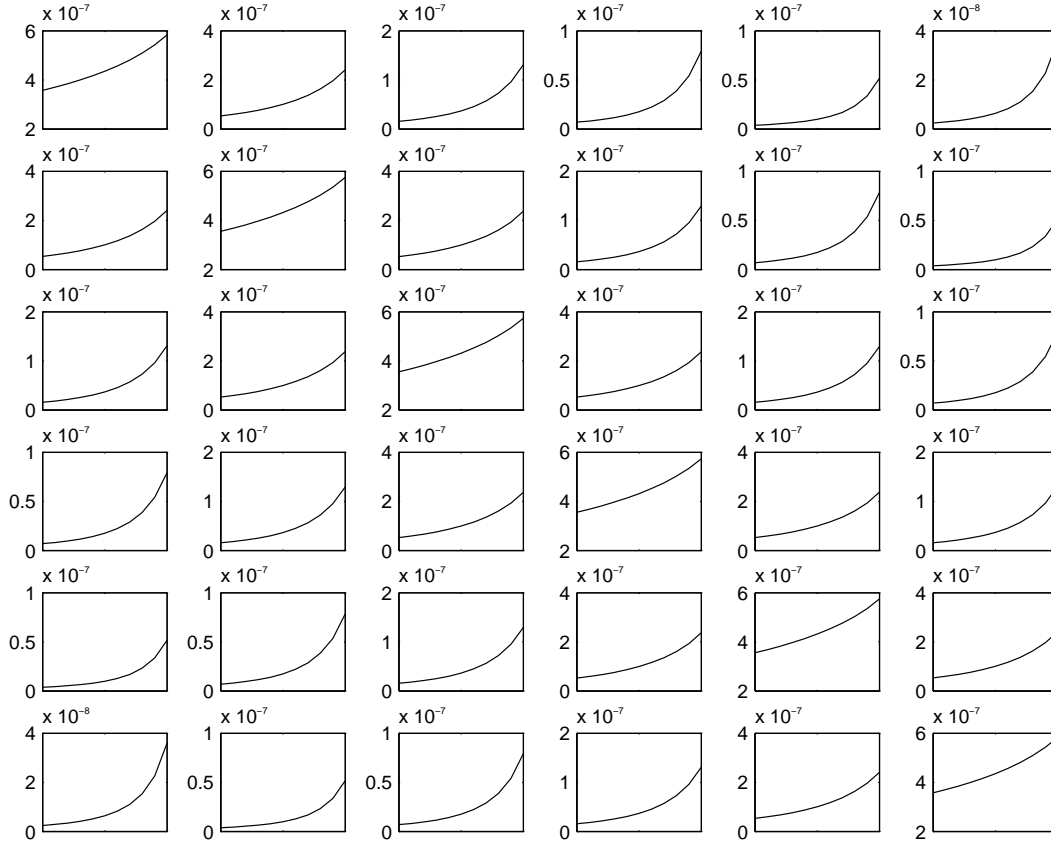
where the unknowns represent the expansion coefficients into the modified hierarchical basis functions. The load equations can now be directly included in these equations by following the same procedure of Sec. 6.3.1. The advantage is that the coefficients  $\lambda_n^0$  and  $\lambda_n^1$  are now all vanishing except one. In fact, only one trial function is nonvanishing at each border. The resulting system of ODE's is now equivalent to the one obtained with the canonical bases, and has the same order of accuracy. The structure of the system matrix  $\hat{\Theta}$  is similar to the one reported in Fig. 6.13. This is the formulation of the adaptive TDSE method that we will adopt in the applications.

### 6.3.3 Applications

#### High-speed packaging interconnect

The TDSE method is applied here to the electrical simulation of the structure depicted in Fig. 6.14. It consists of an array of six conductors providing the electrical connection between components of possibly different nature, like an electrical driver on the left and an optical interconnect module on the right. The conductors are  $20 \mu\text{m}$  thick. Their widths and separations are equal to  $1 \text{ mm}$  at the left termination and to  $0.125 \text{ mm}$  at the right termination. The substrate is  $400 \mu\text{m}$  thick, with a dielectric constant  $\epsilon_r = 4.5$ . The length of the interconnect is  $\mathcal{L} = 5 \text{ mm}$ . The per-unit-length inductance and capacitance matrices have been computed with a commercially available 2D field solver [9] based on the method of moments (MOM). The results, shown in Fig. 6.15 and 6.16 show that the per-unit-length parameters suffer of significant longitudinal variations. Therefore, a significant influence of this nonuniformity is to be expected.

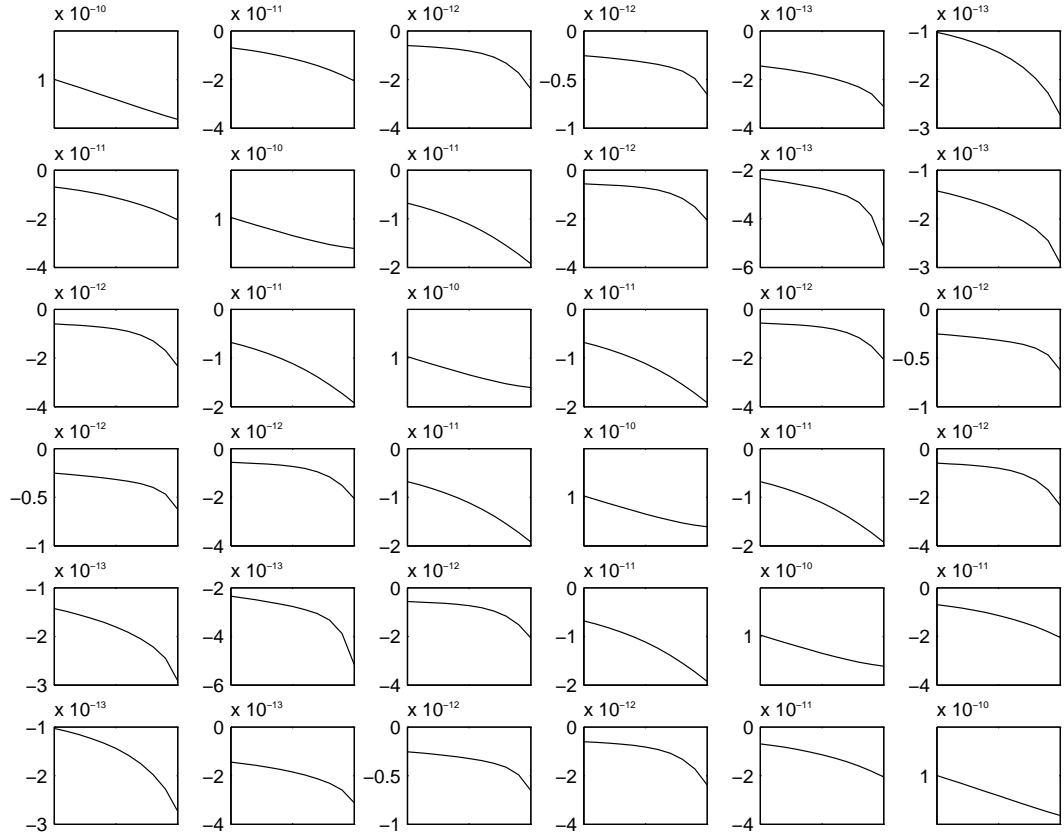
All the conductors are terminated with  $50 \Omega$  resistances, and a  $1 \text{ V}$  step voltage source with a  $20 \text{ ps}$  rise time is applied to one of the middle conductors,



**Figure 6.15:** Per-unit-length inductance matrix entries as functions of the normalized longitudinal coordinate. The units are H/m.

indexed with the subscript  $_3$ . The voltage on this and on the adjacent conductor (indexed with the subscript  $_4$ ) is computed with the TDSE method in two different situations. First, the cross-sectional parameters are evaluated in the middle of the structure and the uniform MTL model is used. Second, the cross-sectional parameters are evaluated section by section and the full NMTL model is used.

The results are plotted in Fig. 6.17, where the dashed lines refer to the uniform case and the continuous lines to the nonuniform case. It should be noted that the maximum crosstalk noise levels on the conductor 4 (bottom panels) are larger in the nonuniform than in the uniform case. This demonstrates that neglecting the nonuniformity of this interconnect in the simulation process leads to underestimate the crosstalk noise level and produces inaccurate predictions for the behavior of the structure.



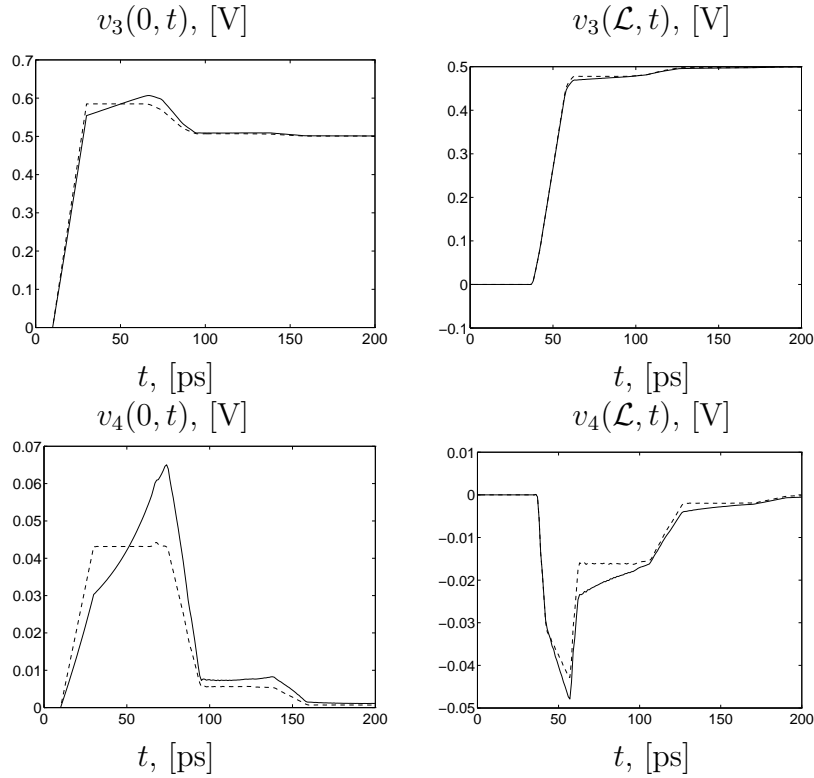
**Figure 6.16:** As in Fig. 6.15 but for the per-unit-length capacitance matrix entries in units of F/m.

### Line with nonuniform phase speed

This section shows the simulations of a line with nonuniform propagation speed. Lines of this type can be obtained when the surrounding medium presents longitudinal variation in the permittivity. We will use a slight modification of the exponential line already analyzed in Sec. 1.3.1, by fixing the per-unit-length capacitance at a constant value along the line, i.e.  $C(z) = C^0 = 1$  F/m. The per-unit-length inductance increases exponentially along the line from 1 H/m up to 4 H/m. These parameters lead to an exponentially increasing nominal characteristic impedance (from 1  $\Omega$  up to 2  $\Omega$ ) and to an exponentially decreasing nominal phase speed,

$$\nu(z) = \frac{1}{\sqrt{L(z)C}}.$$

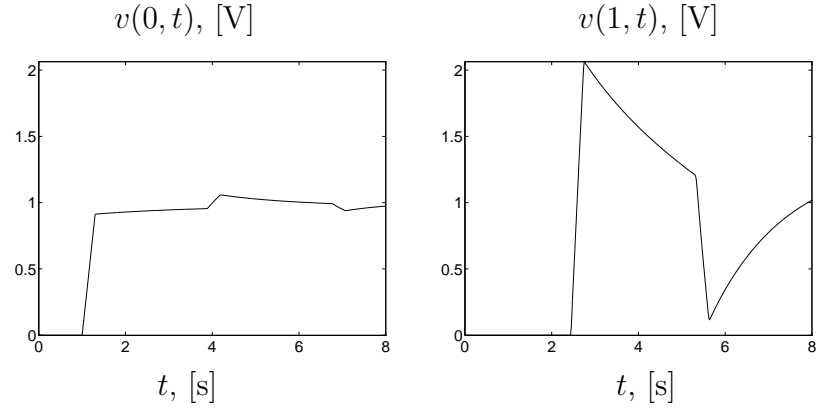
Even if this example is purely academical, it allows to show that the adaptive TDSE method is capable of accurate simulations even when the phase speed



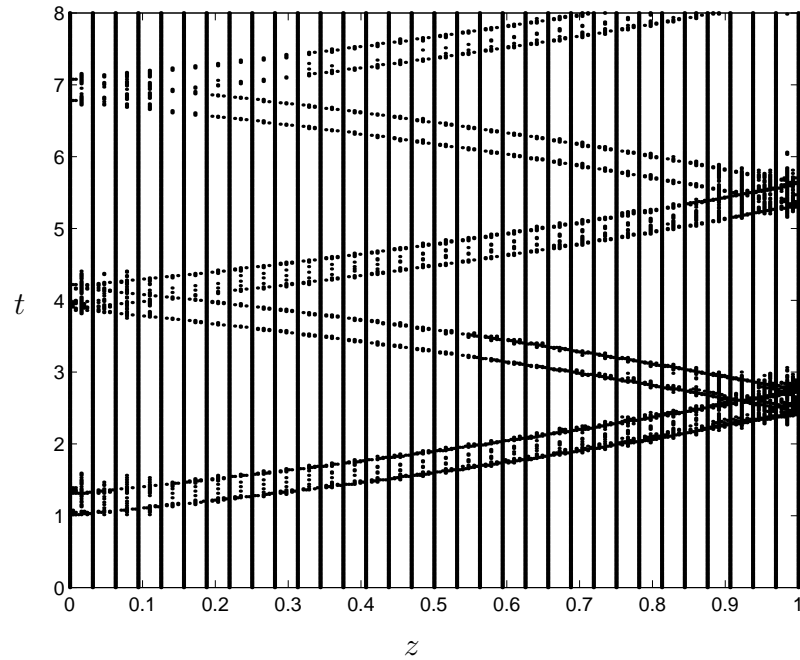
**Figure 6.17:** Voltage on the generator (top row) and receptor (bottom row) conductors of the structure in Fig. 6.14. The continuous and dashed lines indicate the voltages obtained by considering and neglecting, respectively, the longitudinal variation of the per-unit-length parameters.

is not constant. This is obviously not possible when more standard finite difference schemes are used, because these schemes strongly rely on the Courant condition [7], which requires that the time step be matched to the space discretization step through the propagation speed.

We consider a non-matched line with nominal reflection coefficients at the left and right ends equal to  $\Gamma_S = -9/11$  and  $\Gamma_L = 2/3$ , respectively. With these load conditions, the input voltage pulse undergoes significant reflections at the line ends. The voltage waveform used in the following is a 1 V step function with rise time equal to 0.3 s. The resulting voltages at the left and right terminations are plotted in the left and right panels of Fig. 6.18, respectively, while the location of the significant wavelet coefficients (using a threshold  $\varepsilon = 10^{-4}$ ) is plotted in Fig. 6.19. It should be noted that these coefficients trace the characteristic curves of the transmission line equations, tracking the location of the travelling singularities. These curves are no longer straight lines, but are significantly bended, with a tangent at a fixed  $z$  equal to  $\pm 1/\nu(z)$ .



**Figure 6.18:** Voltage at the left and right line ends of the line with nonuniform phase speed of Sec. 6.3.3.



**Figure 6.19:** Location of the significant voltage wavelet coefficients for the line with nonuniform phase speed of Sec. 6.3.3.

# Conclusions

A class of numerical schemes for the transient simulation of the Nonuniform Multiconductor Transmission Lines has been presented. The underlying method, denoted as Time-Domain Space Expansion (TDSE), is based on a weak formulation of the NMTL equations obtained through spatial expansion of the solution into some trial functions and testing of the equations with suitable basis functions. The results presented in this work show that numerical schemes of any fixed approximation order can be obtained by selecting appropriate trial and test functions.

A significant part of this work has been devoted to the use of wavelet bases in the TDSE method. Wavelets allow sparse representations of both regular and singular waveforms and can be used to build adaptive integration schemes for partial differential equations. Therefore, the solution can be obtained at a low computational cost even when the signals travelling along the transmission line present sharp singularities.

The design of such adaptive schemes has required the construction of modified wavelet bases defined on bounded domains, starting from wavelet bases defined on the real line. Indeed, as wavelet bases are intrinsically translation invariant, they cannot be used directly to represent functions defined on bounded domains. This is the case of the solution of NMTL equations. In addition, some optimized algorithms for the computation of integrals of wavelets and their derivatives have been developed. These integrals stem from the discretization of the NMTL equations through a weak formulation. The presented algorithms allow the computation of such integrals at the machine precision without use of quadrature formulas.

The main result of this work is the adaptive TDSE method, which employs adapted approximation spaces based on wavelet expansions for the solution of the NMTL equations. This method uses a time-varying sparse representation of the solution and allows to compute voltage and current along the line at a small computational cost. Several examples illustrate the advantages of the adaptive wavelet expansion with respect to the canonical non-adaptive representation. The results obtained in this work confirm that wavelets constitute a new and effective tool for the numerical solution of the equations commonly encountered in the field of Electromagnetic Compatibility.



# Appendix A

## Index of symbols

We summarize here the symbols used throughout this work. Only the symbols associated to quantities recurring in more than one section are listed here. Therefore, the symbols used in the derivations and not going beyond the scope of one section are not listed. The following tables group the symbols according to the nature of the quantities they refer to.

### NMTL description and TDSE method

$\mathbf{V}(z, t), \mathbf{I}(z, t)$	voltages and currents at location $z$ and time $t$
$\mathbf{L}(z), \mathbf{C}(z), \mathbf{G}(z), \mathbf{R}(z)$	per-unit-length matrices
$P$	number of conductors
$\mathcal{L}$	length of the NMTL
$\mathbf{V}_S(t), \mathbf{V}_L(t)$	Thévenin voltage source vectors
$\mathbf{R}_S, \mathbf{R}_L$	Thévenin load resistance matrices
$\zeta_n$	trial functions for voltage and current
$\phi_k$	expansion functions for the per unit length parameters
$\eta_m$	test functions
$N_\zeta$	number of trial and test functions $\zeta_n$ and $\eta_m$
$N_\phi$	number of expansion functions $\phi_k$
$\Psi, \widehat{\Psi}$	time derivative coefficients in the ODE system
$\Phi, \widehat{\Phi}, \Theta, \widehat{\Theta}$	linear term coefficients in the ODE system
$\Delta_S, \widehat{\Delta}_S, \Delta_L, \widehat{\Delta}_L$	source coefficients in the ODE system
$\Omega_S, \widehat{\Omega}_S, \Omega_L, \widehat{\Omega}_L$	source coefficients in the explicit ODE system
$\Delta_{SD}, \widehat{\Delta}_{SD}, \Delta_{LD}, \widehat{\Delta}_{LD}$	source derivative coefficients in the ODE system
$\mathbf{x}, \widehat{\mathbf{x}}, \overline{\mathbf{x}}, \widehat{\overline{\mathbf{x}}}$	arrays with the unknowns in the ODE system



### Multilevel decompositions: spaces and operators

$V, V(\Omega)$	generic space of functions defined on $\Omega$
$L^2(\Omega)$	space of square integrable functions defined on $\Omega$
$\ell^2$	space of square summable sequences
$H^s$	Sobolev space with regularity $s$
$C^r$	space of functions with $r$ continuous derivatives
$V_j, V_j(\Omega)$	primal approximation spaces at refinement level $j$
$\tilde{V}_j, \tilde{V}_j(\Omega)$	dual approximation spaces at refinement level $j$
$P_j, \tilde{P}_j$	primal and dual projection operators mapping $V$ to $V_j$ and $\tilde{V}_j$
$W_j, W_j(\Omega)$	primal detail spaces at refinement level $j$
$\tilde{W}_j, \tilde{W}_j(\Omega)$	dual detail spaces at refinement level $j$
$Q_j, \tilde{Q}_j$	primal and dual detail projection operators mapping $V$ to $W_j$ and $\tilde{W}_j$
$T_j$	dilation operator
$\mathcal{P}_q$	space of polynomials of degree at most $q$
$\rho_\alpha(x)$	generic basis functions for $\mathcal{P}_q$
$p_\alpha(x)$	basis of monomials for $\mathcal{P}_q$
$\mathcal{B}_{r,b}^n$	Bernstein polynomials

### Multilevel decompositions: scaling functions and wavelets

$\varphi, \tilde{\varphi}$	primal and dual scaling functions
$\varphi_{jk}, \tilde{\varphi}_{jk}$	primal and dual biorthogonal scaling functions on bounded or unbounded domains
$\varphi_{jk}^{\mathbb{R}}, \tilde{\varphi}_{jk}^{\mathbb{R}}$	primal and dual biorthogonal scaling functions on the real line
$\theta_\alpha, \tilde{\theta}_\beta$	modified primal and dual scaling functions
$\theta_{jk}, \tilde{\theta}_{jk}$	primal and dual non-biorthogonal scaling functions on bounded domains
$\psi, \tilde{\psi}$	primal and dual wavelets
$\psi_{jk}, \tilde{\psi}_{jk}$	primal and dual biorthogonal wavelets on bounded or unbounded domains
$\psi_{jk}^{\mathbb{R}}, \tilde{\psi}_{jk}^{\mathbb{R}}$	primal and dual biorthogonal wavelets on the real line
$\xi_{jk}, \tilde{\xi}_{jk}$	primal and dual non-biorthogonal wavelets on bounded domains

### Approximation functions and coefficients

$v$	generic function
$v_j$	projection $v_j = P_j v$
$d_j$	detail $d_j = P_{j+1} v - P_j v = Q_j v$
$\check{v}_{jk}, v_{jk}$	expansion coefficients of $P_j v$ into the scaling functions basis
$\hat{v}_{jk}, \tilde{v}_{jk}$	expansion coefficients of $Q_j v$ into the wavelets basis

### Multilevel decompositions: filters

$h, \tilde{h}$	primal and dual scaling function filters on $\mathbb{R}$
$g, \tilde{g}$	primal and dual wavelet filters on $\mathbb{R}$
$n_0, \tilde{n}_0$	index of the first nonzero elements of $h, \tilde{h}$
$n_1, \tilde{n}_1$	index of the last nonzero elements of $h, \tilde{h}$
$L, \tilde{L}$	number of vanishing moments of primal and dual wavelets
$\sigma, \tilde{\sigma}$	regularity of primal and dual scaling functions and wavelets
$\mathcal{H}, \mathcal{H}^\perp, \tilde{\mathcal{H}}, \tilde{\mathcal{H}}^\perp$	filters of biorthogonal scaling functions on bounded domains
$\mathcal{H}^\angle, \tilde{\mathcal{H}}^\angle$	filters of non-biorthogonal scaling functions on bounded domains
$\mathcal{G}, \mathcal{G}^\perp, \tilde{\mathcal{G}}, \tilde{\mathcal{G}}^\perp$	filters of biorthogonal wavelets on bounded domains
$\mathcal{G}^\angle, \tilde{\mathcal{G}}^\angle$	filters of non-biorthogonal wavelets on bounded domains
$\mathcal{M}, \tilde{\mathcal{M}}$	matrices expressing scaling functions on $[0, 1]$ through scaling functions on $\mathbb{R}$ restricted to $[0, 1]$ .
$j_0, j_0^w, J_0$	minimum refinement level for the multiresolution on $[0, 1]$

### Nonlinear approximations

$\mathcal{S}_\varepsilon$	set of retained coefficients in the nonlinear approximation
$P_j^\varepsilon$	nonlinear projection operator with wavelet thresholding
$S_j^I(\varepsilon)$	sparsity index of a nonlinear approximation



# Bibliography

- [1] C. E. Baum, J. B. Nitsch, and R. J. Sturm, “Analytical Solution for Uniform and Nonuniform Multiconductor Transmission Lines with Sources”, *The Review of Radio Science* 1993-96, URSI-Oxford University Press, 1996, NY.
- [2] P. Besnier, and P. Degauque, “Electromagnetic Topology: Investigations of Nonuniform Transmission Line Networks”, *IEEE Trans. Electromagn. Compat.*, vol. 37, 1995, 227-233.
- [3] C. R. Burrows, “The exponential transmission line”, *Bell System Tech. J.*, vol. 17, 1938, 555-573p.
- [4] F. Chang, “Transient Simulation of Nonuniform Coupled Lossy Transmission Lines Characterized with Frequency-Dependent Parameters-Part I: Waveform Relaxation Analysis”, *IEEE Trans. Circuits Syst. I*, vol. 39, 1992, 585-603.
- [5] T. Dhaene, L. Martens, and D. De Zutter, “Transient Simulation of Arbitrary Nonuniform Interconnection Structures Characterized by Scattering Parameters”, *IEEE Trans. Circuits Syst. I*, vol. 39, 1992, 928-937.
- [6] J. Nitsch, C. E. Baum, “Splitting of degenerate natural frequencies in coupled two-conductor lines by distance variation”, Interaction Notes, Note 477, July 1989.
- [7] C. R. Paul, *Analysis of Multiconductor Transmission Lines*, John Wiley and Sons, NY, 1994.
- [8] V. K. Tripathi, N. Orhanovic, “Time-Domain Characterization and Analysis of Dispersive Dissipative Interconnects”, *IEEE Trans. Circuits Syst. I*, vol. 39, 1992, 938-945.
- [9] A. R. Djordjevic, M. B. Bazdar, R. F. Harrington, T. K. Sarkar, LINPAR, v1.0, Academic Press, 1995.
- [10] IMSL, *IMSL MATH/LIBRARY User's Manual, Version 2.0*, 1991, IMSL, Houston.

- [11] A. Ralston, P. Rabinowitz, *A first course in numerical analysis*, McGraw-Hill, 1982.
- [12] R. A. Adams, *Sobolev Spaces*, Academic Press, 1978.
- [13] R. T. Farouki, C. A. Neff, On the numerical condition of Bernstein-Bézier subdivision processes, *Math. of Comp.*, **55**, 1990, 637-647.
- [14] R. T. Farouki, V. T. Rajan, Algorithms for polynomials in Bernstein form, *Comput. Aided Geom. Des.*, **5**, 1988, 1-26.
- [15] W. Rudin, *Functional Analysis*, McGraw-Hill, New York, 1982.
- [16] C. Canuto, A. Tabacco, Multilevel Decompositions of Functional Spaces, *J. Fourier Anal. and Appl.*, **3**, 1997, 715-742.
- [17] C. Canuto, A. Tabacco, *Ondine biortogonali: teoria ed applicazioni*, Quaderni UMI, 1998, submitted
- [18] A. S. Cavaretta, W. Dahmen, C. A. Micchelli, Stationary Subdivision, *Memoirs of the Amer. Math. Society*, **453**, Providence, Rhode Island, 1991.
- [19] C. K. Chui, *An introduction to Wavelets*, Academic Press, Boston, 1992.
- [20] A. Cohen, I. Daubechies, J. Feauveau, Biorthogonal bases of compactly supported wavelets, *Comm. Pure Appl. Math.*, **45**, 1992, 485-560.
- [21] W. Dahmen, C. Micchelli, Using the refinement equation for evaluating integrals of wavelets, *SIAM J. Num. Anal.*, **30**, 1993, 507-537.
- [22] I. Daubechies, *Ten Lectures on Wavelets*, CBMS-NSF Series in Applied Mathematics 61, SIAM, Philadelphia, 1992.
- [23] I. Daubechies, Orthonormal bases of compactly supported wavelets, *Comm. Pure Appl. Math.*, **41**, 1988, 909-996.
- [24] A. Grossmann e J. Morlet, Decompositions of Hardy functions into square integrable wavelets of constant shape, *SIAM. J. Math. Anal.*, **15**, 1984, 723-736.
- [25] L. Levaggi, A. Tabacco, *Wavelets on the Interval and Related Topics*, Rapporto interno n. 11-1997, Dipartimento di Matematica, Politecnico di Torino.
- [26] A. Haar, Zur Theorie der orthogonalen Funktionen-Systeme, *Math. Ann.*, **69**, 1910, 331-371.
- [27] S. Mallat, Multiresolution approximation and wavelet orthonormal bases of  $l^2$ , *Trans Amer. Math. Soc.*, **315**, 1989, 69-88.

- [28] Y. Meyer, *Ondelettes, vol.I: Ondelettes et Opérateurs*, Hermann, Paris, 1990.
- [29] J. Morlet, G. Arens, E. Fourgeau, D. Giard, *Geophysics*, Vol. 47, N. 2, 1982, 203-236.
- [30] E. Hernández, G. Weiss, *A First Course on Wavelets*, CRC Press, Boca Raton, 1996.
- [31] L. Anderson, H. Hall, B. Jawerth, G. Peters, Wavelets on the closed subsets of the real line, in *Topics in the Theory and Applications of Wavelets*, L. Schumacher and G. Webb (eds.), Academic Press, Boston, 1994, 1-61.
- [32] G. Chiavassa, J. Liandriat, On the effective construction of compactly supported wavelets satisfying homogeneous boundary conditions on the interval, to appear in *Appl. Comput. Harm. Anal.*.
- [33] C. K. Chui and E. Quack, Wavelets on a bounded interval, in *Numerical Methods of Approximation Theory*, D. Braess and L. L. Schumacher (eds.), Birkhäuser, 1992, 1-24.
- [34] A. Cohen, W. Dahmen, R. A. DeVore, Multiscale Decompositions on Bounded Domains, Bericht Nr. 113, IGPM-RWTH-Aachen, 1995, to appear in *Trans. Amer. Math. Soc.*.
- [35] A. Cohen, I. Daubechies, P. Vial, Wavelets on the interval and fast wavelet transform, *Appl. Comp. Harm. Anal.*, **1**, 1993, 54-81.
- [36] W. Dahmen, A. Kunoth, K. Urban, Biorthogonal Spline-Wavelets on the Interval - Stability and Moment Condition, Bericht Nr. 129, IGPM RWTH Aachen, 1996.
- [37] W. Dahmen, A. Kunoth, K. Urban, Wavelets in numerical analysis and their quantitative properties, Bericht Nr. 134, IGPM RWTH Aachen, 1997.
- [38] Y. Meyer, Ondelettes de l'intervalle, *Revista Mat. Iberoamericana*, **7**, 1992, 115-133.
- [39] S. Bertoluzza, Y. Maday, J. C. Ravel, A dynamically adaptive wavelet method for solving partial differential equations, *Comput. Meths. Appl. Mech. Engng.*, **116**, 1994.
- [40] S. Bertoluzza, G. Naldi, A wavelet collocation method for the numerical solution of partial differential equations, *Appl. Comput. Harm. Anal.*, **3**, 1996, 1-22.

- [41] C. Canuto, I. Cravero, A wavelet-based adaptive finite element method for the advection-diffusion equations, *Math. Mod. Meths. Appl. Sci.*, **7**, 1997, 265-289.
- [42] Ph. Charton, V. Perrier, A pseudo-wavelet method for solving the two-dimensional Navier-Stokes equations, *Comp. Appl. Math.*, **15**, 1996, 139-160.
- [43] W. Dahmen, A. Kunoth, K. Urban, A wavelet-Galerkin method for the Stokes equations, *Computing*, **56**, 1996, 259-302.
- [44] W. Dahmen, S. Prössdorf, R. Schneider, Wavelet approximation methods for periodic pseudodifferential equations. Part I: convergence analysis, *Math. Zeit.*, **215**, 1994, 583-620.
- [45] W. Dahmen, S. Prössdorf, R. Schneider, Wavelet approximation methods for periodic pseudodifferential equations. Part II: Fast solution and matrix compression, *Adv. Comput. Math.*, **1**, 1993, 259-335.
- [46] M. Dorobantu, *Wavelet-based Algorithms for Fast PDE Solvers*, PhD. Thesis, Stockholm University, 1995.
- [47] S. Jaffard, Wavelet methods for fast resolution of elliptic problems, *SIAM J. Numer. Anal.*, **29**, 1992, 965-986.
- [48] J. Liandrat, Ph. Tchamitchian, Resolution of the 1D Regularized Burgers Equation using a Spatial Wavelet Approximation, ICASE Report Nr. 90-83, NASA Langley, 1990.
- [49] Y. Maday, V. Perrier, J. C. Ravel, Adaptativité dynamique sur bases d'ondelettes pour l'approximation d'équations aux dérivées partielles, *C. R. Acad. Sci. Paris*, **312** Série I, 1991, 405-410.
- [50] T. von Petersdorff, C. Schwab, Fully Discrete Multiscale Galerkin BEM, Report Nr. 95-08, ETH Zürich, 1995.
- [51] T. von Petersdorff, C. Schwab, Boundary Element Methods with Wavelets and Mesh Refinement, Report Nr. 95-10, ETH Zürich, 1995.
- [52] G. Beylkin, R. Coifman, V. Rokhlin, Fast wavelet transforms and numerical algorithms, *Comm. Pure Appl. Math.*, **44**, 1991, 141-183.
- [53] C. Canuto, A. Tabacco, Absolute and relative cut-off in adaptive approximation by wavelets, Rapporto n. 5-1996, Dipartimento di Matematica, Politecnico di Torino.
- [54] R. A. DeVore, B. Jawerth, V. A. Popov, Compression of wavelet decompositions, *Amer. J. Math.*, **114**, 1992, 737-785.

- 
- [55] Y. Meyer, *Wavelets: Algorithms and Applications*, SIAM, Philadelphia, 1993.
  - [56] G. Strang, T. Nguyen, *Wavelets and Filter Banks*, Wellesley Cambridge Press, Cambridge, 1996.
  - [57] M. V. Wickerhauser, *Adapted Wavelet Analysis from Theory to Software*, AK Peters, Wellesley, 1994.