

2008

# Approximating Solutions to Differential Equations via Fixed Point Theory

Douglas Rizzolo  
*Harvey Mudd College*

---

## Recommended Citation

Rizzolo, Douglas, "Approximating Solutions to Differential Equations via Fixed Point Theory" (2008). *HMC Senior Theses*. 213.  
[https://scholarship.claremont.edu/hmc\\_theses/213](https://scholarship.claremont.edu/hmc_theses/213)

This Open Access Senior Thesis is brought to you for free and open access by the HMC Student Scholarship at Scholarship @ Claremont. It has been accepted for inclusion in HMC Senior Theses by an authorized administrator of Scholarship @ Claremont. For more information, please contact [scholarship@cuc.claremont.edu](mailto:scholarship@cuc.claremont.edu).



# Approximating Solutions to Differential Equations via Fixed Point Theory

**Douglas Rizzolo**

Jon Jacobsen, Advisor

Francis E. Su, Reader

May, 2008

**HARVEY MUDD**  
COLLEGE

Department of Mathematics

Copyright © 2008 Douglas Rizzolo.

The author grants Harvey Mudd College the nonexclusive right to make this work available for noncommercial, educational purposes, provided that this copyright statement appears on the reproduced materials and notice is given that the copying is by permission of the author. To disseminate otherwise or to republish requires written permission from the author.

# Abstract

In the study of differential equations there are two fundamental questions: is there a solution? and what is it? One of the most elegant ways to prove that an equation has a solution is to pose it as a fixed point problem, that is, to find a function  $f$  such that  $x$  is a solution if and only if  $f(x) = x$ . Results from fixed point theory can then be employed to show that  $f$  has a fixed point. However, the results of fixed point theory are often nonconstructive: they guarantee that a fixed point exists but do not help in finding the fixed point. Thus these methods tend to answer the first question, but not the second. One such result is Schauder's fixed point theorem. This theorem is broadly applicable in proving the existence of solutions to differential equations, including the Navier-Stokes equations under certain conditions. Recently a semi-constructive proof of Schauder's theorem was developed in Rizzolo and Su (2007). In this thesis we go through the construction in detail and show how it can be used to search for multiple solutions. We then apply the method to a selection of differential equations.



# Contents

<b>Abstract</b>	<b>iii</b>
<b>Acknowledgments</b>	<b>xi</b>
<b>1 Sperner's Lemma and Search Algorithms</b>	<b>1</b>
1.1 Sperner's Lemma . . . . .	1
1.2 A Basic Search Algorithm . . . . .	6
1.3 Van der Laan and Talman's Basic Algorithm . . . . .	9
1.4 Multiple Completely Labeled Simplices I . . . . .	11
<b>2 Schauder's Fixed Point Theorem</b>	<b>13</b>
2.1 The Original Construction . . . . .	13
2.2 The Construction on the Hilbert Cube . . . . .	19
2.3 Multiple Completely Labeled Simplices II . . . . .	20
<b>3 Hammerstein Integral Equations</b>	<b>23</b>
3.1 Underlying Theory . . . . .	23
3.2 The Second Order Initial Value Problem . . . . .	26
3.3 The Two Point Boundary Value Problem . . . . .	32
3.4 A Note on Uniform Convergence . . . . .	34
<b>4 Future Research</b>	<b>37</b>
<b>Bibliography</b>	<b>41</b>



# List of Figures

- 1.1 An example of the graph  $\Gamma_2$  . . . . . 4
  
- 3.1 (a) A 4 term approximation (blue) plotted against the solution. (b) A 4 term approximation (blue) plotted against the 4 term Fourier expansion to the solution. (c) A 10 term approximation (blue) against the solution. (d) A 10 term approximation (blue) against the 10 term Fourier expansion of the solution. . . . . 29
- 3.2 A Sequence of Ten Term Approximations . . . . . 31
- 3.3 Approximate solutions to the generalized Duffing's equation 34
- 3.4 Approximate solutions (pink) versus actual solutions . . . . 35





# List of Tables

1.1	Pivot Rules for $K_2(m)$ . . . . .	8
1.2	Pivot Rules for van der Laan and Talman's Algorithm . . . .	10



# Acknowledgments

I would like to thank my advisor, Jon Jacobsen, for his help throughout this project, especially in keeping me focused. I would also like to thank my second reader, Francis E. Su, without whom the theoretical backbone of this project might never have been developed.



# Chapter 1

## Sperner's Lemma and Search Algorithms

In this chapter we will go through a constructive proof of Schauder's fixed point theorem. Before we can get into the construction, we need some terminology. This terminology will be used to build up Sperner's Lemma, which forms the basis for our construction. In the following let  $X$  be a real vector space.

### 1.1 Sperner's Lemma

The first result we will need is a result called Sperner's Lemma, which is a lemma concerning the properties of labelings of triangulations of simplices. Below we go through the terminology needed to make these concepts rigorous and we end with a proof of Sperner's Lemma. The first idea we will need is that of *affine* combinations.

**Definition 1.1.** *An affine combination of  $\{x_1, \dots, x_n\} \subset X$  is a sum of the form  $a_1x_1 + \dots + a_nx_n$  such that  $a_1 + \dots + a_n = 1$ . The set  $A \subset X$  is affinely independent if no element in  $A$  is an affine combination of other elements in  $A$ .*

From this definition it is easy to see that a linearly independent set is also affinely independent since all affine combinations are linear combinations. Let  $\text{aff}(A)$  to denote the affine span of  $A$ , the set of all affine combinations of elements of  $A$ . We use  $\text{aff}(A)$  to define the *relative boundary* of  $A$  by  $\partial A = \bar{A} \cap \overline{(\text{aff}(A) \setminus A)}$ . Intuitively the relative boundary of a set  $A$  is the boundary of  $A$  with respect to the lowest dimensional affine space

## 2 Sperner's Lemma and Search Algorithms

---

containing  $A$ . It is fairly easy to see that if  $\text{aff}(A) = X$  then  $\partial A$  is equal to the boundary of  $A$  with the usual definition.

**Definition 1.2.** *The convex hull of the set  $A \subset X$  is intersection of all convex subsets of  $X$  containing  $A$ . This is denoted by  $\text{conv}(A)$ . It is easy to show that*

$$\text{conv}(\{x_1, \dots, x_n\}) = \left\{ \sum_{i=1}^n a_i x_i \mid a_i \geq 0 \text{ and } \sum_{i=1}^n a_i = 1 \right\}.$$

With these two definitions in hand we can now define a simplex — a structure that is vital to our construction of fixed points.

**Definition 1.3.** *Let  $\{x_1, \dots, x_{n+1}\}$  be an affinely independent set of vectors in  $\mathbb{R}^m$  (note this implies that  $m \geq n + 1$ ). The set  $\text{conv}(\{x_1, \dots, x_{n+1}\})$  is called an  $n$ -simplex, often denoted by  $\sigma = \langle x_1, \dots, x_{n+1} \rangle$ . The points  $\{x_1, \dots, x_{n+1}\}$  are the vertices of  $\sigma$ . If  $m = n + 1$  and  $x_i = e_i$  for all  $i$  (where  $e_i$  is the  $i^{\text{th}}$  standard basis vector in  $\mathbb{R}^{n+1}$ ) we call  $\text{conv}(\{x_1, \dots, x_{n+1}\})$  the standard  $n$ -simplex, denoted  $\Delta^n$ .*

A simplex  $\tau$  is a *face* of the simplex  $\sigma$  if the vertices of  $\tau$  are a subset of the vertices of  $\sigma$ ;  $\tau$  is also referred to as an  $i$ -face, where  $i$  is the number of vertices of  $\tau$ . Furthermore, the  $(n - 1)$ -faces of  $\sigma$  are called *facets* of  $\sigma$ . If  $\tau$  is a facet of  $\sigma$  whose vertex set is missing the vertex  $x_i$  we say that  $\tau$  is the facet opposite  $x_i$  and write  $\tau = \sigma_i$ . In order to approximate fixed points we will need to use the concept of a *triangulation* of an  $n$ -simplex, defined as follows:

**Definition 1.4.** *A triangulation (or subdivision) of  $\sigma$  is a finite collection  $\mathcal{T}$  of  $n$ -simplices such that the following two conditions hold:*

- i)  $\bigcup_{\tau \in \mathcal{T}} \tau = \sigma$ .
- ii) If  $\tau_1, \tau_2 \in \mathcal{T}$  then either  $\tau_1 \cap \tau_2 = \emptyset$  or  $\tau_1 \cap \tau_2$  is a facet of both  $\tau_1$  and  $\tau_2$ .

We will often need to refer to faces of simplices in  $\mathcal{T}$ , so we define  $\mathcal{T}^i$  to be the collection of  $i$ -faces of simplices in  $\mathcal{T}$  and  $\mathcal{T}^+$  to be the collection of all faces of simplices in  $\mathcal{T}$ . It is worth noting that  $\mathcal{T}^0$  is the collection of all vertices of simplices in  $\mathcal{T}$ . There are several important properties of triangulations that we summarize in the following lemma:

**Lemma 1.1.** *Let  $\mathcal{T}$  be a triangulation of the simplex  $C$ . Then*

1. *If  $\tau$  is a facet of  $\sigma_1 \in \mathcal{T}$  then either  $\tau \in \partial C$  or there exists exactly one other simplex  $\sigma_2 \in \mathcal{T}$  that has  $\tau$  as a facet.*

2. If  $D$  is a facet of  $C$  then the collection

$$\mathcal{T}_D = \{\tau \mid \tau \subseteq D \text{ and } \tau \in \mathcal{T}^{n-1}\}$$

is a triangulation of  $D$ .

We will not prove this lemma here, but a proof using the equivalent definition of a triangulation by open simplices can be found in Todd (1976) and a proof using subdivisions of abstract simplicial complexes can be found in Spanier (1966).

Let  $\mathcal{T}$  be a triangulation of the simplex  $\sigma$ . A *labeling* of  $\mathcal{T}$  is a map  $\ell : \mathcal{T}^0 \rightarrow \mathbb{N}$ . With this we define a special type of labeling:

**Definition 1.5.** Let  $\mathcal{T}$  be a triangulation of  $\sigma = \text{conv}(\{x_1, \dots, x_{n+1}\})$ . A *Sperner labeling*  $\ell$  of  $\mathcal{T}$  is a labeling with the property that if  $x \in \mathcal{T}^0 \cap \text{conv}(\{x_i \mid i \in J\})$  for some  $J \subseteq \{1, \dots, n+1\}$  then  $\ell(x) \in J$ .

Notice that an immediate consequence of the definition is that  $\ell(x_i) = i$ . Intuitively what this definition says is that a vertex  $x$  of a simplex in  $\mathcal{T}$  carries that same label as one of the vertices of  $\sigma$  on the lowest dimensional face of  $\sigma$  containing  $x$ . A simplex  $\tau \in \mathcal{T}$  is called *completely labeled* (c.l.) if

$$\{\ell(x) \mid x \in \mathcal{T}^0 \cap \tau\} = \{1, \dots, n+1\},$$

that is, if the vertices of  $\tau$  carry all of the labels of the vertices of  $\sigma$ . Furthermore a simplex in  $\mathcal{T}$  or  $\mathcal{T}^{n-1}$  is called *almost completely labeled* (a.c.l.) if its vertices carry the labels  $\{1, \dots, n\}$ . Notice that if a simplex is c.l. then it is also a.c.l.

From here we will follow the exposition in Todd (1976) to prove Sperner's Lemma. The proof of Sperner's Lemma will be based on the properties of the following graph:

**Definition 1.6.** Let  $\mathcal{T}$  be a triangulation of the  $n$ -simplex  $C$  and  $\ell$  a Sperner labeling of  $\mathcal{T}$ . The nodes of the graph  $\Gamma_n$  are the c.l.  $n$ -simplices of  $\mathcal{T}$ , the a.c.l.  $n$ -simplices of  $\mathcal{T}$ , and the a.c.l.  $(n-1)$ -simplices of  $\mathcal{T}$  in  $\partial C$ . Furthermore two nodes  $x$  and  $y$  are adjacent in  $\Gamma_n$  if as simplices one is a face of the other or if they share an a.c.l. face.

Figure 1.1 gives an example of the graph  $\Gamma_2$  for a particular Sperner labeling. The red circles are the nodes and the red lines are the edges. Looking at this example, we see that the connected components of  $\Gamma_2$  are paths either connecting two a.c.l. simplices in the boundary, two c.l. simplices, or an a.c.l. simplex in the boundary to a c.l. simplex. In the following lemma we show that these, in addition to a special type of cycle, are the only behaviors that the components of  $\Gamma_n$  can display.



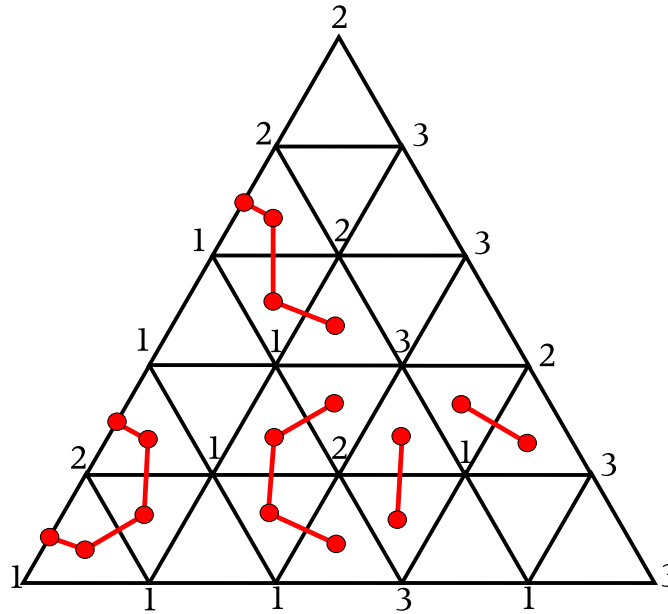


Figure 1.1: An example of the graph  $\Gamma_2$

**Lemma 1.2.** *Each connected component of  $\Gamma_n$  has one of the following forms:*

1. *A cycle whose nodes are a.c.l. but not c.l.  $n$ -simplices.*
2. *A path whose intermediate nodes are a.c.l. but not c.l.  $n$ -simplices and each of whose endpoints is either*
  - (a) *an a.c.l.  $(n - 1)$ -simplex in  $\partial C$  or*
  - (b) *a c.l.  $n$ -simplex.*

*Proof.* To prove this lemma it is clearly sufficient to prove that a node in  $\Gamma_n$  has degree 1 if it is a c.l.  $n$ -simplex or an a.c.l.  $(n - 1)$ -simplex in  $\partial C$  and degree 2 if it is an a.c.l., but not c.l.,  $n$ -simplex. Thus we consider these three cases.

1. Suppose that the node  $x$  is a c.l.  $n$ -simplex. Since  $x$  has  $n + 1$  vertices and  $n + 1$  distinct labels, there is exactly one facet  $\tau$  of  $x$  that is a.c.l. If  $\tau$  is in  $\partial C$  then  $x$  is adjacent to only the node  $\tau$  in  $\Gamma_n$ . Otherwise, there is exactly one other simplex in  $\mathcal{T}$  with  $\tau$  as a facet. Clearly this simplex is a.c.l., and  $x$  is adjacent to only this node. In either case,  $x$  has degree 1.

2. Suppose that the node  $x$  is an a.c.l. but not c.l.  $n$ -simplex. Since  $x$  is a.c.l., it has an a.c.l. facet  $\tau$ . If  $\tau$  is in  $\partial C$ , then  $x$  is adjacent to the node  $\tau$ . Otherwise here is exactly one other simplex in  $\mathcal{T}$  with  $\tau$  as a facet. Clearly this simplex is a.c.l., and  $x$  is adjacent to this simplex. Let  $v$  be the vertex in  $x$  opposite the a.c.l. facet  $\tau$ . Since  $x$  is not c.l. the label of  $v$  is equal to the label of exactly one of the vertices of  $\tau$ . Let  $\tau_v$  be the facet of  $x$  attained when the vertex in  $\tau$  with the same label as  $v$  is replaced by  $v$ . Then, using the same argument as above,  $x$  is either adjacent to the node  $\tau_v$  or to the other simplex in  $\mathcal{T}$  with  $\tau_v$  as a facet. Since  $x$  has no other a.c.l. facets besides  $\tau$  and  $\tau_v$ , it follows that  $x$  is not adjacent to any other nodes of  $\Gamma_n$ . Therefore  $x$  has degree 2.
3. Suppose that  $x$  is an a.c.l.  $(n - 1)$ -simplex in  $\partial C$ . Then  $x$  is a facet of exactly one simplex in  $\mathcal{T}$  and  $x$  is adjacent only to this simplex, and thus has degree 1.

□

This Lemma will be the workhorse in our proof of Sperner's Lemma, which we are now prepared to state and prove.

**Theorem 1.1** (Sperner's Lemma). *Let  $\mathcal{T}$  be a triangulation of  $\Delta^n$  and  $\ell$  a Sperner labeling of  $\mathcal{T}$ . Then  $\mathcal{T}$  contains an odd number of c.l. simplices. In particular,  $\mathcal{T}$  contains at least one c.l. simplex.*

*Proof.* We prove this by induction on  $n$ . When  $n = 0$ ,  $\Delta^0$  is just a point, and thus the theorem is trivially true. Suppose that the theorem holds for  $\Delta^{n-1}$ . Let  $\mathcal{T}$  be a triangulation of  $\Delta^n$  and  $\ell$  a Sperner labeling of  $\mathcal{T}$ . By Lemma 1.1 we know that  $\mathcal{T}_{\Delta^n}$  is a triangulation of  $\Delta^n$ . Furthermore,  $\ell$  is a Sperner labeling of  $\mathcal{T}_{\Delta^n}$  and it is easy to see that there is a linear homeomorphism  $h$  from  $\Delta^n$  to  $\Delta^{n-1}$ . Since linear homeomorphisms preserve simplices, we have a triangulation  $\mathcal{T}_1$  of  $\Delta^{n-1}$  given by

$$\mathcal{T}_1 = \{h(\sigma) \mid \sigma \in \mathcal{T}_{\Delta^n}\},$$

together with the Sperner labeling  $\ell_1 = \ell \circ h^{-1}$ . By the induction hypothesis there are an odd number of c.l. simplices in  $\mathcal{T}_1$ , and thus an odd number of a.c.l.  $(n - 1)$ -simplices in the triangulation of  $\Delta^n$ . By the definition of a Sperner labeling, every a.c.l.  $(n - 1)$ -simplex of  $\mathcal{T}$  in  $\partial\Delta^n$  is in  $\Delta^n$ .

Consider the connected components of the graph  $\Gamma_n$  that are paths. There are three choices for where its endpoints might be. Each path that has both endpoints in  $\partial\Delta^n$  accounts for two of the a.c.l.  $(n - 1)$ -simplices

of  $\mathcal{T}$ . Thus the total number of a.c.l.  $(n - 1)$ -simplices in  $\partial\Delta^n$  accounted for by components that are paths with both endpoints in  $\partial\Delta^n$  is even. Since there are an odd number of a.c.l.  $(n - 1)$ -simplices in  $\partial\Delta^n$  and each one is an endpoint of a path in  $\Gamma_n$ , an odd number of them must be endpoints of paths whose other endpoint is a c.l.  $n$ -simplex. Hence there are an odd number of c.l.  $n$ -simplices in  $\mathcal{T}$  that are endpoints of paths in  $\Gamma_n$  whose other endpoint is in  $\partial\Delta^n$ .

Now, suppose that there is a connected component in  $\Gamma$  that is a path that has neither endpoint in  $\partial\Delta^n$ . Then both endpoints of this path are c.l.  $n$ -simplices. Hence there are an even number of c.l.  $n$ -simplices that are endpoints of paths in  $\Gamma_n$  that have neither endpoint in  $\partial\Delta^n$ . Since the number of c.l.  $n$ -simplices in  $\mathcal{T}$  is equal to the number of c.l.  $n$ -simplices that are endpoints of components that are paths that have exactly one endpoint in  $\partial\Delta^n$  plus the number of c.l.  $n$ -simplices that are endpoints of components that are paths with no endpoints in  $\partial\Delta^n$ , and this sum is an odd number plus an even number, we conclude that there are an odd number of c.l.  $n$ -simplices in  $\mathcal{T}$ .  $\square$

## 1.2 A Basic Search Algorithm

Notice that in the proof of Sperner's Lemma we were able to find an odd number of c.l. simplices that were endpoints of paths that had one endpoint in  $\partial\Delta^n$ . Thus, if we could find the a.c.l. simplices on the boundary of  $\Delta^n$ , we could then trace the paths in  $\Gamma_n$  with them as endpoints to find a c.l. simplex in  $\Delta^n$ . Hence this proof suggests how to design an algorithm to find c.l. simplices in  $\Delta^n$ . Indeed many algorithms have been constructed in this manner; for surveys see Todd (1976) and Talman (1980). Below, we present one of the most basic algorithms, which is the one that follows the method suggested by the proof of Sperner's Lemma.

In order to introduce the search algorithm, we first need to introduce a particular triangulation of  $\Delta^n$ . While the algorithm can be given abstractly, independent of the particular triangulation, it is easier to understand (and the notation is easier) if a particular triangulation is chosen. The triangulation we choose is called Kuhn's triangulation; we use the definition given in Todd (1976) and refer to this text for the proofs of the properties of this triangulation.

**Definition 1.7.** Let  $Q$  be the  $(n + 1) \times n$  matrix

$$Q = \begin{bmatrix} -1 & & & \mathbf{0} \\ 1 & \ddots & & \\ & & \ddots & -1 \\ \mathbf{0} & & & 1 \end{bmatrix}$$

and let  $q^j$  be its  $j$ 'th column. Define  $K_0^2(m) = \{y \in \Delta^n \mid my_i \in \mathbb{Z}\}$  for  $m \in \mathbb{Z}$ . Let  $\pi$  be a permutation of  $\{1, \dots, n\}$ . For  $y^0 \in K_0^2(m)$  we define  $\sigma$  to be the simplex with vertices  $\{y^0, \dots, y^n\}$  where  $y^i = y^{i-1} + m^{-1}q^{\pi(i)}$  for  $i \geq 1$ . If  $\sigma \subseteq \Delta^n$ , we define  $k_2(y^0, \pi) = \sigma$ . The Kuhn triangulation of  $\Delta^n$ , denoted  $K_2(m)$ , is defined to be the collection of all such  $k_2(y^0, \pi)$ .

Notice that the permutation  $\pi$  can be treated like a vector in  $\mathbb{R}^n$ , i.e.,  $\pi^* = (\pi_1, \dots, \pi_n)$  with  $\pi(i) = \pi_i$ . Henceforth we will make no distinction between a permutation and its vector representation.

Given a triangulation  $\mathcal{T}$  of the simplex  $C$ , the mesh size of  $\mathcal{T}$  is defined by

$$\text{mesh}_p(\mathcal{T}) = \sup_{\sigma \in \mathcal{T}} \sup_{x, y \in \sigma} \|x - y\|_p.$$

Intuitively, the mesh size of  $\mathcal{T}$  is furthest apart two elements in  $\Delta^n$  can be given that they are in the same simplex in  $\mathcal{T}$ . Of course, the distance between two points depends on the norm being used, hence the  $p$ -subscript. Given that all norms on  $\mathbb{R}^n$  are equivalent, one might think that the  $p$ -subscript could be omitted entirely, but which  $p$  we are using will become important in later computations. The triangulation  $K_2(m)$  has the property that  $\text{mesh}_\infty(K_2(m)) = m^{-1}$  and  $\text{mesh}_2(K_2(m)) = m^{-1}\sqrt{n+1}$ . Thus the mesh size of the Kuhn triangulation can be made arbitrarily small with respect to any metric whose topology is equivalent to the norm-induced topology on  $\mathbb{R}^n$ .

In the proof of Sperner's Lemma the paths in  $\Gamma_n$  that we would need to follow in order to find a c.l. simplex involve a process called pivoting: that is, given a simplex  $\sigma$  and a facet  $\gamma$  of  $\sigma$ , we must find the other simplex with  $\gamma$  as a facet. Precisely, let  $\sigma = \langle x^0, \dots, x^n \rangle = k_2(x^0, \pi)$  be given and suppose we wish to obtain a simplex  $\tau = \langle y^0, \dots, y^n \rangle = k_2(y^0, \rho)$  such that  $\tau$  has all of the vertices of  $\sigma$  except  $x^i$ . One of the benefits of Kuhn's triangulation is that the rules for this type of pivoting are remarkably simple, and are summarized in Table 2.1.

Let  $\Delta_{(j)}^n = \{x \in \Delta^n \mid x_j = \dots = x_{n+1} = 0\}$  for  $j = 2, \dots, n+1$ . Notice that  $\Delta_{(j)}^n$  naturally corresponds to the simplex  $\Delta^{j-2}$ . The idea of the algo-

	$y^0$	$\pi$
$i = 0$	$y^0 + m^{-1}q^{\pi(1)}$	$(\pi(2), \dots, \pi(n), \pi(1))$
$0 < i < n$	$y^0$	$(\pi(1), \dots, \pi(i+1), \pi(i), \dots, \pi(n))$
$i = n$	$y^0 - m^{-1}q^{\pi(n)}$	$(\pi(n), \pi(1), \dots, \pi(n-1))$

 Table 1.1: Pivot Rules for  $K_2(m)$ 

rithm will be to take the graphs  $\Gamma_j$  as defined above for all of the  $\Delta_{(j)}^n$  and to connect them in a fashion such that there is a path with  $e_1$  as one endpoint and a c.l.  $n$ -simplex as the other endpoint. The construction of this graph and the proof that it has the above property can be found in Kuhn (1969). While the algorithm is described abstractly in Kuhn (1969), we give the concrete interpretation of the algorithm given in Todd (1976).

**Algorithm 1.1.** *Let  $\ell$  be a Sperner Labeling of the triangulation  $K_2(m)$  of  $\Delta^n$ .*

*Step 1 Let  $L = 2$ ,  $y^0 = e_1$ ,  $\sigma_1 = k_2(y^0, \pi)$  with  $\pi = (1)$ , the permutation of  $\{1\}$ . Let  $y^+ = y^1 = e_1 + m^{-1}q^1$ , and set  $j = 1$ .*

*Step 2 If  $\ell(y^+) = L$ , go to Step 4. Otherwise, the label of  $y^+$  duplicates the label of some vertex  $y^-$  of  $\sigma_j$ .*

*Step 3 Let  $\tau$  be the facet of  $\sigma_j$  opposite  $y^-$ . If  $\tau \subseteq \Delta_{(L)}^n$ , go to Step 5. Otherwise, let  $\sigma_{j+1}$  be the unique  $(L-1)$ -simplex in  $\Delta_{(L+1)}^n$  sharing the facet  $\tau$  with  $\sigma_j$ . Set  $j \rightarrow j+1$  and return to Step 2.*

*Step 4 (Increasing Dimension) We have the  $(L-1)$ -simplex  $\sigma_j = k_2(y^0, \pi)$ , say, with  $\pi$  a permutation of  $\{1, \dots, L-1\}$  and the vertices of  $\sigma_j$  have all of the labels  $\{1, \dots, L\}$ . If  $L = n+1$ , stop since  $\sigma_j$  is a c.l. simplex. Otherwise, let  $\sigma_{j+1} = k_2(y^0, \pi')$  where  $\pi' = (\pi(1), \dots, \pi(L-1), L)$ . Let  $y^+$  be the new vertex of  $\sigma_{j+1}$ . Set  $j \rightarrow j+1$ ,  $L \rightarrow L+1$ , and return to Step 2.*

*Step 5 (Decreasing Dimension) We have that  $(L-1)$ -simplex  $\sigma_j = k_2(y^0, \pi) = \langle y^0, \dots, y^{L-1} \rangle$  with  $\pi$  a permutation of  $\{1, \dots, L-1\}$ . We know that  $\sigma_j$  has a facet in  $\Delta_{(L)}^n$  and it is clear that  $\tau$  must be  $\langle y^0, \dots, y^{L-2} \rangle$  and  $\pi(L-1) = L-1$ . The vertices of  $\tau$  have the labels  $1, \dots, L-1$ . Let  $\sigma_{j+1} = k_2(y_0, \pi')$  with  $\pi' = (\pi(1), \dots, \pi(L-2))$ . Set  $j \rightarrow j+1$ ,  $L \rightarrow L-1$ , and return to Step 3.*

This is one of the most basic algorithms and, computationally speaking, it is very inefficient. This is due to the fact that the size of the triangulation

is set before the algorithm starts and the starting point for the algorithm is independent of the problem under consideration. Therefore, if we start with a small mesh size it is potentially very time consuming for the algorithm to terminate since the algorithm is taking small steps and might start far away from a c.l. simplex. However, this algorithm is not without its advantages. One of the primary advantages, as we will discuss in Section 1.4, is that this algorithm can easily be modified to search for multiple c.l. simplices.

### 1.3 Van der Laan and Talman's Basic Algorithm

The problems with the algorithm in Section 1.2, noted at the end of that section, were widely recognized and several algorithms were designed to circumvent them (see e.g. Todd (1976) and Talman (1980)). In this section we present one such algorithm. This algorithm was originally presented in van der Laan and Talman (1979), but we will follow the exposition in Talman (1980). All of the notation from Section 1.2 is carried over into this section with the change that  $\pi$  is now a permutation of a subset of  $\{1, \dots, n + 1\}$  rather than a permutation of a subset of  $\{1, \dots, n\}$  and

$$q^{n+1} = - \sum_{i=1}^n q^i,$$

( $q^{n+1}$  was previously undefined). Additionally, we need the following definition:

**Definition 1.8.** *Let  $\mathcal{T}$  be a triangulation of  $\Delta^n$  with Sperner labeling  $\ell$ . Let  $P$  be a subset of  $\{1, \dots, n + 1\}$  with  $|P| = p$ . A  $(s - 1)$ -simplex  $\sigma$  is called  $P$ -complete if  $\{\ell(x) \mid x \in \mathcal{T}_0 \cap \sigma\} = P$ .*

In this algorithm, we will need a slightly different table for determining how to exchange simplices because we will index differently, and we also need to track a vector  $R \in \mathbb{R}^{n+1}$ . The new table is Table 2.2.

Let  $\mathcal{T} = K_2(m)$  and fix  $y \in \mathcal{T}_0$ . We are now prepared to give the algorithm.

**Algorithm 1.2.**

*Step 1* Set  $p = 0$ ,  $P = \emptyset$ ,  $\pi = \emptyset$ ,  $y^1 = y$ ,  $\sigma = k_2(y^1, \pi)$ ,  $\bar{y} = y^1$ , and  $R = 0 \in \mathbb{R}^{n+1}$ .

	$y^1$	$\pi$	$R$
$i = 1$	$y^1 + m^{-1}q^{\pi(1)}$	$(\pi(2), \dots, \pi(p), \pi(1))$	$R + e_{\pi(1)}$
$2 \leq i \leq p$	$y^1$	$(\pi(1), \dots, \pi(i+1), \pi(i), \dots, \pi(p))$	$R$
$i = p + 1$	$y^1 - m^{-1}q^{\pi(p)}$	$(\pi(p), \pi(1), \dots, \pi(p-1))$	$R - e_{\pi(p)}$

Table 1.2: Pivot Rules for van der Laan and Talman's Algorithm

- Step 2* If  $\ell(\bar{y}) \notin P$ , go to step 4. Otherwise,  $\ell(\bar{y}) = \ell(y^s)$  for exactly one vertex  $y^s \neq \bar{y}$  of  $\sigma$ . The facet  $\tau$  opposite  $y^s$  is  $P$ -complete.
- Step 3* If  $s = p + 1$  and  $R_{\pi(p)} = 0$  go to step 5. Otherwise adapt  $\sigma$  according to Table 2.2 replacing the vertex  $y^s$ . Return to step 2 with  $\bar{y}$  the new vertex of  $\sigma$ .
- Step 4* If  $p = n$ ,  $\sigma$  is c.l. and the algorithm terminates. Otherwise, set  $P = P \cup \{\ell(\bar{y})\}$ ,  $\pi = (\pi, \ell(\bar{y}))$ ,  $\sigma = k_2(y^1, \pi)$ , and  $p = p + 1$ , and  $\bar{y} = y^{p+1}$ . Return to step 2.
- Step 5* Set  $P = P \setminus \{\pi(p)\}$ ,  $\pi = (\pi(1), \dots, \pi(p-1))$ ,  $\sigma = k_2(y^1, \pi)$ , and  $t = t - 1$ . Return to step 3 with  $y^s$  being the vertex of  $\sigma$  with label equal to the integer removed from  $P$  are the beginning of this step.

While the steps that this algorithm goes through, pivoting between simplices and increasing and decreasing dimension until a c.l. simplex is found, are very similar to the previous algorithm, this one is in fact very different. One major difference is that this algorithm can start at an arbitrary point in  $\mathcal{T}_0$ , but the trade off for this freedom is the  $R$  vector that we need to keep track of. The purpose of this vector is to track where the algorithm is relative to the boundary. If the algorithm gets too close to the boundary, in a vague sense of the word close, this vector forces the algorithm to move away from the boundary. Thus, in some non-rigorous sense, this algorithm selects for c.l. simplices that are further from the boundary. Because this algorithm can start at an arbitrary point in  $\mathcal{T}_0$ , it can be used as a restart algorithm. This means that we can run the algorithm with a coarse mesh size to find a c.l. simplex and then use a vertex of that simplex as the starting point for the algorithm on a more refined mesh size, which is a tremendous advantage. Moreover, this restart procedure can be applied iteratively until the mesh size is as small as desired. This method allows the algorithm to move rapidly to parts of  $\Delta^n$  where c.l. simplices are likely to be and then restart the search in that area with a finer mesh. Thus we

expect this algorithm to be substantially faster than the previous one and, indeed, this expectation was realized in our computational experience.

## 1.4 Multiple Completely Labeled Simplices I

Notice that, in the proof of Sperner's Lemma, we actually proved that a subdivision of a Sperner Labeled simplex contains an odd number of completely labeled simplices. However, the search algorithms that we have presented so far search only for a single completely labeled simplex. Thus, it is natural to wonder whether or not these algorithms can be extended to search for multiple completely labeled simplices. Furthermore, this is more than just an idle curiosity, such an extension would have important ramifications for approximating solutions to differential equations.

For example, if we knew that an equation had multiple solutions such an algorithm could potentially be used to approximate more than one solution. Also, perhaps more importantly, if we are working with an equation where the number of solutions is unknown, such an algorithm could be employed to gather experimental evidence as to whether or not multiple solutions exist. Of course, as mathematicians, we deal in the market of proofs, not evidence, but experimental evidence can be extremely useful for formulating conjectures and seeing which direction the theory should go.

One way to accomplish this is to modify the search algorithms. In the VT-algorithm this can be accomplished by changing the initial point and hoping that this yields an approximation to a different solution. This, however, is rather haphazard and we would like a more systematic approach. For this we go to the basic algorithm. Recall that one of the properties of the graph  $\Gamma_n$  is that if you started at a c.l. simplex, then the other end of the component is either another c.l. simplex or an a.c.l. simplex on the boundary. Now, the key insight is to notice that any permutation of the labels in a Sperner labeling is again a Sperner labeling, and this permutation clearly gives rise to a new graph  $\tilde{\Gamma}_n$ . Hence, in order to look for multiple c.l. simplices, we use the basic algorithm to follow a path in  $\Gamma_n$  to a c.l. simplex and then permute the labels and follow that path in  $\tilde{\Gamma}_n$  starting at the c.l. simplex we found until we either find another c.l. simplex or end back at the boundary. This can be repeated at each new c.l. simplex found and for each permutation of the Sperner labeling. Unfortunately, we are not the first to notice this approach; it was first introduced in Jeppson (1972) and was also presented Allgower (1977), though the latter restricts to cyclic



permutations. However, what these two papers overlook is that the permutation can be applied before the initial search is begun. The algorithm can then be started at the vertex with permuted label 1 and move towards the one with permuted label 2. By making this modification one makes the search even more thorough, though efficiency is sacrificed.

For this thesis, we do not implement the most thorough search algorithm. Rather, we limit ourselves to searching through one cyclic permutation from each c.l. simplex found. There is no theoretical advantage to this, but it does save computing time and is easier to implement. In addition we also use another approach that is not based on altering the search algorithm intrinsically. The details of this method are presented in Section 2.3.

## Chapter 2

# Schauder's Fixed Point Theorem

### 2.1 The Original Construction

In this section we will show how Sperner's Lemma can be used to give a semi-constructive proof of Schauder's fixed point theorem. Our construction follows that given in Rizzolo and Su (2007). As a result, we will be able to use the algorithm in Section 1.2 to approximate the fixed points guaranteed by Schauder's theorem. Let us first recall the statement of Schauder's fixed point theorem:

**Theorem 2.1** (Schauder's Fixed Point Theorem). *Let  $B$  be a compact convex subset of the normed space  $X$  and  $f : B \rightarrow B$  a continuous function. Then  $f$  has a fixed point.*

Giving a constructive proof of this theorem in a general normed space is a fairly intractable problem, so the first thing we must do is reduce the problem to a more reasonable setting. It seems to be a general fact about mathematics that it is hard to construct things exactly. The same is true for fixed point problems. Fortunately, the following lemma tells us that all will be right in the world if we can construct a point that is "almost" fixed.

**Lemma 2.1.** *Let  $(M, d)$  be a compact metric space. Suppose that  $f : M \rightarrow M$  is continuous and that for every  $\epsilon > 0$  there exists  $x_\epsilon \in M$  such that  $d(f(x_\epsilon), x_\epsilon) < \epsilon$ . Then  $f$  has a fixed point.*

The proof of this is straightforward and can be found in a number of places (e.g., Rizzolo and Su (2007) or Smart (1974)). The basic idea is to take

a sequence  $\epsilon_n \rightarrow 0$  and for each  $n$  choose  $x_n$  such that  $d(f(x_n), x_n) < \epsilon_n$ . Since  $M$  is compact,  $\{x_n\}$  has a subsequence that converges to the fixed point. This lemma provides the first important reduction. Rather than constructing a fixed point exactly, it will be sufficient to construct points that are approximately fixed. Rigorously, a point  $x$  is an  $\epsilon$ -fixed point of  $f$  if  $d(f(x), x) < \epsilon$ ; we will construct an  $\epsilon$ -fixed point for arbitrary  $\epsilon$ .

Now, a general normed space is a fairly general set to be working in, and in order to create a reasonable construction, we must work in a more specific setting. It turns out that one of the most natural settings to work in is  $\mathbb{R}^\infty$  with the product topology. It is well known that this topology can be metrized with the metric

$$\bar{d}(x, y) = \sum_{i=1}^{\infty} \frac{|x_i - y_i|}{2^i(1 + |x_i - y_i|)},$$

and that  $(\mathbb{R}, \bar{d})$  is a complete metric space. Since Sperner's Lemma applies to finite dimensional simplices, we will need to extend this idea into infinite dimensions. It seems that the most natural extension would be to define the standard infinite dimension simplex as the convex hull of the standard basis vectors in  $\mathbb{R}^\infty$ . Unfortunately, the set that results from this is not compact. Hence we define the standard (closed) infinite dimensional simplex  $\Delta_0^\infty$  to be the closure of this set. We can write this explicitly as

$$\Delta_0^\infty = \left\{ x \in \mathbb{R}^\infty \mid \sum_{i=1}^{\infty} x_i \leq 1 \text{ and } x_i \geq 0 \right\}.$$

While it might seem that restricting our considerations to  $\Delta_0^\infty$  limits the scope of our construction, it turns out that it is theoretically easy to convert fixed points in  $\Delta_0^\infty$  into fixed points in an arbitrary compact convex subset of normed space. We summarize this result in the following theorem:

**Theorem 2.2.** *Suppose that  $B$  is a compact convex subset of a normed space. Then there exists a homeomorphism  $h : B \rightarrow \Delta_0^\infty$ . Furthermore, if  $f : B \rightarrow B$  is continuous, then  $g = h \circ f \circ h^{-1} : \Delta_0^\infty \rightarrow \Delta_0^\infty$  is continuous and if  $x$  is a fixed point of  $g$  then  $h^{-1}(x)$  is a fixed point of  $f$ .*

The only part of this theorem that is nontrivial is the existence of  $h$ . This part is proved in Rizzolo and Su (2007), and it is a minor generalization of a theorem proved in Klee (1955). Hence by only considering the case that  $f : \Delta_0^\infty \rightarrow \Delta_0^\infty$  we are not losing any theoretical generality. However, in order to apply our construction, it does mean that we will need to find the

homeomorphism  $h$ . We are now ready to carry out the construction. The remainder of this section is largely lifted from Rizzolo and Su (2007), with a few corrections and some clarifications.

We will use  $F^n$  to denote the convex hull of the first  $n + 1$  standard basis vectors in  $\mathbb{R}^\infty$ , and we will call  $F^n$  a face of  $\Delta_0^\infty$ . It is clear that there is a natural linear bijection between  $F^n$  and  $\Delta^n$ . The following lemma has as a corollary the fact that this bijection is a homeomorphism.

**Lemma 2.2.** *Let  $A$  be a bounded subset of the normed space  $(\mathbb{R}^n, \|\cdot\|_\infty)$ . On  $A$ , the metric*

$$\bar{d}_n = \sum_{i=1}^n \frac{|x_i - y_i|}{2^i(1 + |x_i - y_i|)}$$

*is equivalent to the metric induced by the norm  $\|\cdot\|_\infty$ .*

*Proof.* Suppose that  $x, y \in \mathbb{R}^n$ . We see that

$$\bar{d}_n(x, y) = \sum_{i=1}^n \frac{|x_i - y_i|}{2^i(1 + |x_i - y_i|)} \leq n\|x - y\|_\infty.$$

Now, since  $A$  is bounded, there is some  $M$  such  $\|x - y\|_\infty \leq M$  for  $x, y \in A$ . Thus we see that

$$\frac{\|x - y\|_\infty}{2^n(1 + M)} \leq \frac{\|x - y\|_\infty}{2^n(1 + \|x - y\|_\infty)} \leq \bar{d}_n(x, y),$$

which implies that

$$\|x - y\|_\infty \leq 2^n(1 + M)\bar{d}_n(x, y). \quad (2.1)$$

Thus  $\bar{d}_n$  is equivalent to the metric induced by the norm on  $A$ .  $\square$

Though technical, the proof of this lemma is important because the bound established in Equation 2.1 will be needed in the proof of Schauder's theorem, which we are now ready to present.

*Proof of Schauder's Theorem.* From Lemma 2.1 it is sufficient to prove that  $f$  has an  $\epsilon$ -fixed point for arbitrary  $\epsilon > 0$ . Let  $\epsilon > 0$  be given and choose  $N \geq \log_2(2/\epsilon) + 1$ . Notice that for  $x, y \in \Delta_0^\infty$ , this implies that

$$\sum_{i=N+1}^{\infty} \frac{|x_i - y_i|}{2^i(1 + |x_i - y_i|)} \leq \sum_{i=N+1}^{\infty} \frac{1}{2^i} < \frac{\epsilon}{2}. \quad (2.2)$$

Since  $f$  maps between countably infinite-dimensional spaces, we can write  $f$  in terms of its components:  $f(x) = (f_1(x), f_2(x), \dots)$ . Since  $f$  is continuous,  $f_i$  is continuous for each  $i$ . Notice that for each  $x \in \mathbb{R}^n$  we can write  $x = \sum_{i=1}^{\infty} x_i e_i$ . We define the function  $P_N(x)$  by  $P_N(x) = \sum_{i=1}^N x_i e_i$ , that is,  $P_N$  is the projection onto the first  $N$  coordinates, which is clearly continuous. Now, consider the function

$$G(x) = (G_1(x), G_2(x), \dots) = (f_1(x), f_2(x), \dots, f_N(x), 1 - \sum_{i=1}^N f_i(x), 0, 0, 0, \dots),$$

and let

$$g(x) = (g_1(x), g_2(x), \dots) = (G \circ P_N)(x).$$

Since each  $f_i$  is continuous and finite sums of continuous function are continuous,  $g_i$  is continuous for each  $i$ . Furthermore, we see that  $g : F^N \rightarrow F^N$ . Consequently,  $g$  is continuous.

Let  $\epsilon_0 = \frac{\epsilon}{8(N+1)}$  and  $\epsilon_1 = \frac{\epsilon}{2^{N+5}(N+1)}$ . Since  $g$  is continuous on a compact set, it is uniformly continuous. Thus there exists  $\delta_1 > 0$  such that  $\bar{d}(x, y) < \delta_1$  implies that  $\bar{d}(g(x), g(y)) < \epsilon_1$ . Let  $\delta = \min(\delta_1, \epsilon_1)$ . Since  $F^N$  can be triangulated with an arbitrarily small triangulation, let  $\mathcal{T}$  be a triangulation with  $\text{mesh}(\mathcal{T}) < \delta$ . Label the vertices of  $\mathcal{T}$  with the map

$$\ell(x) = \operatorname{argmax}_{\{i | x_i \neq 0\}} (x_i - g_i(x)).$$

Recall that the *argmax* function returns the index of the largest element of the argument, and if there are multiple indices that give the maximum value, the *argmax* function returns the least of these indices.

Observe that  $\ell(x)$  produces a Sperner labeling on the vertices of  $\mathcal{T}$ . Thus by Sperner's Lemma, there exists a c.l. simplex in  $\mathcal{T}$ . This simplex can be found using the algorithm in Section 1.2. Let  $\{x^1, x^2, \dots, x^{N+1}\}$  be the vertices of this simplex where the index of each vertex is its Sperner label. From this, we see that for all  $j$ ,

$$x_i^i - g_i(x^i) \geq x_j^i - g_j(x^i).$$

Furthermore, since for each  $x$  in  $F^N$ , we have

$$\sum_{j=1}^{N+1} x_j = \sum_{j=1}^{N+1} g_j(x) = 1,$$

there is at least one  $j$  such that  $g_j(x) \leq x_j$ . In particular, since  $\ell(x^i) = i$ , this implies that for each  $x^i$ ,

$$x_i^i - g_i(x^i) = \max_j (x_j^i - g_j(x^i)) \geq 0.$$

Since  $\text{mesh}(\mathcal{T}) < \delta$  we have that, for all  $i$ ,  $\bar{d}(x^1, x^i) < \delta$ . From the bound (2.1) in Lemma 2.2 (note that in this case  $M = 1$  and  $n = N + 1$ ), we find that for all  $i, j$ ,

$$|x_j^1 - x_j^i| < 2^{N+2}\delta \leq 2^{N+2}\epsilon_1 \leq \epsilon_0. \quad (2.3)$$

By the same logic, we have that for all  $i, j$ ,

$$|g_j(x^1) - g_j(x^i)| < 2^{N+2}\epsilon_1 \leq \epsilon_0. \quad (2.4)$$

Consequently, we have that

$$x_j^1 + \epsilon_0 > x_j^i \quad \text{and} \quad -g_j(x^i) < \epsilon_0 - g_j(x^1)$$

which, in turn, implies that

$$2\epsilon_0 + x_j^1 - g_j(x^1) > x_j^i - g_j(x^i)$$

for all  $i$  and  $j$ . In particular, this implies that the following list of inequalities hold (simply let  $i = j$  and run through all  $i$ ):

$$\begin{array}{rcl} 2\epsilon_0 + x_1^1 - g_1(x^1) & > & x_1^1 - g_1(x^1) & \geq 0, \\ 2\epsilon_0 + x_2^1 - g_2(x^1) & > & x_2^2 - g_2(x^2) & \geq 0, \\ \vdots & & \vdots & \\ 2\epsilon_0 + x_{N+1}^1 - g_{N+1}(x^1) & > & x_{N+1}^{N+1} - g_{N+1}(x^{N+1}) & \geq 0. \end{array}$$

Summing down each column yields the following inequality.

$$2\epsilon_0(N+1) + \sum_{i=1}^{N+1} x_i^1 - \sum_{i=1}^{N+1} g_i(x^1) > \sum_{i=1}^{N+1} (x_i^i - g_i(x^i)) \geq 0.$$

Now we recall that for all  $i$ ,  $x_i^i - g_i(x^i) \geq 0$  and

$$\sum_{i=1}^{N+1} x_i^1 - \sum_{i=1}^{N+1} g_i(x^1) = 1 - 1 = 0.$$

Consequently,

$$\begin{aligned} 2\epsilon_0(N+1) &= 2\epsilon_0(N+1) + \sum_{i=1}^{N+1} x_i^1 - \sum_{i=1}^{N+1} g_i(x^1) \\ &> \sum_{i=1}^{N+1} (x_i^i - g_i(x^i)) \\ &= \sum_{i=1}^{N+1} |x_i^i - g_i(x^i)|. \end{aligned}$$

Using (2.3) and (2.4) and the continuity of  $g$ , for all  $i$ , we have that:  $|x_i^1 - g_i(x^1)| \leq |x_i^1 - x_i^i| + |x_i^i - g_i(x^i)| + |g_i(x^i) - g_i(x^1)| < 2\epsilon_0 + |x_i^i - g_i(x^i)|$ . Hence,

$$\begin{aligned} \bar{d}(x^1, g(x^1)) &= \sum_{i=1}^{N+1} \frac{|x_i^1 - g_i(x^1)|}{2^i(1 + |x_i^1 - g_i(x^1)|)} \leq \sum_{i=1}^{N+1} |x_i^1 - g_i(x^1)| \\ &< \sum_{i=1}^{N+1} (2\epsilon_0 + |x_i^i - g_i(x^i)|) \\ &< 4(N+1)\epsilon_0 \\ &= \frac{\epsilon}{2}. \end{aligned}$$

Let  $y = (x_1^1, x_2^1, \dots, x_N^1, 0, 0, 0, \dots)$ . We see that

$$\begin{aligned} \sum_{i=1}^N \frac{|y_i - f_i(y)|}{2^i(1 + |y_i - f_i(y)|)} &= \sum_{i=1}^N \frac{|y_i - g_i(y)|}{2^i(1 + |y_i - g_i(y)|)} \\ &= \sum_{i=1}^N \frac{|x_i^1 - g_i(x^1)|}{2^i(1 + |x_i^1 - g_i(x^1)|)} \\ &\leq \sum_{i=1}^{N+1} \frac{|x_i^1 - g_i(x^1)|}{2^i(1 + |x_i^1 - g_i(x^1)|)} \\ &< \frac{\epsilon}{2}. \end{aligned} \tag{2.5}$$

From (2.2) and (2.5), we have

$$\begin{aligned} \bar{d}(y, f(y)) &= \sum_{i=1}^{\infty} \frac{|y_i - f_i(y)|}{2^i(1 + |y_i - f_i(y)|)} \\ &= \sum_{i=1}^N \frac{|y_i - f_i(y)|}{2^i(1 + |y_i - f_i(y)|)} + \sum_{i=N+1}^{\infty} \frac{|y_i - f_i(y)|}{2^i(1 + |y_i - f_i(y)|)} \\ &< \frac{\epsilon}{2} + \frac{\epsilon}{2} \\ &= \epsilon. \end{aligned}$$

Therefore,  $y$  is the desired  $\epsilon$ -fixed point.  $\square$

This proof provides a clear framework for approximating fixed points of functions defined in infinite dimensional normed spaces. However, while the framework is clear, there are still several obstacles that must be overcome in order to use this method of approximation. Primarily, given a

compact convex subset  $B$  of a normed space, we still need to find a homeomorphism  $h : B \rightarrow \Delta_0^\infty$ . Unfortunately, such homeomorphisms are not the most natural maps to construct. Thus, in practice, we would like the construction to take place on a nicer set than  $\Delta_0^\infty$ . The set we will use in practice is  $I^\infty = \prod_{i=1}^\infty [-1, 1] \subset \mathbb{R}^\infty$ , which is commonly referred to as the Hilbert Cube. We also introduce the definitions  $I^n = \prod_{i=1}^n [-1, 1]$  and  $I_+^n = \prod_{i=1}^n [0, 1]$ .

## 2.2 The Construction on the Hilbert Cube

Before proving that the construction can be done on  $I^\infty$ , we will show explicitly how to construct the homeomorphism  $h : B \rightarrow I^\infty$  in the case where  $B$  has a special form. To do this we will make use of the following result:

**Theorem 2.3.** *Let  $H$  be a Hilbert Space with inner product  $\langle \cdot, \cdot \rangle$  and maximal orthonormal set  $\{u_\alpha \mid \alpha \in A\}$ . Define  $\hat{x}(\alpha) = \langle x, u_\alpha \rangle$ . Then  $\mathcal{F}(x) = \hat{x}$  is a linear isometry of  $H$  onto  $\ell^2(A)$ .*

We will not prove this here, but it is a classical result and the interested reader can find a proof in Rudin (1987). In particular, we will later be interested in the case where  $H = L^2[a, b]$  and the  $\{u_\alpha\}$  are a Fourier basis.

**Theorem 2.4.** *Let  $\{\delta_n\}$  be a sequence of real numbers such that  $\delta_n \geq 0$  and  $\sum \delta_n^2 < \infty$ . Define  $B_{\{\delta_n\}} \subset \ell^2(\mathbb{N})$  by*

$$B_{\{\delta_n\}} = \{x \in \ell^2(\mathbb{N}) \mid |x_n| \leq \delta_n\}.$$

*Then the function  $K = (K_1, K_2, \dots)$  from  $B$  to  $I^\infty$  defined by  $K_n(x) = \frac{1}{\delta_n} x_n$  is a homeomorphism.*

*Proof.* That  $K$  is bijective is trivial and that  $K$  is continuous follows immediately from the facts that  $I^\infty$  has the product topology and the projection of  $K$  onto each factor is continuous. Furthermore,  $I^\infty$  is Hausdorff and it is an elementary exercise in Hilbert space theory to prove that  $B_{\{\delta_n\}}$  is compact (see e.g. Rudin (1987) Exercise 4.6 or Lang (1993) Exercise 5.3). Hence  $K$  is a continuous bijection from a compact space to a Hausdorff space and is therefore a homeomorphism.  $\square$

Combining these two theorems yields the following proposition:

**Proposition 2.1.** *Suppose that  $H$  is a Hilbert space with maximal orthonormal set  $\{u_i\}_{i=1}^\infty$ . Let  $\{\delta_n\}$  be as above and  $A = \mathcal{F}^{-1}(B_{\{\delta_n\}})$ . Then  $K \circ \mathcal{F} : A \rightarrow I^\infty$  is a homeomorphism.*



The situation covered by this proposition might seem contrived, but it will be central to the results in Chapter 3. Now that we have shown the ease with which we can construct homeomorphisms to  $I^\infty$ , at least in special cases, let us see how the construction on  $\Delta_0^\infty$  can be modified for  $I^\infty$ .

**Theorem 2.5.** *Let  $f : I^\infty \rightarrow I^\infty$  be continuous and  $\epsilon > 0$ . Then there exists  $x \in I^\infty$  such that  $\bar{d}(x, f(x)) < \epsilon$ .*

*Proof.* As before we choose  $N$  such that

$$\sum_{i=N+1}^{\infty} \frac{|x_i - y_i|}{2^i(1 + |x_i - y_i|)} \leq \sum_{i=N+1}^{\infty} \frac{1}{2^i} < \frac{\epsilon}{2},$$

for all  $x, y \in I^\infty$ . Using this, we need only construct an  $\epsilon/2$  fixed point of the restriction of  $f$  to  $I^\infty \cap \text{span}\{e_1, \dots, e_N\}$ . Notice that there is a natural homeomorphism from  $I^\infty \cap \text{span}\{e_1, \dots, e_N\}$  to  $I^N$  given by restriction to the first  $N$  coordinates, i.e.,  $(x_1, \dots, x_N, 0, 0, \dots) \mapsto (x_1, \dots, x_N)$ . Thus it is sufficient to provide a construction on  $I^N$ . This construction is equivalent to providing an explicit homeomorphism from  $I^n$  to  $\Delta^n$  that is valid for all  $n \in \mathbb{N}$ . Fix  $n$  and define  $v \in I^n$  by  $v = (1, \dots, 1)$ . Let  $S(x) = \frac{1}{2}(x + v)$ , which is clearly a homeomorphism from  $I^n$  to  $I_+^n$ . Let  $A = \{x \in \prod_{i=1}^n [0, 1] \mid \sum x_i \leq 1\}$  and define  $h : I_+^n \rightarrow A$  by

$$h(x) = \begin{cases} \frac{\max x_i}{\sum x_i} x & x \neq 0, \\ 0 & x = 0. \end{cases}$$

It is straightforward to show that  $h$  is a homeomorphism. Finally, define  $T : A \rightarrow \Delta^n$  by

$$T(x) = \left( x_1, \dots, x_n, 1 - \sum_{i=1}^n x_i \right).$$

Again,  $T$  is easily seen to be a homeomorphism. Thus we have that  $T \circ h \circ S$  is a homeomorphism from  $I^n$  to  $\Delta^n$ , as was to be shown. Using this homeomorphism, the approximate fixed points we can construct in  $\Delta^n$  using Sperner's Lemma translate to approximate fixed points in  $I^n$ .  $\square$

### 2.3 Multiple Completely Labeled Simplices II

When we previously discussed searching for multiple c.l. simplices we looked at intrinsically changing the search algorithm. In this section, we discuss a different approach; namely we actually change the function we

are looking for a fixed point of. At first this might seem like a strange idea — after all, we are assuming that we are given a function  $f : I^\infty \rightarrow I^\infty$  to look for a fixed point of, how can we change it? Notice though that the search actually takes place on  $\Delta^n$  and we make use of a homeomorphism  $F : I^n \rightarrow \Delta^n$ . However, there are many such homeomorphisms  $F$  and in choosing these we can actually change the function that we are searching for a fixed point of without changing the resulting point's property as an approximate fixed point of  $f$ .

Recall in the proof of Theorem 2.5 we constructed the homeomorphism from  $I^n \rightarrow \Delta^n$  by composing several intermediary homeomorphisms, one of which was  $S : I^n \rightarrow I_+^n$  given by  $S(x) = \frac{1}{2}(x + v)$  with  $v = (1, 1, \dots, 1) \in I^n$ . This is the homeomorphism that we will modify. Let  $\sigma : \{1, 2, \dots, n\} \rightarrow \{0, 1\}$  and define  $v^\sigma$  by  $v_i^\sigma = (-1)^{\sigma(i)}$ . We then define  $S^\sigma = (S_1^\sigma, \dots, S_n^\sigma) : I^n \rightarrow I_+^n$  by  $S_i^\sigma(x) = \frac{(-1)^{\sigma(i)}}{2}(x_i + v_i^\sigma)$ . Thus for each choice of  $\sigma$  we get a distinct homeomorphism  $S^\sigma : I^n \rightarrow I_+^n$ . Using the notation from Theorem 2.5, we get the distinct homeomorphisms  $T \circ h \circ S^\sigma : I^n \rightarrow \Delta^n$ . Since these homeomorphisms are distinct, running the search algorithm with different choices of  $S^\sigma$  can give rise to approximations of different fixed points. We will see an example of this in Section 3.3.



## Chapter 3

# Hammerstein Integral Equations

In this chapter we consider equations of the form

$$u(t) = f(t) + \int_a^b H(t, y)F(y, u(y)) dy, \quad (3.1)$$

which are known as Hammerstein integral equations. In this equation the functions  $f : [a, b] \rightarrow \mathbb{R}$ ,  $H : [a, b] \times [a, b] \rightarrow \mathbb{R}$  and  $F : [a, b] \times \mathbb{R} \rightarrow \mathbb{R}$  are given and we are trying to find the function  $u : [a, b] \rightarrow \mathbb{R}$  satisfying (3.1). The operator corresponding to this problem is

$$\Phi(u)(t) = f(t) + \int_a^b H(t, y)F(y, u(y)) dy. \quad (3.2)$$

The strategy for solving these equations is to show that  $\Phi$  has a fixed point.

### 3.1 Underlying Theory

Fixed point approximation methods have been applied to this type of problem before, notably in Jeppson (1972), Chen (1977), and Allgower (1977). In all of these cases the assumptions have essentially been that  $f = 0$ ,  $F$  and  $H$  are continuous, and  $F$  is bounded. Under these assumptions the existence of a solution to (3.1) is a consequence of Theorem ???. In each of these papers much use was made of the particular form of Hammerstein integral equations and there is not a clear way to generalize the methods they use. In contrast, we have developed a broad framework for approximating

fixed points and are treating this type of equation as a special case. As it happens, upon restriction to these equations our method becomes almost indistinguishable from that in Chen (1977). There are, however, several differences. First, in Chen (1977),  $H$  must be continuous (in fact it must satisfy a weak differentiability condition as well), while we will develop the method assuming only that  $H \in L^2$ . Furthermore, in Chen (1977) the proof of convergence reads as a happy accident with no deep reason for why the method works. Indeed, the link to Schauder's theorem goes completely unmentioned, only Brouwer's theorem is used. As we will show, however, convergence of the method is an almost trivial corollary to the constructive proof of Schauder's theorem.

Before going into specific examples, let us see how we can use Schauder's theorem to guarantee solutions to Hammerstein integral equations.

**Definition 3.1.** *Let  $M_1$  and  $M_2$  be metric spaces and  $f : M_1 \rightarrow M_2$  a continuous function. We say that  $f$  is compact if  $f(\overline{A})$  is compact whenever  $A$  is bounded.*

This condition is linked to Schauder's theorem by the following result.

**Theorem 3.1.** *Suppose that  $B$  is a closed, bounded, convex subset of a normed space  $X$  and  $f : X \rightarrow X$  is compact. Further suppose that  $f(B) \subset B$ . Then there is a compact convex subset  $A$  of  $B$  such that  $f : A \rightarrow A$ .*

This result is a consequence of a classical result, known as Mazur's theorem, which says that the intersection of all closed convex sets containing a given compact set is compact. Thus, if we have a compact map and a closed bounded set that gets mapped to itself, then there is a compact convex set that gets mapped to itself and Schauder's theorem can be applied. We now build up the proof that the operator in Equation 3.2 is compact.

**Theorem 3.2.** *Suppose that  $H \in L^2([a, b] \times [a, b])$ . Then the function*

$$\Psi(u)(x) = \int_a^b H(x, y)u(y) dy$$

*defines a compact operator  $L^2[a, b] \rightarrow L^2[a, b]$ .*

One of the most direct ways to prove this is to approximate  $H$  with its Fourier expansion. That is, we let

$$\Psi_n(u)(x) = \int_a^b H_n(x, y)u(y) dy,$$

where  $H_n$  is the  $n^{\text{th}}$  partial sum of the Fourier series of  $H$ . These define finite dimensional, and thus compact, operators. Furthermore, the convergence of  $H_n \rightarrow H$  implies that  $\Psi_n \rightarrow \Psi$ . It follows that  $\Psi$  is compact because the limit of a sequence of compact operators is compact (i.e., the set of compact operators is a closed subspace of the space of bounded linear operators). However, for the proof I refer the reader to Banach (1987), a classical text that provides one of the early proofs of this theorem. We will also need the following theorem:

**Theorem 3.3.** *Suppose that  $F : [a, b] \times \mathbb{R} \rightarrow \mathbb{R}$  is bounded (or Lipschitz in the second variable) and continuous. Then the function  $u(\cdot) \mapsto F(\cdot, u(\cdot))$  from  $L^2[a, b] \rightarrow L^2[a, b]$  is continuous.*

This theorem is a special case of Vainberg's Lemma, which completely characterizes the functions  $F$  such that  $u(\cdot) \mapsto F(\cdot, u(\cdot))$  from  $L^p(\Omega) \rightarrow L^q(\Omega)$  is continuous for  $\Omega \subset \mathbb{R}^n$  and  $p^{-1} + q^{-1} = 1$ . Vainberg's Lemma, in all its glory, is presented in Vainberg (1953). Since our special case can be proved much more easily than the general result, we provide the proof below.

*Proof.* Since  $F$  is bounded and  $[a, b]$  has finite measure this function clearly maps into  $L^2[a, b]$ . To prove continuity it is sufficient to prove that for any sequence  $\{u_n\}$  with  $u_n \rightarrow u$  there is a subsequence  $\{u_{n_k}\}$  such that  $F(\cdot, u_{n_k}(\cdot)) \rightarrow F(\cdot, u(\cdot))$ . Since  $u_n \rightarrow u$  in  $L^2$ , there is a subsequence  $\{u_{n_k}\}$  such that  $u_{n_k} \rightarrow u$  (a.e.). The continuity of  $F$  thus implies that  $F(\cdot, u_{n_k}(\cdot)) \rightarrow F(\cdot, u(\cdot))$  (a.e.). Since  $F$  is bounded and  $[a, b]$  has finite measure, the Dominated Convergence Theorem (for  $L^2$ ) implies that  $F(\cdot, u_{n_k}(\cdot)) \rightarrow F(\cdot, u(\cdot))$  in  $L^2$ , and thus the theorem is proved.  $\square$

Combining the two results above, we get the following:

**Corollary 3.1.** *The function*

$$\Phi(u)(t) = f(t) + \int_a^b H(t, y)F(y, u(y)) dy$$

(with  $f \in L^2$ ) is a compact function from  $L^2[a, b] \rightarrow L^2[a, b]$ .

From this point forward, we will only deal with the case where  $F$  is bounded. Basic integral estimates then establish that

$$\|\Phi(u)\|_{L^2} \leq \|f\|_{L^2} + C\|F\|_{\infty}\|H\|_{L^2},$$

for some constant  $C$  independent of  $u$ . Letting  $R = \|f\|_{L^2} + C\|F\|_\infty\|H\|_{L^2}$ , we have that  $\Phi$  maps  $B_R(0) \subset L^2$  into itself (where  $B_r(x)$  denotes the ball of radius  $r$  about the point  $x$ ). Combining all of our work above we find that there is a compact convex subset of  $B_R(0)$  that  $\Phi$  maps to itself. The last step we need to complete before our method can be applied is to find this set. It is at this point that we leave the general theory behind and proceed to consider several specific examples.

### 3.2 The Second Order Initial Value Problem

One of the reasons that Hammerstein equations are interesting is that many physically motivated differential equations can be converted into Hammerstein integral equations and it is from these equations that we draw our examples. Let us first consider the following general second order nonlinear initial value problem

$$\begin{cases} u''(t) = f(u(t)), & t > 0, \\ u(0) = \alpha, \\ u'(0) = \beta. \end{cases} \quad (3.3)$$

This class of equations includes many problems of physical interest, including the one which we will concern ourselves with — the pendulum equation. However, before we get too specific, let us see how to solve this equation in the general case. The first thing we do is consider the corresponding linear problem:

$$\begin{cases} u''(t) = f(t), & t > 0, \\ u(0) = \alpha, \\ u'(0) = \beta. \end{cases} \quad (3.4)$$

It is well known that this equation can be solved via integration against the Green's function

$$G(t, \tau) = \begin{cases} 0 & 0 \leq t < \tau, \\ t - \tau, & t > \tau. \end{cases}$$

That is to say, the solution to (3.4) is given by

$$\begin{aligned} u(t) &= \int_0^\infty G(t, \tau) f(\tau) d\tau + \beta t + \alpha \\ &= \int_0^t (t - \tau) f(\tau) d\tau + \beta t + \alpha. \end{aligned} \quad (3.5)$$

The derivation of the Green's function can be found in Stakgold (1979), but given the Green's function it is an elementary exercise to verify that  $u$  defined by equation (3.5) is indeed a solution to (3.4). Returning to the nonlinear case, we define the operator

$$\Phi(u)(t) = \int_0^t (t - \tau)f(u(\tau)) d\tau + \beta t + \alpha.$$

From our discussion above, it is clear that  $u$  is a solution to (3.3) if and only if  $u$  is a fixed point of  $\Phi$ . The fact that  $[0, \infty)$  is a set of infinite measure is something of a problem because our results on the compactness of operators relied on the integral being over a set of finite measure. However, we can avoid this by instead solving the equation on the time interval  $[0, T]$ . Assuming  $f$  is bounded and continuous we see that  $\Phi$  satisfies the hypotheses of Corollary 3.1, so our method applies. The only remaining task is to find the compact convex set that  $\Phi$  maps to itself.

Recall that the cosine functions  $\{\cos(\pi kt/T)\}_{k=0}^{\infty}$  form a basis for  $L^2[0, T]$ , commonly called the Fourier cosine basis. Now, fix  $u \in L^2[0, T]$  and let  $F(t) = \Phi(u)(t)$ . The Fourier expansion of  $F(t)$  is given by

$$F(t) = a_0 + \sum_{k=1}^{\infty} a_k \cos\left(\frac{k\pi t}{T}\right),$$

where

$$a_0 = \frac{1}{T} \int_0^T F(t) dt \quad \text{and} \quad a_k = \frac{2}{T} \int_0^T F(t) \cos\left(\frac{k\pi t}{T}\right) dt, \quad k > 0.$$

In order to get estimates on  $a_k$ , we simplify the problem slightly by assuming that  $\beta = 0$ . Using the orthogonality properties of the cosine basis we then have that, for  $k > 0$ ,

$$\begin{aligned} a_k &= \frac{2}{T} \int_{t=0}^T \left( \int_{\tau=0}^t (t - \tau)f(u(\tau)) d\tau \right) \cos\left(\frac{k\pi t}{T}\right) dt \\ &= \frac{2}{T} \int_{t=0}^T \int_{\tau=0}^t (t - \tau)f(u(\tau)) \cos\left(\frac{k\pi t}{T}\right) d\tau dt \\ &= \frac{2}{T} \int_{\tau=0}^T \int_{t=\tau}^T (t - \tau)f(u(\tau)) \cos\left(\frac{k\pi t}{T}\right) dt d\tau \\ &= \frac{2T}{k^2\pi^2} \int_{\tau=0}^T \left[ \cos(k\pi) - \cos\left(\frac{k\pi\tau}{T}\right) \right] f(u(\tau)) d\tau. \end{aligned}$$



Letting  $M = \|f\|_\infty$ , we have the bound  $|a_k| \leq \frac{4TM}{k^2\pi^2}$ . For  $a_0$ , it is fairly straightforward to derive the bound  $|a_0| \leq T^2/6 + |\alpha|$ . We define the sequence  $\{\delta_n\}$  by  $\delta_1 = T^3/6 + |\alpha|$  and  $\delta_i = \frac{4TM}{(i-1)^2\pi^2}$  for  $i \geq 2$ . Defining  $B_{\{\delta_n\}}$  as in Theorem 2.4, we have that  $B_{\{\delta_n\}}$  is compact, and it is easily seen to be convex. Furthermore, by construction we have that  $A = \mathcal{F}^{-1}(B_{\{\delta_n\}})$  is a compact convex set such that  $\Phi : A \rightarrow A$ . Furthermore, Proposition 2.1 applies to give us a homeomorphism from  $A$  to  $I^\infty$ . Thus we have determined everything necessary to apply our search method to equations of the form

$$\begin{cases} u''(t) = f(u(t)), & t > 0, \\ u(0) = \alpha, \\ u'(0) = 0. \end{cases} \quad (3.6)$$

It is not much more difficult to handle the  $u'(0) \neq 0$  case, but we will not consider that situation here. One interesting problem that has this form is the pendulum equation. To demonstrate our method we consider the equation

$$\begin{cases} u''(t) = -\sin(u(t)), & t \in (0, 1), \\ u(0) = \pi/4, \\ u'(0) = 0. \end{cases} \quad (3.7)$$

For this problem we first ran the basic search algorithm searching for 4 and 10 term approximations, using a mesh size of  $1/2^{10}$  in both cases. The results are shown in Figure 3.1. The algorithm was implemented in Mathematica and all runs were done on a Macintosh Powerbook G4 with a 1.5 GHz PowerPC processor. When run times are mentioned they are only useful for comparisons to each other and are in no way rigorous assessments of algorithmic efficiency.

The  $L^2$  error of this approximation is approximately 0.01 and the  $L^\infty$  error is approximately .045. Since the mesh size of our triangulation is  $\frac{1}{2^{10}} \approx 0.001$ , this error is relatively small. In order to get an idea of how good our approximation is, we have plotted it against the four term Fourier approximation to the solution. As this plot (top right Figure 3.1) shows, our approximation is fairly close to the best possible approximation using only the first four basis functions. However, it is also clear that we could benefit from using a smaller mesh size for our triangulation.

The case for the ten term approximation is similar. This time the  $L^2$  error of the approximation is approximately 0.0075 and the  $L^\infty$  error is approximately .02. In this case increasing the the number of terms improved the

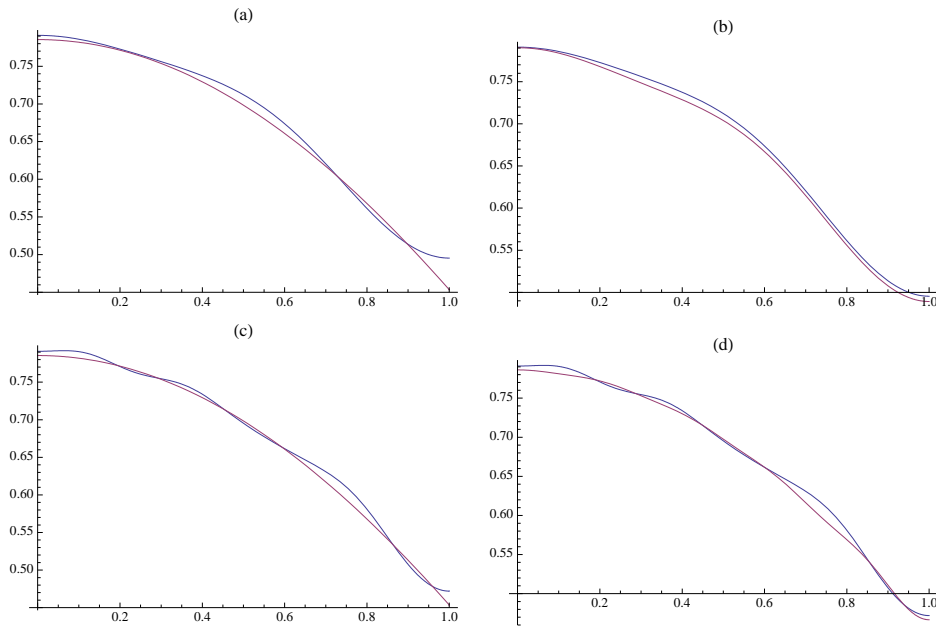


Figure 3.1: (a) A 4 term approximation (blue) plotted against the solution. (b) A 4 term approximation (blue) plotted against the 4 term Fourier expansion to the solution. (c) A 10 term approximation (blue) against the solution. (d) A 10 term approximation (blue) against the 10 term Fourier expansion of the solution.

$L^\infty$  error much more than the  $L^2$  error, and there is some “moral” justification for this. As we saw when introducing  $K_2(m)$ , the  $L^2$  mesh size grows with the dimension of the simplex being triangulated while the  $L^\infty$  mesh size does not. Clearly, however, it would still be desirable to use a finer triangulation. Unfortunately, using a much finer triangulation is computational infeasible with the present algorithm. With the current triangulation size, several hours were needed to compute the four term approximation and the algorithm needed to be run overnight to produce the ten term approximation.

Therefore, in order to get better approximations we turn to van der Laan and Talman’s basic algorithm, which we will henceforth refer to as the VT-algorithm. In addition to speed, one of the nice features of the VT-algorithm is that it is a restart algorithm. That is, it runs a search for a given mesh size, finds a c.l. simplex, and restarts at a vertex of the c.l. simplex with a finer

mesh size. Consequently, one run of this algorithm produces a sequence of approximations. This has the advantage of giving us a look at how such sequences converge. Such a sequence is produced in Figure 3.2. In this figure the graphs are indexed by  $M$ , where  $M$  means that the approximation in the graph is for a triangulation with mesh size  $1/2^M$ . All of the approximations in this figure are ten term approximations. As we would expect, the first several terms of the sequence look nothing like the solution, but as early as  $M = 3$  the approximation has roughly the same shape as the solution and beyond  $M = 10$  the sequence is fairly constant. Also, in terms of time, this sequence took less than an hour to produce, less time than it took to produce the four term approximation using the slower algorithm.

We now analyze the error for the  $M = 20$  approximation. The  $L^2$  error is approximately 0.0021, a bit better than our previous approximation. Furthermore, the ten term Fourier approximation to the solution has  $L^2$  error approximately 0.0017, so our approximation is almost as good as the best ten term approximation in the  $L^2$  sense. Additionally, the  $L^\infty$  error of our approximation is about 0.006, substantially better than our previous approximation.

It is worth noting that, while we have been comparing our approximations to the Fourier approximations, we should not expect the two to be equal. If we look back at the construction, we are essentially producing a fixed point of a finite dimensional approximation to the function  $\Phi$ , and the Fourier approximation to the fixed point of  $\Phi$  is not necessarily a fixed point of the finite dimensional approximation. Indeed, in the current case, the distance from the ten term Fourier approximation to its image under the finite dimensional approximation of  $\Phi$  is about  $5 \times 10^{-7}$  (the distance from our approximation to its image is about  $3 \times 10^{-6}$ ). Thus, while comparison to the Fourier approximation gives a comparison to the theoretical best approximation with a given number of terms, we should not expect our sequence of approximations to be partial sums of the Fourier series for the solution.

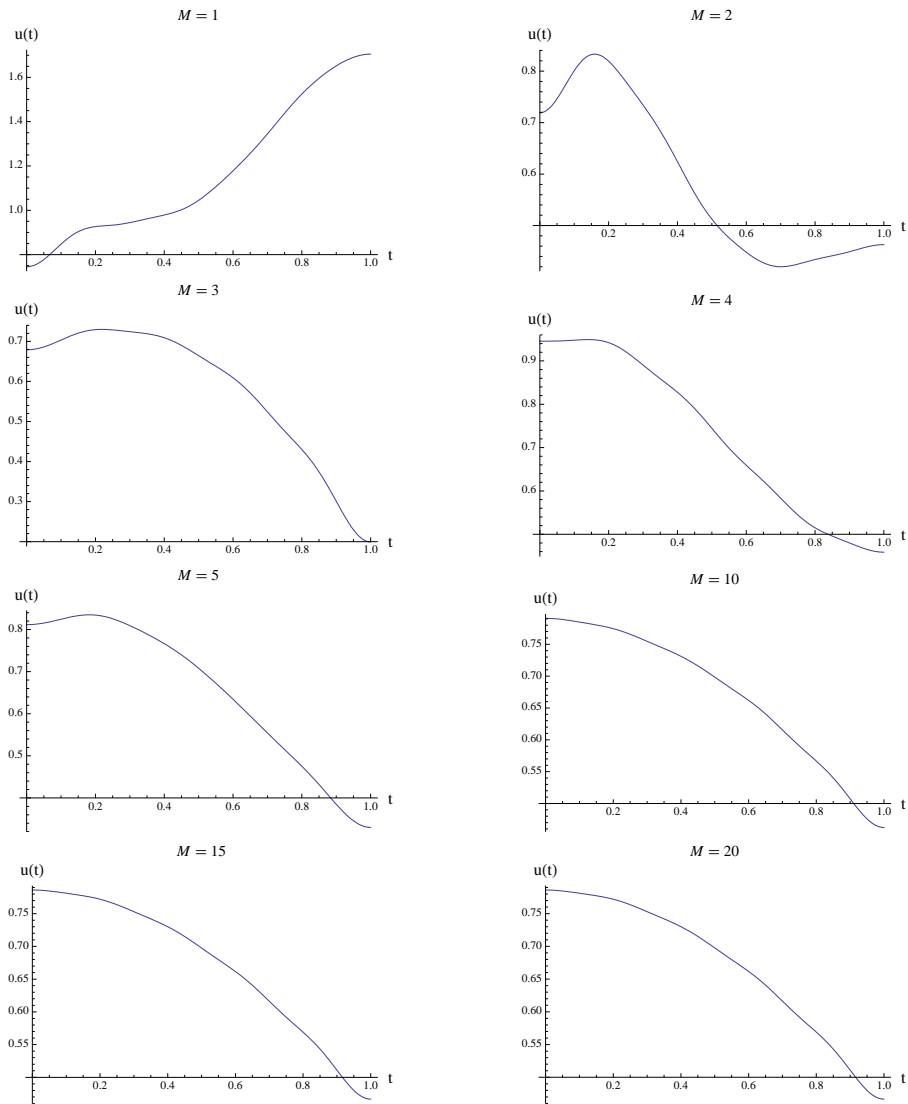


Figure 3.2: A Sequence of Ten Term Approximations

### 3.3 The Two Point Boundary Value Problem

In this section we move from considering initial value problems to considering the boundary value problem

$$\begin{cases} -u''(x) = f(u(x)), & x \in (a, b), \\ u(a) = \alpha, \\ u(b) = \beta. \end{cases} \quad (3.8)$$

While equations of this form are quite similar in appearance to the initial value problems we were considering earlier, the theory regarding the two types of equations is quite different. This is because, in general, one expects initial value problems to have a unique local solution if they have a solution at all, while boundary value problems often have multiple solutions. Nonetheless, for our approach the differences between the two types of problem are minimal. Just as when we solved the initial value problem, we attack this problem by first considering the linear case:

$$\begin{cases} -u''(x) = f(x), & x \in (a, b), \\ u(a) = \alpha, \\ u(b) = \beta. \end{cases} \quad (3.9)$$

Again we can solve this using a Green's function, but this time the Green's function is

$$G(x, y) = \begin{cases} \frac{x(b-y)+a(y-b)}{b-a} & x \leq y \\ \frac{y(b-x)+a(x-b)}{b-a} & x \geq y. \end{cases}$$

This is a special case of derivations in Stakgold (1979), though again given the Green's function it is an elementary exercise to verify that

$$u(x) = \int_a^b G(x, y)f(u(y)) dy + \frac{(b-x)}{b-a}\alpha + \frac{x-a}{b-a}\beta$$

is a solution to equation (3.9). Indeed, results in Stakgold (1979) show that this is the unique solution. The fact that the Green's function is different and the representation of the solution to the linear problem is slightly different are the only substantial differences between our approach to the boundary value problem and the initial value problem. This representation of the solution inspires us to define the operator

$$\Psi(u)(t) = \int_a^b G(x, y)f(y) dy + \frac{(b-x)}{b-a}\alpha + \frac{x-a}{b-a}\beta.$$

Consequently, we have that  $u$  is a solution of equation (3.8) if and only if it is a fixed point of  $\Psi$ . Rather than go through the general computations as we did for the initial value problem, we will go straight to the results. The computations are all similar to those for the initial value problem.

### 3.3.1 The Generalized Duffing's Equation

The first boundary value problem we consider is the generalized Duffing's equation, sometimes also referred to as the boundary value problem for the pendulum. The equation is

$$\begin{cases} -u''(x) = \lambda^2 \sin(u(x)), & x \in (0, \pi), \\ u(0) = 0, \\ u(\pi) = 0. \end{cases} \quad (3.10)$$

In fact, we will only consider the case where  $\lambda = 2.5$ . For this value of  $\lambda$  it is well known that Equation 3.10 has five solutions: A positive solution  $u_1$ , a solution  $u_2$  that is positive on  $(0, \pi/2)$  and negative on  $(\pi/2, \pi)$ , the zero solution  $u_0$ ,  $-u_1$  and  $-u_2$ . This problem was considered both in Allgower (1977) and Chen (1977). In these papers the solutions  $u_1$ ,  $u_2$ , and  $-u_2$  were approximate using a method similar to ours. Their method for searching for multiple solutions consists of the method described in Section 1.4. The way we searched for multiple solutions was to combine a partial implementation of the method in Section 1.4 with the method in 2.3. Doing this we were able to approximate the solutions  $u_1$ ,  $u_2$ ,  $-u_1$ , and  $-u_2$ . Since we only did a partial implementation of the methods described, it is entirely possible that a full implementation will yield an approximation to  $u_0$  as well

Since our primary method for searching for multiple solutions requires using the slower search algorithm, we used a relatively coarse triangulation; namely  $K_2(2^8)$ , and we only looked for a five term approximation. The approximations produced in this fashion appear in Figure 3.3.

It is worth pointing out that, while we know that if  $u$  is a solution to Equation 3.10 then so is  $-u$ , the search algorithm does not know this. That is, the search algorithm arrived at these four approximate solutions without knowing the relationship between them.

Notice that using such a coarse mesh calls into question the accuracy of our approximations. In Figure 3.4, we have plotted our approximations (pink) against approximations obtained using a modified shooting method.

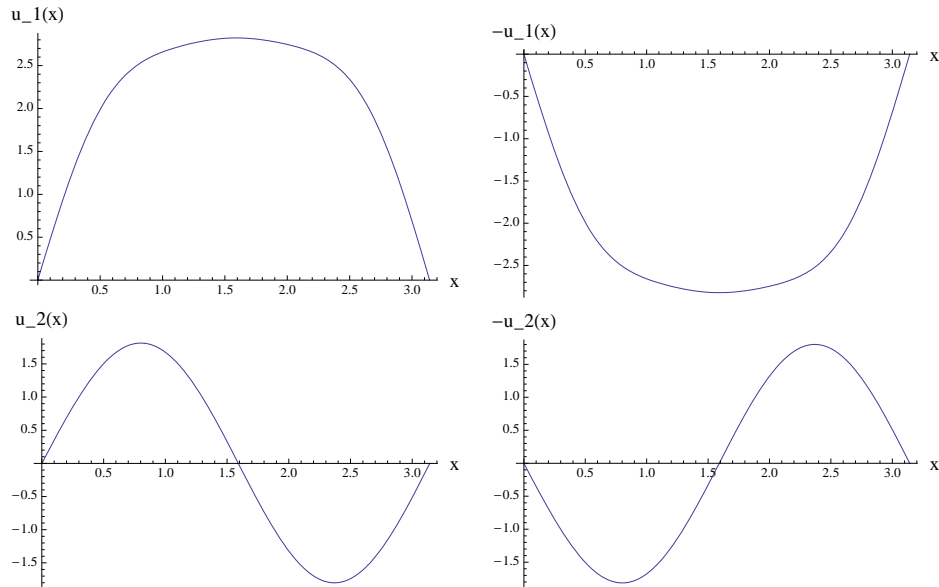


Figure 3.3: Approximate solutions to the generalized Duffing's equation

These graphs show that, while there is certainly room for our approximations to be improved, they are actually good approximations — especially when the large mesh size and small number of terms in the approximation are taken into account. Unfortunately, we are not able to use the VT-algorithm to improve all of these approximations. The VT-algorithm can be used for a less extensive search for multiple solutions by varying any number of the arbitrarily chosen initial parameters or by using the method from Section 2.3. Using the VT-algorithm, the only solutions that have been approximated have been  $u_1$  and  $-u_1$ . Even when approximations to  $u_2$  are used as the initial point the algorithm moves towards  $u_1$  or  $-u_1$ , and a reason for this has yet to be discovered. We do not include images here, but as was expected the VT-algorithm approximations to  $u_1$  and  $-u_1$  were close enough that their graphs were indistinguishable from the approximations attained using the modified shooting method.

### 3.4 A Note on Uniform Convergence

To this point, we have been constructing  $L^2$  approximations to solutions of differential equations. A natural question to ask is whether or not our

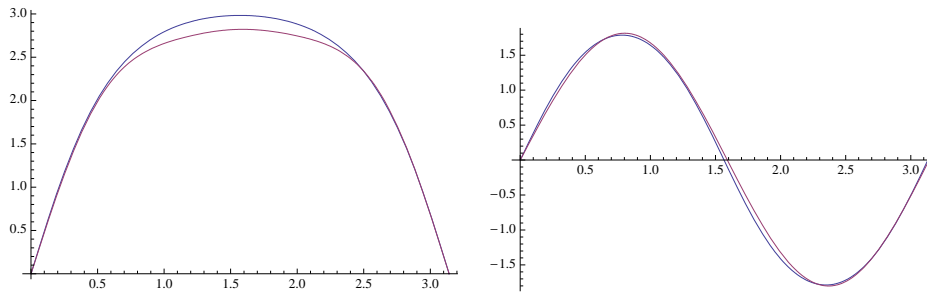


Figure 3.4: Approximate solutions (pink) versus actual solutions

method can produce approximations in a stronger norm. Because of the generality in which we have proven our method, the answer is theoretically yes. But what about practically? As we mentioned before, one of the key difficulties in applying our method is finding the homeomorphism to  $I^\infty$  — how to do this in general is not clear. However, by modifying our method slightly we can, in certain cases, convert the  $L^2$  approximations into uniform approximations. There are several ways of doing this and in this section we outline one of the easiest ways to obtain such approximations.

Let  $\Phi$  be as in Corollary 3.1 and take  $f(t) = 0$ . Further suppose that  $H(t, y)$  is continuous and let  $E$  be the set of continuous functions on  $[a, b]$  with the  $L^2$  norm. Then, not only is  $\Phi : L^2[a, b] \rightarrow L^2[a, b]$  compact, but  $\Phi : E \rightarrow C[a, b]$  is compact as well. Notice that our method, as applied in the previous sections, produces a sequence  $\{u_n\}$  of continuous functions that converge in  $L^2$  to a fixed point  $u$  of  $\Phi$ . Consequently,  $\{u_n\}$  is a bounded sequence in  $E$ . The compactness of  $\Phi$  then implies that the sequence  $\{\Phi(u_n)\}$  has a subsequence  $\{\Phi(u_{n_i})\}$  such that  $\Phi(u_{n_i}) \rightarrow g$  uniformly for some  $g \in C[a, b]$ . This implies that  $\Phi(u_{n_i}) \rightarrow g$  in  $L^2$ . However, in  $L^2$  we know that  $\Phi(u_{n_i}) \rightarrow \Phi(u) = u$  since  $\Phi$  is continuous on  $L^2[a, b]$ . Therefore we have that  $u = g$ . Hence the sequence  $\{\Phi(u_n)\}$  has a subsequence that converges uniformly to a fixed point of  $\Phi$  and this is the same fixed point that  $\{u_n\}$  converges to in  $L^2$ . This procedure thus allows us to convert our  $L^2$  approximations into uniform approximations.





## Chapter 4

# Future Research

This thesis only scratches the surface of the work that needs to be done to fully develop this method of approximation. Indeed, the work here serves mostly as a proof of concept. In this chapter we will introduce several relatively disjoint directions that merit further research. We speak in broad terms with no attempt at self-containment.

The most basic extension would be to generalize our work with Hammerstein integral equations to higher dimensions. Seeing as the theory of Green's functions is well developed for  $n$ -dimensional spaces (see Stakgold (1979) or Evans (1998)), it is clear that the results we developed in Section 3.1 can be extended to the case where the integral in equation (3.1) is taken over region in  $\mathbb{R}^n$  rather than an interval. In addition to generalizing to higher dimensions, we could also generalize to nonstandard domains. One developing field of mathematics is analysis, and differential equations, on fractal domains. Remarkably, Green's functions can be developed on these domain (see Strichartz (2006)). Our method should be easily adaptable to these domains. Furthermore, this could be particularly interesting because our method would provide a method of approximation that is internal to the fractal domain. Indeed, this extension would provide a method to approximate solutions to equations for which no other method of approximation is proven to work.

Furthermore, there are several other well established methods for representing solutions as integrals (see Evans (1998)) and our method can potentially be applied to these equations as well. One example of this would be the inhomogeneous initial value problem for the heat equation in  $\mathbb{R}^n$ , which can be solved via integration against a heat kernel (see Evans (1998)). Thus our method could potentially be used to solve the nonlinear version

of this. However, to do so would be nontrivial because of issues surrounding the integrability and singularity of the heat kernel. Determining if and how these issues could be worked around would be interesting and applicable. Furthermore there is no reason to restrict to equations that can be represented by integration against a kernel. The same basic methods should be easily adaptable to the general case covered by Theorem ??.

Another direction for this research is to generalize the approach to other fixed point theorems. An interesting and straight forward generalization would be to a theorem known both as Schaefer's fixed point theorem and the Leray-Schauder fixed point theorem (see Evans (1998) and McOwen (2003)). The theorem goes as follows:

**Theorem 4.1.** *Suppose that  $X$  is a Banach space and  $\Phi : X \rightarrow X$  is compact. Assume that the set*

$$A = \{u \in X \mid u = \lambda\Phi(u) \text{ for some } 0 \leq \lambda \leq 1\}$$

*is bounded. Then  $A$  has a fixed point.*

The proof of this theorem uses Schauder's fixed point theorem to guarantee the existence of a fixed point of

$$\Psi(u) = \begin{cases} \Phi(u) & \|\Phi(u)\| \leq M \\ \frac{M\Phi(u)}{\|\Phi(u)\|} & \|\Phi(u)\| \geq M, \end{cases}$$

for suitably chosen  $M$ . The fixed point of  $\Psi$  is then shown to be a fixed point of  $\Phi$ . Our method can be used to approximate the fixed points of  $\Psi$ , and thus of  $\Phi$ . This direction is particularly interesting because it is highly applicable. Under certain assumptions (incompressibility, non-slip boundary, etc.) Theorem 4.1 can be used to prove the existence of a steady state solution to the Navier-Stokes equations in two and three dimensions (see McOwen (2003) for details). Hence the method described here presents a way to approximate the steady state solution to the Navier-Stokes equations. Thus this direction should be pursued if for no other reason than the high applicability of the Navier-Stokes equations.

Notice that in this thesis we only applied the approximation method to operators defined on  $L^2[a, b]$ , but there is no theoretical reason for this restriction. Even if we maintain the restriction to Hilbert spaces so that the homeomorphism to  $I^\infty$  is potentially easy to compute we could consider problems on other Hilbert spaces. This has several potential advantages. For example, if we consider operators on the Sobolev space  $H_0^1$  we get an

$L^2$  approximation to the derivative of the function as well. Thus, in one fell swoop, we could get an approximation  $u_n$  to the function  $u$  such that  $u'_n$  is an approximation to  $u'$ . This might not seem interesting at first, but notice that the derivative is, in general, a discontinuous operator. Hence, in general,  $u_n$  being an approximation to  $u$  does not imply that  $u'_n$  is an approximation to  $u'$ . However, approximations in  $H_0^1$  do have this property. This would give our method a leg up on other methods, such as finite elements and finite differences, that do not have this property.

The list of extensions we have mentioned here is far from exhaustive. Other potential extensions include  $k$ -Hessian problems, periodic solutions to systems of equations, and many others. Indeed, the list is almost endless. Though we have tried to highlight several of the most novel and interesting directions for further research, ultimately this section is trying to convey that rather than being a capstone work, this thesis opens a Pandora's box of potential further research.



# Bibliography

- Allgower, E. (1977). Application of a fixed point search algorithm to nonlinear boundary value problems having several solutions. In *Fixed points: algorithms and applications (Proc. First Internat. Conf., Clemson Univ., Clemson, S. C., 1974)*, pages 87–111. Academic Press, New York.
- Banach, S. (1987). *Theory of Linear Operations*, volume 38 of *North-Holland Mathematical Library*. North-Holland Publishing Co., Amsterdam. Translated from the French by F. Jellett, With comments by A. Pełczyński and Cz. Bessaga.
- Chen, H. C. (1977). A constructive existence method for nonlinear boundary value problems. *J. Math. Anal. Appl.*, 59(3):454–468.
- Evans, L. C. (1998). *Partial Differential Equations*, volume 19 of *Graduate Studies in Mathematics*. American Mathematical Society, Providence, RI.
- Jeppson, M. (1972). A search for the fixed points of a continuous mapping. In *Mathematical topics in economic theory and computation (Sympos. Math. Econom., SIAM Fall Meeting, Univ. Wisconsin, Madison, Wis., 1971)*, pages 122–129. Soc. Indust. Appl. Math., Philadelphia, Pa.
- Klee, Jr., V. (1955). Some topological properties of convex sets. *Trans. Amer. Math. Soc.*, 78:30–45.
- Kuhn, H. W. (1969). Approximate search for fixed points. In *Computing Methods in Optimization Problems, 2 (Proc. Conf., San Remo, 1968)*, pages 199–211. Academic Press, New York.
- Lang, S. (1993). *Real and Functional Analysis*, volume 142 of *Graduate Texts in Mathematics*. Springer-Verlag, New York, third edition.
- Marsden, J. E. and Hoffman, M. J. (1993). *Elementary Classical Analysis, second edition*. W.H. Freeman and Company, New York.

- McOwen, R. (2003). *Partial Differential Equations: Methods and Applications*. Prentice Hall, Upper Saddle, NJ.
- Precup, R. (2002). *Methods in Nonlinear Integral Equations*. Kluwer Academic Publishers, Dordrecht.
- Rizzolo, D. and Su, F. E. (2007). A fixed point theorem for the infinite-dimensional simplex. *J. Math. Anal. Appl.*, 332(2):1063–1070.
- Rudin, W. (1987). *Real and Complex Analysis*. McGraw-Hill Book Co., New York, third edition.
- Smart, D. (1974). *Fixed Point Theorems*. Cambridge University Press, London. Cambridge Tracts in Mathematics, No. 66.
- Spanier, E. H. (1966). *Algebraic Topology*. McGraw-Hill Book Co., New York.
- Stakgold, I. (1979). *Green's Functions and Boundary Value Problems*. John Wiley & Sons, New York-Chichester-Brisbane. A Wiley-Interscience Publication, Pure and Applied Mathematics.
- Strichartz, R. S. (2006). *Differential Equations on Fractals*. Princeton University Press, Princeton, NJ. A tutorial.
- Talman, A. (1980). *Variable Dimension Fixed Point Algorithms and Triangulations*, volume 128 of *Mathematical Centre Tracts*. Mathematisch Centrum, Amsterdam. With the collaboration of G. van der Laan.
- Todd, M. J. (1976). *The Computation of Fixed Points and Applications*. Springer-Verlag, Berlin. Lecture Notes in Economics and Mathematical Systems, Vol. 124.
- Vainberg, M. M. (1953). On the structure of an operator. *Doklady Akad. Nauk SSSR (N.S.)*, 92:213–216.
- van der Laan, G. and Talman, A. (1979). A restart algorithm for computing fixed points without an extra dimension. *Math. Programming*, 17(1):74–84.