

A MULTIREOLUTION STRATEGY FOR  
HOMOGENIZATION OF PARTIAL DIFFERENTIAL  
EQUATIONS

by

NICHOLAS ASHTON COULT

B. Sc. in Mathematics

This thesis for the Doctor of Philosophy degree by  
Nicholas Ashton Coult  
has been approved for the  
Department of  
Applied Mathematics  
by

---

Gregory Beylkin

---

Stephen McCormick

Date \_\_\_\_\_

Coult, Nicholas Ashton (Ph. D., Applied Mathematics)

A Multiresolution Strategy for Homogenization of Partial Differential Equations

Thesis directed by Professor Gregory Beylkin

Coefficients of PDE's are often changing across many spatial or temporal scales, whereas we might be interested in the behavior of the solution only on some relatively coarse scale. The multiresolution strategy for reduction and homogenization provides a method for finding an equation for the projection of the solution to a coarse scale. This equation explicitly incorporates the fine-scale behavior of the coefficients.

We present the multiresolution strategy for reduction and homogenization of differential equations, and apply it to linear wave equations in which the coefficients describe a layered medium (the problem reduces to a system of ordinary differential equations in this case) and to elliptic partial differential equations. For the layered-medium wave equations, we discuss and compare the multiresolution approach with classical techniques. For elliptic operators, it is known that the non-standard form has fast off-diagonal decay and the rate of decay is controlled by the number of vanishing moments of the wavelet basis. We prove that if an appropriate (e.g. high order) basis is used, the reduction procedure preserves the rate of decay over any finite number of scales and, therefore, results in sparse matrices for computational purposes. Furthermore, the reduction procedure approximately preserves small eigenvalues of strictly elliptic operators. We also introduce a modified reduction procedure which preserves the small eigenvalues with greater accuracy than the usual y

To Katie.

## ACKNOWLEDGMENTS

I would like to give special thanks to my advisor, Gregory Beylkin, for offering me support and opportunities on all scales; Steve McCormick, for his many useful suggestions throughout my dissertation work; Lucas Monzon, for putting up with my frequent sanity checks and panicked proof-verification; David Trubatch, for his eager participation in spontaneous discussions on research topics; and the faculty, staff, and students of the Department of Applied Mathematics. Thanks also to the University of Colorado, for their partial financial support of my studies, and the National Science Foundation, which supported this work through a National Science Foundation Graduate Research Fellowship.

# CONTENTS

INTRODUCTION	1
1 MULTIREOLUTION REDUCTION AND HOMOGENIZATION	4
1.1 Classical Homogenization	5
1.1.1 Weak Limit Method	6
1.1.2 Asymptotic Method of Multiple Scales	8
1.2 Multiresolution Reduction	10
1.3 Multiresolution Homogenization of Linear ODE's	12
1.3.1 The Homogenization Procedure	13
1.3.2 An Example: MRA Homogenization of Second-Order Ordinary Differential Equations	18
1.4 Homogenization of Acoustic Wave Equations	21
1.4.1 Multiresolution Homogenization of the Acoustic Equations	23
1.4.2 Multiresolution Homogenization of a Two-Layer Medium	26
1.4.3 Comparison to Existing Theory	30
2 REDUCTION OF ELLIPTIC PARTIAL DIFFERENTIAL EQUATIONS	32
2.1 Preservation of Form Under Reduction	34
2.1.1 Reduction of Differential and Convolution Operators	34
2.1.2 Preservation of Divergence Form	37
2.2 Multiresolution Reduction of Elliptic Equations Using High Order Wavelets	40
2.2.1 Preservation of Spectral Bounds	41
2.2.2 Rate of Off-diagonal Decay and Sparsity of Reduced Operators	43
2.2.3 A Fast Method for Computing the Reduced Operator	47

2.2.4	Eigenvalues and Eigenvectors of the Reduced Operators . . . . .	49
2.2.5	Numerical Experiments . . . . .	57
3	CLASSES OF MATRICES . . . . .	64
3.1	Preliminary Considerations . . . . .	64
3.2	Matrices with Exponential Decay . . . . .	65
3.3	Matrices with Polynomial Decay . . . . .	70
4	CONCLUSIONS AND FURTHER DIRECTIONS . . . . .	86
4.1	Multiresolution Reduction of Hyperbolic and Parabolic Partial Differential Equations . . . . .	86
4.1.1	Reduction of Hyperbolic PDE's . . . . .	87
4.1.2	Reduction of Parabolic PDE's . . . . .	88
	BIBLIOGRAPHY . . . . .	90
A	WAVELETS, MULTIREOLUTION ANALYSES, AND OPERATORS . . . . .	93
A.1	Wavelets and Multiresolution Analyses . . . . .	93
A.1.1	Notation and Preliminary Considerations . . . . .	93
A.1.2	Fourier Analysis . . . . .	96
A.2	Operators in the Wavelet Basis . . . . .	99
A.2.1	Notation and Preliminary Considerations . . . . .	99
A.2.2	The Standard and Non-Standard Forms . . . . .	100

## LIST OF TABLES

2.1	Condition numbers and lower bounds for the $\mathbf{A}$ -block of the operator $-\nabla \cdot (a(x, y)\nabla)$ on the unit square with periodic boundary conditions. . . . .	57
2.2	Run times for exact versus truncated computation of $\mathbf{R}_s$ , for various scales. . .	58



## LIST OF FIGURES

2.1	$\mathbf{R}_{\mathbf{S}_j}$ after truncation of entries smaller than $0.02 * \ \mathbf{R}_{\mathbf{S}_j}\ _{\infty}$ . $\mathbf{S}_j$ is the projection of $-\nabla \cdot (a(x, y)\nabla)$ on the unit square with periodic boundary conditions into the multiwavelet basis with two vanishing moments. . . . .	48
2.2	Relative error of eigenvalues of the coarse scale operators obtained by different methods compared to the eigenvalues of the original operator. . . . .	60
2.3	Relative error of eigenvalues of the one-dimensional example operator reduced over four scales, using wavelets with four, eight, and twelve vanishing moments. . . . .	61
2.4	Relative error of eigenvalues of the one-dimensional example operator reduced via the modified reduction procedure over four scales, using wavelets with four, eight, and twelve vanishing moments. . . . .	62
2.5	Relative error of eigenvalues of the operator $-\nabla(2 + \cos(32\pi x))\nabla$ using the three techniques. Multiwavelets with two vanishing moments are used. . . . .	63
A.1	Plots on $[\pi, \pi]$ of the modulus of the low-pass filter $m_0(\xi)$ (solid line) and the high-pass filter $m_1(\xi)$ (dashed line). . . . .	99

## INTRODUCTION

The problem of capturing the influence of fine and intermediate scales on a coarse scale is generally known as that of homogenization. For example, the Earth's crust contains complicated layers of rock ranging in size from sub-millimeters to meters or more in thickness, and long-wavelength (relative to the size of the smallest layers) waves exhibit behavior on the macroscale which is influenced by the microscale structure of the layers.

The motivation for studying problems of this kind is that, in the physical world, coarse scales may be easier or less costly to influence and observe than fine and intermediate scales and, in the mathematical models, solutions on coarse scales require significantly fewer operations to compute. Additionally, the parameters of interest may only be observable on the coarse scale, but the interactions which determine the values of these parameters may occur on many scales.

The mathematical difficulty of the homogeniz



constitute an extension of the results obtained by P. Tchamitchian [37] and S. Jaffard [21].

Finally, Chapter 4 provides a summary of the results of this Thesis and describe directions for future work.

## CHAPTER 1

### MULTIRESOLUTION REDUCTION AND HOMOGENIZATION

In this chapter, we start by introducing the homogenization problem in the classical context and presenting some existing approaches to the problem. The classical approaches we consider are limited in their application to PDE's in that they require separation of fine and coarse scales and do not permit intermediate scales in the problem. We then introduce the multiresolution homogenization method, which allows for coefficients which vary on intermediate scales, and compare it to the existing classical approaches on a one-dimensional example. Since the multiresolution homogenization scheme allows one to incorporate intermediate scales and in this sense is more general than the classical approaches, the purpose of this comparison is to connect the multiresolution and classical approaches in parameter regimes where both are valid and show that classical results may be achieved using the multiresolution approach.

We demonstrate the classical techniques on a simple one-dimensional problem. In the parameter ranges where the classical methods are valid, the multiresolution approach gives the same results for the simple one-dimensional problem, but does not have the restrictions on separation of scales that the classical methods have. Thus, the result of the multiresolution method is more general in that we may apply the multiresolution technique to problems with a wider class of coefficients. We emphasize that the one-dimensional example is meant to

W

## 1.1 Classical Homogenization

Homogenization in the context of differential equations refers to the problem posed by the presence of disparate scales in coefficients of the equation and its solutions. Computationally, this presents a difficulty because the microstructure of the coefficients

homogenization problem the oscillatory component of the solution is not of interest; it is only the gross or coarse-scale features of the solution which are. The goal is to find an effective coefficient  $a_0$  so that, as  $\epsilon \rightarrow 0$ , the sequence  $u^\epsilon$  will have a limit (in some sense) given by  $u^0$  which solves the equation

$$-\frac{d}{dx}(a_0 \frac{d}{dx} u^0) = f(x). \quad (1.1.3)$$

In the following subsections, we present two typical methods for accomplishing this goal.

### 1.1.1 Weak Limit Method

One of the typical approaches (found in e.g. [8], [24]) to homogenization of (1.1.2) is to consider the w

**Proof:** We may restrict the proof to the case when  $\Omega$  is a dilation of the unit cube  $C$  with ratio  $s \geq 1$ . Note that if  $f \in L^p(C)$  and  $\epsilon \leq 1$ , then

$$\int_{\Omega} |f(\frac{x}{\epsilon})|^p dx = \epsilon^n \int_{s\epsilon^{-1}C} |f(x)|^p dx \leq \epsilon^n ([s\epsilon^{-1}] + 1)^n \langle |f|^p \rangle \leq c_0(\Omega) \langle |f|^p \rangle \quad (1.1.6)$$

where  $[s\epsilon^{-1}]$  denotes the largest integer not larger than  $s\epsilon^{-1}$ . Given  $\delta$ , choose a trigonometric polynomial  $w$  such that  $\langle w \rangle = \langle g \rangle$  and  $\langle |g - w|^p \rangle \leq \delta$ . Then for  $\epsilon \leq 1$ , we see that

$$\int_{\Omega} |g(\frac{x}{\epsilon}) - w(\frac{x}{\epsilon})|^p dx \leq c_0(\Omega) \delta. \quad (1.1.7)$$

Of course, for the trigonometric polynomial  $w$ , Lemma 1.1.1 is true by the Riemann-Lebesgue Lemma. The estimate (1.1.7) shows that we may then extend the result to the function  $g$  since  $g$  is arbitrarily well-approximated by a trigonometric polynomial.

From (1.1.2), we extract the relation

$$-\int_0^1 \frac{du^\epsilon}{dx} dx = \int_0^1 a(\frac{x}{\epsilon})^{-1} (F(x) - c_\epsilon) dx = 0 \quad (1.1.8)$$

where  $F(x) = \int_0^x f(x) dx$  and  $c_\epsilon$  is constant in  $x$ . By Lemma 1.1.1, we determine that  $\lim_{\epsilon \rightarrow 0} c_\epsilon = \int_0^1 F(x) dx$ . Noting that

$$p^\epsilon(x) = -a(\frac{x}{\epsilon}) \frac{du^\epsilon}{dx} = F(x) - c_\epsilon \quad (1.1.9)$$

and

$$\frac{du^\epsilon}{dx} = -a(\frac{x}{\epsilon})^{-1} (F(x) - c_\epsilon), \quad (1.1.10)$$

we compute the weak limits of  $p^\epsilon$  and  $\frac{du^\epsilon}{dx}$  by



and, therefore,  $u^0$  solves the equation

$$-\frac{d}{dx}(\langle a^{-1} \rangle^{-1} \frac{d}{dx} u^0) = f(x). \quad (1.1.14)$$

The quantity  $\langle a^{-1} \rangle^{-1}$  is known as the harmonic mean (see e.g. [3]) of the function  $a$ . This quantity shows up quite often in everyday calculations. For example, if a particle is traveling in the positive direction on the real line with a velocity which depends on its position only, i.e.  $v = v(x) > 0$ , then its total travel time is given by  $\langle v^{-1} \rangle = \int_0^1 (v(x))^{-1} dx$ , and we compute its average velocity over this interval as  $\langle v^{-1} \rangle^{-1}$ , the spatial harmonic mean of its velocity.

To summarize, we use a weak limit analysis to obtain the homogenized equation (1.1.14) as an equation on the coarse scale. This equation is a constant-coefficient equation and its coefficients are found by computing the harmonic mean of the function  $a(x)$ .

### 1.1.2 Asymptotic Method of Multiple Scales

The method of multiple scales (see e.g. [7] for details) presents an alternative approach to the homogenization problem, and for (1.1.2) produces the same result as the weak limit method of the previous subsection.

The first step of the multiple-scales method is to identify the variables associated with different powers of the small parameter  $\epsilon$  and treat them as independent variables. There are two scales present in the problem (1.1.2), identified with  $x$  and  $\frac{x}{\epsilon}$ . We consider the variable  $x$  as is and define a new variable  $y = \frac{x}{\epsilon}$  which represents the fine-scale variable. We look for solutions of (1.1.2) of the form

$$u^\epsilon(x, y) = u^0(x) + \epsilon u^1(x, y). \quad (1.1.15)$$

We see that  $\frac{d}{dx} = \frac{\partial}{\partial x} + \frac{1}{\epsilon} \frac{\partial}{\partial y}$ . Thus,

$$\begin{aligned} \frac{d}{dx} \left( a(y) \frac{d}{dx} u^\epsilon(x) \right) &= \epsilon^{-1} \left( \frac{\partial}{\partial y} \left( a(y) \left( \frac{\partial}{\partial x} u^0 \right) \right) + \frac{\partial}{\partial y} \left( a(y) \left( \frac{\partial}{\partial y} u^1 \right) \right) \right) \\ &+ a(y) \left( \frac{\partial}{\partial x} \right)^2 u^0 + \frac{\partial}{\partial x} \left( a(y) \left( \frac{\partial}{\partial y} u^1 \right) \right) + \frac{\partial}{\partial y} \left( a(y) \left( \frac{\partial}{\partial x} u^1 \right) \right) \\ &+ \epsilon \left( a(y) \left( \frac{\partial}{\partial x} \right)^2 u^1 \right) \end{aligned} \quad (1.1.16)$$

$$= f(x). \quad (1.1.17)$$

We now collect terms in powers of  $\epsilon$

## 1.2 Multiresolution Reduction

In contrast to the classical approach to the homogenization problem, the multiresolution approach uses the algebraic transformation between scales provided by the multiresolution analysis to solve for the fine-scale behavior and explicitly eliminate it from the equation. This approach has the advantage that the coefficients may vary on arbitrarily many scales. The chain of subspaces

$$\dots \subset \mathbf{V}_2 \subset \mathbf{V}_1 \subset \mathbf{V}_0 \subset \mathbf{V}_{-1} \subset \mathbf{V}_{-2} \subset \dots \quad (1.2.1)$$

defines the hierarchy of scales that the multiresolution scheme uses. This chain of subspaces is defined in such a way that the space  $\mathbf{V}_j$  is “finer” than the space  $\mathbf{V}_{j+1}$ , in the sense that (1) all of  $\mathbf{V}_{j+1}$  is contained in  $\mathbf{V}_j$ , and (2) the component of  $\mathbf{V}_j$  which is not in  $\mathbf{V}_{j+1}$  consists of functions which resolve features on a scale finer than any function in  $\mathbf{V}_{j+1}$  may resolve. The difference between successive spaces in this chain is captured by the so-called wavelet space  $\mathbf{W}_{j+1}$ , defined to be the orthogonal complement of  $\mathbf{V}_{j+1}$  in  $\mathbf{V}_j$ . An orthogonal basis for the wavelet space  $\mathbf{W}_{j+1}$  is constructed which has vanishing moments, i.e. the basis elements are  $L^2$ -orthogonal to low-degree polynomials (see Appendix A for details). The existence of orthogonal wavelet bases with vanishing moments distinguishes the multiresolution approach from typical multi-scale discretizations provided by finite-element or hierarchical bases (see [5] for a description). If we are considering a multiresolution analysis defined on a bounded domain, then the hierarchy of scales defined by (1.2.1) has a coarsest scale (which we may call  $\mathbf{V}_0$ ), and we write instead

$$\mathbf{V}_0 \subset \mathbf{V}_{-1} \subset \mathbf{V}_{-2} \subset \dots \quad (1.2.2)$$

For more details, see Appendix A.

The multiresolution strategy for the reduction and homogenization of linear problems has been proposed in [10]. Let us briefly review here the reduction procedure (in its general form). Consider a bounded linear operator  $\mathbf{S}_j : \mathbf{V}_j \rightarrow \mathbf{V}_j$ . Since  $\mathbf{V}_j$  is spanned by translations of the function  $\phi(2^j x - k)$ , we know that the operator  $\mathbf{S}_j$  may be written as a matrix. If the multiresolution analysis is defined on a bounded domain, then this matrix is finite; otherwise

it is an infinite matrix which we consider as an operator on  $l^2$ . Let us consider the equation

$$\mathbf{S}_j x = f. \quad (1.2.3)$$

The decomposition  $\mathbf{V}_j = \mathbf{V}_{j+1} \oplus \mathbf{W}_{j+1}$  allows us to split the operator  $\mathbf{S}_j$  into four pieces (recall that  $\mathbf{W}_{j+1}$  is called the wavelet space and is the “detail” or fine-scale component of  $\mathbf{V}_j$ ) and write

$$\begin{pmatrix} \mathbf{A}_{\mathbf{S}_j} & \mathbf{B}_{\mathbf{S}_j} \\ \mathbf{C}_{\mathbf{S}_j} & \mathbf{T}_{\mathbf{S}_j} \end{pmatrix} \begin{pmatrix} d_x \\ s_x \end{pmatrix} = \begin{pmatrix} d_f \\ s_f \end{pmatrix}, \quad (1.2.4)$$

where we have

$$\mathbf{A}_{\mathbf{S}_j} : \mathbf{W}_{j+1} \rightarrow \mathbf{W}_{j+1} \quad (1.2.5)$$

$$\mathbf{B}_{\mathbf{S}_j} : \mathbf{V}_{j+1} \rightarrow \mathbf{W}_{j+1} \quad (1.2.6)$$

$$\mathbf{C}_{\mathbf{S}_j} : \mathbf{W}_{j+1} \rightarrow \mathbf{V}_{j+1} \quad (1.2.7)$$

$$\mathbf{T}_{\mathbf{S}_j} : \mathbf{V}_{j+1} \rightarrow \mathbf{V}_{j+1}, \quad (1.2.8)$$

and  $d_x, d_f \in \mathbf{W}_{j+1}$ ,  $s_x, s_f \in \mathbf{V}_{j+1}$  are the  $L^2$ -orthogonal projections of  $x$  and  $f$  onto the  $\mathbf{W}_{j+1}$  and  $\mathbf{V}_{j+1}$  spaces. The projection  $s_x$  is thus the coarse-scale component of the solution  $x$ , and  $d_x$  is its fine-scale com

Once we have obtained the reduced equation, it may formally be reduced again to produce an equation on  $\mathbf{V}_{j+2}$ , and the solution of this equation is just the  $\mathbf{V}_{j+2}$ -component of the solution of (1).

of equations.

In this Section we describe the MRA homogenization procedure of [10] as applied to linear ODE's, and give an example of the the procedure applied to the one-dimensional version of (1.3.1).

### **1.3.1 The Homogenization Procedure**

The MRA homogeniza i



where  $\delta_j = 2^{-j}$ ,  $\mathbf{I}$  is the  $n \times n$  identity matrix, and  $(B_j)_i$  and  $(A_j)_i$  are the  $i$ -th Haar coefficients on scale  $V_j$  of the  $n \times n$  matrix-valued functions  $B(x)$  and  $A(x)$ .

For equation (1.3.3), the recursion relations are given by

$$(A_{k+1}^{(j)})_i = (S_A)_i - (D_A)_i F^{-1}((D_B)_i + \frac{\delta_k}{2}(S_A)_i) \quad (1.3.12)$$

$$(B_{k+1}^{(j)})_i = (S_B)_i - \frac{\delta_k}{2}(D_A)_i - ((D_B)_i - \frac{\delta_k}{2}(S_A)_i) F^{-1}((D_B)_i + \frac{\delta_k}{2}(S_A)_i) \quad (1.3.13)$$

$$(p_{k+1}^{(j)})_i = (S_p)_i - \frac{\delta_k}{2}(D_A)_i F^{-1}((D_q)_i + (S_p)_i) \quad (1.3.14)$$

$$(q_{k+1}^{(j)})_i = (S_q)_i - \frac{\delta_k}{2}(D_p)_i - \frac{\delta_k}{2}((D_B)_i - (S_A)_i) F^{-1}((D_q)_i + (S_p)_i), \quad (1.3.15)$$

where

$$(S_A)_i = \frac{1}{2} \left( (A_k^{(j)})_{2i} + (A_k^{(j)})_{2i+1} \right) \quad (1.3.16)$$

$$(D_A)_i = \frac{1}{2} \left( (A_k^{(j)})_{2i} - (A_k^{(j)})_{2i+1} \right) \quad (1.3.17)$$

$$(S_B)_i = \frac{1}{2} \left( (B_k^{(j)})_{2i} + (B_k^{(j)})_{2i+1} \right) \quad (1.3.1)$$



(1.3.25) as  $j \rightarrow -\infty$ :

$$\mathbf{B}_0^{(-\infty)} x_0^{(-\infty)} + q_0^{(-\infty)} + \lambda = \mathbf{K}_0(\mathbf{A}_0^{(-\infty)} x_0^{(-\infty)} + p_0^{(-\infty)}). \quad (1.3.26)$$

This amounts to eliminating infinitely many fine scales from the equation. We call the matrices  $\mathbf{B}_0^{(-\infty)}$  and  $\mathbf{A}_0^{(-\infty)}$  the reduced coefficients of the equation (1.3.3).

We then look for the operators and forcing terms  $B^h(t)$ ,  $A^h(t)$ ,  $q^h(t)$ , and  $p^h(t)$  with certain desired qualities (e.g. constant values) such that the equation,

$$(I + B^h(t))x(t) + q^h(t) + \lambda = \int_0^t (A^h(s)x(s) + p^h(s))ds, \quad t \in (0, 1), \quad (1.3.27)$$

when subject to the same reduction and limit procedure as (1.3.3), yields on  $\mathbf{V}_0$  the same equation as in (1.3.26).

For (1.3.3), we usually require that  $A^h$ ,  $B^h$ ,  $p^h$ , and  $q^h$  be constant. The results of homogenization in this case are summarized as follows:

**Proposition 1.3.1** *Given the equation (1.3.3), if the limits which determine the matrices  $\mathbf{B}_0^{(-\infty)}$  and  $\mathbf{A}_0^{(-\infty)}$  exist, then there exist constant matrices  $B^h$ ,  $A^h$  and forcing terms  $p^h$ ,  $q^h$ , such that the reduced coefficients and forcing terms of (1.3.27) are given by  $\mathbf{B}_0^{(-\infty)}$ ,  $\mathbf{A}_0^{(-\infty)}$ ,  $p_0^{(-\infty)}$ ,  $q_0^{(-\infty)}$ . The homogenized coefficients  $B^h$  and  $A^h$  and forcing terms  $p^h$  and  $q^h$  are defined by*

$$A^h = A_0^{(-\infty)} \quad (1.3.28)$$

$$B^h = A^h \tilde{A}^{-1} - \mathbf{I} \quad (1.3.29)$$

$$p^h = p_0^{(-\infty)} \quad (1.3.30)$$

$$q^h = q_0^{(-\infty)} + (\mathbf{I} - \frac{1}{2}\tilde{A} - \tilde{A}(\exp(\tilde{A} - \mathbf{I})^{-1}A^h)^{-1})p^h, \quad (1.3.31)$$

where

$$\tilde{A} = \log(\mathbf{I} + \left(\mathbf{I} + B_0^{(-\infty)} - \frac{1}{2}A^h\right)^{-1}A^h). \quad (1.3.3)$$

Likewise, the recurrence relations for the forcing terms simplify to

$$p_{k+1}^h = p_k^h \tag{1.3.35}$$

$$q_{k+1}^h = q_k^h - \frac{\delta_k}{2} A_k^h (\mathbf{I} + B_k^h)^{-1} p_k^h. \tag{1.3.36}$$

Since the term  $A^h$  is unchanged by reduction, it is clear that  $A^h = A_0^{(-\infty)}$ . Similarly  $p^h$  is unchanged by reduction, so  $p^h = p_0^{(-\infty)}$ . The situation for  $B^h$  and  $q^h$  is more complicated. We solve for them analytically using the solution of (1.3.27).

Consider the case  $p_0^{(-\infty)} = 0$ . Clearly, then, it is the case that  $q^h = q_0^{(-\infty)}$ . The solution of (1.3.27) is therefore given by

$$x(t) = \sum_{k=0}^{\infty} A_k^h \frac{\delta_k}{2} p_k^h$$

As formulated above, the multiresolution approach to homogenization requires the computation of  $A_0^{(-\infty)}$  and  $B_0^{(-\infty)}$ , i.e. a limit over infinitely many scales. In practice, the multiresolution reduction algorithm is applied numerically over only finitely many scales. The typical practice is to compute successive  $A_0^{(-J)}$  and  $B_0^{(-J)}$  terms until finer approximations vary by less than some specified tolerance, and use these matrices as approximations to  $A_0^{(-\infty)}$  and  $B_0^{(-\infty)}$ .

Besides establishing the general framework for multiresolution reduction and homogenization, it is observed in [10] that, for systems of linear ordinary differential equations, using the Haar basis (or a multiwavelet basis) provides a

mogenization strategy of [10]. Equation (1.3.40) may be written as a first

As a first-order system of ordinary differential equations, the homogenized equation yields

$$\begin{cases} \frac{d}{dx}u(x) = (M_1 - 2M_2)v(x) \\ \frac{d}{dx}v(x) = f^h \end{cases} . \quad (1.3.46)$$

The result is somewhat different than the classical result. We trace this difference to the fact that the multiresolution homogenization procedure allows the coefficients  $a(x)$  to vary on arbitrarily many scales, whereas the classical approach in Section 1.1 allows only for coefficients of the form  $a(x/\epsilon)$ . In the multiresolution context this amounts to restricting the coefficients to an asymptotically fine scale. Let us apply

the reduction procedure, and the order of the wavelets has important implications for the eigenvalue problem as well. We study these issues in Chapter 2.

#### 1.4 Homogenization of Acoustic Wave Equations

The linear acoustic equations in a layered medium present an ideal model problem for comparison. These equations have been studied extensively for more than fifty years, allowing the multiresolution method for homogenization to be firmly placed in context and compared with existing approaches.

These equations essentially form a subset of the full linear equations of small motions of an elastic solid, in the sense that the wave solutions of the acoustic equations may be thought of as P-wave solutions of the elasticity equations. A wide variety of phenomena can be observed in the motion of waves through layered media, including dispersion, attenuation, and long-wavelength anisotropy. For a complete discussion of the existing approaches to these phenomena, see e.g. [11], [4], [29], [30], and references therein.

The general linear acoustic equations are given by

$$\rho(\mathbf{x})v_t(\mathbf{x}, t) + \nabla p(\mathbf{x}, t) = 0 \tag{1.4.1}$$

$$\frac{1}{\rho(\mathbf{x})c(\mathbf{x})}p_t(\mathbf{x}, t) + \nabla \cdot v(\mathbf{x}, t) = 0,$$

where  $p(\mathbf{x}, t)$  is pressure,  $v(\mathbf{x}, t)$  is particle velocity,  $\rho(\mathbf{x})$  is density, and  $c(\mathbf{x})$  is sound speed. Assume that the medium varies only in the vertical direction, i.e. that  $\rho(\mathbf{x}) = \rho(x_3)$ ,  $c(\mathbf{x}) = c(x_3)$ . We write (1.4.1) in component form:

$$\rho(x_3) \frac{\partial}{\partial t} v^{(1)}(\mathbf{x}, t) + \nabla_{\mathbf{x}} \cdot \mathbf{I} \mathbf{x} = 0$$

To reduce (1.4.2)

propagating with the same angle  $\theta$  but in the opposite direction. For plane waves propagating with vertical incidence, we set  $\eta = 0$ . In the case of a non-constant medium, such solutions are referred to (see e.g. [29]) as “pseudo plane-wave” solutions, the idea being that in a layered medium, the solution forms a plane wave locally inside each layer.

Now that we have introduced the notation and written the acoustic equations as a system of ODE’s, we apply the multiresolution homogenization technique to this system and compare the results to existing approaches.

#### **1.4.1 Multiresolution Homogenization of the Acoustic Equations**

The homogenization procedure of [1



where  $K^h(\omega, \eta) = (\mathbf{I} + B^h(\omega, \eta))^{-1} A^h(\omega, \eta) = \tilde{\phantom{K}}$

**Proposition 1.4.1** *The matrix  $K^h(\omega, \eta)$  which results from multiresolution homogenization applied to (1.4.4) has eigenvalues  $\lambda_1, \lambda_2$  such that  $\lambda_1 = \bar{\lambda}_2$ .*

**Pro**



To compute  $K^h$ , we use the formula

$$\exp(K^h) = \mathbf{I} + \left( \mathbf{I} + B^h \right)$$

This series of steps is summarized in the following theorem:

**Theorem 1.4.1** *The homogenized matrix  $K^h$  for the two-layer wave equation given by (1.4.14) is defined by*

$$K^h = \log \left( \mathbf{I} + \left( \mathbf{I} + B_0^{(-\infty)} - \frac{1}{2}A_0^{(-\infty)} \right)^{-1} A_0^{-\infty} \right), \quad (1.4.34)$$

where  $B_0^{(-\infty)}$  and  $A_0^{(-\infty)}$  are defined by (1.4.27)-(1.4.33).

Thus, equations (1.4.27)-(1.4.33), together with equation (1.4.19), completely define the steps necessary to compute an exact formula for  $K^h$ . For any given equation, the matrix  $K^h$  may be a very complicated expression. However, the purpose of deriving this exact formula is not for computation but rather to demonstrate that the multiresolution approach can be used to compute analytic rather than numeric results, if desired. Towards this end, we will also compute the first few terms in the small  $\omega$  expansion of  $K^h$  in a more general form in terms of the matrices  $M_0$  and  $M_1$ .

First, we compute  $(B_1^{(-\infty)})_j$  and  $(A_1^{(-\infty)})_j$  in terms of  $M_0$  and  $M_1$ . We see that if  $A = i\omega M$ , then

$$\exp \left( \frac{1}{2}A \right) - \mathbf{I} = \exp \left( \frac{i\omega}{2}M \right) - \mathbf{I} \quad (1.4.35)$$

$$= \frac{i\omega}{2}M - \frac{\omega^2}{8}M^2 - i\frac{\omega^3}{48}M^3 + \mathcal{O}(\omega^4) \quad (1.4.36)$$

$$= \frac{i\omega}{2}M \left( \mathbf{I} + i\frac{\omega}{4}M - \frac{\omega^2}{24} + \mathcal{O}(\omega^3) \right). \quad (1.4.37)$$

Thus,

$$\left( \exp \left( \frac{1}{2}A \right) - \mathbf{I} \right)^{-1} = \left( \mathbf{I} + i\frac{\omega}{4}M - \frac{\omega^2}{24}M^2 \right)^{-1}$$



This leads to

$$S'_A = \omega \bar{S}_A = \omega \left( \frac{1}{\omega} \right)$$

see, however, that the leading order behavior in  $\omega$  of  $K^h$  may be found simply by averaging  $M_0$  and  $M_1$ .

Also, since the multiresolution homogenization procedure gives the same result for all values of the initial conditions, we may extend the exact-homogenization formula for the two-layer problem to a periodically layered medium with two distinct, repeated layers of equal size. In the case where we have many non-periodic layers of equal size, the leading order behavior in  $\omega$  of  $K^h$  may be found simply by averaging  $\hat{M}$  over all the layers. This can be seen by starting with the two-layer asymptotic expansion and applying the formula recursively.

#### 1.4

technique, will give the same results for a periodic layered medium in one dimension. In this context, the multiresolution approach may be viewed as an alternative to computing the exact solution directly.

As far as we are aware, the propagator matrix technique has not been generalized to fully two or three dimensional problems. The multiresolution approach, on the other hand, does generalize to higher dimensional problems. This generalization is the subject of study in Chapter 2.



## CHAPTER 2

REDUCT

- a symmetric partial differential operator of the form  $-\nabla(a(\mathbf{x})\nabla)$ , where  $a(\mathbf{x}) > 0$

or, more generally,

- a symmetric second-order elliptic pseudodiff

elliptic operators in one and two dimensions.

## 2.1 Preservation of Form Under Reduction

As we demonstrated in Chapter 1, the multiresolution strategy for homogenization of ODE's preserves an explicit form of the equations under the reduction procedure, provided a basis with non-overlapping supports is used for the reduction step.

In this section we consider second-order elliptic operators, and in particular operators of the form

$$-\nabla \cdot (a(\mathbf{x})\nabla) \tag{2.1.1}$$

on  $L^2(\mathbf{R}^1)$  and  $L^2(\mathbf{R}^2)$ . We would like to find out whether the form of the matrix  $\mathbf{S}_j$ , which is the projection of the operator (2.1.1) into a multiresolution analysis, is preserved under reduction. This would provide us with an analogy to the recursion relation for the coefficients in the multiresolution homogenization scheme for ODE's.

### 2.1.1 Reduction of Differential and Convolution Operators

For matrices  $\mathbf{S}_j$  which are projections of operators on  $L^2(\mathbf{R})$ , the reduction procedure simplifies greatly if the matrix  $\mathbf{S}_j$  is a convolution matrix. We show that if  $\mathbf{S}_j$  is a convolution matrix which approximates the second derivative to second order, then  $\mathbf{S}_j$  is a convolution

as multiplication by a function in the Fourier domain. Furthermore, if a convolution has an inverse then we may represent application of its inverse as division by a function in the Fourier domain.

The  $\mathbf{A}_{\mathfrak{S}_j}$ ,  $\mathbf{B}_{\mathfrak{S}_j}$ ,  $\mathbf{C}_{\mathfrak{S}_j}$ , and  $\mathbf{T}_{\mathfrak{S}_j}$  blocks are also

for all  $\xi \in [-\pi, \pi]$ .

Additionally, we see that if  $\mathbf{S}_j$  is symmetric and positive, then  $s_j(\xi)$  is real and non-negative. For the second-derivative operator, then, in order to ensure the existence of  $\mathbf{A}_{\mathbf{S}_j}^{-1}$ , we construct  $\mathbf{S}_j$  so that it is symmetric and positive, and 0 is an element of the spectrum of  $\mathbf{S}_j$  only if  $s_j(\xi)$  is non-zero everywhere except at  $\xi = 0$ .

Now, assume that we have constructed  $\mathbf{S}_j$  as a discretization of  $-\frac{d^2}{dx^2}$  in such a way that we assure the existence of  $\mathbf{A}_{\mathbf{S}_j}^{-1}$  (as above). We use properties of  $m_0$  and  $m_1$  (see Appendix A) to arrive at

$$r_j(2\xi) = t_j(2\xi) - \frac{b_j(2\xi)c_j(2\xi)}{a_j(2\xi)} = \dots = \frac{s_j(\xi)s_j(\xi + \pi)}{|m_1(\xi)|^2 s_j(\xi) + |m_0(\xi)|^2 s(\xi + \pi)}. \quad (2.1.8)$$

(In fact the above equation holds for any  $s_j(\xi)$  as long as  $\mathbf{A}_{\mathbf{S}_j}^{-1}$  exists.)

We summarize our results as

**Proposition 2.1.1** *Let  $\mathbf{V}_j$  and  $\mathbf{V}_{j+1}$  be successive subspaces of a multiresolution analysis with  $m - 1$  vanishing moments. Suppose  $\mathbf{S}_j$  is a symmetric and positive convolution matrix  $\mathfrak{B}$*

from which we deduce

$$(\mathbf{R}_{\mathfrak{S}_j})_{0,-k} = \frac{1}{2\pi} \int_{-\pi}^{\pi} r_j(-\xi) \cos(-k\xi) d\xi = \frac{1}{2\pi} \int_{-\pi}^{\pi} r_j(\xi) \cos(k\xi) d\xi = (\mathbf{R}_{\mathfrak{S}_j})_{0,k}. \quad (2.1.12)$$

Therefore,  $\mathbf{R}_{\mathfrak{S}_j}$  is a symmetric convolution operator. Using (2.1.8), we see that  $r_j(\xi) \neq 0$  except at  $\xi = 0$ .

Furthermore, if we are using wavelets with  $m - 1$  vanishing moments, then for small  $\xi$  we know that

$$|m_0(\xi)|^2 \sim 1 + \mathcal{O}(\xi^{2m}) \quad (2.1.13)$$

and

$$|m_1(\xi)|^2 \sim \mathcal{O}(\xi^{2m}). \quad (2.1.14)$$

We can rewrite (2.1.8) as

$$r_j(2\xi) = \frac{s_j(\xi)}{|m_1(\xi)|^2 \frac{s_j(\xi)}{s_j(\xi+\pi)} + |m_0(\xi)|^2}.$$

Assuming  $\xi \ll 1$  gives us

$$r_j(2\xi) \sim \frac{s_j(\xi)}{\mathcal{O}(\xi^{2m})(\mathcal{O}(\xi^2) + \mathcal{O}(\xi^q)) + 1 + \mathcal{O}(\xi^{2m})} \sim s_j(\xi) + \mathcal{O}(\xi^q) \quad (2.1.15)$$

which yields

$$r_j(\xi) \sim s_j\left(\frac{\xi}{2}\right)$$

In one dimension, this operator may be written in the Haar basis as

$$\mathbf{S}_j$$

The analysis in [17] does not make it clear if in general the homogenized coefficient matrix  $H$  has a form more particular than that given by (2.1.22). Numerical experiments in [17] indicate that  $H$  has a faster rate of decay in the entries away from the diagonal than the matrix  $\mathbf{R}_{\mathbf{S}_j}$ , which may have quite slow decay. When  $H$  has this property, it may be truncated and compressed more effectively than  $\mathbf{R}_{\mathbf{S}_j}$ . Additionally, since the matrix  $H$  is not diagonal, it does not appear that the form given by (2.1.21) may be used to find a recursion relation for  $H$  across many scales.

However, we may say more about  $H$  if we make some restrictions on the coefficients, as is done in [17]. The analysis of [17] yields:

**Proposition 2.1.2** *Suppose  $a(x) = \bar{a} + \tilde{a}(x)$ , where  $\tilde{a}(x) \in \mathbf{W}_{j+1}$ , the finest wavelet space in the domain of  $\mathbf{S}_j$ , and  $\tilde{a}(x)$  has constant amplitude such that  $|\tilde{a}(x)| < \bar{a}$ . Then, for any function  $v(x) \in L^2([0, 1])$  such that  $v(x)$  has a continuous and bounded fourth derivative, we have*

$$\|\mathbf{R}_{\mathbf{S}_j}(\mathbf{P}_j v(x)) - \alpha \frac{1}{h_j^2} \Delta_+ \Delta_- (\mathbf{P}_j v(x))\|_\infty \leq C h_j^2 \|v^{(4)}(x)\|_\infty, \quad (2.1.24)$$

where  $\alpha = \langle a^{-1} \rangle^{-1}$  is the harmonic average of  $a(x)$  on  $[0, 1]$ .

We do not prove this proposition here; for the complete proof, see [17].

The basic idea of this proposition is that, for highly oscillatory coefficients which are resolved only by the finest scale of  $\mathbf{V}_j$ , the reduction procedure applied to projections of operators of the form  $-\frac{d}{dx}(a(x)\frac{d}{dx})$  yields the same result as classical homogenization, plus a small perturbation term. (Given the results of Section 1.3, this is not unexpected.) We note that this is the case even if the oscillatory part of  $a(x)$  is set to zero (i.e.  $a(x)$  is a constant).

This corresponds with our analysis in [17].



*y* directions

appear to present som

Since  $\mathbf{S}_j$  is positive definite, so is  $\begin{pmatrix} \mathbf{A}_{\mathbf{S}_j} & \mathbf{B}_{\mathbf{S}_j} \\ \mathbf{B}_{\mathbf{S}_j}^* & \mathbf{T}_{\mathbf{S}_j} \end{pmatrix}$  and, thus, it follows that  $\mathbf{A}_{\mathbf{S}_j}$  is positive definite and  $\mathbf{A}_{\mathbf{S}_j}^{-1}$  exists. Let us consider the operator

$$\mathbf{Z} = \begin{pmatrix} \mathbf{I} & -\mathbf{A}_{\mathbf{S}_j}^{-1}\mathbf{B}_{\mathbf{S}_j} \\ 0 & \mathbf{I} \end{pmatrix}.$$

Then we have

$$\mathbf{Z}^* \begin{pmatrix} \mathbf{A}_{\mathbf{S}_j} & \mathbf{B}_{\mathbf{S}_j} \\ \mathbf{B}_{\mathbf{S}_j}^* & \mathbf{T}_{\mathbf{S}_j} \end{pmatrix} \mathbf{Z} = \begin{pmatrix} \mathbf{A}_{\mathbf{S}_j} & 0 \\ 0 & \mathbf{R}_{\mathbf{S}_j} \end{pmatrix},$$

and

$$(\mathbf{R}_{\mathbf{S}_j} x, x) = \begin{pmatrix} x & \nabla \\ \mathbf{B}_j & \mathbf{R} \end{pmatrix}$$

# R

The estimate of (2.2.3) raises the important question of whether it is possible (and under which conditions) to have exactly or approximately the lower eigenvalues of  $\mathbf{S}_j$  as eigenvalues of  $\mathbf{R}_{\mathbf{S}_j}$ . We will consider these questions in Section 2.2.4 below.

### **2.2.2 Rate of Off-diagonal Decay and**

“weak cancellation condition”)

$$\left| \int_{I \times I} K(x, y) dx dy \right| \leq C|I|, \quad (2.2.15)$$

for all dyadic intervals  $I$ . Under these conditions, we have (see [9])

**Theorem 2.2.2** *If the wavelet basis has  $M$  vanishing moments, then, for any kernel  $K$  satisfying the conditions (2.2.13), (2.2.14), and (2.2.15), the matrices  $\alpha^j$ ,  $\beta^j$ ,  $\gamma^j$  satisfy the estimate*

$$|\alpha_{k,l}^j| + |\beta_{k,l}^j| + |\gamma_{k,l}^j| \leq C_M^j (1 + |k - l|)^{-M-1}, \quad (2.2.16)$$

for all in

**Theorem 2.2.4** *If a matrix  $\{m_{k,k',l,l'}\}_{k,k',l,l' \in \mathbf{Z}}$  satisfies*

$$|m_{k,k',l,l'}| < C(1 + |k - k'| + |l - l'|)^{-2-\alpha} \quad (2.2.19)$$

*(where  $\alpha \in \mathbf{Z}, \alpha \geq 2$ ) and if the matrix is invertible on  $l^2$ , then*

$$|m_{k,k',l,l'}^{-1}| < C''(1 + |k - k'| + |l - l'|)^{-2-\alpha}. \quad (2.2.20)$$

See Figure 2.1 for an example of such a matrix after truncation of elements below a given threshold. Matrices which satisfy (2.2.19) also form an algebra under multiplication; for a proof of this see Chapter 3.

We use Theorems 2.2.3 and 2.2.4 to show that, at all stages of the reduction procedure in both one and two dimensions, the matrices representing the **A**, **B**, and **C** blocks of the reduced operators (1.2.10) satisfy the same off-diagonal decay estimate (2.2.16) as the blocks of the non-standard form in Theorem 2.2.2 and its two-dimensional analogue. In other words, the reduction procedure preserves sparsity for a wide class of operators. In this sense, the form (or structure) is preserved under the reduction procedure which allows us to apply it over a finite number of scales. The following theorem applies to the one-dimensional case, but analogous results for two dimensions can be proved using Theorem 2.2.4.

**Theorem 2.2.5 (Preservation of structure over finitely many scales)** *Assume that the operator **S** and the wavelet basis satisfy the conditions of Theorem 2.2.2. Let  $\mathbf{R}_j$  be the reduced operator on some scale  $j$ , where reduction started at some scale  $n$ ,  $n \leq j$ ,  $n, j \in \mathbf{Z}$ , and let  $\mathbf{A}_{\mathbf{R}_j}$ ,  $\mathbf{B}_{\mathbf{R}_j}$  and  $\mathbf{C}_{\mathbf{R}_j}$  be its blocks. Then the bi-infinite matrices  $\alpha^{r,j}$ ,  $\beta^{r,j}$  and  $\gamma^{r,j}$  representing these blocks satisfy*

$$|\alpha_{k,l}^{r,j}| + |\beta_{k,l}^{r,j}| + |\gamma_{k,l}^{r,j}| \leq C_M^{n,j} (1 + |k - l|)^{-M-1}, \quad (2.2.21)$$

*for all integers  $k, l$ .*

**Proof:** Our starting point is the operator  $\mathbf{S}_n$  and its blocks,  $\mathbf{A}_{\mathbf{S}_n}$ ,  $\mathbf{B}_{\mathbf{S}_n}$ ,  $\mathbf{C}_{\mathbf{S}_n}$  and  $\mathbf{T}_{\mathbf{S}_n} = \mathbf{S}_{n+1}$ . Matrices representing these blocks satisfy the estimate of Theorem 2.2.2. Since  $\mathbf{S}_n$  is positive definite, so is  $\mathbf{A}_{\mathbf{S}_n}$  (see Section 2.2.1) and, thus,  $\mathbf{A}_{\mathbf{S}_n}^{-1}$  exists and, according to Theorem 2.2.3,

satisfies the estimate in (2.2.16). Since  $\mathbf{B}_{\mathfrak{g}_n}$  and

for all integer  $k, l$  and  $m$ . Matrices  $\{m_{kl}\}_{k,l \in \mathbf{Z}}$  which are invertible on  $l^2$  and which satisfy for all integer  $m$  the inequality

$$|m_{k,l}^{-1}| < \tilde{C}_m (1 + |k - l|)^{-m}, \quad (2.2.27)$$

form an algebra (see [37]). We may thus repeat the above considerations to prove a version of Theorem 2.2.5 with the decay condition replaced by a decay condition of the form of (2.2.27).

**Remark 2.** It is clear that Theorem 2.2.5 may be view



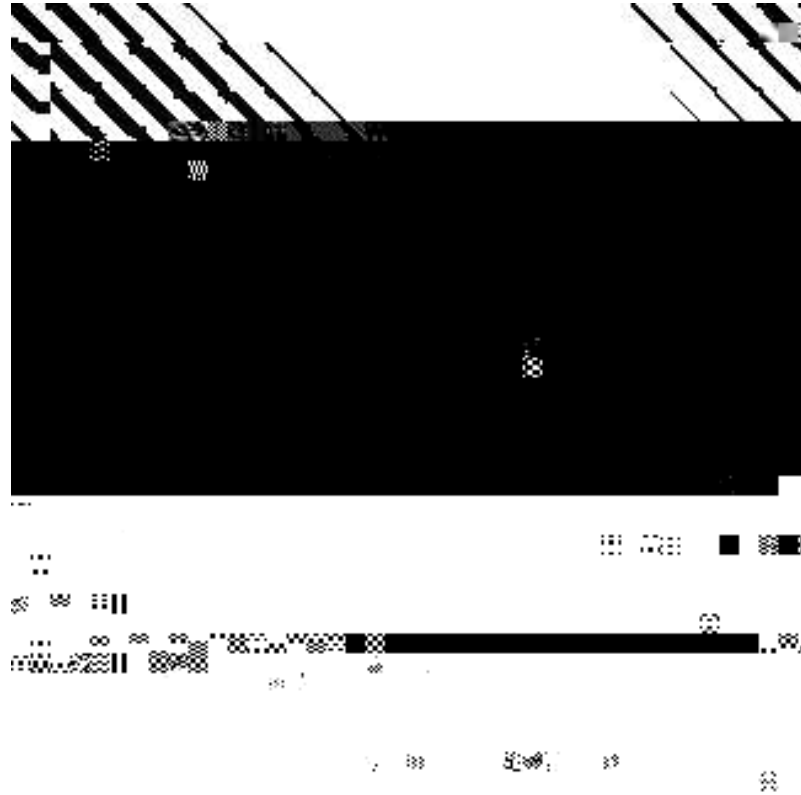


Figure 2.1:  $\mathbf{R}_{\mathbf{S}_j}$  after truncation of entries smaller than  $0.02 * \|\mathbf{R}_{\mathbf{S}_j}\|_{\infty}$ .  $\mathbf{S}_j$  is the projection of  $-\nabla \cdot (a(x, y) \nabla)$  on the unit square with periodic boundary conditions into the multiwI H H H mg H

In the one-dimensional case, if  $\mathbf{A}_{\mathbf{S}_j}$  is banded with bandwidth  $m$ , then its LU-factors will also be banded with bandwidth  $m$ , so they may be computed in  $\mathcal{O}(Nm^2)$ . If  $\mathbf{B}_{\mathbf{S}_j}$  is also banded with this same bandwidth, then we may solve for  $\tilde{\mathbf{B}}_{\mathbf{S}_j}$  in  $\mathcal{O}(Nm^2)$ ; and similarly for  $\hat{\mathbf{C}}_{\mathbf{S}_j}$ . For fixed relative accuracy  $\epsilon$  (and, hence, fixed bandwidth  $m$ ), this leads directly to the  $\mathcal{O}(N)$  procedure for computing  $\mathbf{R}_{\mathbf{S}_j}$  via the sparse incomplete block LU decomposition given by (2.2.28).

The two-dimensional case is more complicated. Each of the blocks on the left-hand side of (2.2.28) will in general exhibit a multi-banded structure, as seen in Figure 2.1. Thus, one may expect the LU-factors of  $\mathbf{A}_{\mathbf{S}_j}$  to fill in between the bands. Indeed, this is the case, but the fill-in which occurs is observed in practice to be fill-in with rapid decay, so that truncating to a given accuracy as we compute the LU factors results in a fast method for computing the reduction (as in the one-dimensional case).

There are many details involved in the description of the multiresolution LU decomposition, and we refer to [20] for a full treatment of them. We note finally that, by virtue of this algorithm the reduction procedure requires  $\mathcal{O}(N)$  operations.

#### 2.2.4 Eigenvalues and Eigenvectors of the Reduced Operators

In this section, we further investigate the relations between the spectra of the operators  $\mathbf{S}_j$  and  $\mathbf{R}_{\mathbf{S}_j}$ . In Section 2.2.1, we established relations between the spectral bounds of these operators. Here, we consider relations between the small eigenvalues and corresponding eigenvectors of the operators  $\mathbf{S}_j$  and  $\mathbf{R}_{\mathbf{S}_j}$ .

We view  $\mathbf{S}_j$  as the projection of a positive definite self-adjoint elliptic operators with a compact inverse; this class includes variable-coefficient elliptic operators. For such an operator  $\mathbf{S}$ , the spectrum consists of isolated eigenvalues with finite multiplicity and the only accumulating point is at infinity. The eigenvalues may be ordered according to

$$0 < \lambda_0 \leq \lambda_1 \leq \lambda_2 \dots$$

The eigenvectors of such operators form an orthonormal basis in the Hilbert space  $\mathcal{H}$ , and each eigenspace is a finite-dimensional subspace. Heuristically, e.g. in numerical literature,

it is always assumed for elliptic operators that the eigenvectors which correspond to small eigenvalues are less oscillatory than those which correspond to large eigenvalues and the number of oscillations increases as  $\lambda_n \rightarrow \infty$

approximate the eigenspaces in a natural sequence, proceeding from less oscillatory to more oscillatory. For practical purposes, however, we have to construct the MRA very carefully if we want to achieve this property for the first few scales that are involved. For example, it is clear that in order to have good approximating properties, the basis functions of the MRA have to satisfy  $\| \hat{f}_i \|_{\infty} \rightarrow 0$  as  $i \rightarrow \infty$ .

Given the eigenvalue problem

$$\begin{pmatrix} \mathbf{A}_{\mathfrak{S}_j} & \mathbf{B}_{\mathfrak{S}_j} \\ \mathbf{B}_{\mathfrak{S}_j}^* & \mathbf{T}_{\mathfrak{S}_j} \end{pmatrix} \begin{pmatrix} d \\ s \end{pmatrix} = \lambda \begin{pmatrix} d \\ s \end{pmatrix}, \quad (2.2.33)$$

we use the same approach as in deriving (1.2.9). Solving for  $d$  in terms of  $s$  and assuming that  $(\mathbf{A}_{\mathfrak{S}_j} - \lambda \mathbf{I})^{-1}$  exists, we obtain

$$(\mathbf{T}_{\mathfrak{S}_j} - \mathbf{B}_{\mathfrak{S}_j}^* (\mathbf{A}_{\mathfrak{S}_j} - \lambda \mathbf{I})^{-1} \mathbf{B}_{\mathfrak{S}_j}) s = \lambda s. \quad (2.2.34)$$

The existence of

$$\mathbf{G}(\lambda) = (\mathbf{A}_{\mathfrak{S}_j} - \lambda \mathbf{I})^{-1} \quad (2.2.35)$$

is assured if we consider (2.2.34) for  $\lambda$  smaller than the lower bound of  $\mathbf{A}_{\mathfrak{S}_j}$ .

We now consider approximations of the left-hand side of (2.2.34) and the accuracy of solutions based on these approximations. We will use the following simple lemma.

**Lemma 2.2.1** *For a normal matrix  $\mathbf{M}$ , if*

$$\mathbf{M}x = \lambda x + \xi, \quad (2.2.36)$$

*then there exists an eigenvalue  $\lambda_{\mathbf{M}}$  of  $\mathbf{M}$  such that*

$$|\lambda - \lambda_{\mathbf{M}}| \leq \frac{\|\xi\|}{\|x\|}. \quad (2.2.37)$$

**Proof:** The proof is straightforward. Let  $\mathbf{G} = \mathbf{M} - \lambda \mathbf{I}$ . Then there is a singular value  $\sigma_0$  of  $\mathbf{G}$  such that

$$\sigma_0 = \inf_{\|y\| \neq 0} \frac{(\mathbf{G}^* \mathbf{G} y, y)^{\frac{1}{2}}}{\|y\|} \leq \frac{\|\mathbf{G}x\|}{\|x\|} = \frac{\|\xi\|}{\|x\|}. \quad (2.2.38)$$

Since  $\mathbf{G}$  is normal, it is diagonalizable by a unitary matrix  $\mathbf{Q}$ . Therefore, the singular values of  $\mathbf{G}$  are given by the absolute values of its eigenvalues. Since at least one singular value of  $\mathbf{G}$  satisfies (2.2.38), the estimate (2.2.37) follows.  $\square$

From (2.2.33), it is clear that

$$d = -\mathbf{G}(\lambda) \mathbf{B}_{\mathfrak{S}_j} s. \quad (2.2.39)$$

We rewrite (2.2.34) as

$$\mathbf{T}_{\mathfrak{S}_j} s = \lambda s + \mathbf{B}_{\mathfrak{S}_j}^* d. \quad (2.2.40)$$

Using

$$\mathbf{G}(\lambda) - \mathbf{G}(0) = \lambda \mathbf{G}(\lambda) \mathbf{G}(0), \quad (2.2.41)$$

where  $\mathbf{G}(0) = \mathbf{A}$

The equations (2.2.48), (2.2.50), and (2.2.51) represent the modified reduction procedure. The operator  $\mathbf{Y}_{\mathfrak{g}}$  is self-adjoint and  $\mathfrak{g}$

**Theorem 2.2.6** *Given an eigenvector  $x$  of  $\mathbf{S}_j$  such that  $\mathbf{S}_j x = \lambda x$ ,  $\|x\|_2 = 1$ ,  $d = \mathbf{Q}_{j+1} x$ , and  $\|d\|_2^2 \ll \frac{1}{2}$ , there exist real  $\lambda_{\mathbf{T}}$ ,  $\lambda_{\mathbf{R}}$ , and  $\lambda_{\mathbf{Y}}$  which solve (2.2.45), (2.2.46), and (2.2.47), respectively, such that*

$$|\lambda_{\mathbf{T}} - \lambda| \leq C_d \|\mathbf{B}_{\mathfrak{S}_j}\|_2 \|d\|_2 \quad (2.2.57)$$

$$|\lambda_{\mathbf{R}} - \lambda| \leq C_d \|\mathbf{B}_{\mathfrak{S}_j}\|_2 \|d\|_2 \left( \frac{\lambda}{m_{\mathbf{A}}^j} \right) \quad (2.2.58)$$

$$|\lambda_{\mathbf{Y}} - \lambda| \leq C_d \|\mathbf{B}_{\mathfrak{S}_j}\|_2 \|d\|_2 \left( \frac{\lambda}{m_{\mathbf{A}}^j} \right)^2, \quad (2.2.59)$$

where  $1 \leq C_d \leq \sqrt{2}$ .

We may now identify two factors that affect the estimate:  $\|d\|_2$  and the ratio  $\lambda/m_{\mathbf{A}}^j$ . In order for  $\|d\|_2$  to be small, we have to assume **th**



**Remark.** In Section 2.2.3, we outlined an  $\mathcal{O}(N)$  procedure for computing the reduced operator to relative accuracy  $\epsilon$ . For small eigenvalues, however, it might be necessary to maintain absolute rather than relative accuracy while performing the reduction. This puts an additional computational burden on the reduction procedure in the case of ill-conditioned operators.

In particular, if we compute  $\hat{\mathbf{A}}_{\mathbf{S}_j}$  and  $\hat{\mathbf{C}}_{\mathbf{S}_j} = \tilde{\mathbf{B}}_{\mathbf{S}_j}^*$  to some absolute accuracy  $\delta$ , and from this compute  $\mathbf{R}'_{\mathbf{S}_j}$ , it is clear (from (2.2.29)) that

$$\|\mathbf{R}'_{\mathbf{S}_j} - \mathbf{R}_{\mathbf{S}_j}\| < \delta \|\hat{\mathbf{C}}_{\mathbf{S}_j}\|. \quad (2.2.60)$$

In the worst case, the eigenvalues of  $\mathbf{R}'_{\mathbf{S}_j}$  will approximate the eigenvalues of  $\mathbf{R}_{\mathbf{S}_j}$  with accuracy no better than  $\delta \|\hat{\mathbf{C}}_{\mathbf{S}_j}\|$  (see e.g. [25]). For a typical second-order elliptic operator  $\mathbf{S}$ , the norms of each of the blocks  $\mathbf{A}_{\mathbf{S}_j}$  and  $\mathbf{C}_{\mathbf{S}_j} = \mathbf{B}_{\mathbf{S}_j}^*$  behave like  $\mathcal{O}(h_j^{-2})$  (where  $h_j$  is the step size of the discretization). Furthermore, in the Cholesky decomposition, the norm of the lower triangular factor is equal to the square root of the norm of the matrix. Therefore, if we compute the LU factorization defined in Section 2.2.3 to absolute accuracy  $\delta$ , then the resulting matrix  $\mathbf{R}'_{\mathbf{S}_j}$  approximates  $\mathbf{R}_{\mathbf{S}_j}$  to absolute accuracy  $\delta h_j^{-1}$ , as can easily be seen from (2.2.60).

In other words, to compute  $\mathbf{R}'_{\mathbf{S}_j}$  so that its eigenvalues approximate the small eigenvalues of  $\mathbf{R}_{\mathbf{S}_j}$  with absolute accuracy  $\epsilon$ , it is necessary to compute the multiresolution LU decomposition with working precision  $\epsilon \cdot h_j$ . For a given accuracy  $\delta$ , the bandwidth  $m$  of matrices which satisfy (2.2.16) (or its two-dimensional analogue) is given by  $m \sim (C\delta)^{-\frac{1}{M}}$ , where  $M$  is the number of vanishing moments of the wavelet basis (see e.g. [9] for details). Thus, as  $h_j$  decreases (and the scale becomes finer (

Table 2.1: Condition numbers and lower bounds of  $\mathbf{A}_{\mathbf{S}_j}$ .

$N$	$\kappa(\mathbf{A}_j)$	$m_{\mathbf{A}_j}$	$\kappa(\mathbf{S}_j)$
256	6.18	$1.4577 \times 10^3$	$1.16 \times 10^2$
1024	8.06	$5.0363 \times 10^3$	$6.26 \times 10^2$
2304	10.48	$1.0043 \times 10^4$	$1.53 \times 10^3$
4096	11.66	$1.5948 \times 10^4$	$2.86 \times 10^3$
5184	13.06	$1.9077 \times 10^4$	$3.66 \times 10^3$

Table 2.1: Condition numbers and lower bounds for the  $\mathbf{A}$ -block of the operator  $-\nabla \cdot (a(x, y)\nabla)$  on the unit square with periodic boundary conditions. Here,  $N$  is the number of unknowns in the two-dimensional spatial grid. Multiwavelets with two vanishing moments are used, and the coefficients  $a(x, y)$  are set to  $a(x, y) = 2 + \cos(16\pi x) \cos(16\pi y)$ , which provides a moderate amount of oscillation in the coefficients. The condition number depends only weakly on the scale, unlike the condition number of the original matrix (denoted in the table as  $\kappa(\mathbf{S}_j)$ ), which for second-order elliptic operators scales as  $h^{-2}$  (where  $h$  is the step-size of the discretization). Note that  $m_{\mathbf{A}_j}$  also scales as  $h^{-2}$ .

that the multiresolution LU decomposition requires  $\mathcal{O}(NN^{\frac{4n}{M}}) = \mathcal{O}(N^{1+\frac{4}{M}})$  operations to compute the matrix  $\mathbf{R}'_{\mathbf{S}_j}$ , so that its eigenvalues approximate the eigenvalues of  $\mathbf{R}_{\mathbf{S}_j}$  to absolute accuracy  $\epsilon$ . This means, for example, that when  $M = 2$ , the computational complexity could be  $\mathcal{O}(N^3)$ , which is as bad as computing the Cholesky decomposition exactly. However, we see in Table 2.2 that in practice, even with  $M = 2$  things may

It in praO a g QH Rgo



Table 2.2: Run times for the sparse versus full reduction procedure.

$N$	$T_{\text{exact}}$	$T_{\epsilon h_j}$	$\ \mathbf{R}_{\mathbf{S}_j} - \mathbf{R}'_{\mathbf{S}_j}\ _{\infty}$
576	5.21	4.33	$3.3 \times 10^{-3}$
1296	51.69	25.54	$4.1 \times 10^{-3}$
2304	287.52	66.04	$4.9 \times 10^{-3}$
3600	18.3 min*	121.67	$5.2 \times 10^{-3}$
5184	54.6 min*	195.74	$5.9 \times 10^{-3}$
9216	5.1 hrs*	417.47	$7.8 \times 10^{-3*}$
16384	28.7 hrs*	901.21	†

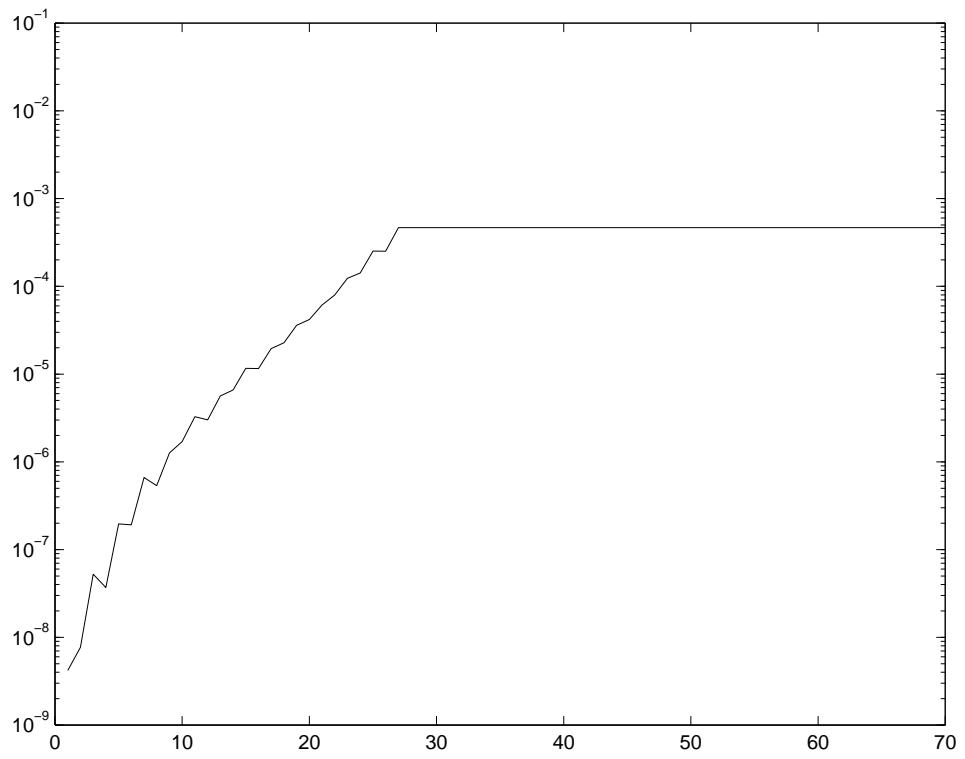
Table 2.2: Run times for exact versus truncated computation of  $\mathbf{R}_{\mathbf{S}_j}$  for various scales. The operator  $\mathbf{S}_j$  is the projection of  $-\nabla \cdot (a(x, y)\nabla)$  on the unit square with periodic boundary conditions.

) ON fpe

difference approximation to the first derivative,  $M$  is a diagonal matrix with uniform samples

- Relative error for the eigenvalues of the T-block
- - Relative error for the eigenvalues of the reduced operator
- ..... Relative error for the eigenvalues of the modified reduced operator









## CHAPTER 3

### CLASSES OF MATRICES

In practice, one of the chief benefits of using wavelets for numerical analysis of differential and integral equations is that the matrix representations of the operators have sparse approximations. The structure of these sparse approximations typically takes the form of fast decay in the magnitude of the elements away from the diagonal, so that most of the elements of the matrix may be truncated to zero. For matrices of finite size, it is impossible to quantify the rate of decay of the elements away from the diagonal for a class of matrices since such a statement is inherently a limit statement. Statements about the rate of decay apply only to the infinite-dimensional (which we also call bi-infinite) matrices. In this chapter, we consider several classes of bi-infinite matrices, and prove some results which were used in Chapter 2. As the experimental results of Chapter 2 demonstrate, the decay rate for matrices of finite size is sufficient to render the algorithms presented in that chapter useful in practice.

#### 3.1 Prelimi

$\mathbf{m} \stackrel{\epsilon}{\mathcal{L}} \mathbf{C}$

$\mathbf{m}$

$\mathbf{m} \stackrel{\epsilon}{d}$

where the number  $m_{\mathbf{k},\mathbf{l}}$  is the  $\mathbf{k},\mathbf{l}$  entry of the matrix  $\mathbf{M}$ . We may also denote this entry by  $(\mathbf{M})_{\mathbf{k},\mathbf{l}}$ , so that we may write  $(\mathbf{M}_1\mathbf{M}_2)_{\mathbf{k},\mathbf{l}}$  for the  $\mathbf{k},\mathbf{l}$  entry of the matrix defined by the matrix-matrix product of  $\mathbf{M}_1$  and  $\mathbf{M}_2$ .

In the following sections, we study matrices which have various types of decay in the magnitude of the elements as one moves away from the diagonal. We define the distance between multi-indices in terms of the absolute values of the differences of their components, i.e.

$$|\mathbf{k} - \mathbf{l}| = \sum_{n=1}^d |k_n - l_n|. \quad (3.1.2)$$

As defined above,  $|\mathbf{k} - \mathbf{l}|$  satisfies the triangle inequality, as well as other properties of the usual notion of distance. The diagonal of a matrix is defined by  $|\mathbf{k} - \mathbf{l}| = 0$ , or  $\mathbf{k} = \mathbf{l}$ . Operations on matrices and vectors, such as transposes and multiplication, are defined using this notation in the usual

for  $1 \leq p < \infty$ , and

$$\|x\|_\infty = \max_{1 \leq i \leq n} |x_i|$$

We start by showing the following:

**Lemma 3.2.1** *If  $\mathbf{M} \in \mathcal{X}$ , then  $\mathbf{M}$  is a bounded operator on  $l^1$  and  $l^\infty$ .*

**Proof:** Via Young's Inequality, we see that

$$\begin{aligned}\|\mathbf{M}\|_{1,1} &= \sup_{\|x\|_1=1} \|\mathbf{M}x\|_1 \\ &= \sup_{\|x\|_1=1} \sum_{\mathbf{k} \in \mathbf{Z}}\end{aligned}$$

where  $A(\epsilon) = \sum_{\mathbf{j} \in \mathbb{Z}^2} e^{-\epsilon|\mathbf{j}|} < \infty$ . Thus the function  $A(\epsilon)$  depends only on  $\epsilon$ , and, furthermore,  $A(\epsilon) > 1$ . We see that if  $\epsilon_1 < \epsilon_2$ , then (3.2.11) implies that  $\mathbf{M}_1 \mathbf{M}_2 \in \mathcal{X}$ . If  $\epsilon_1 = \epsilon_2$  we note that we may always adjust  $\epsilon_1$  so that it is smaller than  $\epsilon_2$ . Thus, for any values of  $\epsilon_1, \epsilon_2 > 0$ , we have  $\mathbf{M}_1 \mathbf{M}_2 \in \mathcal{X}$ .  $\square$

It is shown in [37] and [22] if  $d = 1$ , then invertible elements of the class  $\mathcal{X}$  have inverses which are elements of  $\mathcal{X}$ . We use Lemma 3.2.2 to prove that the same is true if  $d = 2$ . Our proof follows that of [37] and [22] very closely.

**Theorem 3.2.3** *If  $\mathbf{M} \in \mathcal{X}$ , and  $\mathbf{M}$  is invertible on  $l^2$ , then  $\mathbf{M}^{-1} \in \mathcal{X}$ .*

**Proof:** First, consider the case where  $\mathbf{M} = \mathbf{I} - \mathbf{U}$  and  $\|\mathbf{U}\|_{2,2} < 1$ . Then we have  $\mathbf{M}^{-1} = \sum_{n=0}^{\infty} \mathbf{U}^n$ . We denote the  $\mathbf{k}, \mathbf{l}$  entry of  $\mathbf{U}^n$  by  $(u^{(n)})_{\mathbf{k}, \mathbf{l}}$ . Since  $\mathbf{U} \in \mathcal{X}$ , we write  $|(u^{(1)})_{\mathbf{k}, \mathbf{l}}| \leq C e^{-\epsilon|\mathbf{k}-\mathbf{l}|}$ . Using Lemma 3.2.2 and noting that  $|(u^{(1)})_{\mathbf{k}, \mathbf{l}}| \leq C e^{-\epsilon|\mathbf{k}-\mathbf{l}|} \leq C e^{-\frac{\epsilon}{2}|\mathbf{k}-\mathbf{l}|}$ , we estimate

$$|(u^{(n)})_{\mathbf{k}, \mathbf{l}}| \leq C^n A\left(\frac{\epsilon}{2}\right)^{n-1} e^{-\frac{\epsilon}{2}|\mathbf{k}-\mathbf{l}|}. \quad (3.2.18)$$

Since  $A > 1$ , we then estimate

$$\left| \sum_{n=0}^N (\mathbf{U}^n)_{\mathbf{k}, \mathbf{l}} \right| = \left| \sum_{n=0}^N (u^{(n)})_{\mathbf{k}, \mathbf{l}} \right| \leq \sum_{n=0}^N (AC)^n e^{-\frac{\epsilon}{2}(|\mathbf{k}-\mathbf{l}|)}. \quad (3.2.19)$$

For the remainder, we write

$$\sum_{n=N+1}^{\infty} \|\mathbf{U}^n\|_{2,2} \leq \sum_{n=N+1}^{\infty} \|\mathbf{U}\|_{2,2}^n \frac{\|\mathbf{U}\|_{2,2}^{N+1}}{1 - \|\mathbf{U}\|_{2,2}}. \quad (3.2.20)$$

Since we have

$$\sum_{n=0}^N (AC)^n = \frac{1 - (AC)^{N+1}}{1 - AC}, \quad (3.2.21)$$

and the norm  $\|\mathbf{U}^n\|_{2,2}$  provides a uniform bound on the elements of the matrix  $\mathbf{U}^n$ , we obtain

$$\left| \left( \sum_{n=0}^{\infty} \mathbf{U}^n \right)_{\mathbf{k}, \mathbf{l}} \right| \leq \left| \sum_{n=0}^N (\mathbf{U}^n)_{\mathbf{k}, \mathbf{l}} \right| + \left| \left( \sum_{n=N+1}^{\infty} \mathbf{U}^n \right)_{\mathbf{k}, \mathbf{l}} \right| \quad (3.2.22)$$

$$\begin{aligned} &\leq \frac{1 - (AC)^{N+1}}{1 - AC} e^{-\frac{\epsilon}{2}|\mathbf{k}-\mathbf{l}|} + \frac{\|\mathbf{U}\|_{2,2}^{N+1}}{1 - \|\mathbf{U}\|_{2,2}} \\ &= \frac{1 - e^{(N+1)\log(AC)}}{1 - AC} e^{-\frac{\epsilon}{2}|\mathbf{k}-\mathbf{l}|} + e^{\log(\|\mathbf{U}\|_{2,2})} \frac{1 - \|\mathbf{U}\|_{2,2}^{N+1}}{1 - \|\mathbf{U}\|_{2,2}} \end{aligned} \quad (3.2.23)$$



In

where  $\|\mathbf{U}\|_{2,2} \leq \frac{B-A}{B+A} < 1$ . Thus, if  $\mathbf{L} \in \mathcal{X}$  is symmetric and positive-definite, it may b

$$= C_1 C_2 \sum_{\mathbf{j} \in \mathbf{Z}^2} (1 + |(\mathbf{k} - \mathbf{1}) - \mathbf{j}|)^{-2-\alpha} (1 + |\mathbf{j}|)^{-2-\alpha}. \quad (3.3.5)$$

We set  $\tilde{\mathbf{k}} = \mathbf{k} - \mathbf{1}$  and split  $\mathbf{Z}^2$  into two sets:

$$R_1 = \left\{ \mathbf{j} \in \mathbf{Z}^2 \mid m(\tilde{\mathbf{k}}, \mathbf{j}) \geq m(\mathbf{j}, 0) \right\} \quad (3.3.6)$$

and

$$R_2 = \mathbf{Z}^2 \setminus R_1, \quad (3.3.7)$$

where  $m$  is the usual euclidean distance  $m(\mathbf{p}, \mathbf{q}) = \sqrt{(p_1 - q_1)^2 + (p_2 - q_2)^2}$ . We note that  $\frac{1}{\sqrt{2}}|\mathbf{p} - \mathbf{q}| \leq m(\mathbf{p}, \mathbf{q}) \leq |\mathbf{p} - \mathbf{q}|$ , which implies that

$$(1 + |\mathbf{p} - \mathbf{q}|)^{-2-\alpha} \leq (1 + m(\mathbf{p}, \mathbf{q}))^{-2-\alpha} \leq 2^{1+\frac{\alpha}{2}} (1 + |\mathbf{p} - \mathbf{q}|)^{-2-\alpha}. \quad (3.3.8)$$

Additionally, it is clear that if  $\mathbf{j} \in R_1$ , then

$$m(0, \tilde{\mathbf{k}}) \leq m(0, \mathbf{j}) + m(\tilde{\mathbf{k}}, \mathbf{j}) \leq 2m(\tilde{\mathbf{k}}, \mathbf{j}), \quad (3.3.9)$$

so, for the sum over  $R_1$ , we may write

$$\sum_{\mathbf{j} \in R_1} (1 + |\tilde{\mathbf{k}} - \mathbf{j}|)^{-2-\alpha} (1 + |\mathbf{j}|)^{-2-\alpha} \leq C(1 + |\tilde{\mathbf{k}}|)^{-2-\alpha} \sum_{\mathbf{j} \in R_1} (1 + |\mathbf{j}|)^{-2-\alpha} \leq C'(1 + |\tilde{\mathbf{k}}|)^{-2-\alpha} \quad (3.3.10)$$

for  $|\tilde{\mathbf{k}}|$  large enough.

We note that  $\mathbf{j} \in R_2$  implies that  $m(\mathbf{j}, 0) > m(\tilde{\mathbf{k}}, \mathbf{j})$ , from which we deduce  $m(\mathbf{j} - \tilde{\mathbf{k}}, -\tilde{\mathbf{k}}) > m(0, \mathbf{j} - \tilde{\mathbf{k}})$ , which means that  $\tilde{\mathbf{k}} - \mathbf{j} \in R_1$ . For the sum over  $R_2$ , we write

$$\sum_{\mathbf{j} \in R_2} (1 + |\tilde{\mathbf{k}} - \mathbf{j}|)^{-2-\alpha} (1 + |\mathbf{j}|)^{-2-\alpha} = \sum_{\mathbf{j}' \in R_1} (1 + |\mathbf{j}'|)^{-2-\alpha} (1 + |\tilde{\mathbf{k}} - \mathbf{j}'|)^{-2-\alpha} \leq C' \mathbf{W}^{-\alpha} \quad (3.3.1) \quad 6 \quad 6 \quad \mathbf{4} \quad 6 \quad 6$$





Since  $(|\mathbf{k} - \mathbf{l}|)^{2+\alpha} \leq (1 + |\mathbf{k} - \mathbf{l}|)^{2+\alpha}$ , we may divide both sides of the above inequality by  $(1 + |\mathbf{k} - \mathbf{l}|)^{2+\alpha}$  to obtain

$$|m_{\mathbf{k}, \mathbf{l}}^{-1}| \leq C(\mathbf{M}, \alpha)(1 + |\mathbf{k} - \mathbf{l}|)^{-2-\alpha}, \quad (3.3.22)$$

which implies that  $\mathbf{M}^{-1} \in \mathcal{I}_\alpha$ . In order to complete the proof of Theorem 3.3.2, we need to show that

**Proposition 3.3.1** *If  $\mathbf{M} \in \mathcal{I}_\alpha$  is invertible on  $l^2$  and  $\alpha \in \mathbf{Z}$ ,  $\alpha \geq 2$ , then the commutator  $[\mathbf{X}_i, \mathbf{M}]_{2+\alpha}$  is a bounded linear operator from  $l^1$  to  $l^\infty$ .*

**Proof:** The proof of this proposition hinges on the following:

**Proposition 3.3.2** *If  $\mathbf{M} \in \mathcal{I}_\alpha$  is invertible on  $l^2$  and  $\alpha \geq 2$ , then  $\mathbf{M}$  is invertible on  $l^p$ ,  $1 \leq p \leq \infty$ .*

We prove Proposition 3.3.2 at the end of this chapter.

The proof of Proposition 3.3.1 also requires the following lemma (based on a lemma from [37]).

**Lemma 3.3.1** *If  $\mathbf{A}$  is a bounded linear operator on  $l^2$  such that*

$$|a_{\mathbf{k}, \mathbf{l}}| \leq C(1 + |\mathbf{k} - \mathbf{l}|)^{-s}, \quad s \in (0, 2], \quad (3.3.23)$$

*and  $\frac{1}{p} - \frac{1}{q} > 1 - \frac{s}{2}$ , then  $\mathbf{A}$  is continuous from  $l^p$  to  $l^q$ .*

**Proof:** We define the sequence  $\mathcal{P}(\gamma)$  by  $(\mathcal{P}(\gamma))_{\mathbf{k}} = (1 + |\mathbf{k}|)^{-\gamma}$ . Suppose  $x \in l^p$  and  $\|x\|_p = 1$ .

Then

$$\|\mathbf{A}x\|_q \leq C\|\mathcal{P}(s) * x\|_q \leq C\|\mathcal{P}(s)\|_r \|x\|_p \quad (3.3.24)$$

if  $\frac{1}{q} = \frac{1}{r} + \frac{1}{p} - 1$  (by Young's inequality). But  $\|\mathcal{P}(s)\|_r < \infty$  iff  $sr > 2$ , so if  $\frac{1}{p} - \frac{1}{q} > 1 - \frac{s}{2}$  we can choose an  $r$  such that  $sr > 2$  and  $\frac{1}{q} = \frac{1}{r} + \frac{1}{p} - 1$ .  $\square$

Now, we continue with the proof of Proposition 3.3.1. We use the commutators of  $\mathbf{X}_i$  with  $\mathbf{M}^{-1}$  to obtain all but the last two degrees of decay in the entries of  $\mathbf{M}^{-1}$ . We make use of the commutator identities

$$[\mathbf{X}_i, \mathbf{M}^{-1}] = -\mathbf{M}^{-1}[\mathbf{X}_i, \mathbf{M}]\mathbf{M}^{-1} \quad (3.3.25)$$

and

$$[\mathbf{X}_i, [\mathbf{X}_i, \mathbf{M}^{-1}]] = \mathbf{M}^{-1}(-[\mathbf{X}_i, [\mathbf{X}_i, \mathbf{M}]] - 2[\mathbf{X}_i, \mathbf{M}]\mathbf{M}^{-1}[\mathbf{X}_i, \mathbf{M}])\mathbf{M}^{-1} \quad (3.3.26)$$

and their higher-order analogues. That is,  $[\mathbf{X}_i, \mathbf{M}^{-1}]_\alpha$  is a sum of products of  $\mathbf{M}^{-1}$  and commutators of  $\mathbf{X}_i$  and  $\mathbf{M}$  of order no greater than  $\alpha$  (see [37], [6]). By Lemma 3.3.1, it is clear that each of these commutators is a bounded linear operator from  $l^p$  to  $l^{p+\epsilon}$  for any  $p > 1$ ,  $\epsilon > 0$ , so we may string together these commutators with  $\mathbf{M}^{-1}$ , which by Proposition 3.3.2 is bounded on all  $l^p$  spaces, in any combination of products to obtain a bounded linear operator from  $l^1$  to  $l^\infty$ . This implies that  $\mathbf{M}^{-1} \in \mathcal{I}_{\alpha-2}$ .

We obtain the last two degrees of decay using a bootstrapping proof following [37]. The following lemma provides the first step:

**Lemma 3.3.2** *If  $\mathbf{M} \in \mathcal{I}_\alpha$  ( $\alpha \in \mathbf{Z}$ ,  $\alpha \geq 2$ ) and  $\mathbf{M}^{-1} \in \mathcal{I}_{\alpha+\gamma-2}$ , where  $0 \leq \gamma \leq 1$ , then*

$$[\mathbf{X}_i, \mathbf{M}^{-1}]_\alpha = -\mathbf{M}^{-1}([\mathbf{X}_i, \mathbf{M}]_\alpha + [\mathbf{X}_i, \mathbf{M}]_{\alpha-1}\mathbf{M}^{-1}[\mathbf{X}_i, \mathbf{M}] + [\mathbf{X}_i, \mathbf{M}]\mathbf{M}^{-1}[\mathbf{X}_i, \mathbf{M}]_{\alpha-1} + \mathbf{K}_\alpha)\mathbf{M}^{-1} \quad (3.3.27)$$

where  $\mathbf{K}_\alpha \in \mathcal{I}_\gamma$ .

**Proof:** First, consider the case where  $\alpha = 2$ . We see that the identity (3.3.26) implies that  $\mathbf{K}_2 = 0$  and so the result is trivially true.

Now consider the case where  $\alpha > 2$ , which we prove in this case by induction. As the first step in the proof, we use the identity

of  $\mathcal{I}_{\alpha-2+\gamma}$ . The commutator of  $\mathbf{X}_i$  with any element of  $\mathcal{I}_{\alpha-2+\gamma}$  is an element of  $\mathcal{I}_{\alpha-3+\gamma}$ , so from this we see that  $\mathbf{K}_3 \in \mathcal{I}_{\alpha-3+\gamma}$ .

Now, assume that

$$[\mathbf{X}_i, \mathbf{M}^{-1}]_n = -\mathbf{M}^{-1}([\mathbf{X}_i, \mathbf{M}]_n + [\mathbf{X}_i, \mathbf{M}]_{n-1} \mathbf{M}^{-1}[\mathbf{X}_i, \mathbf{M}] + [\mathbf{X}_i, \mathbf{M}] \mathbf{M}^{-1}[\mathbf{X}_i, \mathbf{M}]_{n-1} + \mathbf{K}_n) \mathbf{M}^{-1} \quad (3.3.30)$$

where  $\mathbf{K}_n \in \mathcal{I}_{\alpha-n+\gamma}$  and  $n$  is an integer strictly less than  $\alpha$ . We wish to show that  $[\mathbf{X}_i, \mathbf{M}^{-1}]_{n+1}$  has this same form and  $\mathbf{K}_{n+1} \in \mathcal{I}_{\alpha-(n+1)+\gamma}$ .

By the definition of the commutator, we have

$$[\mathbf{X}_i, \mathbf{M}^{-1}]_{n+1} = \mathbf{X}_i [\mathbf{X}_i, \mathbf{M}^{-1}]_n - [\mathbf{X}_i, \mathbf{M}^{-1}]_n \mathbf{X}_i, \quad (3.3.31)$$

which leads us to the relation

$$\begin{aligned} [\mathbf{X}_i, \mathbf{M}^{-1}]_{n+1} &= \mathbf{M}^{-1}(-[\mathbf{X}_i, \mathbf{M}]_{n+1} \\ &\quad - [\mathbf{X}_i, \mathbf{M}]_n \mathbf{M}^{-1}[\mathbf{X}_i, \mathbf{M}] \\ &\quad - [\mathbf{X}_i, \mathbf{M}] \mathbf{M}^{-1}[\mathbf{X}_i, \mathbf{M}]_n + \mathbf{K}_{n+1}) \mathbf{M}^{-1} \end{aligned} \quad (3.3.32)$$

where

$$\mathbf{K}_{n+1} = [\mathbf{X}_i, \mathbf{M}]_{n-1} \mathbf{M}^{-1}[\mathbf{X}_i, \mathbf{M}] \mathbf{M}^{-1}[\mathbf{X}_i, \mathbf{M}] \quad (3.3.33)$$

$$+ 2[\mathbf{X}_i, \mathbf{M}] \mathbf{M}^{-1}[\mathbf{X}_i, \mathbf{M}]_{n-1} \mathbf{M}^{-1}[\mathbf{X}_i, \mathbf{M}] \quad (3.3.34)$$

$$+ [\mathbf{X}_i, \mathbf{M}] \mathbf{M}^{-1}[\mathbf{X}_i, \mathbf{M}] \mathbf{M}^{-1}[\mathbf{X}_i, \mathbf{M}]_{n-1} \quad (3.3.35)$$

$$+ -[\mathbf{X}_i, [\mathbf{X}_i, \mathbf{M}]_{n-1} \mathbf{M}^{-1}[\mathbf{X}_i, \mathbf{M}]] - [\mathbf{X}_i, [\mathbf{X}_i, \mathbf{M}] \mathbf{M}^{-1}[\mathbf{X}_i, \mathbf{M}]_{n-1}] \quad (3.3.36)$$

$$+ \mathbf{M}[\mathbf{M}^{-1} \mathbf{K}_n \mathbf{M}^{-1}, \mathbf{X}_i] \mathbf{M}. \quad (3.3.37)$$

It remains to show that  $\mathbf{K}_{n+1} \in \mathcal{I}_{\alpha-(n+1)+\gamma}$ . Since  $[\mathbf{X}_i, \mathbf{M}]_{n-1} \in \mathcal{I}_{\alpha-n+1}$ , it is clear that each of the three terms (3.3.33), (3.3.34), and (3.3.35) is an element of  $\mathcal{I}_{\alpha-n+1}$ . Additionally, we see that the term (3.3.36) is an element of  $\mathcal{I}_{\alpha-n}$ . The last term, (3.3.37), is an element of  $\mathcal{I}_{\alpha-n+\gamma-1}$  since  $\mathbf{K}_n \in \mathcal{I}_{\alpha-n+\gamma}$  and  $\mathbf{M}^{-1} \in \mathcal{I}_{\alpha-2}$ . Thus, we have shown that  $\mathbf{K}_{n+1} \in \mathcal{I}_{\alpha-(n+1)+\gamma}$  for  $n < \alpha$ . In particular, we may set  $n = \alpha - 1$  to obtain  $\mathbf{K}_\alpha \in \mathcal{I}_\gamma$ , and we have completed the proof of Lemma 3.3.2.  $\square$

In the next step of the bootstrap proof of Proposition 3.3.1, we use Lemma 3.3.2. The technique we use is a slight extension of that of [37].

**Lemma 3.3.3** *If  $\mathbf{M} \in \mathcal{I}_\alpha$  ( $\alpha \in \mathbf{Z}$ ,  $\alpha \geq 2$ ) and  $\mathbf{M}^{-1} \in \mathcal{I}_{\alpha+\gamma-2}$ , where  $0 \leq \gamma \leq 1$ , then the matrix whose  $\mathbf{k}, \mathbf{l}$  element is defined by*

$$|k_i - l_i|^{\beta+\alpha} |m_{\mathbf{k}, \mathbf{l}}^{-1}| \tag{3.3.38}$$

*is a bounded linear operator from  $l^1$  to  $l^\infty$  if  $\beta \leq 2$  and  $\beta < 2 + \gamma$ .*

**Proof:** We note that  $|k_i - l_i|^{\beta+\alpha} |m_{\mathbf{k}, \mathbf{l}}^{-1}| = |k_i - l_i|^\beta ([\mathbf{X}_i, \mathbf{M}^{-1}]_\alpha)_{\mathbf{k}, \mathbf{l}}$ . We use equation (3.3.27) to break the right-hand side of this equation into three parts:

$$|k_i - l_i|^\beta (\mathbf{M}^{-1} [\mathbf{X}_i, \mathbf{M}]_\alpha \mathbf{M}^{-1})_{\mathbf{k}, \mathbf{l}} = |k_i - l_i|^\beta \left( \mathbf{M}^{-1} [\mathbf{X}_i, \mathbf{M}]_\alpha \mathbf{M}^{-1} \right)_{\mathbf{k}, \mathbf{l}} + |k_i - l_i|^\beta \left( \mathbf{M}^{-1} [\mathbf{X}_i, \mathbf{M}]_\alpha \right)_{\mathbf{k}, \mathbf{l}} + |k_i - l_i|^\beta \left( [\mathbf{X}_i, \mathbf{M}]_\alpha \mathbf{M}^{-1} \right)_{\mathbf{k}, \mathbf{l}}$$

We show (under certain restrictions on  $\beta$ ) that each of the matrices  $A = \left\{ \sum_{\mathbf{p}, \mathbf{q}} a_{\mathbf{p}, \mathbf{q}} \right\}$ ,  $B = \left\{ \sum_{\mathbf{p}, \mathbf{q}} b_{\mathbf{p}, \mathbf{q}} \right\}$ , and  $C = \left\{ \sum_{\mathbf{p}, \mathbf{q}} c_{\mathbf{p}, \mathbf{q}} \right\}$  is a bounded linear operator from  $l^1$  to  $l^\infty$ . Thus, we may use the inequality (3.3.43) to show that the matrix defined by (3.3.39) is a bounded linear operator from  $l^1$  to  $l^\infty$ .

Each of the three matrices  $A$ ,  $B$ , and  $C$  is a product of three of the following matrices. The matrix defined by

- $|m_{\mathbf{k}, \mathbf{p}}^{-1}|$  is bounded on  $l^1$  and  $l^\infty$ .
- $|k_i - p_i|^\beta |m_{\mathbf{k}, \mathbf{p}}^{-1}|$  is an element of  $\mathcal{I}_{\alpha-2+\gamma-\beta}$  and therefore (1) is bounded from  $l^1$  to  $l^r$  if  $(\alpha + \gamma - \beta) > \frac{2}{r}$ , by Young's inequality; and (2) is bounded from  $l^r$  to  $l^\infty$  if  $\frac{2}{r} > 2 - (\alpha + \gamma - \beta)$ , again by Young's inequality.
- $|p_i - q_i|^{\beta+\alpha} |m_{\mathbf{p}, \mathbf{q}}|$  is an element of  $\mathcal{I}_{-\beta}$  and if  $\beta \leq 2$  it is bounded from  $l^1$  to  $l^\infty$ .
- $|p_i - q_i|^\alpha |m_{\mathbf{p}, \mathbf{q}}|$  is an element of  $\mathcal{I}_0$  and so is bounded from  $l^r$  to  $l^{r+\epsilon}$  for any  $\epsilon > 0$ .

Therefore, we see that the matrix  $A$  is bounded from  $l^1$  to  $l^\infty$  if there exists an  $r > 1$  such that  $\frac{2}{r} > 2 - (\alpha + \gamma - \beta)$ . Since  $\alpha \geq 2$ , we can only guarantee this if  $\frac{1}{r} > \frac{\beta-\gamma}{2}$ . Since  $r > 1$ , this implies that we require  $\beta < 2 + \gamma$ . Likewise, the matrix  $C$  is bounded from  $l^1$  to  $l^\infty$  if there exists an  $r \geq 1$  such that  $\frac{2}{r} < \alpha + \gamma - \beta$ . Since  $\gamma \geq 0$ , we can guarantee this if  $\beta < \alpha + \gamma$ . Since  $\alpha \geq 2$ , this is already covered by the condition for the matrix  $A$ . Additionally, the matrix

$$\tilde{b}_{\mathbf{p},\mathbf{q}} = |m_{\mathbf{k},\mathbf{p}}^{-1}| |p_i - q_i|^\beta |(\mathbf{K}_\alpha)_{\mathbf{p},\mathbf{q}}| |m_{\mathbf{q},\mathbf{1}}^{-1}|, \quad (3.3.49)$$

and

$$\tilde{c}_{\mathbf{p},\mathbf{q}} = |m_{\mathbf{k},\mathbf{p}}^{-1}| |(\mathbf{K}_\alpha)_{\mathbf{p},\mathbf{q}}| |q_i - l_i|^\beta |m_{\mathbf{q},\mathbf{1}}^{-1}|. \quad (3.3.50)$$

We use the fact that  $\mathbf{K}_\alpha \in \mathcal{I}_\gamma$  to derive the same restrictions on  $\beta$  that were derived for the matrix (3.3.39). The details are very similar to those of the previous case, so we do not display them here.

The matrix defined by (3.3.40) is dealt in two different ways depending on the value of  $\alpha$ . First, consider the case where  $\alpha > 2$ . If  $\beta \leq 2$ , then we see that the matrix defined by (3.3.40) is at worst an element of  $\mathcal{I}_{-1}$ , so there exists a uniform bound on its elements. Any s ele

**Proof of Proposition 3.3.2.** The main idea of our proof comes from [37]. First we define the matrix  $\mathbf{M}_{(N)}$  by

$$(\mathbf{M}_{(N)})_{\mathbf{k},\mathbf{l}} = \begin{cases} 0 & |\mathbf{k} - \mathbf{l}| > N \\ m_{\mathbf{k},\mathbf{l}} & |\mathbf{k} - \mathbf{l}| \leq N \end{cases}. \quad (3.3.52)$$

We define the remainder as  $\mathbf{R} = \mathbf{M} - \mathbf{M}_{(N)}$ . Then, by making estimates on the norm of  $\mathbf{M}_{(N)}$ ,  $\mathbf{M}_{(N)}^{-1}$ , and  $\mathbf{M} - \mathbf{M}_{(N)}$ , we will show that, by making the bandwidth of  $\mathbf{M}_{(N)}$  large enough we can make the perturbation  $\mathbf{M} - \mathbf{M}_{(N)}$  small enough so that the inverse of  $\mathbf{M}$  can be shown to exist. We construct the proof in a series of lemmas.

First, we claim that

**Lemma 3.3.4** *If  $\mathbf{M} \in \mathcal{I}_\alpha$ , then  $\|\mathbf{R}\|_{1,1} + \|\mathbf{R}\|_{\infty,\infty} \leq C(\mathbf{M}, \alpha)N^{-\alpha}$ .*

To prove this, we see that

$$\|\mathbf{R}\|_{1,1} \leq C \sup_{\|x\|_1=1} \sum_{\mathbf{k} \in \mathbf{Z}^2} \sum_{|\mathbf{k}-\mathbf{l}| > N} (1 + |\mathbf{k} - \mathbf{l}|)^{-2-\alpha} |x_1| \quad (3.3.53)$$

$$\leq \sum_{\mathbf{l} \in \mathbf{Z}, |\mathbf{l}| > N} (1 + |\mathbf{l}|)^{-2-\alpha} \quad (3.3.54)$$

by Young's inequality. We then write

$$\sum_{\mathbf{l} \in \mathbf{Z}, |\mathbf{l}| > N} (1 + |\mathbf{l}|)^{-2-\alpha} = \sum_{|l_1| > N} \sum_{l_2 \in \mathbf{Z}} (1 + |l_1| + |l_2|)^{-2-\alpha} + \quad (3.3.55)$$

$$\sum_{|l_1| \leq N} \sum_{|l_2| > N - |l_1|} (1 + |l_1| + |l_2|)^{-2-\alpha}. \quad (3.3.56)$$

We estimate

$$\begin{aligned} \sum_{|l_1| > N} \sum_{l_2 \in \mathbf{Z}} (1 + |l_1| + |l_2|)^{-2-\alpha} &= \sum_{|l_1| > N} \left( (1 + |l_1|)^{-2-\alpha} + 2 \sum_{l_2=1}^{\infty} (1 + |l_1| + |l_2|)^{-2-\alpha} \right) \\ &\leq \end{aligned}$$



$$\begin{aligned}
&= \frac{2}{1+\alpha}(1+N)^{-1-\alpha} + \frac{4}{\alpha+\alpha^2}(1+N)^{-\alpha} \\
&\leq CN^{-\alpha},
\end{aligned}$$

and

$$\begin{aligned}
\sum_{|l_1| \leq N} \sum_{|l_2| > N - |l_1|} (1 + |l_1| + |l_2|)^{-2-\alpha} &= 4 \sum_{0 \leq l_1 \leq N} \sum_{l_2 < l_1 - N} (1 + l_1 - l_2)^{-2-\alpha} \\
&\leq 4 \sum_{0 \leq l_1 \leq N} \int_{-\infty}^{l_1 - N} (1 + l_1 - y)^{-2-\alpha} dy \\
&= 4 \sum_{0 \leq l_1 \leq N} \frac{1}{1+\alpha} (1+N)^{-1-\alpha} \\
&\leq CN^{-\alpha},
\end{aligned}$$

and so we arrive at

$$\|\mathbf{R}\|_{1,1} \leq CN^{-\alpha}. \quad (3.3.57)$$

Additionally,  $\|\mathbf{R}\|_{\infty,\infty} = \|\mathbf{R}^*\|_{1,1}$ , and  $\|\mathbf{R}^*\|_{1,1}$  satisfies an inequality of the same form since each step in the derivation of the inequality is also valid for  $\mathbf{R}^*$ , so the result of Lemma 3.3.4 is proved.

We use the above Lemma and Theorem 3.2.2 to prove the following lemma:

**Lemma 3.3.5** *If  $\mathbf{M} \in \mathcal{I}_\alpha$  is invertible on  $l^p$  and  $l^{p'}$ , where  $1/p + 1/p' = 1$ , and  $p \leq p'$ , then for  $N$  large enough,  $\mathbf{M}_{(N)}$  is invertible on  $l^p$  and  $l^{p'}$ .*

**Proof:** By Theorem 3.2.2 and Lemma 3.3.4,

$$\|\mathbf{R}\|_{q,q} \leq \|\mathbf{R}\|_{1,1}^{1/q} \text{ and } L$$



so combining (3.3.68) and (3.3.69) yields

$$|(k_i - k'_i)^2 (\mathbf{M}_{(N)})_{k_1, k'_1, k_2, k'_2}| \leq C(p) N^{3 - \frac{2}{p}} (1 + |k_1 - k'_1| + |k_2 - k'_2|)^{-s} \quad (3.3.70)$$

provided  $k_1 \neq k'_1$ ,  $k_2 \neq k'_2$ ; this case may be handled by adjusting the constant  $C(p)$ . Thus, we have proved the inequality (3.3.62); the inequality (3.3.63) is proved similarly.  $\square$

**Lemma 3.3.8** *If  $\mathbf{M} \in \mathcal{I}_\alpha$  and  $\mathbf{M}$  is invertible*

Lemma 3.3.8 shows that  $[\mathbf{X}_i, [\mathbf{X}_i, \mathbf{M}_{(N)}]]$  is continuous from  $l^p$  to

finite sum of products between  $\mathbf{M}_{(N)}^{-1}$  and the commutators  $[\mathbf{X}_i, \mathbf{M}_{(N)}]$ ,  $[\mathbf{X}_i, [\mathbf{X}_i, \mathbf{M}_{(N)}]]$ , and  $[\mathbf{X}_i, [\mathbf{X}_i, [\mathbf{X}_i, \mathbf{M}_{(N)}]]]$ , it is continuous on  $l^2$  and its norm is less than  $C(\mathbf{M})N^J$ , where  $J \geq 3$ . Since the norm of an operator on  $l^2$  is also a bound on the elements of the matrix which represents the operator, we know that

$$|(\mathbf{M}_{(N)})_{\mathbf{k}, \mathbf{l}}^{-1}| \leq C(\mathbf{M}, \alpha)N^{27}(1 + |\mathbf{k} - \mathbf{l}|)^{-3}. \quad (3.3.81)$$

Now we can combine the inequalities (3.3.72) and (3.3.81) by raising one to the  $1 - \epsilon$  power and the other to the  $\epsilon$  power and taking their product; this yields

$$|(\mathbf{M}_{(N)})_{\mathbf{k}, \mathbf{l}}^{-1}| \leq C(\mathbf{M}, \alpha)N^{(3-\frac{2}{p})(1-\epsilon)}N^{J\epsilon}(1 + |\mathbf{k} - \mathbf{l}|)^{-2(1-\epsilon)-3\epsilon}, \quad (3.3.82)$$

which after simplification becomes the inequality (3.3.79).  $\square$

where  $\frac{1}{p_1} + \frac{1}{p_1'} = 1$ .

**Proof:** By Theorem 3.2.2,

$$\|\mathbf{M}_{(N)}^{-1}\|_{p_1, p_1} \leq \|\mathbf{M}_{(N)}^{-1}\|_{p, p}^t \|\mathbf{M}_{(N)}^{-1}\|_{1, 1}^t. \quad (3.3.87)$$

However,  $\|\mathbf{M}_{(N)}^{-1}\|_{p, p}$  is uniformly bounded in  $N$ , so

$$\|\mathbf{M}_{(N)}^{-1}\|_{p_1, p_1} \leq C(\mathbf{M}, \alpha) \|\mathbf{M}_{(N)}^{-1}\|_{1, 1}^t. \quad (3.3.88)$$

Now,

$$\|\mathbf{M}_{(N)}^{-1}\|_{p_1', p_1'} \leq C(\mathbf{M}, \alpha) \|\mathbf{M}_{(N)}^{-1}\|_{\infty, \infty}^t \quad (3.3.89)$$

by Theorem 3.2.2, so combining Proposition 3.3.3 with the inequalities (3.3.88) and (3.3.89) yields the result of this lemma.  $\square$

**Lemma 3.3.12**  $\mathbf{M}$  is invertible on  $l^{p_1}$  and  $l^{p_1'}$  as soon as  $(3 - \frac{2}{p})t < \alpha$ .

**Proof:** The inequality

$$\|\mathbf{M}_{(N)}^{-1} \mathbf{R}\|_{p_1, p_1} + \|\mathbf{M}_{(N)}^{-1} \mathbf{R}\|_{p_1', p_1'} \leq C(\mathbf{M}, \alpha) N^{(3 - \frac{2}{p})t - \alpha} (\log N)^t \quad (3.3.90)$$

follows immediately from the inequality (3.3.58) and Lemma 3.3.11. Thus, if  $(3 - \frac{2}{p})t < \alpha$

## CHAPTER 4

### CONCLUSIONS AND FURTHER DIRECTIONS

In Chapter 1, we described the multiresolution homogenization approach in the context of linear ODE's. We compared the multiresolution approach with existing approaches and found that classical results may be obtained using the multiresolution technique.

In Chapter 2, we discussed the generalization of the multiresolution approach to partial differential equations. We showed that the multiresolution reduction procedure preserves the sparsity of operators which are compressible in the wavelet basis, and also approximately preserves small eigenvalues of elliptic operators.

In Chapter 3, we proved results concerning algebras of bi-infinite matrices and tensors. We showed that bi-infinite matrices or tensors with exponential or polynomial decay away from the diagonal form an algebra under inversion, and that, if a matrix in either class is invertible, then its inverse is in that class as well. These results were used in Chapter 2.

In this chapter, we describe directions for future research.

#### 4.1 Multiresolu . Th u H H basis, H H h H H H g Ry CH H UG. H





we obtain

$$u^j(x, t) = \sum_n a_n \cos((\lambda_n^j)^1)$$





[16] M. Dorobantu. *Wavelet-based Algorithms for*

- [34] E.M. Stein and G. Weiss. *Introduction to Fourier Analysis on Euclidean Spaces*. Princeton University Press, 1971.
- [35] B. Z. Steinberg and J. J. McCoy. Towards effective parameter theories using multiresolution decomposition. *J. Acoust. Soc. Am.*, 96:1130–1143, 1994.
- [36] L. Tartar. *Compensated Compactness and Applications to Partial Differential Equations*. Heriot-Watt Sympos. vol. IV. Pitman, NY, 1979.
- [37] P. Tchamitchian. *Wavelets: T*

## Appendix A

# WAVELETS, MULTIREOLUTION ANALYSES, AND OPERATORS

### A.1 Wavelets and Multiresolution Analyses

In this section, we set our notation and give a brief description of the concept of multiresolution analysis (MRA) and wavelets. For details, we refer to e.g. [15].

#### A.1.1 Notation and Preliminary Considerations

As usual, we consider a chain of subspaces

$$\dots \subset \mathbf{V}_2 \subset \mathbf{V}_1 \subset \mathbf{V}_0 \subset \mathbf{V}_{-1} \subset \mathbf{V}_{-2} \subset \dots \quad (\text{A.1.1})$$

such that

$$\bigcap_j \mathbf{V}_j = \{0\} \text{ and } \overline{\bigcup_j \mathbf{V}_j} = \mathbf{L}^2(\Omega) \quad (\text{A.1.2})$$

where  $\Omega$  is some domain in  $\mathbf{R}^d$ . If the domain  $\Omega$  is bounded, then there is a coarsest space  $\mathbf{V}_0$  and, instead of (A.1.1), we write

$$\mathbf{V}_0 \subset \mathbf{V}_{-1} \subset \mathbf{V}_{-2} \subset \dots \quad (\text{A.1.3})$$

The subspace  $\mathbf{V}_j$  is spanned by an orthonormal basis  $\{\phi_k^j(x) = 2^{-j/2} \phi(2^{-j}x - k)\}_{k \in \mathbf{Z}}$ . The function  $\phi$  is called the scaling function, and it satisfies the two-scale difference equation

$$\phi(x/2) = \sqrt{2} \sum_k h_k \phi(x - k). \quad (\text{A.1.4})$$

We consider the space  $V_j$  to specify a *scale* or *resolution* of the space of  $\mathbf{L}^2$  functions on  $\Omega$ , and use the index  $j$  to identify the scale. As  $j \rightarrow \infty$ , the scale grows “coarser,” and as

$j \rightarrow -\infty$ , the scale grows “finer.” Functions in  $\mathbf{L}^2(\Omega)$  which are smooth or slowly-varying may be represented on a coarser scale of the MRA than those which are highly-oscillatory or have steep gradients.

We denote by  $\mathbf{W}_j$  the orthogonal complement of  $\mathbf{V}_j$  in  $\mathbf{V}_{j-1}$ ,  $\mathbf{V}_{j-1} = \mathbf{V}_j \oplus \mathbf{W}_j$  and use  $\mathbf{P}_j$  and  $\mathbf{Q}_j$  to denote the orthogonal projection operators onto  $\mathbf{V}_j$  and  $\mathbf{W}_j$ . Note that  $\mathbf{Q}_{j+1} = \mathbf{P}_j - \mathbf{P}_{j+1}$ . If  $x \in \mathbf{V}_j$ , we write  $s_x = \mathbf{P}_{j+1}x$  and  $d_x = \mathbf{Q}_{j+1}x$ , where  $s_x \in \mathbf{V}_{j+1}$  and  $d_x \in \mathbf{W}_{j+1}$ . If  $d = 1$ , then the subspace  $\mathbf{W}_j$  is spanned by an orthonormal basis  $\{\psi_k^j(x) = 2^{-j/2}\psi(2^{-j}x - k)\}_{k \in \mathbf{Z}}$ . The function  $\psi$  is called the wavelet, and it may be computed using the scaling function  $\phi$  via the two-scale relation

$$\psi(x/2) = \sqrt{2} \sum_k g_k \phi(x - k). \tag{A.1.5}$$

The space  $\mathbf{W}_{j+1}$  represents the “detail” component of the space  $\mathbf{V}_j$ , and the function  $\psi_k^j(x)$  captures the highly-oscillatory, quickly-varying component of functions in  $\mathbf{V}_j$ .

From (A.1.2), we see that

$$\mathbf{L}^2(\Omega) = \bigoplus_j \mathbf{W}_j. \tag{A.1.6}$$

If  $d \geq 2$ , then (for rectangular domains) the basis in the subspace  $\mathbf{W}_j$  may be constructed using products of wa

and

$$\psi(x) = \begin{cases} 1 & \text{if } 0 \leq x < \frac{1}{2} \\ -1 & \text{if } \frac{1}{2} \leq x < 1 \\ 0 & \text{otherwise.} \end{cases} \quad (\text{A.1.9})$$

There are many examples of wavelets with more vanishing moments. The Haar basis is the only (anti) symmetric orthogonal wavelet basis with compact support. Daubechies in [15] constructed compactly supported wavelets with more vanishing moments.

The sequences  $h$  and  $g$  in (A.1.4) and (A.1.5) may be either finite or infinite. Given a function

$$f^j(x) = \sum_k f_k^j \phi_k^j(x) \quad (\text{A.1.10})$$

in  $\mathbf{V}_j$ , we may compute the coefficients of  $s_f$  and  $d_f$  in the bases  $\phi^{j+1}$  and  $\psi^{j+1}$ . Using the relations (A.1.4) and (



In this context, we represent the projection  $s_f = \mathbf{P}_{j+1} f$  as application of the matrix  $H$  to the vector of coefficients of  $f$ .

Just as we have defined the wavelet decomposition, we may also define the wavelet reconstruction (or synthesis). Given functions  $s(x) \in \mathbf{V}_{j+1}$  and  $d(x) \in \mathbf{W}_{j+1}$  with coefficients  $s_k$  and  $d_k$ , we may write

$$f(x) = \sum_k f_k \phi_k^j(x) = s(x) + d(x), \quad (\text{A.1.14})$$

where

$$f_k = \sum_l (h_{k-2l} s_l + g_{k-2l} d_l). \quad (\text{A.1.15})$$

### A.1.2 Fourier Analysis

The sequences  $h$  and  $g$  are *filters* which are applied to sequences. In the Fourier domain, we represent these filters as trigonometric polynomials defined by

$$m_0(\xi) = \frac{1}{\sqrt{2}} \sum_k h_k e^{-ik\xi} \quad (\text{A.1.16})$$

and

$$m_1(\xi) = \frac{1}{\sqrt{2}} \sum_k g_k e^{-ik\xi}. \quad (\text{A.1.17})$$

We define  $\hat{g}$  to be the Fourier transform of the function  $g$  as follows:

$$\hat{g}(\xi) = \frac{1}{2\pi} \int e^{-i\xi x} g(x) dx. \quad (\text{A.1.18})$$

The relations (A.1.4) and (A.1.5) may be recast in the Fourier domain as

$$\hat{\phi}(\xi) = m_0(\xi/2) \hat{\phi}(\xi/2) \quad (\text{A.1.19})$$

and

$$\hat{\psi}(\xi) = m_1(\xi/2) \hat{\phi}(\xi/2). \quad (\text{A.1.20})$$

The functions  $m_0$  and  $m_1$  (and hence the sequences  $h$  and  $g$ ) define a *quadrature-mirror filter*. The function  $m_0$  is a “low-pass” filter and captures low-frequency components in the Fourier domain;  $m_1$  is a “high-pass” filter and captures the high-frequency components in the Fourier domain. We list here some properties of the functions  $m_0$  and  $m_1$ :

1.  $|m_0(\xi)|^2 + |m_1(\xi)|^2 = 1$
2.  $m_0(\xi) = 1$
3. If the wavelet has  $M$  vanishing moments, then

$$\left. \left( \frac{d}{d\xi} \right)^k m_1(\xi) \right|_{\xi=0} = 0 \quad \text{for } k = 0, \dots, M-1 \quad (\text{A.1.21})$$

$$\left. \left( \frac{d}{d\xi} \right)^k m_0(\xi) \right|_{\xi=0} = 0 \quad \text{for } k = 1, \dots, M-1 \quad (\text{A.1.22})$$

The convolution-decimation operations in (A.1.11) and (A.1.12) may be represented in the Fourier domain. Note that, via (A.1.10), we derive

$$\hat{f}(\xi) = \tilde{f}(\xi)\hat{\phi}(\xi). \quad (\text{A.1.23})$$

where  $\tilde{f}(\xi) = \sum_k f_k^j e^{-ik\xi}$ . From this we obtain

$$\hat{s}_f(\xi) = \tilde{s}_f(\xi)\hat{\phi}(\xi/2) \quad (\text{A.1.24})$$

and

$$\hat{d}_f(\xi) = \tilde{d}_f(\xi)\hat{\phi}(\xi/2), \quad (\text{A.1.25})$$

where

$$\tilde{s}_f(2\xi) = m_0(\xi)\tilde{f}(\xi) + m_1(\xi)\tilde{f}(\xi)$$

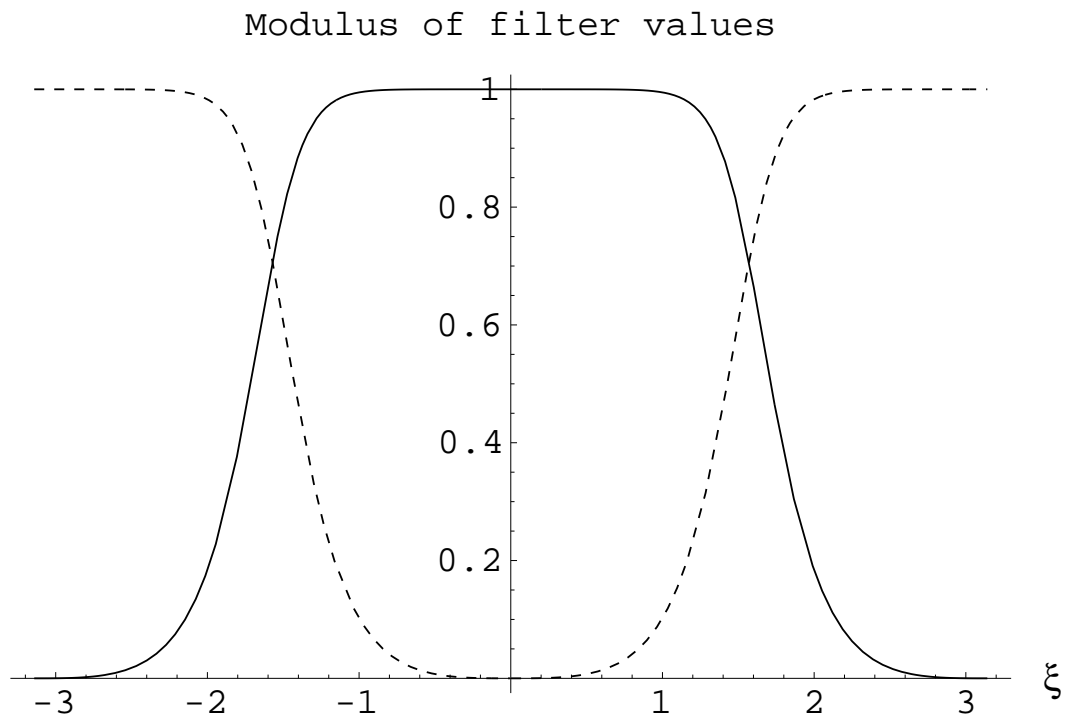


Figure A.1: Plots on  $[-\pi, \pi]$  of the modulus of the low-pass filter  $m_0(\xi)$  (solid line) and the high-pass filter  $m_1(\xi)$  (dashed line). In this example the filters were derived from non-compact orthogonal spline wavelets of degree 2. Clearly,  $m_0$  low-frequency and  $m_1$  is high-frequency.

## A.2 Operators in the Wavelet Basis

In the wavelet basis, certain classes of operators may be represented using matrices with relatively few significant coefficients, so that fast application of matrices and matrices and matrices to vectors may be achieved (see [9] for more details).

In this section we set the notation for representation of operators in the wavelet basis and describe some features of the “non-standard form.”

### A.2.1 Notation and Preliminary Considerations

Given a bounded linear operator  $\mathbf{S}$  on  $\mathbf{L}^2(\mathbf{R}^d)$ , consider its projection  $\mathbf{S}_j$  on  $\mathbf{V}_j$ ,  $\mathbf{S}_j = \mathbf{P}_j \mathbf{S} \mathbf{P}_j$ . Since  $\mathbf{V}_j$  is a subspace spanned by translations of  $\phi^j$ , we may represent the operator  $\mathbf{S}_j$  as a (possibly infinite) matrix in that basis. With a slight abuse of notation, we will use the same symbol  $\mathbf{S}_j$  to represent both the operator and its matrix. Since  $\mathbf{V}_j = \mathbf{V}_{j+1} \oplus \mathbf{W}_{j+1}$ , we may also write  $\mathbf{S}_j : \mathbf{V}_j \rightarrow \mathbf{V}_j$  in a block form

$$\mathbf{S}_j = \begin{pmatrix} \mathbf{A}_{\mathbf{S}_j} & \mathbf{B}_{\mathbf{S}_j} \\ \mathbf{C}_{\mathbf{S}_j} & \mathbf{T}_{\mathbf{S}_j} \end{pmatrix} : \mathbf{V}_{j+1} \oplus \mathbf{W}_{j+1} \rightarrow \mathbf{V}_{j+1} \oplus \mathbf{W}_{j+1}, \quad (\text{A.2.1})$$

where

$$\begin{aligned} \mathbf{A}_{\mathbf{S}_j} &= \mathbf{Q}_{j+1} \mathbf{S}_j \mathbf{Q}_{j+1}, \\ \mathbf{B}_{\mathbf{S}_j} &= \mathbf{Q}_{j+1} \mathbf{S}_j \mathbf{P}_{j+1}, \\ \mathbf{C}_{\mathbf{S}_j} &= \mathbf{P}_{j+1} \mathbf{S}_j \mathbf{Q}_{j+1}, \\ \mathbf{T}_{\mathbf{S}_j} &= \mathbf{P}_{j+1} \mathbf{S}_j \mathbf{P}_{j+1}. \end{aligned} \quad (\text{A.2.2})$$

We note that  $\mathbf{T}_{\mathbf{S}_j} = \mathbf{S}_{j+1}$ . E

and then apply the wavelet decomposition again to the columns of that matrix.

The operators (and their matrix representations) in (A.2.2) are referred to as the  $\mathbf{A}$ ,  $\mathbf{B}$ ,  $\mathbf{C}$ , and  $\mathbf{T}$  *blocks* of  $\mathbf{S}_j$ . Also, for an operator  $\mathbf{Z}$ , we use notation  $\mathbf{A}_Z$ ,  $\mathbf{B}_Z$ ,  $\mathbf{C}_Z$ , and  $\mathbf{T}_Z$  to indicate its blocks.

### A.2.2 The Standard and Non-Standard Forms

In the standard form, the idea is to represent the operator  $\mathbf{S}_j$  in the coordinates of the decomposition (A.1.13). If  $d = 1$ , this is done by applying the wavelet decomposition over several scales to the rows and columns of the matrix  $\mathbf{S}_j$ . Thus, the standard form is simply a representation of the original matrix  $\mathbf{S}_j$  in a different basis.

The non-standard form (see e.g. [9]) is an alternative representation of the operator which is not (strictly speaking) a representation of the matrix in a different basis. We start with the telescoping series

$$\mathbf{S} = \sum_j \mathbf{P}_{j-1} \mathbf{S} \mathbf{P}_{j-1} - \mathbf{P}_j \mathbf{S} \mathbf{P}_j. \quad (\text{A.2.3})$$

We note that  $\mathbf{Q}_j = \mathbf{P}_{j-1} - \mathbf{P}_j$  and rewrite this series as

$$\mathbf{S} = \sum_j \mathbf{A}_j + \mathbf{B}_j + \mathbf{C}_j \quad (\text{A.2.4})$$

where  $\mathbf{A}_j = \mathbf{Q}_j \mathbf{S} \mathbf{Q}_j$ ,  $\mathbf{B}_j = \mathbf{Q}_j \mathbf{S} \mathbf{P}_j$ , and  $\mathbf{C}_j = \mathbf{P}_j \mathbf{S} \mathbf{Q}_j$ . If  $\mathbf{S}_{j_0}$  is a bounded operator on  $\mathbf{V}_{j_0}$ , then we see that

$$\mathbf{S}_{j_0} = \left( \sum_{j=j_0}^{j_1} \mathbf{A}_j + \mathbf{B}_j + \mathbf{C}_j \right) + \mathbf{P}_{j_1} \mathbf{S}_{j_0} \mathbf{P}_{j_1}. \quad (\text{A.2.5})$$

This equation shows that we may think of the operator  $\mathbf{S}_{j_0}$  as a coarse-scale component  $\mathbf{P}_{j_1} \mathbf{S}_{j_0} \mathbf{P}_{j_1}$  together with interactions between successively finer scales. Slightly more work is required to apply the operator in this form to another operator or a vector.

If the operator  $\mathbf{S}$  is an integral operator with kernel  $K(x, y)$ , then the entries of the matrices which represent the operators  $\mathbf{A}_j$ ,  $\mathbf{B}_j$ , and  $\mathbf{C}_j$  are given by

$$a_{k, k'}^j = \int_{\Omega} \int_{\Omega} K(x, y) \psi_k^j(x) \psi_{k'}^j(y) dx$$

$$b_{k,k'}^j = \int_{\Omega} \int_{\Omega} K(x,y) \psi_k^j(x) \phi_{k'}^j(y) dx dy \quad (\text{A.2.7})$$

$$c_{k,k'}^j = \int_{\Omega} \int_{\Omega} K(x,y) \phi_k^j(x) \psi_{k'}^j(y) dx dy \quad (\text{A.2.8})$$