

Some pages of this thesis may have been removed for copyright restrictions.

If you have discovered material in AURA which is unlawful e.g. breaches copyright, (either yours or that of a third party) or any other law, including but not limited to those relating to patent, trademark, confidentiality, data protection, obscenity, defamation, libel, then please read our [Takedown Policy](#) and [contact the service](#) immediately

The University of Aston in Birmingham

**APPROXIMATING DIFFERENTIABLE RELATIONSHIPS
BETWEEN DELAY EMBEDDED DYNAMICAL SYSTEMS
WITH RADIAL BASIS FUNCTIONS**

Michael Alan Sherred Potts
Doctor of Philosophy

December 1996

This copy of the thesis has been supplied on condition that anyone who consults it is understood to recognise that its copyright rests with its author and that no quotation from the thesis and no information derived from it may be published without proper acknowledgement.

The University of Aston in Birmingham

**APPROXIMATING DIFFERENTIABLE RELATIONSHIPS
BETWEEN DELAY EMBEDDED DYNAMICAL SYSTEMS
WITH RADIAL BASIS FUNCTIONS**

**Michael Alan Sherred Potts
Doctor of Philosophy**

December 1996

This thesis is about the study of relationships between experimental dynamical systems. The basic approach is to fit radial basis function maps between time delay embeddings of manifolds. We have shown that under certain conditions these maps are generically diffeomorphisms, and can be analysed to determine whether or not the manifolds in question are diffeomorphically related to each other. If not, a study of the distribution of errors may provide information about the lack of equivalence between the two.

The method has applications wherever two or more sensors are used to measure a single system, or where a single sensor can respond on more than one time scale: their respective time series can be tested to determine whether or not they are coupled, and to what degree. One application which we have explored is the determination of a minimum embedding dimension for dynamical system reconstruction. In this special case the diffeomorphism in question is closely related to the predictor for the time series itself.

Linear transformations of delay embedded manifolds can also be shown to have nonlinear inverses under the right conditions, and we have used radial basis functions to approximate these inverse maps in a variety of contexts. This method is particularly useful when the linear transformation corresponds to the delay embedding of a finite impulse response filtered time series. One application of fitting an inverse to this linear map is the detection of periodic orbits in chaotic attractors, using suitably tuned filters. This method has also been used to separate signals with known bandwidths from deterministic noise, by tuning a filter to stop the signal and then recovering the chaos with the nonlinear inverse. The method may have applications to the cancellation of noise generated by mechanical or electrical systems.

In the course of this research a sophisticated piece of software has been developed. The program allows the construction of a hierarchy of delay embeddings from scalar and multi-valued time series. The embedded objects can be analysed graphically, and radial basis function maps can be fitted between them asynchronously, in parallel, on a multi-processor machine. In addition to a graphical user interface, the program can be driven by a batch mode command language, incorporating the concept of parallel and sequential instruction groups and enabling complex sequences of experiments to be performed in parallel in a resource-efficient manner.

Keywords: neural networks, total least squares, time series prediction, inverting filters, signal separation.

To Mary, who kept the Last Homely House west of the Cotswolds,
and to Raymond and Jill, who could usually be found there.

Acknowledgements

I would like to acknowledge the great amount of help and encouragement I have received from members of the Signal Processing Theory group at the Defence Research Agency (formerly the Royal Signals and Radar Establishment) in Malvern: in particular, from David Broomhead and Jerry Huke (now both with the Department of Mathematics at UMIST), Geoff de Villiers, Ian Proudler and Robin Jones; and also from Chris Booth, the resident C guru in the Parallel Processing group. I would also like to thank David Bounds (now with Recognition Systems) and David Lowe, of Aston University, and Lenny Smith, of Oxford University, for their constructive input. The work on separating chaotic signals from transmitted messages, described in the last section of the final chapter, was performed in collaboration with David Broomhead and Jerry Huke, and its results have been independently reported in the literature. This thesis was typeset in Knuth's \TeX , and the figures were generated either directly in Postscript, or through Williams and Kelley's Gnuplot. The software which was written in the course of this work incorporates several routines adapted from Press et al's Numerical Recipes in C.

Contents

1	Introduction	12
1.1	Overview	14
2	Dynamical systems and embeddings	17
2.1	The dynamical system	18
2.1.1	The role of dissipation	19
2.2	The delay embedding	20
2.2.1	Differentiable equivalence	25
2.3	Filtered embedding	28
2.3.1	Sampling	29
2.3.2	Singular subspaces	32
2.3.3	Time series filtering	32
3	Approximation	34
3.1	Characteristics of the RBF map	36
3.2	The method of least squares	37
3.2.1	The least squares solution	38
3.2.2	Forward selection	41
3.2.3	Control of over-fitting	42
3.2.4	Detecting non-invertible maps with least squares	50
3.3	The method of total least squares	52
3.3.1	The total least squares solution	55

3.3.2	Detecting non-invertible maps with total least squares	56
3.3.3	Numerical instability in the total least squares solution	58
4	Maps on manifolds	60
4.1	Embedding a circle in the plane	60
4.1.1	The circle and least squares	63
4.1.1.1	Analysis of the least squares solution	65
4.1.1.2	Approximating the identity map	69
4.1.1.3	Analytical calculation of Lipschitz constants	75
4.1.1.4	Approximation of the Lipschitz constants	76
4.1.2	The circle and total least squares	82
4.2	Embedding a torus in three dimensions	85
4.2.1	The torus and least squares	87
4.2.1.1	Analysis of the least squares solution	89
4.2.1.2	Approximating the identity map	91
4.2.1.3	Analytical calculation of Lipschitz constants	91
4.2.1.4	Approximation of the Lipschitz constants	93
4.2.2	The torus and total least squares	94
5	Maps on dynamical systems	99
5.1	Determination of a minimum embedding dimension	100
5.1.1	Embedding the Ikeda system	102
5.1.2	Embedding the Hénon system	104
5.1.3	Embedding the Lorenz system	104
5.1.4	Embedding the laser system	108
5.2	Constructing a singular subspace	112
5.2.1	Singular subspaces of the embedded Lorenz system	113
5.2.2	Singular subspaces of the embedded laser system	116
5.3	Detecting unstable periodic orbits	116
5.3.1	Filtering the Ikeda attractor	119
5.3.2	Filtering the Hénon attractor	121
5.3.3	Filtering the Lorenz attractor	123
5.4	Signal separation	126
5.4.1	Isolating a sinusoidal message from Ikeda chaos	129
5.4.2	Isolating a phase modulated message from Lorenz chaos	133

6	Conclusions	141
	References	146
A	Inverting radial basis functions	149
	A.1 Proof that φ is an injective immersion	149
	A.2 Approximating the inverse of φ	150
B	Implementation	152
	B.1 Structure	152
	B.2 Process control	154
	B.3 The batch interface	154
	B.4 The graphical user interface	157

Figures

2.1	Attractors of discrete dynamical systems.	19
2.2	Continuous attractor of the Lorenz system.	19
2.3	Time series measured from numerically simulated Ikeda and Hénon systems.	21
2.4	Time series from the first component of the numerically integrated Lorenz system.	22
2.5	Comparing the intensity time series measured from a laser experiment with a time series measured on the Lorenz system.	23
2.6	Embedding the Ikeda and Henon systems.	26
2.7	Embedding the Lorenz and laser systems.	27
2.8	Filtered embedding of the Lorenz and laser systems.	31
3.1	Comparing the repulsive and forward selection methods.	40
3.2	Over-fitting on the laser system prediction problem, for two individual sets of centers, selected by both repulsive and forward selection, respectively.	43
3.3	The effects of rank truncation on the laser time series prediction problem with centers selected by both repulsive and forward selection methods.	46
3.6	Scatter plot of the reordering map relating the blind and targeted rank selection criteria.	47
3.4	Improving generalisation with repulsive centers by rank truncation.	48
3.5	Improving generalisation with forward selection by rank truncation.	49
4.1	The circle embedded in \mathbb{R}^3 and its projections in the plane.	61
4.2	Comparing the mean, normalised test errors $\langle \epsilon_0^{(\phi)} \rangle$ and $\langle \epsilon_\phi^{(0)} \rangle$, versus ϕ , for the LS approximations $\widehat{\mathbf{f}}_\phi$ and $\widehat{\mathbf{f}}_\phi^{-1}$ on the projected circles \mathcal{S}_0 and \mathcal{S}_ϕ .	64
4.3	Approximating the map $\mathbf{f}_\phi: \mathcal{S}_0 \rightarrow \mathcal{S}_\phi$ and its inverse, for $\phi = 20, 40$ degrees.	66
4.4	Approximating the map $\mathbf{f}_\phi: \mathcal{S}_0 \rightarrow \mathcal{S}_\phi$ and its inverse, for $\phi = 26, 27$ degrees.	67

4.5	Colour-coding the projected circles \mathcal{S}_ϕ , for $\phi = 20, 26, 27, 40$ degrees, by the relative magnitudes of the per-point errors $\epsilon_\phi^{(0)}(\mathbf{y})$ to which they give rise under $\widehat{\mathbf{f}}_\phi^{-1}$.	68
4.6	Comparing the mean, normalised test set errors $\langle \eta_0^{(\phi)} \rangle$ and $\langle \eta_\phi^{(0)} \rangle$, versus ϕ , arising from the LS identity approximations $\widehat{\mathbf{I}}_\phi^{(\phi)}$ and $\widehat{\mathbf{I}}_\phi^{(0)}$ on the projected circles \mathcal{S}_0 and \mathcal{S}_ϕ .	70
4.7	Approximating the identity maps between \mathcal{S}_0 and \mathcal{S}_ϕ , for $\phi = 20, 40$ degrees.	71
4.8	Approximating the identity maps between \mathcal{S}_0 and \mathcal{S}_ϕ , for $\phi = 26, 27$ degrees.	72
4.9	Colour-coding the identity maps between \mathcal{S}_0 and \mathcal{S}_ϕ , for $\phi = 20$ and 40 degrees, by the per-point error magnitudes $\ \eta_0^{(\phi)}(\mathbf{x})\ $ and $\ \eta_\phi^{(0)}(\mathbf{y})\ $.	73
4.10	Colour-coding the identity maps between \mathcal{S}_0 and \mathcal{S}_ϕ , for $\phi = 26$ and 27 degrees, by the per-point error magnitudes $\ \eta_0^{(\phi)}(\mathbf{x})\ $ and $\ \eta_\phi^{(0)}(\mathbf{y})\ $.	74
4.11	Numerically simulated, and empirically estimated, upper and lower bounds for the growth of errors under $\mathbf{f}_\phi: \mathcal{S}_0 \rightarrow \mathcal{S}_\phi$ and, where appropriate, its inverse.	77
4.12	Scatter plot of per-point identity errors $\ \eta_\phi^{(0)}(\mathbf{y})\ $ versus $\ \eta_0^{(\phi)}(\mathbf{x})\ $ for $\phi = 20, 40$ degrees, averaged over 500 sets of randomly-seeded repulsive centers.	79
4.13	Scatter plot of per-point identity errors $\ \eta_\phi^{(0)}(\mathbf{y})\ $ versus $\ \eta_0^{(\phi)}(\mathbf{x})\ $ for $\phi = 26, 27$ degrees, averaged over 500 sets of randomly-seeded repulsive centers.	80
4.14	Empirically estimated upper and lower bounds \widehat{U}'_ϕ and \widehat{L}'_ϕ , obtained by moving the average over random seeds inside the calculation of extrema, superimposed on the numerically simulated analytical bounds U_ϕ and L_ϕ to which they correspond.	81
4.15	Comparing the mean, normalised forward and inverse errors and condition numbers for the TLS approximation to $\mathbf{f}_\phi: \mathcal{S}_0 \rightarrow \mathcal{S}_\phi$, calculated over the test set.	83
4.16	Approximating the map $\mathbf{f}_\phi: \mathcal{S}_0 \rightarrow \mathcal{S}_\phi$ and its inverse, for $\phi = 20$, with a symmetrical RBF map.	85
4.17	Illustrating the 2-tori $\mathcal{T}_{3,r}$ for $r = 0.8, 1$ and 1.2 .	86
4.18	Comparing the mean, normalised errors $\langle \epsilon_{4,r} \rangle$ and $\langle \epsilon_{3,r} \rangle$, versus r , for the LS approximations $\widehat{\mathbf{f}}_r$ and $\widehat{\mathbf{f}}_r^{-1}$ on the tori $\mathcal{T}_{4,r}$ and $\mathcal{T}_{3,r}$, calculated over training and test sets.	88
4.19	Illustrating the errors arising from the approximation $\widehat{\mathbf{f}}_r^{-1}: \mathcal{T}_{3,r} \subset \mathbb{R}^3 \rightarrow \mathbb{R}^4$ by colour-coding the elements of $\mathcal{T}_{3,r}$, for $r = 0.8, 1$ and 1.2 .	90
4.20	Comparing the mean, normalised test set errors $\langle \eta_{4,r} \rangle$ and $\langle \eta_{3,r} \rangle$ versus r for the LS identity approximations $\mathbf{I}_{4,r}$ and $\mathbf{I}_{3,r}$ on $\mathcal{T}_{4,r}$ and $\mathcal{T}_{3,r}$.	92
4.21	Empirically estimated upper and lower bounds for the growth of errors under $\mathbf{f}_r: \mathcal{T}_{4,r} \rightarrow \mathcal{T}_{3,r}$.	94
4.22	Scatter plot of per-point identity errors $\ \eta_{3,r}(\mathbf{y})\ $ versus $\ \eta_{4,r}(\mathbf{x})\ $ for $r = 0.8, 1$ and 1.2 , averaged over 500 sets of randomly-seeded repulsive centers.	95
4.23	Empirically estimated upper and lower bounds \widehat{U}'_r and \widehat{L}'_r , obtained by moving the average over random seeds inside the calculation of extrema, superimposed on the analytical bounds U_r and L_r to which they correspond.	96

4.24	Comparing the mean, normalised forward and inverse errors and condition numbers for the TLS approximation to $f_r: \mathcal{T}_{4,r} \rightarrow \mathcal{T}_{3,r}$, calculated over the test set.	97
5.1	Establishing a minimum embedding dimension for the Ikeda system by plotting the mean time series prediction error $\langle \epsilon_m \rangle$, versus m , calculated over both training and test sets.	103
5.2	Establishing a minimum embedding dimension for the Hénon system by plotting the mean, normalised prediction error $\langle \epsilon_m \rangle$, versus m , over training and test sets.	105
5.3	Establishing a minimum embedding dimension for the 0.01-step Lorenz system by plotting the mean, normalised prediction error $\langle \epsilon_m \rangle$, versus m , over training and test sets.	106
5.4	Establishing a minimum embedding dimension for the 0.01-step Lorenz system, generated from a noisy time series, by plotting the mean, normalised prediction error $\langle \epsilon_m \rangle$, versus m , over training and test sets.	107
5.5	Establishing a minimum embedding dimension for the 0.01-step Lorenz system, reconstructed with a lag of $\tau = 10$, by plotting the mean, normalised prediction error $\langle \epsilon_m \rangle$, versus m , over training and test sets.	109
5.6	Establishing a minimum embedding dimension for the 0.01-step Lorenz system, reconstructed with a lag of $\tau = 10$ from a noisy time series, by plotting the mean, normalised prediction error $\langle \epsilon_m \rangle$, versus m , over training and test sets.	110
5.7	Establishing a minimum embedding dimension for the laser system by plotting the mean, normalised prediction error $\langle \epsilon_m \rangle$, versus m , over training and test sets.	111
5.8	Comparing the mean, normalised prediction errors $\langle \epsilon_{m,n} \rangle$ for the 0.01-step Lorenz system, reconstructed in singular subspaces of \mathbb{R}^{50} , \mathbb{R}^{20} and \mathbb{R}^{10} with a lag-10 delay map.	114
5.9	Comparing the mean, normalised prediction errors $\langle \epsilon_{m,n} \rangle$, obtained with $m = 10, 20$ and 50 , for the 0.01-step Lorenz system, generated with a lag-10 delay map from a noisy time series.	115
5.10	Comparing the mean, normalised prediction errors $\langle \epsilon_{m,n} \rangle$, obtained with $m = 10, 20$ and 50 , for the laser system.	117
5.11	Detecting periodic orbits in the Ikeda attractor by plotting the test error obtained by approximating the inverse to a family of FIR filtered delay embeddings into \mathbb{R}^5 .	119
5.12	Colour-coding the first three components of the Ikeda attractor, reconstructed with filter coefficients $a_1 = -1$ and $a_1 = 0$, by the errors to which they give rise in predicting the unfiltered Ikeda time series.	120
5.13	Prediction error for a FIR filtered delay embedding of the Hénon attractor into \mathbb{R}^4 .	121
5.14	Colour-coding the first three components of the Hénon attractor, reconstructed with filter coefficients $a_1 = -1$ and $a_1 = 0$, by the errors to which they give rise in predicting the unfiltered Hénon time series.	122
5.15	Prediction error obtained from a 9-delay filtered embedding of the 0.1-step Lorenz attractor.	123

5.16	Delay reconstructions of the 0.1-step Lorenz system, obtained with various FIR filter coefficients.	124
5.17	Colour-coding the reconstructions of the 0.1-step Lorenz system, obtained with non-generic and generic FIR filters, by the corresponding prediction errors.	125
5.18	Coarsely sampled sinusoid with additive deterministic noise generated from the Ikeda map, superimposed on the sinusoid itself.	129
5.19	Power spectra for a sinusoid with additive Ikeda chaos, before and after filtering.	130
5.20	Isolated chaotic signal obtained through blind prediction from the filtered time series.	131
5.21	Isolating a sinusoidal message by subtracting the blind predicted chaos from the corrupted signal, plotted in both time and frequency domains.	132
5.22	Isolated chaotic signal obtained through targeted prediction from a filtered time series.	133
5.23	Isolating a sinusoidal message by subtracting the targeted prediction of the chaos from the corrupted signal, plotted in both time and frequency domains.	134
5.24	A phase modulated message corrupted by additive Lorenz chaos is superimposed on both its individual message and chaos components.	135
5.25	Power spectra for a phase modulated message with additive Lorenz chaos, before and after filtering.	136
5.26	Chaotic time series estimated by blind prediction, superimposed on the original.	137
5.27	Reconstructing a message by subtracting the chaos estimated by blind prediction, in both time and frequency domains.	138
5.28	Result of demodulating the original and recovered messages.	139
B.1	Example command file for batch-mode time series analysis.	156

Chapter 1

Introduction

When analysing an experimental system, the experimenter does not have direct access to the state space of the system itself, but is instead forced to rely solely on the output of one or more probes, or sensors, incorporated into that system. Such probes might measure, for instance, the temperature somewhere in a human body, the rate of flow at a certain point in a hydrodynamic system or the current at a particular place in an electronic circuit. If the system under investigation is a stochastic one then a statistical approach is clearly indicated. If, on the other hand, it is a deterministic process then we are dealing with a dynamical system, whose evolution is uniquely determined, for all time, by its state at any given instant. As a consequence of dissipation in non-Hamiltonian systems, a dynamical system evolves asymptotically—and frequently chaotically, if the system is a nonlinear one—on a ‘differentiable manifold’, a topological space which is usually only *locally* Euclidean.

The output of a given probe is typically sampled to produce a scalar time series. If the probe in question is sufficiently responsive to the dynamics it is measuring, and the sampling rate is sufficiently fast, then it turns out to be generically possible (in a sense to be defined) to recover all of the important dynamical information in the system under investigation from that single time series, using what is known as the ‘method of delays’ [39]. This astounding result is achieved in practice by passing the time series through a ‘tapped delay line’ to obtain, with m taps, a sequence of vectors in \mathbb{R}^m which is a ‘differentiably equivalent’ copy \mathcal{M}_m of the manifold \mathcal{M} on which the state evolves, meaning that the topological and differentiable structure is preserved. We say that we have ‘embedded’ this system; the relationship between \mathcal{M} and its embedding \mathcal{M}_m is known as a ‘diffeomorphism’, a continuous, invertible and, in both directions, differentiable map.

Having obtained a reconstructed state space in this manner we are free to analyse it in place of the (inaccessible) original system, for instance to calculate dynamical or topological invariants [29], or to

construct predictive models for the time series itself [32]. This approach depends crucially on the fact that \mathcal{M}_m is *embedded* in \mathbb{R}^m , but it is often not clear, a priori, whether or not this is the case. In this thesis we attempt to answer this, and related questions by describing a methodology for determining whether or not a particular relationship $f: \mathcal{M}_m \rightarrow \mathcal{M}_n$ between two submanifolds $\mathcal{M}_m \subset \mathbb{R}^m$ and $\mathcal{M}_n \subset \mathbb{R}^n$ is a diffeomorphism. (It is important to note that we are concerned here only with that particular f which relates specific points in \mathcal{M}_m and \mathcal{M}_n , for instance by a common time index, and not with the more general question of whether or not \mathcal{M}_m and \mathcal{M}_n are diffeomorphic to each other.) This is a desirable test to be able to make, as it applies not only to delay reconstructions of dynamical systems but to diffeomorphisms between arbitrary submanifolds of Euclidean space. For instance, \mathcal{M}_m and \mathcal{M}_n might be delay reconstructions obtained using time series measured from two different probes in some experimental system, such as the rate of flow at two separate locations in a hydrodynamic system, and we might wish to know whether or not those two time series contain mutually independent information about the system in question: this could be thought of as a form of ‘nonlinear correlation’. Alternatively, we might be investigating a family of experimental systems by varying some control parameter, with \mathcal{M}_m and \mathcal{M}_n the embeddings of two different members of that family. We might then try to reveal the presence of a bifurcation at some critical value of that parameter by looking for a breakdown in the diffeomorphism between the two delay reconstructions. (Strictly speaking, this would only work if the particular bifurcation created a difference in the topological or differentiable structure of \mathcal{M}_m and \mathcal{M}_n ; moreover, for the approach adopted in this thesis, we would require additional information to allow us to index the two time series in a consistent way.) Finally, we might use such a test to establish whether or not applying a given filtering operation to a time series effects a qualitative change in the dynamics of the reconstructed system.

A statistical approach to this problem has been described by Pecora, Carroll and Heagy [31]. They propose separate tests for continuity and differentiability, which they apply to both f and its inverse. The test for continuity (for example, in the forward direction) consists of obtaining a set of data points sampled from \mathcal{M}_m and \mathcal{M}_n , defining an open ball $\mathcal{B}_\epsilon \subset \mathbb{R}^n$ of radius ϵ , centered at some $\mathbf{y}_0 \in \mathbb{R}^n$, then calculating the largest radius δ of the open ball $\mathcal{B}_\delta \subset \mathbb{R}^m$, centered at $\mathbf{x}_0 = f^{-1}(\mathbf{y}_0)$, such that every data point within \mathcal{B}_δ is mapped under f to a point inside \mathcal{B}_ϵ . For each ϵ , the number of points found inside \mathcal{B}_δ and \mathcal{B}_ϵ is incorporated into a statistic, averaged over a random set of open balls in \mathbb{R}^n , which reflects the likelihood that the data set is randomly distributed over \mathcal{M}_m and \mathcal{M}_n . A somewhat more complicated statistic for differentiability is based on the construction of linear approximations to the restriction of f to an open ball in \mathbb{R}^m . The authors find that their statistics yield positive results, but note that the success of the statistic for continuity depends upon a sensible choice of ϵ , so as to avoid complications due to the curvature of \mathcal{M}_n , and that the statistic for differentiability depends on prior knowledge of the topological dimension d of the manifold itself.

In contrast to this statistical technique, our approach is to construct empirical models $\widehat{f}: \mathbb{R}^m \rightarrow \mathbb{R}^n$ of

f and $\widehat{f^{-1}}: \mathbb{R}^n \rightarrow \mathbb{R}^m$ of its inverse, optimised over a set of data pairs sampled from \mathcal{M}_m and \mathcal{M}_n such that the restriction of \widehat{f} to \mathcal{M}_m is as good an approximation as possible (in the mean squared sense) to f , and the restriction of $\widehat{f^{-1}}$ to \mathcal{M}_n is a similarly good approximation to f^{-1} . We then base our decision on whether or not f is a diffeomorphism on an analysis of those models. The class of nonlinear models which we consider are the radial basis function (RBF) maps [6]. These maps can be shown, under certain conditions, to be universal approximators [33]. We will also show in this thesis that, under the appropriate conditions, RBF maps are generically diffeomorphisms—a result which has important consequences for the use to which we wish to put them. Although RBF maps can be implemented in the form of ‘feed-forward neural networks’, they also possess, over other such models, the significant advantage of being extremely easy to optimise, comprising an adaptive *linear* transformation composed with a fixed nonlinear one. (Strictly speaking, the nonlinear part can also be adapted—a process which consists of selecting a set of ‘centers’ from the domain—but this is usually not part of the optimisation process.) The RBF map is usually optimised by minimising the least squares (LS) error, but we will also investigate the solution obtained by minimising the total least squares (TLS) error. The former results in a model which simply approximates its range, on average, as well as possible. The latter, implemented in a more symmetrical form, ignores the distinction between domain and range. Both have their advantages and disadvantages, and will be compared in the experiments to follow.

1.1 Overview

We will now give a brief overview of the structure of the thesis, which divides basically into two parts, each containing two chapters. The first part describes the theoretical results, chapter 2 introducing dynamical systems, diffeomorphisms and embeddings and chapter 3 approximation theory, as applied to the RBF map. The second part describes the application of these results to some experimental data sets, consisting in chapter 4 of simple submanifolds of Euclidean space and in chapter 5 of delay embedded dynamical systems. The conclusions form chapter 6. In addition, there are two appendices: appendix A contains a proof of the invertibility of the nonlinear part of the RBF map and appendix B describes a sophisticated time series analysis software package, written in the course of this work.

In chapter 2 we summarise the relevant theory behind the fields of dynamical systems analysis and time delay embedding. We begin by stating a few basic theorems, culminating in the definitions of the diffeomorphism and the differentiable manifold. We then use these definitions to define the dynamical system, in both discrete and continuous form, and in particular define the Ikeda [20], Hénon [17] and Lorenz [24] systems, all of which will be used extensively in the experiments to be described in chapters 4 and 5. We also describe an experimental system, a far infra-red NH_3 laser [19]. After discussing the role of dissipation in restricting the dynamics under observation to the manifolds in which we are interested,

we introduce the time delay embedding and define the differentiable equivalence which makes it such an important tool in the analysis of experimental systems. Finally, we establish the conditions under which a linear transformation can be an embedding of submanifolds of Euclidean space. We then use this result to justify the projection of such manifolds onto principal component, and other, bases as a valid, and useful, experimental technique.

Chapter 3 formally introduces the RBF map, briefly discussing an ad hoc method for the selection of centers from training data and establishing the necessary results on universal approximation and diffeomorphism. We then go on to consider the ways in which a given map can fail to be a diffeomorphism, and how we might go about detecting such a case, either directly, through the LS error, or by composing RBF maps as an approximation to the identity map, finding Lipschitz constants with which to bound the resulting per-point errors. The LS solution is derived from geometrical considerations, for both the standard case of a fixed nonlinear transformation and the more specialised case known as ‘forward selection’, in which a more nearly optimal set of centers is chosen from the training set during adaptation. These two methods are justified experimentally by comparison with a random choice of centers. We then consider the problem of over-fitting, and its control through rank-reduction methods. We also introduce a novel re-ordering of the basis functions designed to achieve an optimal error with the smallest possible rank. This method is again compared with the standard methods on an experimental data set. Finally, we define the symmetrical variant of the RBF map used to calculate the TLS solution, again from geometrical considerations, and make analytic comparisons between the forward and inverse LS errors and both the TLS error and the condition numbers of the TLS submatrices.

In chapter 4 we describe the results of applying LS and TLS RBF maps to the detection of diffeomorphisms between subsets of Euclidean space without any intrinsic dynamical structure. The chapter contains two distinct, but similar experiments, which act as a simple test bed for the techniques involved. The first of these concerns maps between a family of projections into the plane of a topological circle in \mathbb{R}^3 . The images of these projections vary continuously with a control parameter between a circle and a figure-of-eight. The task which we set ourselves is to determine whether or not two such projections are diffeomorphic to each other, and hence to establish a critical value of the control parameter separating circles from figures-of-eight. In the LS case we investigate not only the LS error, as a measure of diffeomorphism, but also the distribution of errors arising from the approximation to the identity discussed in chapter 3. The second experiment concerns maps between a family of 2-tori in \mathbb{R}^3 and \mathbb{R}^4 , obtained by varying the ratio of the two radii. Once again, we attempt to establish the parameter range within which these tori are diffeomorphic to each other, using LS and TLS RBF maps.

Finally, in chapter 5 we begin applying the test for diffeomorphism to maps on delay reconstructed dynamical systems. The chapter deals with three related experiments, variously carried out on manifolds obtained from the numerically simulated Ikeda, Hénon and Lorenz systems introduced in chapter 2, and also from a time series of intensities measured in the laser experiment, all of which concern the effect

of a linear transformation on the manifold in question. The first of these experiments is an investigation into whether or not we can determine a minimum embedding dimension for a given system, based on fitting maps between reconstructions at ever-increasing delay lengths. The second examines the result of projecting an embedded manifold onto a lower-dimensional principal component basis, and again tries to determine a minimum value for the dimension of this basis. The third experiment introduces the concept of a linear transformation which corresponds to an embedding of a finite impulse response (FIR) filtered copy of the original time series, and uses this approach to look for periodic orbits in embedded attractors. Finally, the same technique is applied to the separation of chaotic ‘noise’ from a transmitted signal; the two experiments in this last section were performed in collaboration with David Broomhead and Jerry Huke at the Defence Research Agency, and the results have been independently reported in Broomhead, Huke and Potts [4].

Appendix A presents a proof that the nonlinear transformation from which the RBF map inherits its powerful approximation properties is, under certain circumstances, an embedding of compact subsets of its domain. This result is relied on heavily throughout the thesis, as it enables us to show that RBF maps are generically diffeomorphisms. Having shown that the transformation is invertible, this appendix then briefly goes on to discuss the implementational issues concerned with obtaining a LS approximation to this inverse.

Carrying out the investigations described in this thesis has entailed the implementation of a large, and fairly powerful, time series analysis package, described in some detail in appendix B. This program allows the user to load in one or more (multi-dimensional) time series and, from each one, make any number of delay reconstructions of the underlying system, incorporating variable lags and principal component subspaces. These data sets can then be related by fitting LS and TLS RBF maps between them. The program has an X Window front end for real-time data analysis, but can also be run in batch mode, using a full-featured, purpose-built scripting language. Combined with a sophisticated scheduling system, this batch mode of operation allows the program to make efficient use of its computing resources through the fitting of multiple RBF maps in parallel as subprocesses on multi-processor Unix workstations.

In the remainder of this thesis, where it will aid understanding, figures depicting three-dimensional objects will be displayed as stereo pairs. These can be viewed, without any external apparatus, by holding the figure at arm’s length and allowing the lines of sight from both eyes to converge in a plane a little distance behind the plane of the figure, so that the two images are fused stereoscopically. We will also make use of colour in certain figures, for which we apologise if this copy of the thesis is in black and white.

Chapter 2

Dynamical systems and embeddings

Although we will assume that the reader has a passing familiarity with analysis and differential topology (good introductions may be found in Chillingworth [10], Hirsch [18] and Milnor [28]), for the sake of completeness we begin by briefly stating a few of the more important definitions. Consider a map $f: \mathcal{M} \rightarrow \mathcal{N}$ between subsets $\mathcal{M} \subset \mathbb{R}^m$ and $\mathcal{N} \subset \mathbb{R}^n$, not necessarily open. We say that f is differentiable if it is the restriction to \mathcal{M} of some differentiable map on an open set in \mathbb{R}^m containing \mathcal{M} . We call f a homeomorphism if f is invertible and both f and f^{-1} are continuous. We say that f is a diffeomorphism if f is a homeomorphism and both f and f^{-1} are differentiable. As already mentioned, we will be particularly interested in diffeomorphisms of differentiable manifolds. We say that $\mathcal{M} \subset \mathbb{R}^m$ is a d -dimensional differentiable manifold if it is locally diffeomorphic to \mathbb{R}^d ; loosely speaking, \mathcal{M} is locally, differentiably parameterised by open subsets of \mathbb{R}^d . From now on we will take the differentiability of \mathcal{M} as read, and just call it a manifold.

Another useful map on manifolds, on which we rely strongly, is the embedding. We are interested in embeddings because if f is an embedding then $f: \mathcal{M} \rightarrow f\mathcal{M}$ is a diffeomorphism and its image $f\mathcal{M}$ is itself a manifold. In order to define an embedding we must first make the following definition: f is an immersion if f is differentiable and its derivative Df is everywhere injective. We say that f is an embedding if it is an immersive homeomorphism. Indeed, for an arbitrary subset \mathcal{M} of \mathbb{R}^m we can say, as above, that f is an embedding if it is the restriction to \mathcal{M} of some embedding on a manifold in \mathbb{R}^m containing \mathcal{M} . In fact it turns out that if we restrict our attention to compact subsets then f is an embedding if it is an injective immersion. We make use of this latter result in appendix A:

In the rest of this chapter we define pairs of diffeomorphisms and manifolds which together constitute nonlinear dynamical systems. These will arise frequently in the following chapters. We then go on to introduce the method of delays, with which we can construct empirical models of these systems directly

from scalar time series of measurements made on them. Finally, we discuss linear transformations of embedded manifolds, and the conditions under which they are diffeomorphisms.

2.1 The dynamical system

A discrete dynamical system is the pair (\mathcal{M}, ψ) , consisting of a differentiable manifold \mathcal{M} which is mapped onto itself by a diffeomorphism $\psi: \mathcal{M} \rightarrow \mathcal{M}$. The elements of \mathcal{M} are called the states of the system; they evolve as $\xi \mapsto \psi\xi$ or, more generally, under p iterations as $\xi_p = \psi \circ \dots \circ \psi \xi_0 = \psi^p \xi_0$. We therefore call \mathcal{M} the state, or phase space. Examples of discrete dynamical systems which we will make use of in this thesis are the Ikeda system [20] and the Hénon system [17]. The Ikeda system evolves under a family of diffeomorphisms $\psi: \mathbb{R}^2 \rightarrow \mathbb{R}^2$ defined by

$$\xi \mapsto \begin{pmatrix} 1 - \mu(\xi_1 \cos \theta - \xi_2 \sin \theta) \\ \mu(\xi_1 \sin \theta + \xi_2 \cos \theta) \end{pmatrix} \quad (2.1)$$

where $\theta = \alpha - \beta/(1 + \|\xi\|^2)$ and $\alpha, \beta, \mu \in \mathbb{R}$ are control parameters. We select a particular Ikeda map by taking $\alpha = 0.4$, $\beta = 6$ and $\mu = 0.7$ in all of the experiments described in this thesis. The Hénon system evolves, also in \mathbb{R}^2 , under the family of diffeomorphisms

$$\xi \mapsto \begin{pmatrix} 1 - a\xi_1^2 + \xi_2 \\ b\xi_1 \end{pmatrix} \quad (2.2)$$

with control parameters $a, b \in \mathbb{R}$. We select the Hénon map defined by $a = 1.4$ and $b = 0.3$ in this thesis.

We will also consider continuous dynamical systems, in which \mathcal{M} is mapped under a flow $\{\psi_t\}_{t \in \mathbb{R}}$, which is a group of diffeomorphisms $\psi_t: \mathcal{M} \rightarrow \mathcal{M}$ with $\psi_s \psi_t = \psi_{s+t}$ and identity ψ_0 . In this case the state describes a continuous trajectory $\xi: \mathbb{R} \rightarrow \mathcal{M}$ parameterised by t , so that $\xi(t) = \psi_t \xi(0)$. In numerical simulations a particular member $\psi \equiv \psi_t$ of $\{\psi_t\}$ is typically selected by integrating a set of differential equations with a fixed step size t , for instance by the Runge Kutta method [34], to obtain a set of discrete values $\xi_i \equiv \xi(it)$ sampling the continuous trajectory ξ . We say that the continuous system induces a discrete system (\mathcal{M}, ψ) , and we notionally work directly with this induced system. In physical systems, this sampling is usually imposed by the act of measurement, to be described in the next section. An example of a continuous dynamical system which we will make use of later is the Lorenz system [24]. This may be thought of as a family of flows on the manifold $\mathcal{M} = \mathbb{R}^3$ obtained by solving the ordinary differential equations

$$\dot{\xi} = \begin{pmatrix} -\sigma(\xi_1 - \xi_2) \\ r\xi_1 - \xi_2 - \xi_1\xi_3 \\ \xi_1\xi_2 - b\xi_3 \end{pmatrix} \quad (2.3)$$

with parameters $r, \sigma, b \in \mathbb{R}$, where $\dot{\xi}$ is the derivative of ξ with respect to t . We take $\sigma = 10$, $b = \frac{8}{3}$ and $r = 28$ in this thesis.

2.1.1 The role of dissipation

Since ψ is a diffeomorphism, ξ_{i+1} is uniquely determined by ξ_i (and vice versa), and the state evolves through a set of points in \mathcal{M} called a trajectory. The trajectories obtained from these systems by iteration or integration do not, in general, occupy their entire state space \mathcal{M} but instead evolve asymptotically on subsets of \mathcal{M} called attractors. This contraction of state space is often, in physical systems, the result of so-called dissipative forces, such as friction. To illustrate this principle, we show the attractors corresponding to equations (2.1) and (2.2) in figure 2.1, and that of equation (2.3) in figure 2.2, the latter having been integrated with a step size of 0.01. Figures 2.1(a) and (b) are plotted using separate points, since they are attractors of discrete maps, but figure 2.2 is plotted with lines joining consecutive points to illustrate the continuous nature of the Lorenz flow.

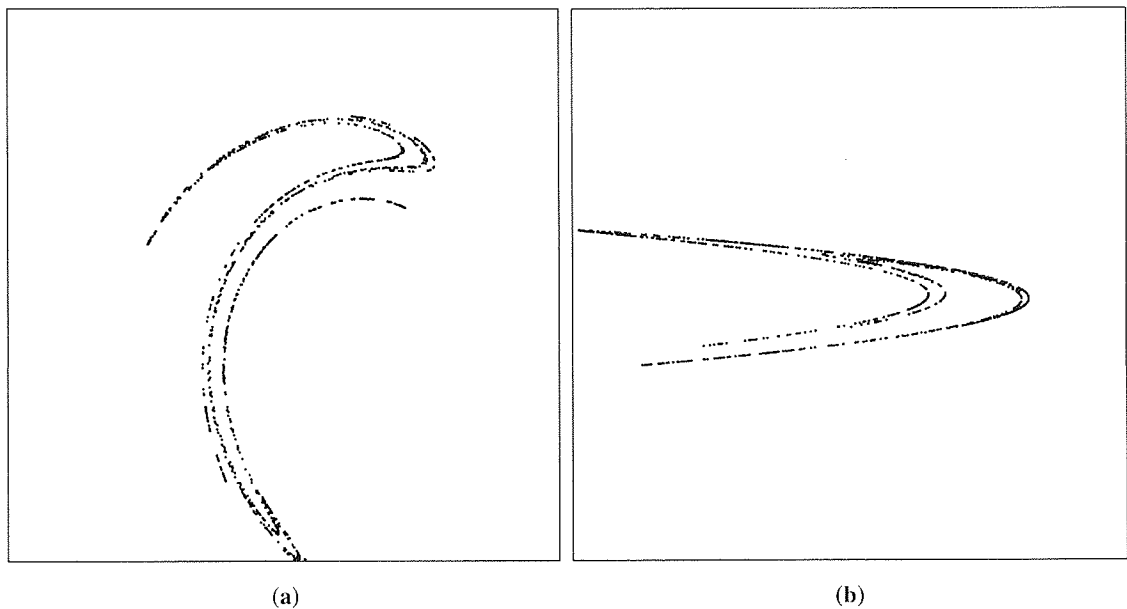


Figure 2.1 Attractors of discrete dynamical systems. Obtained with 1000 iterations each of (a) the Ikeda map and (b) the Hénon map.

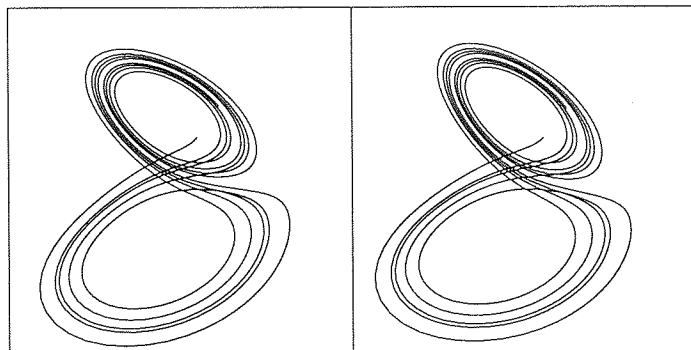


Figure 2.2 Continuous attractor of the Lorenz system. Obtained by integrating the Lorenz system in \mathbb{R}^3 for 1000 steps of size 0.01; plotted as a stereo pair.

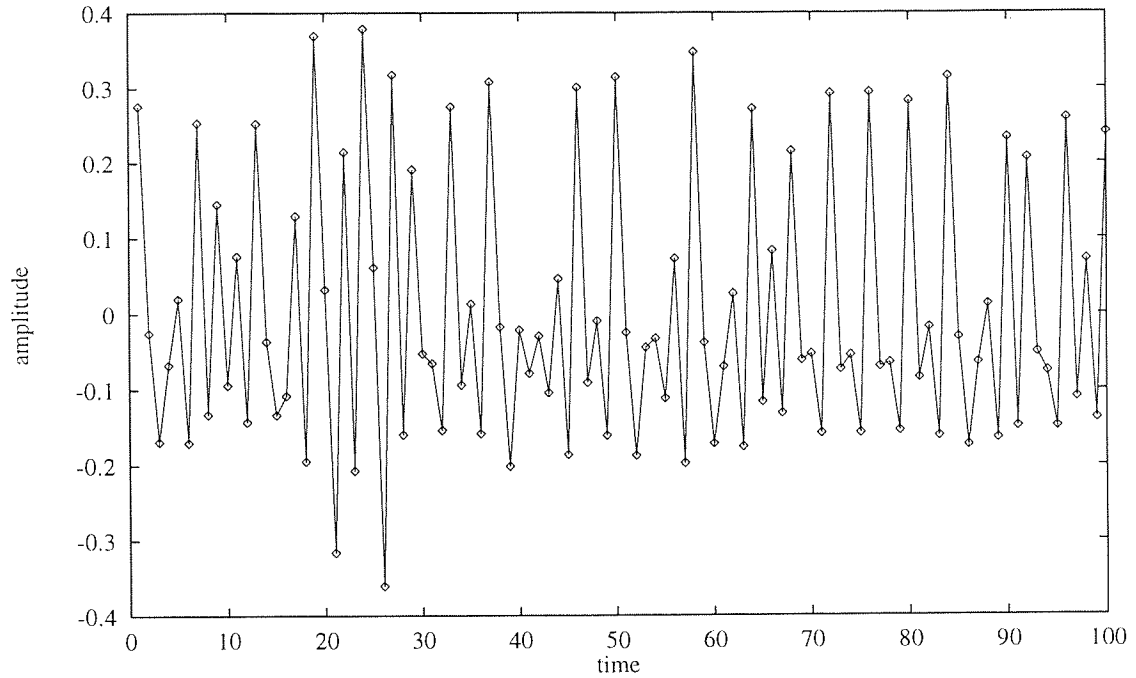
We categorise nonlinear dynamical systems by the manner in which small displacements in the state vector propagate under ψ . In many systems, within certain parameter ranges almost any small error in initial position on \mathcal{M} grows exponentially on small scales. This sensitive dependence on initial conditions is a defining property of chaotic dynamical systems. The multiplicative ergodic theorem of Oseledec [30] establishes the existence, for almost every $\xi \in \mathcal{M}$, of a set of constant ‘characteristic exponents’. The largest of these is an upper bound for the asymptotic, exponential growth rate of tangent vectors under repeated applications of $D\psi$, and so quantifies approximately the growth rate of small errors as ψ is iterated. Oseledec’s theorem is discussed by Eckmann and Ruelle in their review paper [12]; in Potts and Broomhead [32] the characteristic exponents, and their local analogues, of an experimental system are estimated using the techniques to be described in the chapters to follow.

2.2 The delay embedding

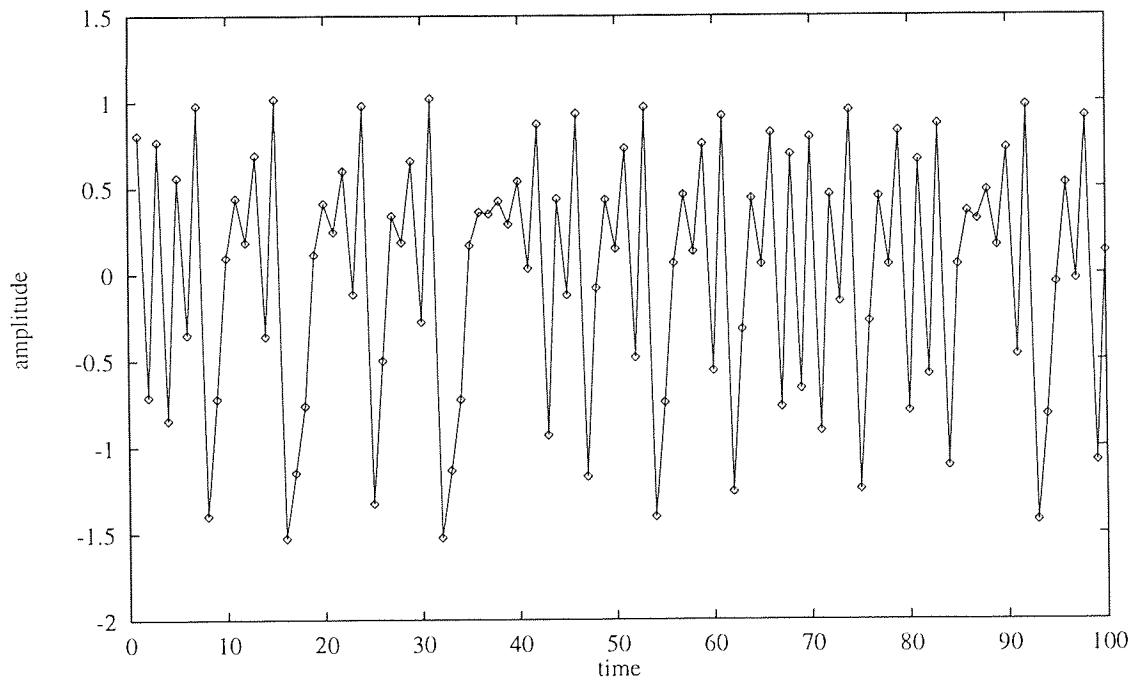
Although in numerical simulations, such as the ones discussed above, we can clearly measure the state of a system directly, under experimental conditions we would usually expect only to have available the output of one or more probes or sensors incorporated into that system. The output of one of these probes will be assumed to take the form of a smooth and, in general, nonlinear measurement function $v: \mathcal{M} \rightarrow \mathbb{R}$. By sampling this measurement function we obtain a time series $\{v_i\}$ of values $v_i \equiv v(\xi_i)$. In this thesis we will use time series constructed from the numerically simulated Ikeda, Hénon and Lorenz systems. For simplicity, we shall use measurement functions which correspond to (linear combinations of) single components of the state vectors ξ_i . Examples of time series obtained from the first component of the Ikeda and Hénon maps are shown in parts (a) and (b), respectively, of figure 2.3.

The effect of the sampling interval on time series obtained from continuous systems is illustrated in figure 2.4, which plots, in parts (a) and (b), the time series obtained from the numerically integrated Lorenz map for integration steps of 0.01 and 0.1, respectively, the latter interval clearly under-sampling the map quite severely. Although it is tempting to think of these time series merely as two measurements on a single system, differing only in their sampling intervals, it is important to note that a process such as Runge Kutta necessarily imposes its own dynamical structure on the system being integrated, so the two time series do not, in fact, correspond to the same Lorenz system.

We will also make use of an experimentally obtained time series of fluctuations in the intensity of a single-mode, far infra-red NH_3 laser, described in Hubner, Weiss, Abraham and Tang [19] and featured in the Santa Fe time series prediction competition [42]. Hubner shows how the semiclassical laser equations can be transformed into the Lorenz model, in which form the intensity can be regarded as the square of the first component of the Lorenz map. We illustrate this correspondence in figure 2.5, plotting the time series of intensity fluctuations in part (a) and the squared Lorenz time series (integration step 0.01) in part (b).

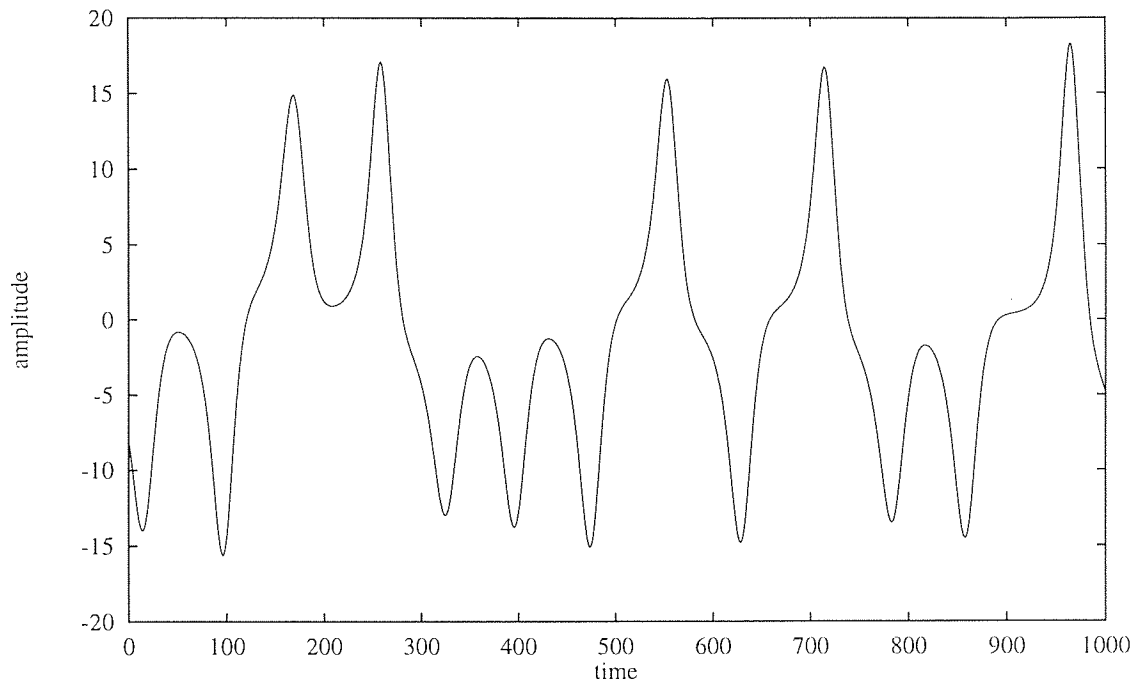


(a)

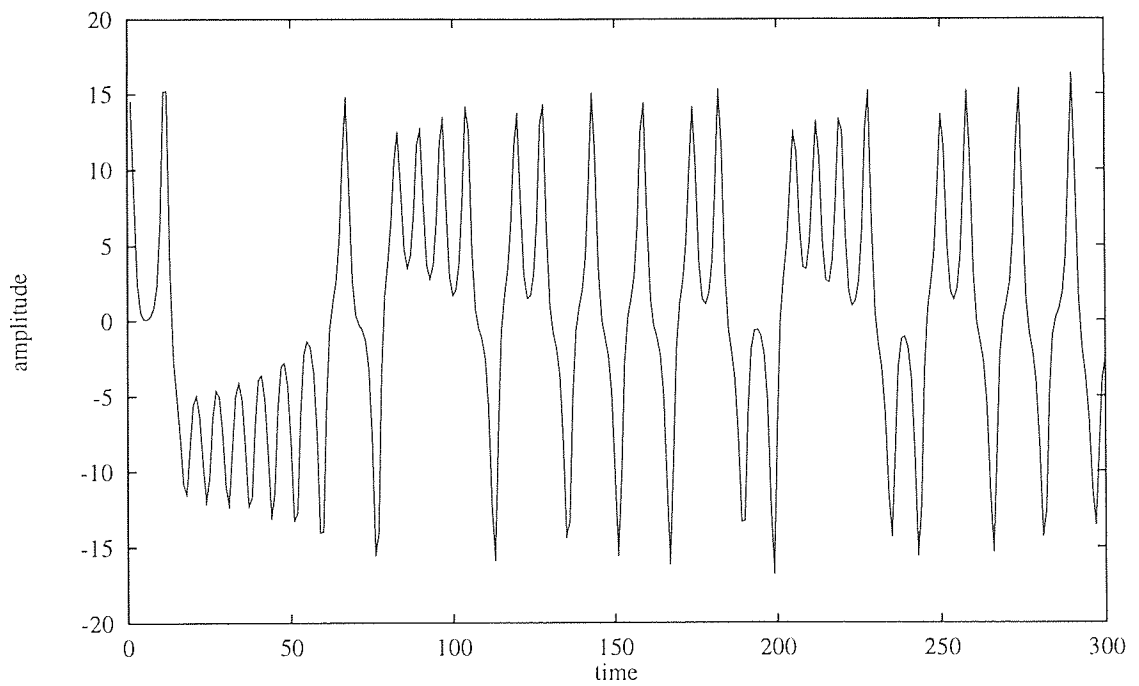


(b)

Figure 2.3 Time series measured from numerically simulated Ikeda and Hénon systems. Plotting 100 samples from the first component of (a) the Ikeda map and (b) the Hénon map.

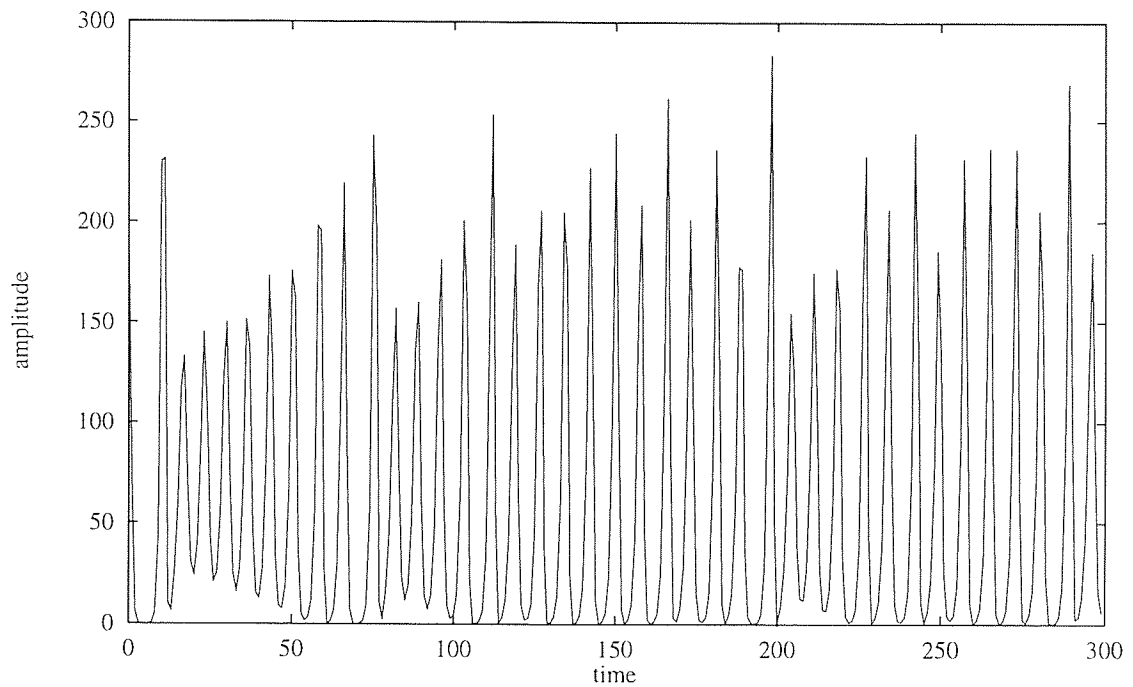


(a)

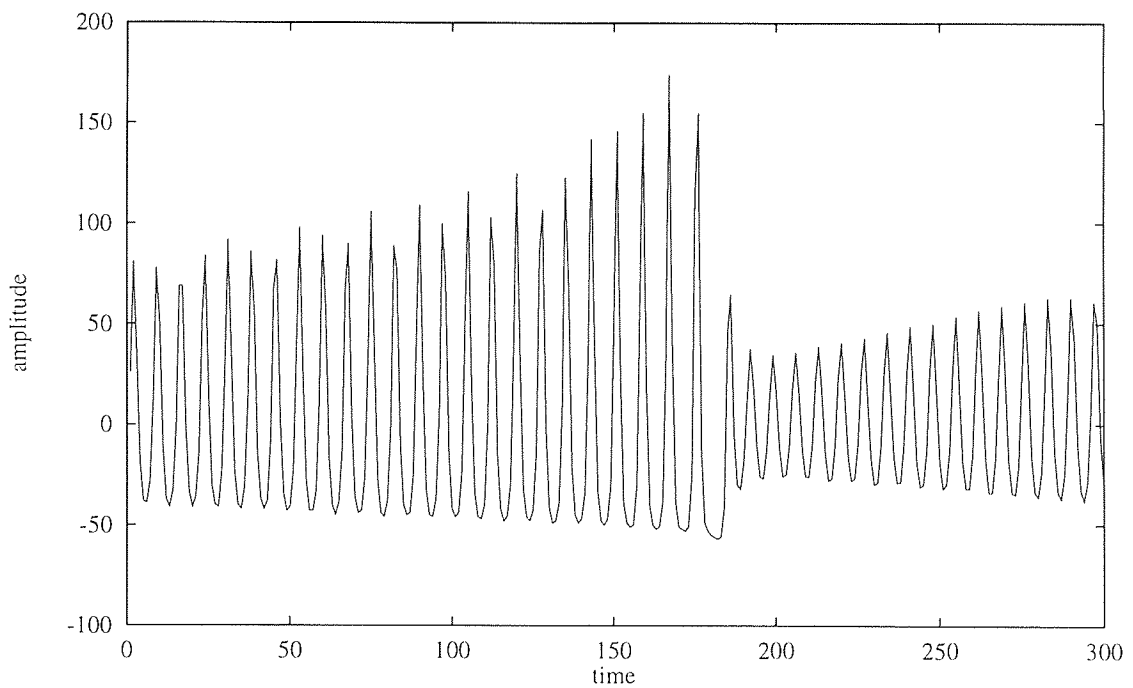


(b)

Figure 2.4 Time series from the first component of the numerically integrated Lorenz system. Obtained by integrating the Lorenz map for (a) 1000 samples with a step size of 0.01 and (b) 300 samples with a step size of 0.1.



(a)



(b)

Figure 2.5 Comparing the intensity time series measured from a laser experiment with a time series measured on the Lorenz system. Part (a) plots 300 samples from the square of the first component of the Lorenz system, integrated with a step size of 0.01, and can be seen to be similar in form to the laser intensity time series, in part (b), obtained by Hubner, shown here renormalised with zero mean.

These time series are, by definition, a function of the dynamical system on which they are measured, but they represent a bottleneck in the flow of information. Is it possible that such a function might, under suitably general conditions, contain enough dynamical information to enable us in some way to reconstruct the system under observation? Takens [39] has shown that this is indeed possible: his work provides us with a method which underpins the entire field of experimental dynamical systems analysis. Following Takens, we define a map $\Phi_{v,m}: \mathcal{M} \rightarrow \mathbb{R}^m$ by

$$\mathbf{x}_i = \Phi_{v,m}(\xi_i) = \begin{pmatrix} v(\xi_i) \\ v(\psi^{-1}\xi_i) \\ \vdots \\ v(\psi^{-(m-1)}\xi_i) \end{pmatrix} = \begin{pmatrix} v_i \\ v_{i-1} \\ \vdots \\ v_{i-(m-1)} \end{pmatrix} \quad (2.4)$$

so that each component of \mathbf{x}_i is a delayed element of the time series $\{v_i\}$, that is,

$$(\mathbf{x}_i)_j = v_{i-(j-1)} \quad (2.5)$$

From an operational viewpoint, this map is constructed by passing the time series through a so-called delay window—specifically, an m -delay window—also known as a tapped delay line, whose contents are the elements of the $\mathbf{x}_i \in \mathbb{R}^m$. We call the \mathbf{x}_i delay vectors. Takens showed that, for generic choices of manifold \mathcal{M} and measurement function v , $\Phi_{v,m}$ is an embedding of \mathcal{M} , provided that the number of delays $m > 2d$, where d is the dimension of \mathcal{M} . (In certain simple cases it may be possible to embed \mathcal{M} for some $d \leq m \leq 2d$.) That is, $\Phi_{v,m}: \mathcal{M} \rightarrow \Phi_{v,m}\mathcal{M}$ is a diffeomorphism. We call $\mathcal{M}_m = \Phi_{v,m}\mathcal{M}$ the pseudo-phase space and call $\Phi_{v,m}$ a delay map, or a delay reconstruction; if $\Phi_{v,m}$ is an embedding then we call it a delay embedding. The process by which $\Phi_{v,m}$ is constructed is called the method of delays, and a more thorough analysis may be found in Sauer, Yorke and Casdagli [36].

The delay structure in the co-domain of $\Phi_{v,m}$ becomes clear when we form the N by m trajectory matrix \mathbf{X} , whose rows are the \mathbf{x}_i , for $i = 1, \dots, N$, with which we perform the experiments in later chapters: the j -th column of \mathbf{X} is just a shifted copy of the time series, with elements $\{v_{j+(N-1)}, \dots, v_j\}$. Matrices of this form are called Hankel. As a trivial consequence of equation (2.5) a sequence of delay vectors \mathbf{x}_i can be seen to obey the relationship

$$(\mathbf{x}_i)_j = (\mathbf{x}_{i+k})_{j+k} \quad (2.6)$$

which we call the shift property of delay embeddings. This constraint becomes an important consideration when we consider approximating diffeomorphisms on delay embeddings in chapter 5.

2.2.1 Differentiable equivalence

Suppose we have a delay embedding $\Phi_{v,m}$, so that $\Phi_{v,m}: \mathcal{M} \rightarrow \mathcal{M}_m$ is a diffeomorphism. We can compose $\Phi_{v,m}$ with $\psi: \mathcal{M} \rightarrow \mathcal{M}$ to induce a new diffeomorphism $\psi_m: \mathcal{M}_m \rightarrow \mathcal{M}_m$, with $\psi_m = \Phi_{v,m} \circ \psi \circ \Phi_{v,m}^{-1}$ or, by composition, $\psi_m^p = \Phi_{v,m} \circ \psi^p \circ \Phi_{v,m}^{-1}$, so that $x_{i+1} = \psi_m x_i$. (The subscript on ψ_m is not to be confused with its earlier use in denoting the flow ψ_t .) The usefulness of the method of delays is a consequence of the following result: it can be shown [16] that the original system (\mathcal{M}, ψ) and its delay embedding (\mathcal{M}_m, ψ_m) are differentially equivalent. For instance, since \mathcal{M} is homeomorphic to \mathcal{M}_m it inherits the same topology, which means that topological features such as periodic (and hence also fixed) points are invariant under $\Phi_{v,m}$. Thus, if $\xi_0 = \psi^p(\xi_0)$ is a periodic point of ψ then $x_0 = \Phi_{v,m}(\xi_0)$ is a periodic point of ψ_m , since

$$\psi_m^p(x_0) = \Phi_{v,m} \circ \psi^p(\xi_0) = \Phi_{v,m}(\xi_0) = x_0 \quad (2.7)$$

That $\Phi_{v,m}$ is a diffeomorphism implies more, however: by differentiating $\psi_m^p \circ \Phi_{v,m} = \Phi_{v,m} \circ \psi^p$ to get

$$D\psi_m^p(x_0) = D\Phi_{v,m}(\xi_0) D\psi^p(\xi_0) [D\Phi_{v,m}(\xi_0)]^{-1} \quad (2.8)$$

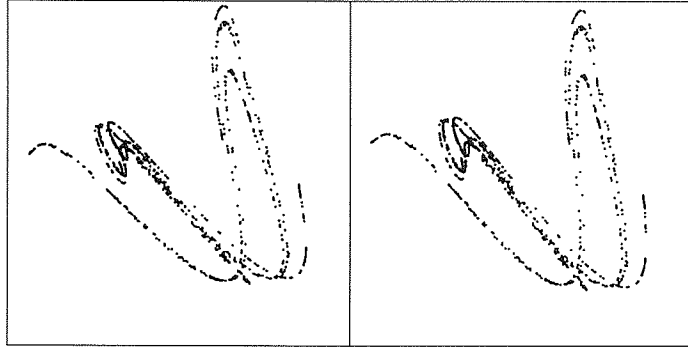
we see that $D\psi_m^p(x_0)$ and $D\psi^p(\xi_0)$ are similar matrices and therefore have the same eigenvalues. In other words, $\Phi_{v,m}$ preserves the eigenvalues at periodic points. This result can seem almost magical at first sight: by recording the rate of flow at a single location in our hypothetical fluid system we have used Takens' theorem to recover a differentially equivalent copy of this system, simply by constructing delay vectors from a scalar time series. We will show in later chapters how we can build models of this embedded system with which to analyse or predict the original.

Figures 2.6(a) and (b) show a three-delay reconstruction of the Ikeda attractor and a two-delay reconstruction of the Hénon attractor, obtained from the time series in figures 2.3(a) and (b), respectively. These may be compared to the attractors shown in figures 2.1(a) and (b) in the previous section. The Ikeda attractor turns out to be impossible to embed with less than four delays, due to self-intersections, but the Hénon attractor is clearly embedded in \mathbb{R}^2 . The reason for this distinction between the two maps is that the latter may be written as

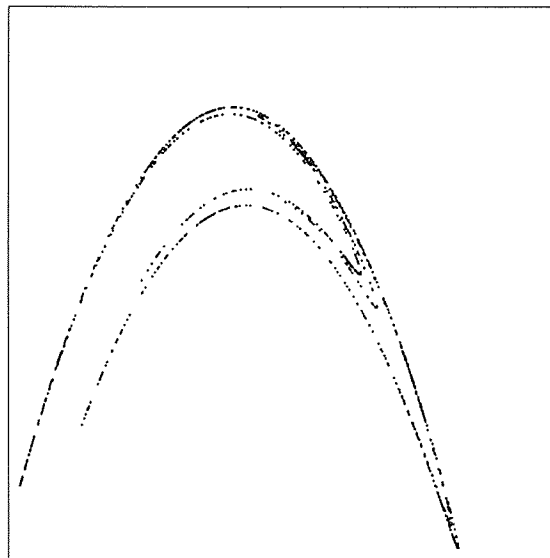
$$(\xi_{i+1})_1 = 1 - a(\xi_i)_1^2 + b(\xi_{i-1})_1 \quad (2.9)$$

by eliminating ξ_2 from equation (2.2), and hence is completely determined by two delays.

Figure 2.7(a) shows a three-delay reconstruction of the Lorenz attractor from the 0.01-step time series of figure 2.4(a), to be compared with the attractor shown in figure 2.2. It is a consequence of the 0.01 step size that this reconstructed attractor is nearly one-dimensional: the components of any given delay vector

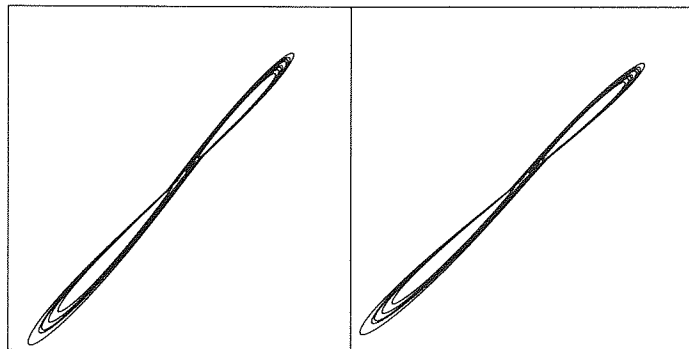


(a)

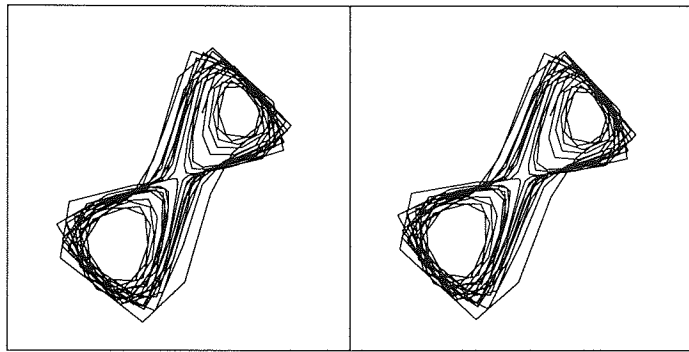


(b)

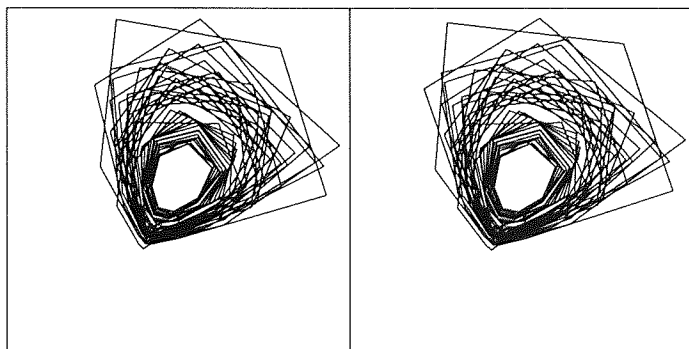
Figure 2.6 Embedding the Ikeda and Henon systems. Showing in (a) a stereo plot of the first three components of a delay embedding from a time series of 1000 samples of the first component of the Ikeda map and in (b) the attractor corresponding to a delay embedding from a time series of 1000 samples of the first component of the Hénon map.



(a)



(b)



(c)

Figure 2.7 Embedding the Lorenz and laser systems. Parts (a) and (b) were obtained from the first component of the Lorenz map, integrated (a) for 1000 steps of size 0.01, and (b) for 300 steps of size 0.1, and reconstructed in \mathbb{R}^3 . The overly small integration step which characterises (a) is revealed by the narrowness of the reconstructed object. Part (c) was obtained with a three-delay reconstruction from 300 samples of the laser intensity time series, illustrating the ‘folding over’ of the laser attractor in comparison to the Lorenz attractor. Plotted in stereo.

are almost equal, so the reconstructed trajectory lies close to the diagonal in \mathbb{R}^3 . In order to combat this effect we may wish to adopt the filtered embedding strategy described in the next section to sub-sample or otherwise transform the time series. In contrast, the attractor in figure 2.7(b), which is a three-delay embedding of the Lorenz attractor obtained with an integration step of 0.1, is far better resolved in \mathbb{R}^3 , but at the price of a very sparse trajectory. In figure 2.7(c) we show a three-delay reconstruction of the attractor corresponding to the time series of laser intensities plotted in figure 2.5(b). Owing to the relationship between the laser and the Lorenz systems it is useful to compare figure 2.7(c) with the attractor in figure 2.7(b). The fact that these two reconstructed attractors seem topologically distinct leads us to the conclusion that intensity is not a generic measurement, in the sense of Takens, for the laser system: the effect of using intensity is to identify the two unstable foci in figure 2.7(b).

There are some practicalities to be taken into account when using the method of delays: most obviously, we would not usually expect to have prior knowledge of d , the manifold dimension. Various methods for empirically determining a minimum embedding dimension m from time series have been proposed [21, 22, 5], and a method based on approximating ψ_m for various values of m will be described in section 5.1.

2.3 Filtered embedding

We will also make use in this thesis of linear transformations of embedded manifolds. The Whitney embedding theorem states that a d -dimensional manifold \mathcal{M} can be embedded in \mathbb{R}^n provided that $n > 2d$. In proving this theorem, Hirsch [18] shows that given a d -dimensional submanifold $\mathcal{M}_m \subset \mathbb{R}^n$ (which in our case will be a delay embedding $\mathcal{M}_m = \Phi_{v,m}\mathcal{M}$) then, generically, a (not necessarily orthogonal) projection $\mathcal{F}: \mathbb{R}^m \rightarrow \mathbb{R}^n$, with $2d < n < m$, is an embedding on \mathcal{M}_m . We define

$$\mathcal{F}(x) = c_n + F(x - c_m) \quad (2.10)$$

where $c_m \in \mathbb{R}^m$ and $c_n \in \mathbb{R}^n$ are translations and F is an n by m matrix, and write $\mathcal{M}_n = \mathcal{F}\mathcal{M}_m$. We can generalise this result by writing F in terms of the n by n matrices U and Σ and the m by n matrix V making up its singular value decomposition (SVD), defined by

$$F = U\Sigma V^T \quad (2.11)$$

where the columns of U are left singular vectors of F —an orthonormal basis for \mathbb{R}^n —the columns of V are right singular vectors of F —an orthonormal basis for the row space of F —and Σ is a diagonal matrix whose diagonal elements are the singular values $\sigma_j \in \mathbb{R}$ of F , ordered such that $\sigma_j \geq \sigma_{j+1} \geq 0$. (We will make extensive use of the SVD, in a more general form, in chapter 3.) Decomposed in this manner we see

that \mathbf{F} just consists of a projection \mathbf{V}^T which, according to Hirsch, is generically an embedding, followed by a scaling by the singular values in Σ , which is clearly an embedding provided that all $\sigma_j > 0$, and finally multiplication by \mathbf{U} , which is just a change of basis. So the only further condition on \mathbf{F} necessary for it to be an embedding of \mathcal{M}_m is that $\text{rank } \mathbf{F} = n$. This is, of course, a generic argument, and we must be careful when carrying it over into the case of specific projections \mathbf{F} to check that it still applies in each case. We can extend the argument in section 2.2.1 to accommodate a linear transformation of this form by composing \mathcal{F} with the delay embedding $\Phi_{v,m}$ to define a new embedding $\mathcal{G} = \mathcal{F} \circ \Phi_{v,m}$, which induces the differentiably equivalent dynamical system (\mathcal{M}_n, ψ_n) , with $\psi_n = \mathcal{G} \circ \psi \circ \mathcal{G}^{-1}$ and $\mathcal{M}_n = \mathcal{G}\mathcal{M}$.

It is instructive to write the image $\mathbf{y}_i = \mathcal{F}\mathbf{x}_i$ of the delay vector $\mathbf{x}_i \in \mathcal{M}_m$ in terms of the time series elements $v_{i-(m-1)}, \dots, v_i$ from which it is constructed. We write

$$\begin{aligned} (\mathbf{y}_i)_j &= (\mathbf{c}_n)_j + \sum_{k=1}^m F_{jk}[(\mathbf{x}_i)_k - (\mathbf{c}_m)_k] \\ &= (\mathbf{c}_n)_j + \sum_{k=1}^m F_{jk}v_{i-k+1} - \sum_{k=1}^m F_{jk}(\mathbf{c}_m)_k \end{aligned} \quad (2.12)$$

from which, neglecting contributions from the constant terms, we see that the j -th component of \mathbf{y}_i generates a time series $\{u_i^{(j)}\}$, with $u_i^{(j)} = (\mathbf{y}_i)_j$, which we can view as the output of one of n finite impulse response (FIR) filters applied to the original time series $\{v_i\}$. We can express this relationship in a more conventional form by making the definition $a_{k-1}^{(j)} = F_{jk}$, so that

$$u_i^{(j)} \propto \sum_{k=0}^{m-1} a_k^{(j)} v_{i-k} \quad (2.13)$$

where $a_0^{(j)}, \dots, a_{m-1}^{(j)}$ are the coefficients of the j -th FIR filter. For this reason \mathcal{F} is frequently referred to as a filter bank when applied to a sequence of delay vectors; if \mathcal{F} embeds \mathcal{M}_m then we call it a filtered embedding. In general, each generated time series $\{u_i^{(j)}\}$ is the output of a distinct FIR filter, with the result that the N by n matrix $\mathbf{Y} = \mathbf{F}\mathbf{X}$, whose rows are the \mathbf{y}_i , is not Hankel and hence \mathcal{F} is not a delay map. Its inverse $\mathcal{F}^{-1}: \mathcal{M}_n \rightarrow \mathcal{M}_m$, on the other hand, is a delay map, since its co-domain is in the image of $\Phi_{v,m}$. We therefore write $\Phi_{w,m} \equiv \mathcal{F}^{-1}$, where $w: \mathcal{M}_n \rightarrow \mathbb{R}$ is the ‘measurement function’ induced by the delay structure in \mathcal{M}_m , and is defined by $w(\mathbf{y}_i) = v_i$ or, if \mathcal{G} is invertible, $w = v \circ \mathcal{G}^{-1}$. We will make use of this definition in chapter 5 when we consider fitting maps between embedded dynamical systems.

2.3.1 Sampling

We have already referred to the use of Hirsch’s result in the previous section, where we discussed the construction of a delay embedding $\Phi_{v,m}: \mathcal{M} \rightarrow \mathbb{R}^m$, with $m \gg 2d$, from a finely-sampled time series.

A method which is used to combat the near-redundancy of neighbouring elements in the resulting delay vectors $\mathbf{x}_i \in \mathbb{R}^m$ is to make the definition $\mathbf{F} = (\mathbf{e}_1, \mathbf{e}_{1+\tau}, \dots, \mathbf{e}_m)^\top$, where $\mathbf{e}_j \in \mathbb{R}^m$ is the unit vector which picks out the j -th dimension of \mathbb{R}^m and $\tau \geq 0$ is an integer number of samples called the ‘lag’; c_m and c_n are just set to the zero vector in \mathbb{R}^m and \mathbb{R}^n , respectively. We can now extend the definition of the delay map, writing the lagged analogue of equation (2.5) as

$$(\mathbf{y}_i)_j = \mathbf{e}_{1+(j-1)\tau} \cdot \mathbf{x}_i = v_{i-(j-1)\tau} \quad (2.14)$$

and the corresponding lagged shift property as

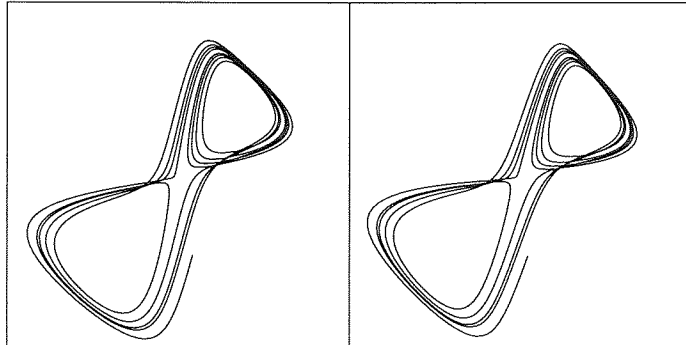
$$(\mathbf{y}_i)_j = (\mathbf{y}_{i+k\tau})_{j+k} \quad (2.15)$$

This projection does not simply sub-sample the time series, as can be seen from an examination of the trajectory matrix \mathbf{Y} , whose j -th column is once more just a shifted copy of the time series, with elements $\{v_{N+(j-1)\tau}, \dots, v_{1+(j-1)\tau}\}$. Instead, we have sub-sampled the m -delay window itself, maintaining its overall width whilst reducing the number of delays within it. We will refer to this new window as an (m, τ) -delay window. It can be shown that Takens’ theorem still holds for lagged delay vectors of this form, so as \mathcal{F} preserves the shift property we are able to write \mathcal{G} in the form of a lagged delay map $\Phi_{v,n,\tau} \equiv \mathcal{G} = \mathcal{F} \circ \Phi_{v,m}$ by defining

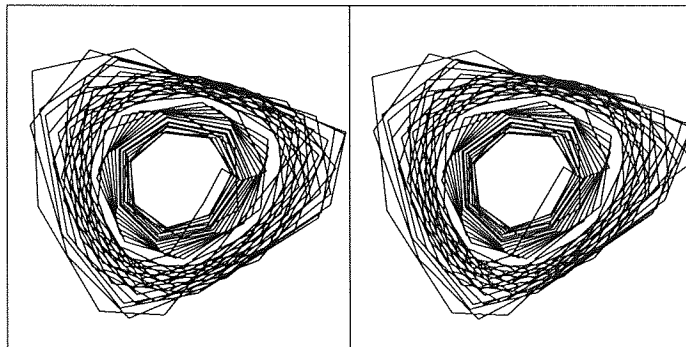
$$\mathbf{y}_i = \Phi_{v,n,\tau}(\xi_i) = \begin{pmatrix} v(\xi_i) \\ v(\psi^{-\tau}\xi_i) \\ \vdots \\ v(\psi^{-(m-1)\tau}\xi_i) \end{pmatrix} = \begin{pmatrix} v_i \\ v_{i-\tau} \\ \vdots \\ v_{i-(m-1)\tau} \end{pmatrix} \quad (2.16)$$

which is an embedding of \mathcal{M} provided that $n > 2d$. Indeed, the theorem holds even if we allow the lag to vary between consecutive delays, forming \mathbf{F} from some arbitrary subset of $\{\mathbf{e}_j\}$. By composition, therefore, we see that this specialised form of projection \mathcal{F} embeds \mathcal{M}_m . We will make explicit use of this result in section 5.1, when we come to consider the determination of a minimum embedding dimension for \mathcal{M} .

As an illustration of the effect of incorporating lags into a delay embedding, consider the reconstructed Lorenz system illustrated in figure 2.8(a), obtained by sub-sampling the 0.01-step Lorenz time series of figure 2.4(a) with a lag of $\tau = 10$. Comparing this plot with the un-lagged reconstruction in figure 2.7(a), the sub-sampling is seen to successfully combat the one-dimensional nature of the un-lagged reconstruction, in a manner almost identical to that of the 10-delay embedding of the 0.1-step time series of figure 2.4(b) but without the inevitable sparseness of the latter. Indeed, if we were to throw away (the appropriate) nine out of every ten points along this new trajectory we would, in fact, be left with almost precisely the trajectory shown in figure 2.4(b); that they are not identical (the difference being due to the latter time series having been generated with an integration step of 0.1, rather than 0.01, as discussed



(a)



(b)

Figure 2.8 Filtered embedding of the Lorenz and laser systems. Part (a) was obtained by sub-sampling the first component of the Lorenz map, integrated with a step size of 0.01, and (b) was obtained by projecting a ten-delay embedding of the laser systems onto its principal components. Plotted in stereo.

earlier), of course means that in this case the two trajectories correspond to completely different dynamical systems.

2.3.2 Singular subspaces

In order to improve on this arbitrary choice of projection, Broomhead and King [7] proposed choosing \mathcal{F} so as to reduce the impact on the delay reconstruction of stochastic noise present in the time series. This is achieved by calculating the SVD of the zero-mean trajectory matrix $\mathbf{X} - \mathbf{1}\bar{x}^T = \mathbf{U}\Sigma\mathbf{V}^T$, where $\mathbf{1} \in \mathbb{R}^N$ has elements identically equal to one and $\bar{x} \in \mathbb{R}^m$ is the mean of the distribution of $x_i \in \mathbb{R}^m$, for $i = 1, \dots, N$. The column vectors \mathbf{v}_k of \mathbf{V} are principal components of this distribution, and the singular values σ_k have the interpretation that σ_k^2 is N times the variance of the projection of \mathbf{X} onto \mathbf{v}_k . This decomposition allows us to define a linear ‘singular’ subspace of \mathbb{R}^m by the n -element subset of principal components whose associated singular values σ_k are greater than some predetermined or empirically calculated level; we call this level the ‘noise floor’. We then make the definition $\mathbf{F} = (\mathbf{v}_1, \dots, \mathbf{v}_n)^T$, setting $\mathbf{c}_m = \bar{x}$ and $\mathbf{c}_n = \mathbf{0}$ in equation (2.10), to define a projection \mathcal{F} whose null-space is the linear subspace of \mathbb{R}^m spanned by those \mathbf{v}_k in whose directions the variance lies below this floor. (As an implementational note, for large N the calculation of \bar{x} can be avoided—neglecting end effects—by arranging for the time series $\{v_1, \dots, v_{N+m-1}\}$ itself to have zero mean.) In this case, however, \mathcal{F} does not, in general, preserve the shift property, as each of the n FIR filters of equation (2.12) has as coefficients the m elements of the singular vector \mathbf{v}_j , so we cannot write \mathcal{G} in the form of a delay map.

In figure 2.8(b) we illustrate a projection of this form, based on the laser system whose time series is shown in figure 2.5(b). In this figure we project the 10-delay reconstruction, whose first three dimensions are shown in figure 2.7(c), into a singular subspace spanned by the first three singular vectors obtained from an SVD of the trajectory in \mathbb{R}^{10} . On comparing the two figures we find, not surprisingly, that the principal component projection does a much better job of embedding the attractor in question than does the arbitrary projection shown in the previous figure, although in practice we are, of course, free to make use of all ten dimensions of the latter. Nevertheless, it seems clear that by using a singular basis we should be able to produce an embedded attractor with substantially fewer delays than would otherwise be the case. In chapter 5 we will analyse this claim in some detail for both the Lorenz and laser systems.

2.3.3 Time series filtering

The third, and final, class of linear transformations which we consider in this chapter is rather more specific: up till now, with the trivial exception of the sampling matrix of section 2.3.1, we have been dealing with filtering operations \mathcal{F} which generically do not impose a delay structure on their images $\mathcal{M}_n = \mathcal{F}\mathcal{M}_m$ because each row of \mathbf{F} contains the coefficients $a_{k-1}^{(j)}$ of a distinct FIR filter. In signal

processing applications, however, we are often interested in the effect of a *particular* FIR filter on the time series under investigation—for instance we might expect some form of linear filtering to be unavoidably performed in the process of measurement. In other words, we wish to work with a single time series $\{u_i\}$ which we obtain from $\{v_i\}$ by applying the FIR filter

$$u_i = \sum_{k=0}^{c-1} a_k v_{i-k} \quad (2.17)$$

with c coefficients $a_0, \dots, a_{c-1} \in \mathbb{R}$.

It is not immediately obvious that a delay reconstruction of the form (2.4), using the output $\{u_i\}$ of such a filter, will generically be diffeomorphic to a delay embedding of the system on which the unfiltered time series $\{v_i\}$ was measured, but Broomhead, Huke and Muldoon have shown [3] that this is generically the case: if $\Phi_{v,m}: \mathcal{M} \rightarrow \mathbb{R}^m$ is a delay embedding of \mathcal{M} then we can also embed \mathcal{M} with the delay map $\Phi_{u,n}: \mathcal{M} \rightarrow \mathbb{R}^n$, provided that $n > 2d$ as usual, where $u: \mathcal{M} \rightarrow \mathbb{R}$ is the measurement function ‘induced’ by the FIR filter. It is easy to show, by analogy with equation (2.12), that $\mathcal{M}_m = \Phi_{v,m}\mathcal{M}$ and $\mathcal{M}_n = \Phi_{u,n}\mathcal{M}$ are related by the linear transformation $\mathcal{F}: \mathbb{R}^m \rightarrow \mathbb{R}^n$ defined, with c_m and c_n set to zero, by the $n = m - c + 1$ by m matrix

$$\mathbf{F} = \begin{pmatrix} a_0 & a_1 & a_2 & \cdots & a_{c-1} & 0 & 0 & \cdots & 0 \\ 0 & a_0 & a_1 & \cdots & a_{c-2} & a_{c-1} & 0 & \cdots & 0 \\ 0 & 0 & a_0 & \cdots & a_{c-3} & a_{c-2} & a_{c-1} & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \cdots & a_0 & a_1 & a_2 & \cdots & a_{c-1} \end{pmatrix} \quad (2.18)$$

such that $\Phi_{u,n} \equiv \mathcal{G} = \mathcal{F} \circ \Phi_{v,m}$, the banding in \mathbf{F} imposing the necessary delay structure in \mathbb{R}^n . We will make use of this form of filtered embedding in sections 5.3 and 5.4 to come.

It is perhaps worth mentioning that the FIR filter (2.17) can be thought of as a special case of the more general infinite impulse response (IIR) filter,

$$u_i = \sum_{k=0}^{c-1} a_k v_{i-k} + \sum_{k=0}^{d-1} b_k u_{i-k} \quad (2.19)$$

which incorporates delays of both unfiltered and filtered time series. An IIR filter therefore constitutes a linear dynamical system in its own right; Badii and Politi have shown [1] that the system obtained from a delay reconstruction using a filtered time series of this form is not, in general, a diffeomorphism of the underlying system, and so we do not consider such filters here.

Chapter 3

Approximation

We now come to the central problem of this thesis: given two compact subsets $\mathcal{M} \subset \mathbb{R}^m$ and $\mathcal{N} \subset \mathbb{R}^n$ and a map $f: \mathcal{M} \rightarrow \mathcal{N}$ through which they are related, is f a diffeomorphism? In chapters 4 and 5 we will investigate this question by considering empirical models of f and f^{-1} . The models which we adopt are based upon the radial basis function (RBF) map $\hat{f}: \mathbb{R}^m \rightarrow \mathbb{R}^n$, described in Broomhead and Lowe [6] and Powell [33], and are derived by fitting data from a ‘training’ set of N vector pairs $\{(x_i, y_i)\}$ sampled from f so that $y_i = f(x_i)$, for $i = 1, \dots, N$. Although \hat{f} is defined on \mathbb{R}^m our task is to make \hat{f} as ‘good’ a fit to f on \mathcal{M} as possible, in a sense which will shortly be defined.

The ‘classical’ RBF map is a linear combination of p nonlinear basis functions $\varphi_j: \mathbb{R}^m \rightarrow \mathbb{R}$. These basis functions are defined with respect to p distinct points $c_j \in \mathbb{R}^m$, known as ‘centers’, by $\varphi_j(x) = \phi(\|x - c_j\|)$, where $\phi: \mathbb{R}^+ \rightarrow \mathbb{R}$ is a nonlinear function and $\|\cdot\|$ denotes the vector 2-norm. Treated as components of a map $\varphi: \mathbb{R}^m \rightarrow \mathbb{R}^p$, together with the linear transformation $\mathcal{W}: \mathbb{R}^p \rightarrow \mathbb{R}^n$, the φ_j form the RBF map

$$\hat{f} = \mathcal{W} \circ \varphi \tag{3.1}$$

Powell also includes polynomial basis functions, for instance a zeroth order term $\varphi_{p+1}(x) = 1$ and first order terms $\varphi_{p+j}(x) = x_{j-1}$, for $j = 1, \dots, m$, and so on, in his models. Although we do not explicitly incorporate such terms into φ we do implicitly include the zeroth order term by writing $\mathcal{W}(\varphi) = \bar{y} + \mathcal{W}^T(\varphi - \bar{\varphi})$, where \mathcal{W} is a p by n matrix which we refer to as the ‘weight matrix’ or just the ‘weights’ and \bar{y} and $\bar{\varphi}$, the latter defined by

$$\bar{\varphi} = \frac{1}{N} \sum_{i=1}^N \varphi(x_i) \tag{3.2}$$

denote a mean over the N element training set (from now on, the precise definition of a training set mean, denoted by an overbar, will be understood from its context).

The question of how good a fit to f is achieved on \mathcal{M} by the RBF approximation \hat{f} is usually answered in terms of the error function $\epsilon: \mathbb{R}^m \rightarrow \mathbb{R}^n$, defined by $\epsilon = \hat{f} - f$. This function enables us to define a normalised mean squared error $\epsilon \in \mathbb{R}$ by

$$\epsilon^2 = \sigma^{-2} \sum_{i=1}^N \|\epsilon(x_i)\|^2 = \sigma^{-2} \sum_{i=1}^N \|\hat{y}_i - y_i\|^2 \quad (3.3)$$

where $\hat{y} = \hat{f}(x)$ is the RBF estimate of y . The normalising scalar σ is defined (for purely notational convenience in the discussion to come) by

$$\sigma^2 = \sum_{i=1}^N \|y_i - \bar{y}\|^2 \quad (3.4)$$

so that σ^2 is N times the variance of the training set in \mathbb{R}^n (the precise definition of a given normaliser will also be understood from its context in future usage.) Normalisation of this form ensures that $\epsilon = 1$ when \hat{f} fits the mean in \mathbb{R}^n —that is, $\hat{f}(x_i) = \bar{y}$ for all i —a sure sign that f is completely unamenable to RBF approximation. The error is normalised in this fashion, rather than by incorporating a diagonal matrix of individual variances into the norm, to avoid scaling each dimension independently: if the y_i are delay vectors then their associated noise components are necessarily isotropic.

In order to approximate a given f we usually fix φ by choosing p centers from the training set and use a fixed nonlinearity ϕ . This restricts the optimisation of \hat{f} to the weight matrix \mathbf{W} , and so avoids the need for time-consuming nonlinear optimisation techniques. A common ad hoc approach is to choose the centers from the training set at random or, in the case of a delay embedded dynamical system, at regular intervals in the time variable. However, for reasons to be explained in section 3.2.1 we define an alternative selection method which consists of choosing the first center c_1 (which we will refer to as the ‘seed’) from the training set, either at random or following some arbitrary criterion (such as maximising $\|x_i\|$), then choosing as the r -th center the training set element x_i which maximises the expression $\min_{1 \leq j \leq r} \|x_i - c_j\|$. We therefore call this method ‘repulsive selection’. This approach is similar to one proposed by Smith [37], which consists of choosing centers uniformly on \mathcal{M} , subject to the constraint that no two centers should be closer together than a given distance; a suitable value for this inter-center distance is iteratively determined by decreasing its value from an initial maximum, until a sufficient number of centers have been so obtained. We will also investigate the nonlinear ‘forward selection’ method in section 3.2.2, for use in circumstances where we are willing to wait a long time for a particularly good fit.

3.1 Characteristics of the RBF map

In an appendix to Powell [33], Brown shows that \hat{f} is capable of universal approximation, under the uniform norm, on the space of continuous real-valued maps on compact subsets of \mathbb{R}^m (such as the manifolds we consider in this thesis). That is, $\sup_{\mathbf{x} \in \mathcal{M}} \|\hat{f}(\mathbf{x}) - \mathbf{f}(\mathbf{x})\|$ can be made arbitrarily small on \mathcal{M} by taking p large enough, provided that ϕ obeys certain constraints. Specifically, either $\phi(r) = e^{-r^2}$ or $\phi(r) = r$ must hold, or $D\phi(r^{1/2})$ must exist and be strictly completely monotonic. Another way of stating this result is that given enough basis functions the relationship between $\varphi\mathcal{M}$ and \mathcal{N} can be made as close to linear as required. Although the existence of an argument such as this is clearly desirable, in practice—working with finite p —we choose the cubic $\phi(r) = r^3$, which has the considerable advantage of being parameterless and yet is of comparable performance on experimental data sets to suitably tuned gaussians or other functions admissible by Brown.

Brown’s result, and empirical observations, lead us to expect to be able to achieve a reasonably good approximation \hat{f} to \mathbf{f} , in terms of a small ϵ , with a finite number of centers. We might, therefore, naively assume that if \mathbf{f} is a diffeomorphism then \hat{f} will also be a diffeomorphism, and conversely, if \mathbf{f} is *not* a diffeomorphism then neither is \hat{f} . If so, we could clearly establish strong evidence for whether or not \mathbf{f} is a diffeomorphism by analysing the injectivity and immersivity of \hat{f} . This approach, however, turns out to be a little too naive. In appendix A we show that, under certain mild conditions on the choice of ϕ and the positions of the centers, φ is an embedding of compact subsets of \mathbb{R}^m , provided that $p > m$. (In fact, it can be shown that with a monotonic increasing basis function, such as a cubic, φ embeds \mathbb{R}^m itself, although this is clearly not the case with a decreasing function, such as a gaussian). From section 2.3 we know that a d -dimensional submanifold $\varphi\mathcal{M} \subset \mathbb{R}^p$ is embedded by a generic choice of \mathcal{W} , provided that $n > 2d$. It follows, by composition, that \hat{f} is generically an embedding of compact \mathcal{M} if $p > m$ and $n > 2d$, where d is the dimension of \mathcal{M} , if a manifold, or of the lowest-dimensional manifold containing \mathcal{M} , if not. (In fact, a finite set of points sampled from \mathcal{M} is, by its very nature, a compact subset of \mathbb{R}^m , irrespective of whether or not \mathcal{M} is a manifold.) In other words, $\hat{f}: \mathcal{M} \rightarrow \hat{f}\mathcal{M}$ is generically a diffeomorphism.

As we can now expect to find, under the appropriate conditions, a diffeomorphism \hat{f} arbitrarily close to \mathbf{f} , we must modify the naive test, described above, to reflect this situation. Specifically, we must now attempt to determine whether or not \mathbf{f} is a diffeomorphism by measuring how close, in a sense to be discussed, \hat{f} is to a map which fails to be a diffeomorphism. To this end, we have investigated two distinct varieties of RBF model, each with its own strategy for the optimisation of \mathcal{W} , minimising the least squares (LS) and total least squares (TLS) errors. These two methods are described in sections 3.2 and 3.3 below.

It is important to note that—practically speaking—a particular error will often only be significant, as an indicator of the existence of a (differentiable) map, when compared with an error calculated for some other map, whose injectivity and immersivity may already have been established by some other means. For the experiments to be described later in this thesis, we will usually consider a family of maps $\{\mathbf{f}_\mu\}$,

where μ is some experimental parameter, and examine the variation of LS and TLS errors as a function of μ . Within a given family of maps the mean squares and/or distributions of these errors will then be used to establish those values of μ for which f_μ is a diffeomorphism, but between different families there may be no realistic comparison.

3.2 The method of least squares

In order to investigate this problem further it is useful to consider the manner in which f itself may fail to be a diffeomorphism by examining its inverse f^{-1} and the conditions under which it can exist. The LS solution is to attempt to model f and f^{-1} separately with RBF maps $\widehat{f}: \mathbb{R}^m \rightarrow \mathbb{R}^n$ and $\widehat{f}^{-1}: \mathbb{R}^n \rightarrow \mathbb{R}^m$, respectively. (Note that even when f has no inverse we can still obtain such an approximation to the relationship $y \mapsto x$, optimised over the training set, which we call \widehat{f}^{-1} for convenience.) To keep our notation consistent we label \widehat{f} and \widehat{f}^{-1} with their respective domains by rewriting (3.1) as $\widehat{f} = \mathcal{W}_M \circ \varphi_M$, where $\varphi_M: \mathbb{R}^m \rightarrow \mathbb{R}^p$ and $\mathcal{W}_M: \mathbb{R}^p \rightarrow \mathbb{R}^n$, and similarly $\widehat{f}^{-1} = \mathcal{W}_N \circ \varphi_N$, where $\varphi_N: \mathbb{R}^n \rightarrow \mathbb{R}^q$ is composed of q basis functions (we will usually just take $q = p$) and $\mathcal{W}_N: \mathbb{R}^q \rightarrow \mathbb{R}^m$. The linear maps are then implemented by $\mathcal{W}_M(\varphi) = \bar{y} + \mathbf{W}_M^T(\varphi - \overline{\varphi_M})$ and $\mathcal{W}_N(\varphi) = \bar{x} + \mathbf{W}_N^T(\varphi - \overline{\varphi_N})$. With error functions $\epsilon_M: \mathbb{R}^m \rightarrow \mathbb{R}^n$ and $\epsilon_N: \mathbb{R}^n \rightarrow \mathbb{R}^m$ defined by $\epsilon_M = \widehat{f} - f$ and $\epsilon_N = \widehat{f}^{-1} - f^{-1}$ the corresponding normalised LS errors ϵ_M and ϵ_N are

$$\epsilon_M^2 = \sigma_N^{-2} \sum_{i=1}^N \|\widehat{y}_i - y_i\|^2, \quad \epsilon_N^2 = \sigma_M^{-2} \sum_{i=1}^N \|\widehat{x}_i - x_i\|^2 \quad (3.5)$$

where $\widehat{x} = \widehat{f}^{-1}(y)$ is the RBF estimate of x and the normalisers σ_M and σ_N are calculated on $\mathcal{M} \subset \mathbb{R}^m$ and $\mathcal{N} \subset \mathbb{R}^n$ in the usual way. The procedure whereby these errors are minimised with respect to a particular training set is described in section 3.2.1 below.

So how do we use these models in practice? For the sake of argument we will phrase our discussion in terms of the relationship $x \mapsto y$, although we could just as well consider the map in the other direction. If f is a diffeomorphism then so, trivially, is f^{-1} . Provided that f is close, in the LS sense, to the space of functions spanned by the components $(\varphi_M)_j$ of φ_M , then we can clearly expect \widehat{f} to be a good fit to f , with uniformly small per-point errors $\|\epsilon_M(x)\|$ giving rise to a similarly small mean squared error ϵ_M . We can apply the same criteria to f^{-1} , in terms of the basis functions $(\varphi_N)_j$, yielding similarly small errors $\|\epsilon_N(y)\|$ and ϵ_N .

In considering how f may fail to be a diffeomorphism we will first examine the case in which f is injective but not immersive. This will occur when there is some subset $\mathcal{U} \subset \mathcal{M}$ for which Df fails to be injective, mapping tangent vectors $\nabla x \in \mathbb{R}^m$ to the zero vector in \mathbb{R}^n , with the result that Df^{-1} is not defined on $\mathcal{V} = f\mathcal{U}$. Now, although we may be able to find an \widehat{f} as close to f as we wish, we might reasonably expect to obtain a higher error ϵ_N in fitting \widehat{f}^{-1} than would otherwise be the case, since we will

be approximating with smooth $(\varphi_{\mathcal{N}})_j$, a map with infinite derivative on \mathcal{V} . This is, however, a marginal effect at best, particularly when we consider the fact that the training set is merely *sampled* from $\mathcal{M} \times \mathcal{N}$ and may not even incorporate those points at which the derivative breaks down.

Alternatively, f may fail to be injective on some subset $\mathcal{U} \subset \mathcal{M}$, whose image under f is a self-intersecting set $\mathcal{V} \subset \mathcal{N}$, in which case f^{-1} will not even be defined. For instance, f might be a delay map with insufficient delays to embed \mathcal{M} . It is instructive to view the loss of immersivity as a limiting case of the loss of injectivity—a transitory state through which f must pass on the way to becoming a many-to-one map. In attempting to model such a one-to-many relationship on \mathcal{V} the method of LS will result in an RBF map $\widehat{f^{-1}}$ where f^{-1} itself does not exist, taking each $\mathbf{y} \in \mathcal{V}$ to a weighted average of those $\mathbf{x} \in \mathcal{U}$ which lie in its pre-image under f . Given the finite nature of the training set, this will typically involve the incorporation of highly oscillatory basis functions, in an attempt to interpolate what is in reality a many-valued relationship. We may now be fairly confident of a large error $\epsilon_{\mathcal{N}}$, particularly if we use the techniques described in section 3.2.3, below, to constrain $\widehat{f^{-1}}$.

It is worth noting that there is another way in which f can fail to be an immersion, which we would not expect to be able to detect with a LS RBF fit, and that is where Df approaches different limits from different directions on some subset of \mathcal{M} . Given the discrete nature of the training set, such points are unlikely to have any effect on the LS error. Indeed, a breakdown of differentiability of this type would also defeat the error distribution analysis to be described in section 3.2.4. Happily, this situation is unlikely to occur when \mathcal{M} is the image of a delay map, given the assumptions made by Takens on the smoothness of both dynamical system and measurement function.

3.2.1 The least squares solution

In this section we will briefly describe the method by which we solve the LS fitting problem. For the time being we will drop the subscripts on φ , \mathcal{W} and \mathbf{W} for notational convenience. We wish to find the linear map $\mathcal{W}: \mathbb{R}^p \rightarrow \mathbb{R}^q$, realised by the p by q matrix \mathbf{W} , which takes each point $\mathbf{a} \in \mathbb{R}^p$ onto its image $\mathbf{b} \in \mathbb{R}^q$. This map might complete the RBF approximation $\widehat{f} = \mathcal{W}_{\mathcal{M}} \circ \varphi_{\mathcal{M}}$, in which case we set $\mathbf{a} = \varphi_{\mathcal{M}}(\mathbf{x}) - \overline{\varphi_{\mathcal{M}}}$, $\mathbf{b} = \mathbf{y} - \overline{\mathbf{y}}$ and $q = n$, or the inverse $\widehat{f^{-1}} = \mathcal{W}_{\mathcal{N}} \circ \varphi_{\mathcal{N}}$, in which case we set $\mathbf{a} = \varphi_{\mathcal{N}}(\mathbf{y}) - \overline{\varphi_{\mathcal{N}}}$, $\mathbf{b} = \mathbf{x} - \overline{\mathbf{x}}$ and $q = m$. Regardless of its interpretation, we optimise this map by minimising the LS error ϵ , defined over the training set by

$$\epsilon^2 = \sigma^{-2} \sum_{i=1}^N \|\widehat{\mathbf{b}}_i - \mathbf{b}_i\|^2 \quad (3.6)$$

where $\widehat{\mathbf{b}} = \mathbf{W}^T \mathbf{a}$ is the LS RBF estimate of \mathbf{b} and either $\sigma = \sigma_{\mathcal{N}}$, in the case of solving for \widehat{f} —in which case $\epsilon = \epsilon_{\mathcal{M}}$ —or, when solving for $\widehat{f^{-1}}$, $\sigma = \sigma_{\mathcal{M}}$ and $\epsilon = \epsilon_{\mathcal{N}}$.

We first recast this error in matrix form, defining an N by p matrix \mathbf{A} whose rows are the transposed \mathbf{a}_i and an N by q matrix \mathbf{B} whose rows are the transposed \mathbf{b}_i . Then the N by q matrix $\widehat{\mathbf{B}}$ of estimates $\widehat{\mathbf{b}}_i = \mathbf{W}^T \mathbf{a}_i$ is defined by $\widehat{\mathbf{B}} = \mathbf{A}\mathbf{W}$ and the LS error becomes

$$\epsilon^2 = \sigma^{-2} \|\widehat{\mathbf{B}} - \mathbf{B}\|_{\mathbb{F}}^2 \quad (3.7)$$

where $\|\cdot\|_{\mathbb{F}}$ is the Frobenius norm (the squared Frobenius norm $\|\mathbf{B}\|_{\mathbb{F}}^2$ of a matrix \mathbf{B} is equal to the trace of $\mathbf{B}^T \mathbf{B}$ or, equivalently, the sum of the squares of the elements of \mathbf{B}). If we write $\beta, \widehat{\beta} \in \mathbb{R}^N$ for the column vectors of \mathbf{B} and $\widehat{\mathbf{B}}$ respectively then we see that the minimum is achieved when each $\widehat{\beta}$ is the orthogonal projection of the corresponding β into the column space $\mathcal{R}(\mathbf{A}) \subset \mathbb{R}^N$. In other words, each component $\widehat{f}_j: \mathbb{R}^m \rightarrow \mathbb{R}$ of the multi-dimensional LS solution $\widehat{\mathbf{f}}$ is found independently of the rest (this is just a restatement of the linearity of $\widehat{\mathbf{f}}$ in terms of its basis functions). To write down this solution it is convenient to find an orthogonal basis for $\mathcal{R}(\mathbf{A})$, spanned by the column vectors $\alpha \in \mathbb{R}^N$. One useful decomposition for this purpose is the SVD, which we first defined in section 2.3. The SVD of \mathbf{A} is written in full as

$$\mathbf{A} = (\mathbf{U}; \overline{\mathbf{U}}) \begin{pmatrix} \Sigma \\ \mathbf{0} \end{pmatrix} \mathbf{V}^T \quad (3.8)$$

where the N by N matrix $(\mathbf{U}; \overline{\mathbf{U}})$ is partitioned into an N by p submatrix \mathbf{U} , whose column vectors $\mathbf{u}_k \in \mathbb{R}^N$ are singular vectors of the distribution of $\alpha \in \mathbb{R}^N$, and an N by $(N - p)$ matrix $\overline{\mathbf{U}}$, whose columns span the orthogonal complement of $\mathcal{R}(\mathbf{A})$. Written more concisely, as in equation (2.11) of the preceding chapter, this becomes $\mathbf{A} = \mathbf{U}\Sigma\mathbf{V}^T$. The \mathbf{u}_k are, as usual, ordered so that their corresponding singular values $\sigma_k \geq 0$, which appear as elements of the p by p diagonal matrix Σ , satisfy $\sigma_k \geq \sigma_{k-1}$. In this basis the orthogonal projection of β into $\mathcal{R}(\mathbf{A})$ is given by $\widehat{\beta}_j = \mathbf{U}\mathbf{U}^T \beta_j$ or, collectively, $\widehat{\mathbf{B}} = \mathbf{U}\mathbf{U}^T \mathbf{B}$, and the solution can be found with

$$\begin{aligned} \mathbf{A}\mathbf{W} &= \widehat{\mathbf{B}} \\ \Rightarrow \mathbf{U}\Sigma\mathbf{V}^T\mathbf{W} &= \mathbf{U}\mathbf{U}^T\mathbf{B} \\ \Rightarrow \mathbf{W} &= \mathbf{V}\Sigma^{-1}\mathbf{U}^T\mathbf{B} \end{aligned} \quad (3.9)$$

The matrix $\mathbf{V}\Sigma^{-1}\mathbf{U}^T$ is the pseudo-inverse \mathbf{A}^\dagger of \mathbf{A} in the special case that the rank of \mathbf{A} is p , that is, when Σ^{-1} exists. In this case \mathbf{A}^\dagger is more usually written as $(\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T$.

Another method of obtaining an orthogonal basis for $\mathcal{R}(\mathbf{A})$ is the QR decomposition [38], written as

$$\mathbf{A} = (\mathbf{Q}; \overline{\mathbf{Q}}) \begin{pmatrix} \mathbf{R} \\ \mathbf{0} \end{pmatrix} \quad (3.10)$$

where the N by N orthogonal matrix $(\mathbf{Q}; \overline{\mathbf{Q}})$, with N by p submatrix \mathbf{Q} and N by $(N - p)$ submatrix $\overline{\mathbf{Q}}$, is chosen so that the p by p matrix \mathbf{R} is upper-triangular. Written more concisely as $\mathbf{A} = \mathbf{Q}\mathbf{R}$, with the

columns of Q an orthogonal basis for $\mathcal{R}(A)$, this decomposition is unique if A is full rank, and can be used in exactly the same manner as we used the SVD to find a LS projection $\hat{B} = QQ^T B$ and solution

$$\begin{aligned} AW &= \hat{B} \\ \Rightarrow QRW &= QQ^T B \\ \Rightarrow W &= R^{-1}Q^T B \end{aligned} \quad (3.11)$$

Both of these decompositions necessarily lead to the same solution, assuming that Σ^{-1} and R^{-1} exist, but we will generally use the SVD approach in this thesis as its interpretation in terms of the singular spectrum of A provides a useful method for controlling ill-conditioning and over-fitting, to be described in section 3.2.3. However we note that, from the operation counts for QR and singular value decompositions given by Lawson and Hanson [23], in the case that $N > \frac{5}{3}p$, rather than find the SVD of the N by p matrix A directly it is computationally more efficient to first write A in the form of equation (3.10) and then find the SVD of the p by p matrix R , say $R = L\Sigma V^T$, so that we can make the straightforward identification

$$A = QL\Sigma V^T = U\Sigma V^T \quad (3.12)$$

and hence $U = QL$.

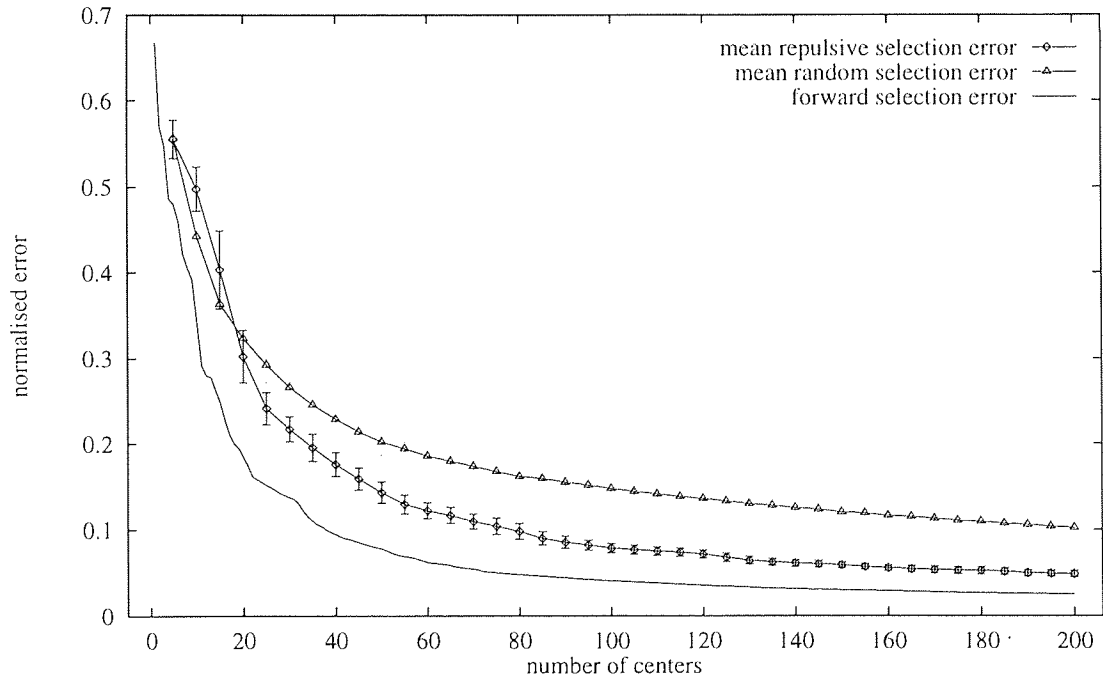


Figure 3.1 Comparing the repulsive and forward selection methods. Plotting the expected value $\langle \epsilon_p \rangle$, and standard deviation (denoted by error bars), of the fitting error ϵ_p versus p , for a function on a five component singular subspace of a ten-delay embedding of the laser system, trained using both randomly-seeded repulsive and purely random center selection methods, with sample sizes of 500 and 1000 centers sets, respectively, together with the error ϵ_p obtained from a single set of centers selected by forward selection. The error bars on the purely random selection error have been lightened for clarity.

To demonstrate the characteristics of the LS solution, we will briefly investigate the laser system of figure 2.8(b), in a five-dimensional singular subspace \mathcal{N} of a ten-delay embedding obtained from the time series in figure 2.5(b). The function $f: \mathcal{N} \rightarrow \mathbb{R}$ which we approximate in this experiment is the one-step generator for the time series; it will be described in more detail in chapter 5, where the choice of embedding dimension will also be explained. In figure 3.1 we plot the fitting error ϵ_p as a function of $p = 5, 10, \dots, 200$ centers chosen using the repulsive selection method from a training set of $N = 2000$ vectors in \mathbb{R}^5 . In order to eliminate—as much as possible—the dependence of ϵ_p on the particular set of centers chosen in this manner, the value plotted is actually the mean error $\langle \epsilon_p \rangle$, obtained by averaging ϵ_p , for each p , over 500 sets of centers obtained by selecting the seed for the repulsive selection algorithm at random (without replacement) from the training set. (We will make use of this procedure—averaging over the model—in most of the experiments to come, so we adopt the notation $\langle \cdot \rangle$ of statistical expectation in order to distinguish this form of expectation from the average incorporated into the calculation of ϵ_p itself.) For the sake of comparison, we also plot the mean error obtained by fitting the same function with 1000 sets of purely random centers selected (also without replacement) from the training set. (It is worth noting that the number of distinct sets of repulsive centers obtainable through random selection of initial center is necessarily no greater than N .) The error bars in both cases represent one standard deviation in each direction. The additional curve, corresponding to the method of forward selection, will be discussed in the following section.

It is clear from this figure that for $p \gtrsim 20$ the method of repulsive selection achieves, on average, an error value several standard deviations below that obtained by purely random selection. An appealing explanation of this effect is that the repulsive method avoids choosing any two centers closer together than necessary in \mathbb{R}^m , which is not generally the case in random selection: in the limit that one such inter-center distance goes to zero there will be two corresponding identical column vectors $\alpha \in \mathbb{R}^N$, resulting in an ill-conditioned data matrix \mathbf{A} . Nevertheless, we see also that the random selection method itself achieves a respectably small error on average, provided that p is sufficiently large.

3.2.2 Forward selection

Having shown how the optimisation of an RBF map may be restricted to a linear LS problem by first fixing the centers in an ad hoc manner, as discussed above—and indeed that the choice of a particular set of centers is not necessarily of crucial importance, provided that an adequate number are chosen—we will now describe an exception to this approach. In [9] Chen, Cowan and Grant describe a method which they call ‘orthogonal least squares’, in which centers are selected from the training set in an incremental fashion, using a recursive form of the QR decomposition.

We will assume that we have already identified p distinct training vectors $\mathbf{x}_i \in \mathbb{R}^m$ as centers, each of which gives rise to a vector $\alpha_j \in \mathbb{R}^N$ in the usual way, where $1 \leq j \leq p$, forming the columns of the

data matrix A_p . We write the QR decomposition of A_p , by analogy with (3.10), as

$$A_p = (Q_p; \overline{Q}_p) \begin{pmatrix} R_p \\ 0 \end{pmatrix} \quad (3.13)$$

leading to a LS solution $W_p = Q_p R_p^{-1} B$ and fitting error $\epsilon_p = \sigma^{-2} \|A_p W_p - B\|_F^2$, by analogy with equations (3.9) and (3.7), respectively. It is easily shown [38] that $(Q_p; \overline{Q}_p)$ can be written as a product of p 'elementary reflectors' P_p , so that $(Q_p; \overline{Q}_p) = P_1 P_2 \dots P_p$. Then, for a new center, and hence column vector $\alpha_{p+1} \in \mathbb{R}^N$, we can find Q_{p+1} from $(Q_{p+1}; \overline{Q}_{p+1}) = (Q_p; \overline{Q}_p) P_{p+1}$ and R_{p+1} from

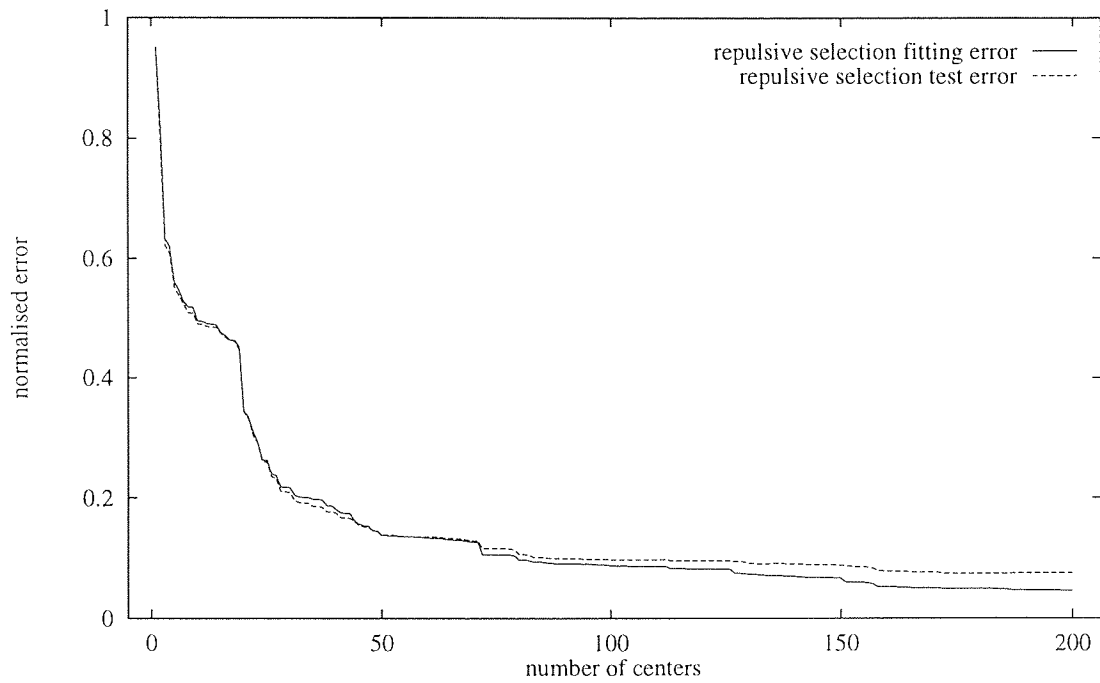
$$\begin{aligned} \begin{pmatrix} R_{p+1} \\ 0 \end{pmatrix} &= (Q_{p+1}; \overline{Q}_{p+1})^T A_{p+1} \\ &= P_{p+1}^T (Q_p; \overline{Q}_p) (A_p; \alpha_{p+1}) \\ &= P_{p+1}^T \left(\begin{pmatrix} R_p \\ 0 \end{pmatrix}; (Q_p; \overline{Q}_p) \alpha_{p+1} \right) \end{aligned} \quad (3.14)$$

Chen's method consists of choosing, at each iteration, that element of the training set the inclusion of whose corresponding column vector α_{p+1} minimises ϵ_{p+1} over the set of all remaining candidates. Of course, this technique will not generally choose the best possible set of p centers, in the sense of a minimum error ϵ_p , since each new center candidate is only considered in the context of the centers already chosen. Moreover, it should be noted that this combinatorial technique has been found to be substantially more computationally intensive than the ad hoc method and direct decomposition techniques combined on the experiments described in this thesis.

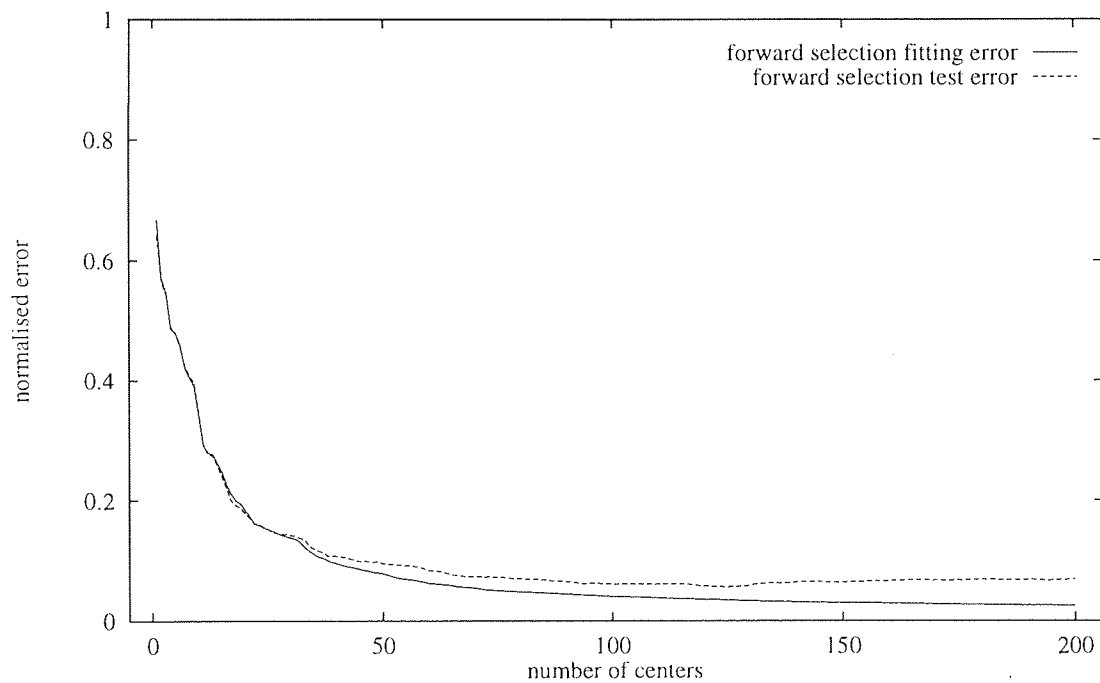
The forward selection method is also illustrated in figure 3.1, in which we plot the fitting error ϵ_p versus $p = 1, 2, \dots, 200$ for the laser prediction experiment. (The curve obtained in this case is necessarily unique, so we plot it with a continuous line for the sake of clarity.) We see in this figure that the forward selection error is substantially smaller than both the mean error obtained through both random selection and randomly-seeded repulsive selection, typically achieving an error some three or four standard deviations below either error. In particular—and in contrast with repulsive selection—there is no value of p for which the model trained by forward selection produces a higher error than the mean error arising from the random selection method.

3.2.3 Control of over-fitting

Although we might, given a suitable choice of ϕ , invoke Brown and Powell [33] to argue that we can make ϵ_p arbitrarily small by taking p large enough, it is important not to lose sight of the fact that \hat{f} has been chosen by minimising ϵ_p on a single training set of vector pairs $\{(x_i, y_i)\}$, sampled from a joint distribution in \mathbb{R}^{p+q} whose elements satisfy $y_i = f(x_i)$. When f is a map between delay embedded manifolds, this sampling notionally takes place in both space and time: the sampling in space is a result



(a)



(b)

Figure 3.2 Over-fitting on the laser system prediction problem, for two individual sets of centers, selected by both repulsive and forward selection, respectively. Plotting ϵ_p versus p , we see that at $p = 200$, in part (a), using repulsive centers, the test error has clearly begun to saturate, as the fitting error continues its monotonic decrease, while in part (b), using forward selection, the test error not only saturates, but actually begins to rise again.

of the quantisation, and any other sources of measurement error arising from the particular measurement function (for the numerically simulated time series in this thesis, these contributions can be assumed to be negligible); the sampling in time corresponds to the sampling interval (or integration step) with which the time series in question was obtained. As a consequence of this sampling process, a particular training set is necessarily consistent, not only with f itself, but also with any number of other functions whose values coincide with those of f on the x_i in question. It is therefore quite likely that a model which is sufficiently powerful (has sufficiently many degrees of freedom) to achieve a small error on the training set may nevertheless give rise to a disproportionately large error on another, unseen ‘test’ set of vector pairs sampled from the same distribution; we call this phenomenon ‘over-fitting’.

To illustrate this point, in figure 3.2 we plot the errors calculated on both the training set and a non-overlapping test set of 2000 vectors for the laser time series prediction problem discussed above, using centers selected by both repulsive and forward selection (the errors in the former case now correspond to a single set of centers whose seed is arbitrarily chosen to maximise $\|x_i\|$ over the training set). Over-fitting appears on these plots as a systematic divergence of fitting and test errors (due to the tendency of the test error to saturate) with increasing p . In part (a), which shows the errors arising from the repulsive center selection method, although this saturation has clearly begun at $p = 200$, both fitting and test errors are still apparently decreasing. In contrast, in part (b), which corresponds to the forward selection method, the test error appears to saturate quite early and even—if anything—to rise slightly as p approaches its limit. The fitting error in both cases decreases monotonically as expected.

We notice, in figure 3.2, that although the fitting error falls monotonically with increasing p , as it must, it does not fall particularly smoothly. Instead, it appears to approach a series of plateaus, within which the addition of each subsequent center makes relatively little difference to the error, but between which the error can drop significantly. However, it is important to remember that the particularly large decrease in fitting error resulting from the inclusion of a given center is a function of the entire basis set thus far selected, and not merely of that center alone. Certain linear combinations of centers, in other words, are particularly effective in ‘explaining’ the variance in B . The effect is most evident for the repulsive selection method illustrated in figure 3.2(a), with ϵ_p exhibiting a substantial decrease over a short interval near $p \approx 20$, but it is observable even in the forward selection error of 3.2(b).

In a process analogous to that described in section 2.3.2, this observation is commonly exploited in terms of the SVD of A , by eliminating those singular vectors v_k of A whose corresponding singular values σ_k lie below the noise floor defined by some $\sigma_{r \leq p}$. This solution is implemented by defining $\mathcal{K} = \{1, 2, \dots, r\}$, so as to make a partition of U by forming an N by r submatrix $U_{\mathcal{K}}$ whose columns u_k , indexed by $k \in \mathcal{K}$, span an r -dimensional subspace of $\mathcal{R}(A)$. Corresponding to $U_{\mathcal{K}}$ there is a p by r submatrix $V_{\mathcal{K}}$ of V , with columns v_k , and an r by r submatrix $\Sigma_{\mathcal{K}}$ of Σ , with diagonal elements σ_k . We are thus able to define the rank- r truncated pseudo-inverse $A_{\mathcal{K}}^{\dagger} = V_{\mathcal{K}} \Sigma_{\mathcal{K}}^{-1} U_{\mathcal{K}}^T$, which we use to form a new weight matrix $W_{\mathcal{K}} = A_{\mathcal{K}}^{\dagger} B$. Clearly, if $r = p$ then $W_{\mathcal{K}} = W$. If A is rank-deficient, however—or

merely ill-conditioned—then by restricting the rank of $\mathbf{A}_{\mathcal{K}}^{\dagger}$ to $r < p$ we obtain a truncated solution which (we assume) more concisely expresses the relationship between \mathbf{A} and \mathbf{B} .

Truncating the pseudo-inverse in the manner outlined above corresponds to constraining the LS problem to use the fullest available representation (in terms of variance) of the basis set $\{\alpha_j\}$, for a given rank r . But it makes no use whatsoever of the information contained in the columns of \mathbf{B} , the approximation of which is the sole purpose for which we constructed $\mathbf{W}_{\mathcal{K}}$ in the first place. To explore this situation in more detail, we can use the Frobenius norm (as defined in section 3.2.1) to rewrite the error ϵ , for a given p , in terms of the projections $(\mathbf{I} - \mathbf{U}\mathbf{U}^T)\boldsymbol{\beta}$ of the column vectors $\boldsymbol{\beta}$ into the orthogonal complement of $\mathcal{R}(\mathbf{A})$, where \mathbf{I} is the N by N identity matrix, giving

$$\epsilon^2 = \sigma^{-2} \|(\mathbf{I} - \mathbf{U}\mathbf{U}^T)\mathbf{B}\|_{\text{F}}^2 = \sigma^{-2} \|\mathbf{B}\|_{\text{F}}^2 - \sigma^{-2} \|\mathbf{U}^T \mathbf{B}\|_{\text{F}}^2 \quad (3.15)$$

and then define a truncated error $\epsilon_{\mathcal{K}}$ by

$$\epsilon_{\mathcal{K}}^2 = \sigma^{-2} \|\mathbf{B}\|_{\text{F}}^2 - \sigma^{-2} \|\mathbf{U}_{\mathcal{K}}^T \mathbf{B}\|_{\text{F}}^2 \quad (3.16)$$

We can now take advantage of the orthogonality of the \mathbf{u}_k to expand $\mathbf{W}_{\mathcal{K}}$ as

$$\mathbf{W}_{\mathcal{K}} = \sum_{k \in \mathcal{K}} \Delta \mathbf{W}_k \quad (3.17)$$

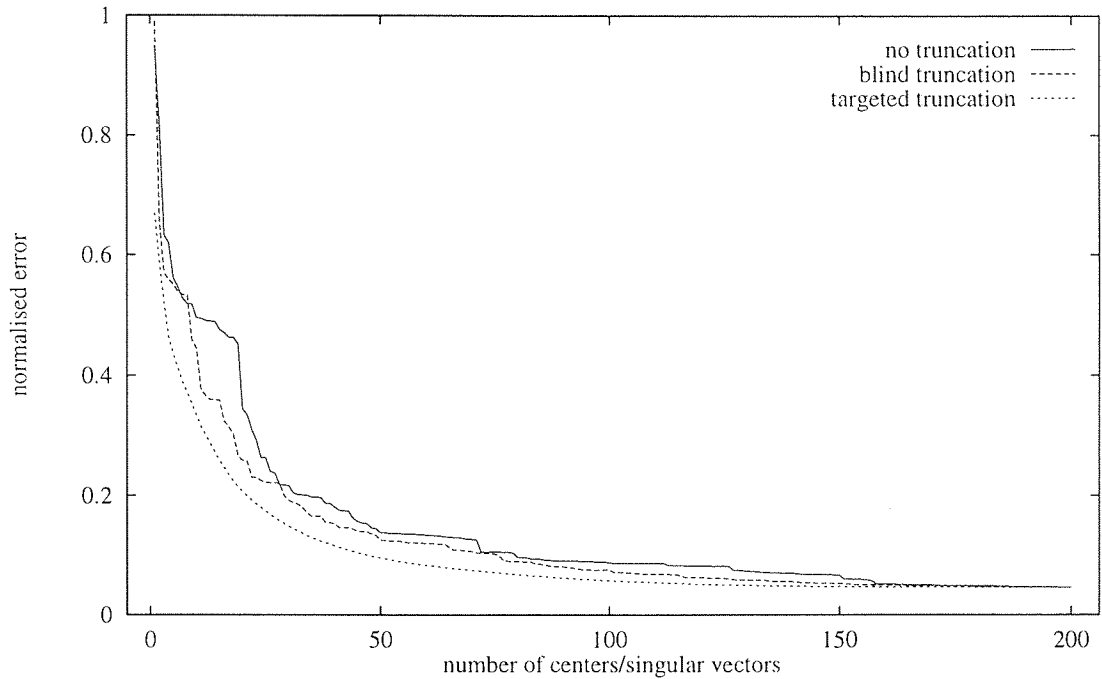
where $\Delta \mathbf{W}_k = \sigma_k^{-1} (\mathbf{u}_k^T \mathbf{B}) \mathbf{v}_k$, and similarly,

$$\epsilon_{\mathcal{K}}^2 = \epsilon_0^2 - \sum_{k \in \mathcal{K}} \Delta \epsilon_k^2 \quad (3.18)$$

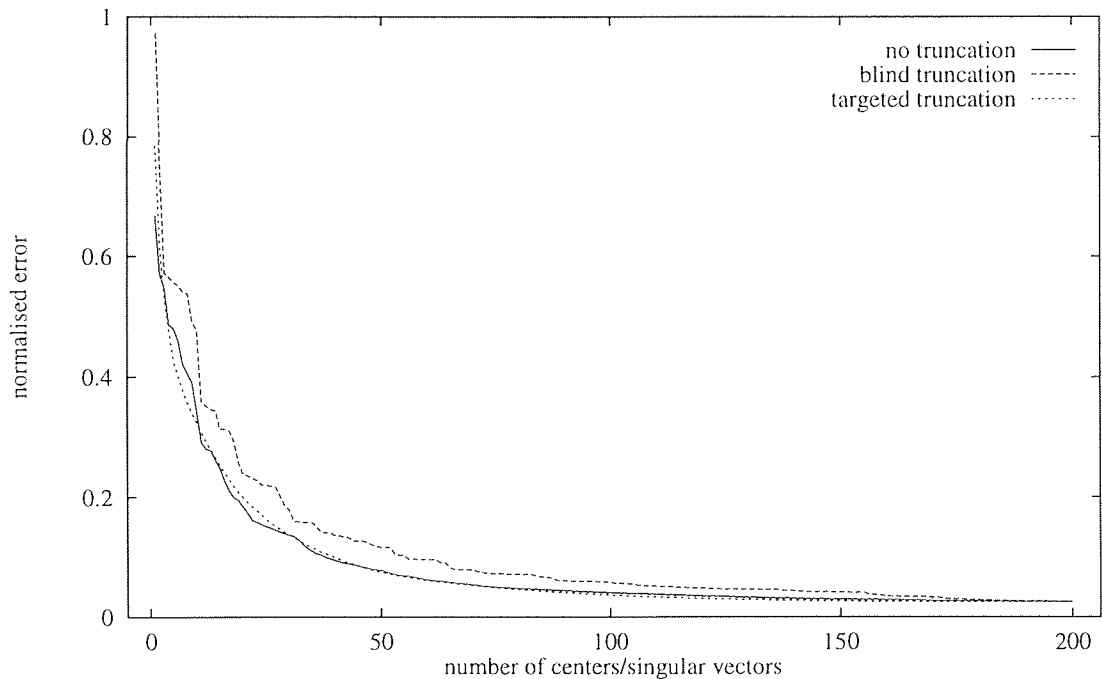
where $\epsilon_0 = \sigma^{-1} \|\mathbf{B}\|_{\text{F}}$ and $\Delta \epsilon_k = \sigma^{-1} \|\mathbf{u}_k^T \mathbf{B}\|$.

It is immediately clear that equation (3.18) provides us with an alternative—and perhaps more effective—truncation criterion: choose \mathcal{K} so as to exclude those \mathbf{u}_k which result in the smallest contributions $\Delta \epsilon_k$ to $\epsilon_{\mathcal{K}}$. By thus making use of the information contained in \mathbf{B} , we are now solving the restricted rank LS problem subject to the potentially more useful constraint that the representation embodied by \mathcal{K} results in the smallest possible error in approximating \mathbf{f} on the training set. For notational convenience, we will therefore call the former (σ_k) criterion ‘blind’ truncation, and the latter ($\Delta \epsilon_k$) criterion will be called ‘targeted’ truncation.

In figures 3.3(a) and (b) we plot the result of truncating the rank of laser system predictors trained with $p = 200$ centers chosen, respectively, by the repulsive and forward selection methods. In both figures we plot the errors $\epsilon_{\mathcal{K}}$, versus the rank $r = \text{card } \mathcal{K}$, for \mathcal{K} ordered according to both truncation criteria and $\text{card } \mathcal{K} = 1, \dots, p$. We also plot, on the same axes, the error ϵ_p , versus $p = 1, \dots, 200$, arising from the un-truncated model trained with p centers obtained by both center selection methods (naturally, in both



(a)



(b)

Figure 3.3 The effects of rank truncation on the laser time series prediction problem with centers selected by both repulsive and forward selection methods. Plotting the fitting errors $\epsilon_{\mathcal{K}}$ versus $r = \text{card } \mathcal{K}$, for both blind and targeted truncations of a $p = 200$ center model, superimposed on the error ϵ_p versus $p = 1, \dots, 200$, for comparison. (a) With repulsive centers, although for a given rank both criteria can be seen to substantially out-perform the error obtained without truncation for a corresponding number of centers, the advantage of using the latter criterion is clear; (b) with centers obtained by forward selection it is difficult to choose between a targeted truncation and the corresponding un-truncated error.

plots all three curves intersect at $p = 200$). In examining figure 3.3(a), corresponding to the repulsive selection method, we see that both truncated $p = 200$ models consistently achieve a smaller error $\epsilon_{\mathcal{K}}$, for a given rank $r = \text{card } \mathcal{K}$, than that of the un-truncated model corresponding to $p = r$. In the case of the blind truncation criterion, this result follows naturally from the ad hoc method by which the centers were chosen. By construction, the targeted truncation criterion provides, for a given rank r , a lower bound for all possible truncations of that rank in the SVD basis, and this can be seen in the figure. In contrast, in part (b) of figure 3.3 we see that blind truncation results in a error which is consistently *larger* than that of the corresponding un-truncated model, while the error for the targeted truncation criterion, although necessarily everywhere lower than the blind truncation error, is virtually indistinguishable from its un-truncated equivalent. In fact, the targeted truncation error can actually be seen to be higher than the un-truncated error for some values of r , in an effective demonstration of the power of the forward selection method.

Having demonstrated the usefulness—at least in the case of a model trained with repulsive centers—of the rank truncation method, it is also interesting to observe its effect on generalisation. To this end we compare, in figures 3.4 and 3.5, the fitting and test errors resulting from each of the four models considered above, figure 3.4 illustrating the repulsive model and figure 3.5 the forward selection case. In both figures, we plot the errors arising from the blind truncation criterion in part (b) and the targeted truncation criterion in part (c), duplicating in part (a) the un-truncated error curves of figure 3.2, for the purposes of comparison. Happily, in the case of all four truncated rank models, we see test error follow fitting error as closely as in the corresponding un-truncated cases.

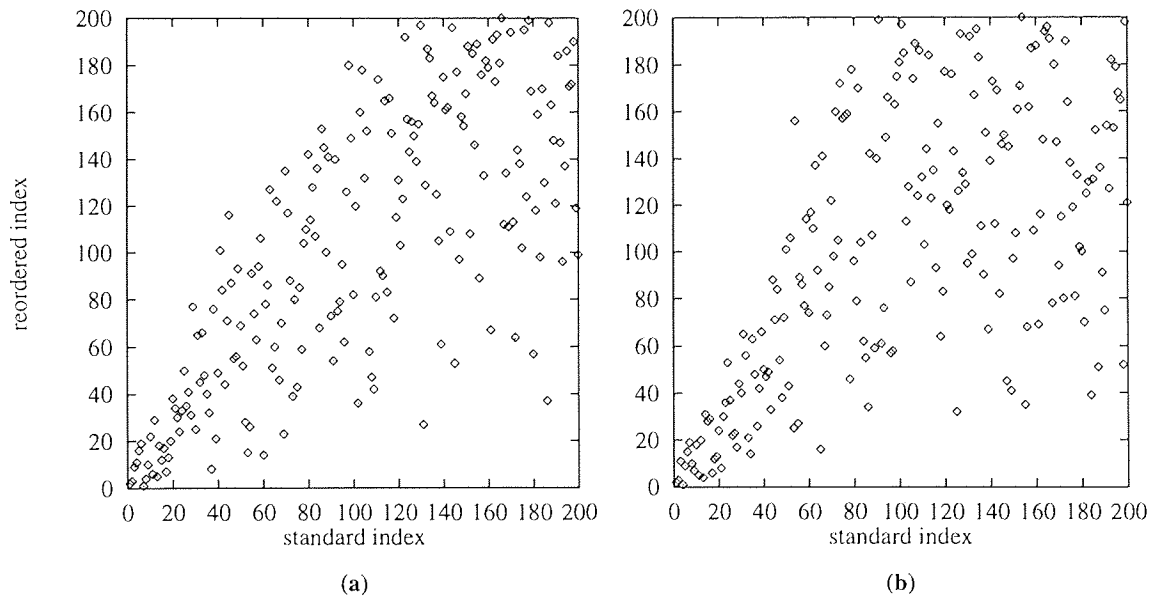
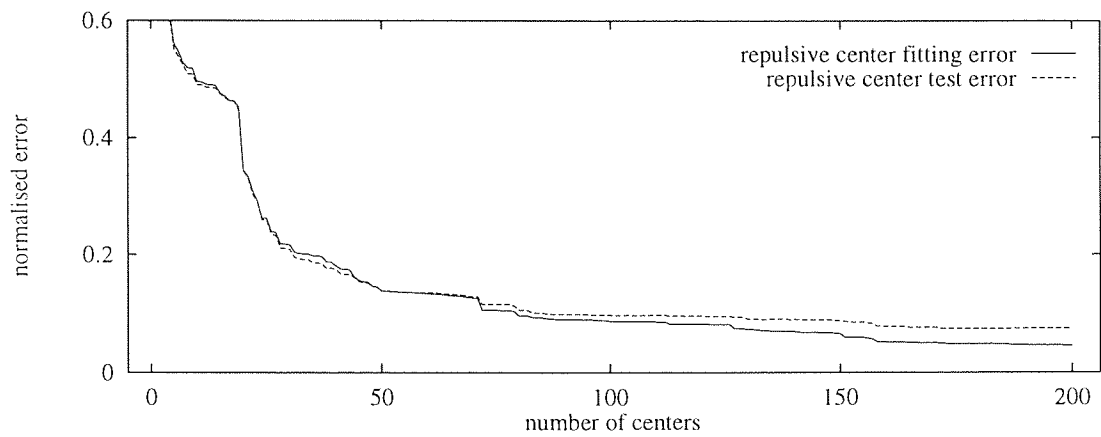
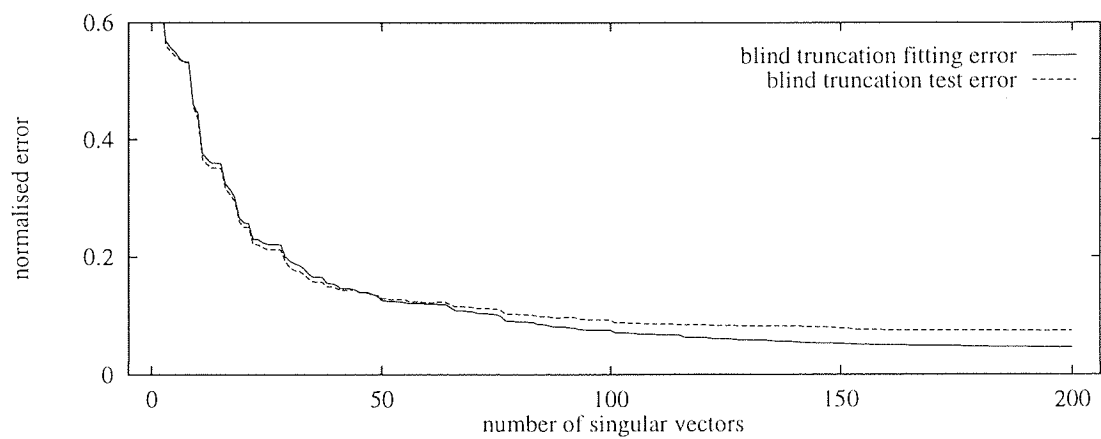


Figure 3.6 Scatter plot of the reordering map relating the blind and targeted rank selection criteria. Centers chosen by (a) repulsive selection and (b) forward selection.

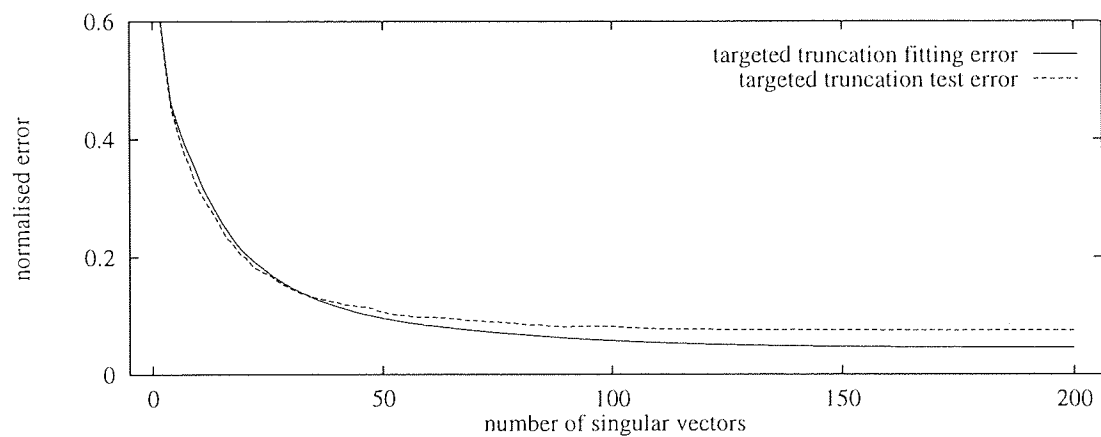
Another way in which to compare the two truncation criteria is to plot the one-to-one mapping between



(a)

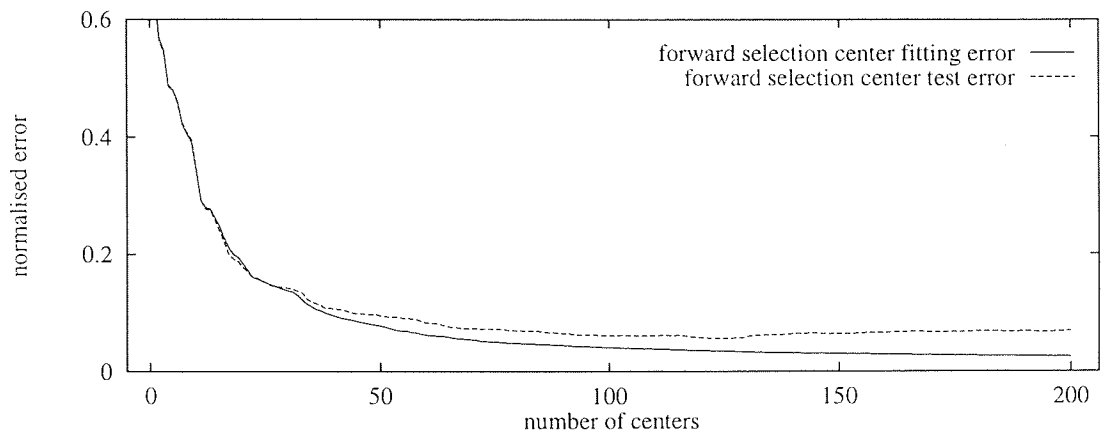


(b)

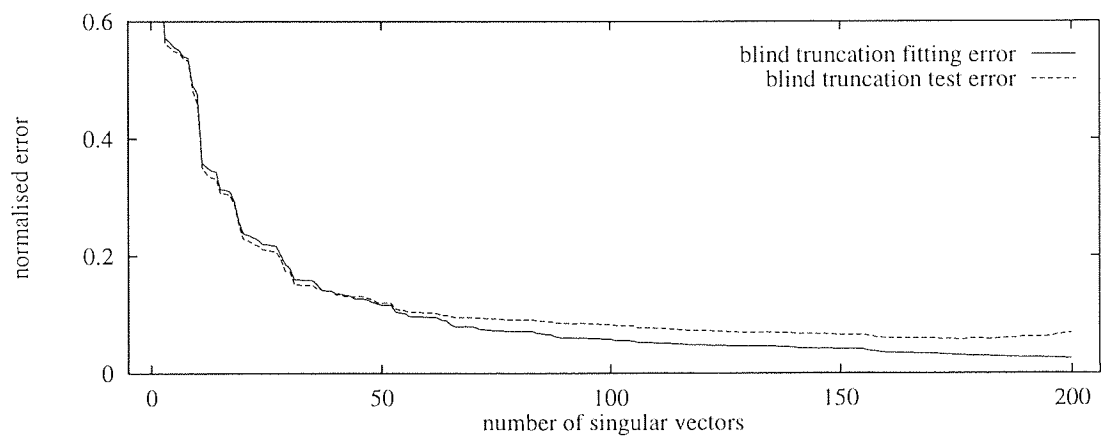


(c)

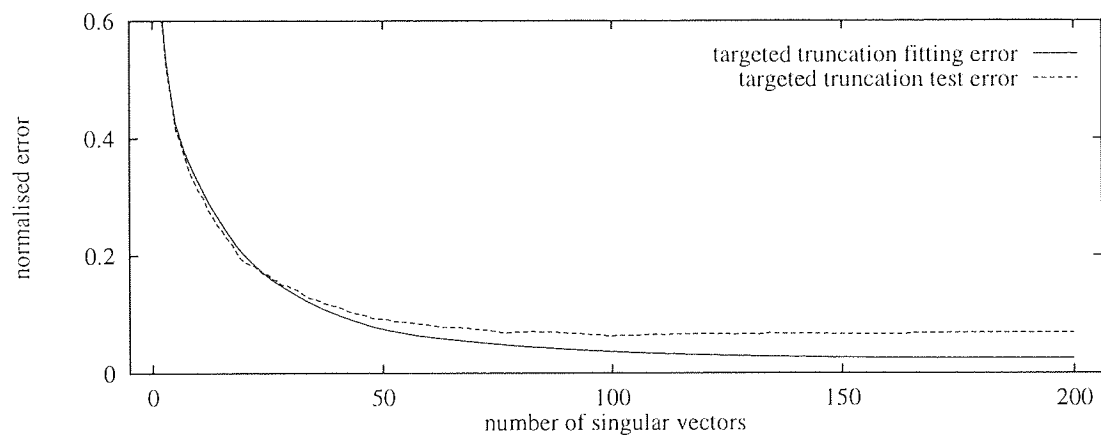
Figure 3.4 Improving generalisation with repulsive centers by rank truncation. Comparing the fitting and test errors $\epsilon_{\mathcal{K}}$, plotted versus $r = \text{card } \mathcal{K}$, for the laser system prediction problem, showing the corresponding un-truncated errors in part (a) as a baseline. In comparison, part (b) shows the effect on generalisation of the blind truncation criterion and part (c) shows the effect of targeted truncation.



(a)



(b)



(c)

Figure 3.5 Improving generalisation with forward selection by rank truncation. Comparing the fitting and test errors $\epsilon_{\mathcal{K}}$, plotted versus $r = \text{card } \mathcal{K}$, for the laser system prediction problem, showing the full-rank errors in part (a) as a baseline. In comparison, part (b) shows the effect on generalisation of the blind truncation criterion and part (c) shows the effect of targeted truncation.

the elements of the set \mathcal{K} chosen by blind truncation and those of the set chosen by targeted truncation. A scatter plot of this relationship is shown in figures 3.6(a) and (b), for repulsive and forward selection methods, respectively, labelling the x-axis by the index of the singular values, in the standard σ_k ordering, and the y-axis by the index of the same singular value after re-ordering by $\Delta\epsilon_k$. The clearly discernible structure in both figures, a ‘spreading out’ of the diagonal with increasing k , may best be thought of as a measure of the closeness of the ‘natural’ σ_k ordering to the reordering by $\Delta\epsilon_k$.

Despite the evident desirability of the rank-reduced models based on targeted truncation we do not make extensive use of this technique in later chapters. This is because we will usually be interested in analysing the variation of the error ϵ_μ or per-point errors $\|\epsilon_\mu(\mathbf{x})\|$ with some parameter μ , for a fixed RBF map architecture, and it is deemed inappropriate to allow the map itself to vary with this parameter.

3.2.4 Detecting non-invertible maps with least squares

Since the LS error (3.6) is an average, over the training or test set on which it is calculated, it is by definition quite insensitive to infrequently occurring per-point errors of large magnitude (as opposed to the supremum error, for instance). If f fails to be a diffeomorphism, through loss of either injectivity or immersivity, on a subset $\mathcal{V} \subset \mathcal{N}$ of relatively low weight in \mathcal{N} then the method of LS is entirely capable of suffering a relatively large $\|\epsilon_{\mathcal{N}}(\mathbf{y})\|$, for $\mathbf{y} \in \mathcal{V}$, if it is enabled thereby to offset that contribution to $\epsilon_{\mathcal{N}}$ by achieving a correspondingly better fit elsewhere on \mathcal{N} . In other words, the scalar quantity $\epsilon_{\mathcal{N}}$ may not convey sufficient information about the *distribution* of per-point errors arising from \widehat{f}^{-1} to establish whether or not f is injective, let alone immersive. Furthermore, since such a model must approximate on \mathcal{V} a map which at best is non-differentiable and at worst plain does not exist, a direct analysis of the per-point errors $\|\epsilon_{\mathcal{N}}(\mathbf{y})\|$ may not even be appropriate.

To overcome this potential limitation we need a measure which is sensitive to the *local* behaviour of \widehat{f} . The approach which we propose is to attempt to bound the errors associated with f by establishing the existence of constants L and U such that

$$L \leq \frac{\|f(\mathbf{x} + \Delta\mathbf{x}) - f(\mathbf{x})\|}{\|\Delta\mathbf{x}\|} \leq U \quad (3.19)$$

for all $\mathbf{x}, \Delta\mathbf{x} \in \mathcal{M}$. Clearly if f is injective then this expression can equivalently be written as

$$U^{-1} \leq \frac{\|f^{-1}(\mathbf{y} + \Delta\mathbf{y}) - f^{-1}(\mathbf{y})\|}{\|\Delta\mathbf{y}\|} \leq L^{-1} \quad (3.20)$$

where $\mathbf{y} = f(\mathbf{x})$. For convenience, we make equations (3.19) and (3.20) concrete by finding a lower bound for U and an upper bound for L such that

$$U \equiv \max_{\mathbf{y}=f(\mathbf{x})} \frac{\|\Delta\mathbf{y}\|}{\|\Delta\mathbf{x}\|}, \quad L \equiv \min_{\mathbf{y}=f(\mathbf{x})} \frac{\|\Delta\mathbf{y}\|}{\|\Delta\mathbf{x}\|} \quad (3.21)$$

where $\Delta \mathbf{y} = \mathbf{f}(\mathbf{x} + \Delta \mathbf{x}) - \mathbf{f}(\mathbf{x})$. If \mathbf{f} fails to be a diffeomorphism—through either non-immersivity or non-injectivity on some $\mathcal{U} \subset \mathcal{M}$ —then, trivially, $L = 0$ in equation (3.21). Similarly, the case where \mathbf{f}^{-1} exists, but is not a diffeomorphism, will result in a value of $U \rightarrow \infty$. If a finite lower bound U can be found then we say that \mathbf{f} is Lipschitz, with Lipschitz constant U , and if an $L > 0$ can be found then we can also say that \mathbf{f}^{-1} is Lipschitz, with Lipschitz constant L^{-1} . The existence of Lipschitz constants is a necessary, but not sufficient, condition for the existence of a diffeomorphism, so if we can show that a finite upper bound for one or both of L^{-1} and U cannot be found then we have ruled out the possibility that \mathbf{f} is a diffeomorphism.

To estimate these bounds we assume that $\widehat{\mathbf{f}}$ is a good fit to \mathbf{f} on \mathcal{M} , and find an expression for the inverse $\widehat{\mathbf{f}}^{-1}$ of $\widehat{\mathbf{f}}$ (as distinct from $\widehat{\mathbf{f}}^{-1}$, the RBF approximation of \mathbf{f}^{-1}). Although $\widehat{\mathbf{f}}$ is generically a diffeomorphism, given the appropriate conditions, it may be very close to an \mathbf{f} which is not a diffeomorphism. If this is the case, the estimate \widehat{L} of L , calculated from $\widehat{\mathbf{f}}$, should be noticeably small if $L = 0$, and similarly \widehat{U} should be correspondingly large if U itself is infinite.

It turns out that we can calculate these bounds if we examine the compositions $\widehat{\mathbf{I}}_{\mathcal{M}} = \widehat{\mathbf{f}}^{-1} \circ \widehat{\mathbf{f}}$ and $\widehat{\mathbf{I}}_{\mathcal{N}} = \widehat{\mathbf{f}} \circ \widehat{\mathbf{f}}^{-1}$, which we treat as approximations to the identity maps $\mathbf{I}_{\mathcal{M}}: \mathcal{M} \rightarrow \mathcal{M}$ and $\mathbf{I}_{\mathcal{N}}: \mathcal{N} \rightarrow \mathcal{N}$. (These should not be confused with the direct approximations to $\mathbf{I}_{\mathcal{M}}$ and $\mathbf{I}_{\mathcal{N}}$, which would be trivial to construct but of no practical use.) To this end, we define two new error functions, $\eta_{\mathcal{M}}: \mathbb{R}^m \rightarrow \mathbb{R}^m$ and $\eta_{\mathcal{N}}: \mathbb{R}^n \rightarrow \mathbb{R}^n$, by $\eta_{\mathcal{M}} = \widehat{\mathbf{I}}_{\mathcal{M}} - \mathbf{I}_{\mathcal{M}}$ and $\eta_{\mathcal{N}} = \widehat{\mathbf{I}}_{\mathcal{N}} - \mathbf{I}_{\mathcal{N}}$, respectively, and their corresponding error measures, $\eta_{\mathcal{M}}, \eta_{\mathcal{N}} \in \mathbb{R}$, by

$$\eta_{\mathcal{M}}^2 = \sigma_{\mathcal{M}}^{-2} \sum_{i=1}^N \|\eta_{\mathcal{M}}(\mathbf{x}_i)\|^2 = \sigma_{\mathcal{M}}^{-2} \sum_{i=1}^N \|\widehat{\mathbf{x}}_i - \mathbf{x}_i\|^2 \quad (3.22)$$

and

$$\eta_{\mathcal{N}}^2 = \sigma_{\mathcal{N}}^{-2} \sum_{i=1}^N \|\eta_{\mathcal{N}}(\mathbf{y}_i)\|^2 = \sigma_{\mathcal{N}}^{-2} \sum_{i=1}^N \|\widehat{\mathbf{y}}_i - \mathbf{y}_i\|^2 \quad (3.23)$$

where $\widehat{\mathbf{x}} = \widehat{\mathbf{I}}_{\mathcal{M}}(\mathbf{x})$ and $\widehat{\mathbf{y}} = \widehat{\mathbf{I}}_{\mathcal{N}}(\mathbf{y})$ are the LS RBF identity estimates of \mathbf{x} and \mathbf{y} , respectively. Then by noting that $\widehat{\mathbf{f}}^{-1} = \mathbf{f}^{-1} + \epsilon_{\mathcal{N}}$ and substituting $\widehat{\mathbf{f}}$ for \mathbf{f} we can expand $\eta_{\mathcal{M}}$, acting on the point $\mathbf{x} \in \mathcal{M}$, as

$$\begin{aligned} \eta_{\mathcal{M}}(\mathbf{x}) &= \widehat{\mathbf{f}}^{-1} \circ \widehat{\mathbf{f}}(\mathbf{x}) - \mathbf{x} \\ &= (\mathbf{f}^{-1} + \epsilon_{\mathcal{N}}) \circ \widehat{\mathbf{f}}(\mathbf{x}) - \mathbf{x} \\ &\approx (\mathbf{f}^{-1} + \epsilon_{\mathcal{N}}) \circ \mathbf{f}(\mathbf{x}) - \mathbf{x} \\ &= \epsilon_{\mathcal{N}} \circ \mathbf{f}(\mathbf{x}) \end{aligned} \quad (3.24)$$

The equivalent expression for $\eta_{\mathcal{N}}$ is less simple, with $\eta_{\mathcal{N}}$, acting on the point $\mathbf{y} = \mathbf{f}(\mathbf{x})$, written as

$$\begin{aligned}
\eta_{\mathcal{N}}(\mathbf{y}) &= \widehat{\mathbf{f}} \circ \widehat{\mathbf{f}}^{-1}(\mathbf{y}) - \mathbf{y} \\
&= \widehat{\mathbf{f}} \circ (\mathbf{f}^{-1} + \epsilon_{\mathcal{N}})(\mathbf{y}) - \mathbf{y} \\
&= \widehat{\mathbf{f}}(\mathbf{x} + \epsilon_{\mathcal{N}} \circ \mathbf{f}(\mathbf{x})) - \mathbf{f}(\mathbf{x}) \\
&\approx \widehat{\mathbf{f}}(\mathbf{x} + \eta_{\mathcal{M}}(\mathbf{x})) - \widehat{\mathbf{f}}(\mathbf{x})
\end{aligned} \tag{3.25}$$

Assuming that $\widehat{\mathbf{f}}^{-1}$ exists, we can also express this relationship as

$$\eta_{\mathcal{M}}(\mathbf{x}) \approx \widehat{\mathbf{f}}^{-1}(\mathbf{y} + \eta_{\mathcal{N}}(\mathbf{y})) - \widehat{\mathbf{f}}^{-1}(\mathbf{y}) \tag{3.26}$$

We can now rearrange these equations to get, in analogy to equations (3.19) and (3.20),

$$\widehat{L} \lesssim \frac{\|\widehat{\mathbf{f}}(\mathbf{x} + \eta_{\mathcal{M}}(\mathbf{x})) - \widehat{\mathbf{f}}(\mathbf{x})\|}{\|\eta_{\mathcal{M}}(\mathbf{x})\|} \lesssim \widehat{U} \tag{3.27}$$

and

$$\widehat{U}^{-1} \lesssim \frac{\|\widehat{\mathbf{f}}^{-1}(\mathbf{y} + \eta_{\mathcal{N}}(\mathbf{y})) - \widehat{\mathbf{f}}^{-1}(\mathbf{y})\|}{\|\eta_{\mathcal{N}}(\mathbf{y})\|} \lesssim \widehat{L}^{-1} \tag{3.28}$$

In practice, given a finite training set, we replace the inequalities in equations (3.27) and (3.28) with definitions analogous to those of equation (3.21),

$$\widehat{L} \equiv \min_{\mathbf{y}=\mathbf{f}(\mathbf{x})} \frac{\|\eta_{\mathcal{N}}(\mathbf{y})\|}{\|\eta_{\mathcal{M}}(\mathbf{x})\|}, \quad \widehat{U} \equiv \max_{\mathbf{y}=\mathbf{f}(\mathbf{x})} \frac{\|\eta_{\mathcal{N}}(\mathbf{y})\|}{\|\eta_{\mathcal{M}}(\mathbf{x})\|} \tag{3.29}$$

If \widehat{L} is sufficiently large then, to the extent that our assumption that \mathbf{f} and $\widehat{\mathbf{f}}$ are interchangeable is justified, we can now claim that $\widehat{\mathbf{f}}^{-1}$ is Lipschitz, with Lipschitz constant \widehat{L}^{-1} , and if \widehat{U} is sufficiently small we can similarly claim that $\widehat{\mathbf{f}}$ is Lipschitz, with Lipschitz constant \widehat{U} . If either of these constants is outside of acceptable parameters then we have established that $\widehat{\mathbf{f}}$ approximates an \mathbf{f} which is not a diffeomorphism. We will investigate the use of this method, as an adjunct to the analysis of the potentially less sensitive LS errors, in the following chapter.

3.3 The method of total least squares

In the previous section we concentrated on an analysis of the errors arising from the LS approximations of \mathbf{f} and \mathbf{f}^{-1} , individually and in composition, as indicators of how nearly \mathbf{f} fails to be a diffeomorphism. Returning to the discussion in section 3.1, we will now consider a slightly different approach, in which we attempt to answer the same question—is \mathbf{f} a diffeomorphism—in terms of a single, ‘symmetrical’ model $\widehat{\mathbf{f}}: \mathbb{R}^m \rightarrow \mathbb{R}^n$. (By symmetrical, we mean that if $\mathbf{g}: \mathbb{R}^m \rightarrow \mathbb{R}^n$ is a model of \mathbf{f} , and $\mathbf{h}: \mathbb{R}^n \rightarrow \mathbb{R}^m$ is a

model of f^{-1} , then $h = g^{-1}$.) With \widehat{f} defined as in equation (3.1), this task is made particularly hard by the fact that φ is a nonlinear transformation. But we show, in appendix A, that φ is generically an embedding of compact sets: we can therefore restrict our attention to the linear part \mathcal{W} of \widehat{f} , for which the properties of injectivity and immersivity are equivalent (ie. an immersive linear map is injective onto its image). This is a reasonable approach because, as we have already stated, with enough centers we can make the relationship between $\varphi\mathcal{M}$ and \mathcal{N} as nearly linear as we wish: this is precisely the source of the approximating power of the RBF map.

In general, of course, \mathcal{W}^{-1} (if it exists) is itself a nonlinear map. For this reason, and also from general considerations of symmetry, we now adapt the classical RBF map, incorporating *two* nonlinear transformations $\varphi_{\mathcal{M}}: \mathbb{R}^m \rightarrow \mathbb{R}^p$ and $\varphi_{\mathcal{N}}: \mathbb{R}^n \rightarrow \mathbb{R}^p$, and writing

$$\widehat{f} = \varphi_{\mathcal{N}}^{-1} \circ \mathcal{W} \circ \varphi_{\mathcal{M}} \quad (3.30)$$

with the caveat that $\varphi_{\mathcal{N}}^{-1}$ is only defined—by virtue of the analysis in appendix A—on the image, under $\varphi_{\mathcal{N}}$, of compact subsets of \mathbb{R}^n . We define $\mathcal{W}: \mathbb{R}^p \rightarrow \mathbb{R}^p$ by $\mathcal{W}(\varphi) = \overline{\varphi_{\mathcal{N}}} + \mathbf{W}^T(\varphi - \overline{\varphi_{\mathcal{M}}})$, with \mathbf{W} a p by p matrix. In this symmetrical form the inverse of \mathcal{W} is itself a linear map, with $\mathcal{W}^{-1}(\varphi) = \overline{\varphi_{\mathcal{M}}} + \mathbf{W}^{-T}(\varphi - \overline{\varphi_{\mathcal{N}}})$, where \mathbf{W}^{-T} represents the transpose of \mathbf{W}^{-1} . Assuming that \mathbf{W}^{-1} exists—that is, $\text{rank } \mathbf{W} = p$ —we can now invert (3.30) directly, asserting $\widehat{f}^{-1} = \widehat{f}^{-1}$, to get

$$\widehat{f}^{-1} = \varphi_{\mathcal{M}}^{-1} \circ \mathcal{W}^{-1} \circ \varphi_{\mathcal{N}} \quad (3.31)$$

with the appropriate condition on the existence of $\varphi_{\mathcal{M}}^{-1}$.

For notational convenience, since we are now interested solely in the invertibility of \mathcal{W} , we make the definitions $\mathcal{A} = \varphi_{\mathcal{M}}\mathcal{M}$ and $\mathcal{B} = \varphi_{\mathcal{N}}\mathcal{N}$, and write $\mathbf{a} = \varphi_{\mathcal{M}}(x) - \overline{\varphi_{\mathcal{M}}}$ and $\mathbf{b} = \varphi_{\mathcal{N}}(y) - \overline{\varphi_{\mathcal{N}}}$. We also define two new error terms, $\epsilon_{\mathcal{A}}, \epsilon_{\mathcal{B}} \in \mathbb{R}$, with

$$\epsilon_{\mathcal{A}}^2 = \sigma_{\mathcal{B}}^{-2} \sum_{i=1}^N \|\widehat{\mathbf{b}}_i - \mathbf{b}_i\|^2, \quad \epsilon_{\mathcal{B}}^2 = \sigma_{\mathcal{A}}^{-2} \sum_{i=1}^N \|\widehat{\mathbf{a}}_i - \mathbf{a}_i\|^2 \quad (3.32)$$

where $\widehat{\mathbf{b}} = \mathbf{W}^T \mathbf{a}$, $\widehat{\mathbf{a}} = \mathbf{W}^{-T} \mathbf{b}$ and the normalisers $\sigma_{\mathcal{A}}$ and $\sigma_{\mathcal{B}}$ are calculated on $\mathcal{A}, \mathcal{B} \subset \mathbb{R}^p$, respectively. We work with these errors, calculated in \mathbb{R}^p , rather than with $\epsilon_{\mathcal{M}}$ and $\epsilon_{\mathcal{N}}$, because we are specifically interested in the linear part of the symmetrical RBF map: the latter two errors now suffer from the effects induced by transformation through an inverted RBF nonlinearity, and are therefore no longer suitable for our purposes.

Having defined the symmetrical RBF map we must now reconsider the manner in which the weight matrix \mathbf{W} —now common to both \widehat{f} and \widehat{f}^{-1} —is calculated from a particular set of training data. Clearly, if we were to use the method of LS we could minimise either the forward error, $\epsilon_{\mathcal{A}}$, or the inverse error, $\epsilon_{\mathcal{B}}$, but not both together: each would be minimised at the expense of the other. We therefore need to find

a new method, one which ignores the distinction between mapping in the forward and inverse directions. For this purpose, it will prove useful to rewrite (3.32) in terms of a norm calculated in the ‘product’ space $\mathbb{R}^{2p} = \mathbb{R}^p \times \mathbb{R}^p$, whose elements are the vectors $\begin{pmatrix} \mathbf{a} \\ \mathbf{b} \end{pmatrix}$, to get

$$\epsilon_{\mathcal{A}}^2 = \sigma_{\mathcal{B}}^{-2} \sum_{i=1}^N \left\| \begin{pmatrix} \mathbf{a}_i \\ \widehat{\mathbf{b}}_i \end{pmatrix} - \begin{pmatrix} \mathbf{a}_i \\ \mathbf{b}_i \end{pmatrix} \right\|^2, \quad \epsilon_{\mathcal{B}}^2 = \sigma_{\mathcal{A}}^{-2} \sum_{i=1}^N \left\| \begin{pmatrix} \widehat{\mathbf{a}}_i \\ \mathbf{b}_i \end{pmatrix} - \begin{pmatrix} \mathbf{a}_i \\ \mathbf{b}_i \end{pmatrix} \right\|^2 \quad (3.33)$$

From a geometrical standpoint, \mathcal{W} is defined by the p -dimensional hyperplane $\mathcal{H} \subset \mathbb{R}^{2p}$ containing the points $\begin{pmatrix} \mathbf{a} \\ \widehat{\mathbf{b}} \end{pmatrix}$ and $\begin{pmatrix} \widehat{\mathbf{a}} \\ \mathbf{b} \end{pmatrix}$. We could clearly specify this hyperplane by independently minimising either $\epsilon_{\mathcal{A}}$ or $\epsilon_{\mathcal{B}}$, as discussed above, but this would result in a distinct hyperplane—and hence RBF map—in each case. To retain the symmetrical aspect of this new model we must instead find that single hyperplane \mathcal{H} which (in some sense) best captures the linear relationship embodied in the *joint distribution* of data pairs $\begin{pmatrix} \mathbf{a} \\ \mathbf{b} \end{pmatrix} \in \mathbb{R}^{2p}$. This interpretation suggests a new error, ϵ_{\perp} , defined by

$$\epsilon_{\perp}^2 = \sigma^{-2} \sum_{i=1}^N \left\| \begin{pmatrix} \widetilde{\mathbf{a}}_i \\ \widetilde{\mathbf{b}}_i \end{pmatrix} - \begin{pmatrix} \mathbf{a}_i \\ \mathbf{b}_i \end{pmatrix} \right\|^2 \quad (3.34)$$

and normalised by

$$\sigma^2 = \sigma_{\mathcal{A}}^2 + \sigma_{\mathcal{B}}^2 = \sum_{i=1}^N \left\| \begin{pmatrix} \mathbf{a}_i \\ \mathbf{b}_i \end{pmatrix} \right\|^2 \quad (3.35)$$

where $\begin{pmatrix} \widetilde{\mathbf{a}} \\ \widetilde{\mathbf{b}} \end{pmatrix} \in \mathbb{R}^{2p}$ is the image of $\begin{pmatrix} \mathbf{a} \\ \mathbf{b} \end{pmatrix}$ under normal projection into \mathcal{H} . We call ϵ_{\perp} the total least squares (TLS) error: its minimisation defines a p -plane \mathcal{H} which is the best rank- p linear approximation to the joint distribution in \mathbb{R}^{2p} . It is to the TLS error that we look when we wish to know how well we have modelled a particular training set: if ϵ_{\perp} is large then we can only state that $\widehat{\mathbf{f}}$ is not a particularly good model of \mathbf{f} —in other words, $\varphi_{\mathcal{N}} \circ \mathbf{f} \circ \varphi_{\mathcal{M}}^{-1}$ is not sufficiently close to a linear map—and continue no further; if, on the other hand, ϵ_{\perp} is deemed to be acceptably small then the determination of whether or not $\widehat{\mathbf{f}}$ is a diffeomorphism boils down to the linear relationship represented by \mathcal{H} .

It is important to note that in choosing to minimise ϵ_{\perp} , rather than the ‘directional’ errors $\epsilon_{\mathcal{A}}$ and $\epsilon_{\mathcal{B}}$, we are no longer able to assume the existence of either \mathcal{W} or \mathcal{W}^{-1} , based purely on the direction of the RBF map in question: whether or not \mathcal{H} represents an *invertible* map \mathcal{W} between the $\mathbf{a}, \mathbf{b} \in \mathbb{R}^p$ depends solely on the orientation of \mathcal{H} in \mathbb{R}^{2p} . Our labelling of \mathcal{W} and \mathcal{W}^{-1} as forward and inverse maps is completely arbitrary in this symmetrical model, and will be effectively circumvented in section 3.3.1 below. In section 3.3.2 we discuss the means by which we investigate the invertibility of \mathcal{W} , using either the directional errors or an equivalent measure of the ill-conditioning of \mathcal{W} . Then, in section 3.3.3, we discuss the origin of a numerical instability which has been observed, in practice, to detract from the overall usability of the TLS methodology. First, however, we describe the method by which we minimise ϵ_{\perp} , and hence obtain an analytic expression for \mathcal{W} . For a more comprehensive description of TLS, see

Van Huffel and Vandewalle [41]; a definitive, but concise treatment can also be found in Golub and Van Loan [15].

3.3.1 The total least squares solution

To find an analytic expression for \mathbf{W} we rely once more on the SVD, writing the composite matrix $(\mathbf{A}; \mathbf{B})$ as $\mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$. This decomposition provides us with a basis for the approximating hyperplane \mathcal{H} , obtained by further decomposing the $2p$ by $2p$ matrix \mathbf{V} of right singular vectors into $\mathbf{V} = (\mathbf{V}_p; \overline{\mathbf{V}}_p)$, where $\mathbf{V}_p, \overline{\mathbf{V}}_p$ have p columns each. In this basis, \mathcal{H} is spanned by the columns of \mathbf{V}_p , and its orthogonal complement in \mathbb{R}^{2p} by those of $\overline{\mathbf{V}}_p$. The projection of $\begin{pmatrix} \mathbf{a} \\ \mathbf{b} \end{pmatrix}$ into \mathcal{H} is thus $\begin{pmatrix} \tilde{\mathbf{a}} \\ \tilde{\mathbf{b}} \end{pmatrix} = \mathbf{V}_p \mathbf{V}_p^T \begin{pmatrix} \mathbf{a} \\ \mathbf{b} \end{pmatrix}$, and equation (3.34) can be rewritten as

$$\begin{aligned} \epsilon_{\perp}^2 &= \sigma^{-2} \sum_{i=1}^N \|(1 - \mathbf{V}_p \mathbf{V}_p^T) \begin{pmatrix} \mathbf{a}_i \\ \mathbf{b}_i \end{pmatrix}\|^2 \\ &= \sigma^{-2} \sum_{i=1}^N \|\overline{\mathbf{V}}_p \overline{\mathbf{V}}_p^T \begin{pmatrix} \mathbf{a}_i \\ \mathbf{b}_i \end{pmatrix}\|^2 \\ &= \sigma^{-2} \sum_{i=1}^N \|\overline{\mathbf{V}}_p^T \begin{pmatrix} \mathbf{a}_i \\ \mathbf{b}_i \end{pmatrix}\|^2 \\ &= \sigma^{-2} \sum_{i=1}^N \|\mathbf{P}^T \mathbf{a}_i + \mathbf{Q}^T \mathbf{b}_i\|^2 \end{aligned} \tag{3.36}$$

where the last line follows from a further partitioning of $\overline{\mathbf{V}}_p$ into p by p submatrices \mathbf{P} and \mathbf{Q} , with $\overline{\mathbf{V}}_p = \begin{pmatrix} \mathbf{P} \\ \mathbf{Q} \end{pmatrix}$. The variance σ^2/N of the distribution in \mathbb{R}^{2p} can also be expressed directly in terms of the singular values σ_j comprising the diagonal elements of $\mathbf{\Sigma}$, with

$$\sigma^2 = \sum_{j=1}^{2p} \sigma_j^2 \tag{3.37}$$

and the TLS error can be similarly written as

$$\epsilon_{\perp}^2 = \sigma^{-2} \sum_{j=p+1}^{2p} \sigma_j^2 \tag{3.38}$$

This form is used, in practice, to detect the onset of numerical error in the calculation of \mathbf{W} and \mathbf{W}^{-1} by comparing the resulting value with that calculated according to equation (3.36).

To find an expression for \mathbf{W} we notice that we can now write, for $\begin{pmatrix} \mathbf{a} \\ \mathbf{b} \end{pmatrix} \in \mathcal{H} \subset \mathbb{R}^{2p}$,

$$\overline{\mathbf{V}}_p^T \begin{pmatrix} \mathbf{a} \\ \mathbf{b} \end{pmatrix} = \mathbf{P}^T \mathbf{a} + \mathbf{Q}^T \mathbf{b} = 0 \tag{3.39}$$

and hence $\mathbf{b} = -\mathbf{Q}^{-\text{T}}\mathbf{P}^{\text{T}}\mathbf{a}$, assuming that \mathbf{Q}^{-1} exists, and $\mathbf{a} = -\mathbf{P}^{-\text{T}}\mathbf{Q}^{\text{T}}\mathbf{b}$, assuming that \mathbf{P}^{-1} exists (bear in mind that \mathbf{P} and \mathbf{Q} are not, in general, orthonormal matrices), enabling us to make—in each case—the definitions

$$\mathbf{W} = -\mathbf{P}\mathbf{Q}^{-1} \quad (3.40)$$

and, equivalently,

$$\mathbf{W}^{-1} = -\mathbf{Q}\mathbf{P}^{-1} \quad (3.41)$$

We write \mathbf{Q}^{-1} and \mathbf{P}^{-1} in equations (3.40) and (3.41) as inverses, rather than pseudo-inverses, because \mathbf{P} and \mathbf{Q} are square and (by assumption) invertible matrices. Nevertheless, in order to calculate \mathbf{W} and/or \mathbf{W}^{-1} for a given training set we use the familiar SVD to get $\mathbf{P} = \mathbf{L}\mathbf{S}\mathbf{T}^{\text{T}}$, where the p by p matrices \mathbf{L} , \mathbf{S} and \mathbf{T} have the usual interpretation. Noticing that $\overline{\mathbf{V}}_p^{\text{T}}\overline{\mathbf{V}}_p = \mathbf{P}^{\text{T}}\mathbf{P} + \mathbf{Q}^{\text{T}}\mathbf{Q} = \mathbf{I}$, and hence $\mathbf{Q}^{\text{T}}\mathbf{Q} = \mathbf{I} - \mathbf{P}^{\text{T}}\mathbf{P}$, we now write

$$\mathbf{Q}^{\text{T}}\mathbf{Q}\mathbf{T} = (\mathbf{I} - \mathbf{P}^{\text{T}}\mathbf{P})\mathbf{T} = \mathbf{T} - \mathbf{T}\mathbf{S}^2 = \mathbf{T}\mathbf{C}^2 \quad (3.42)$$

where

$$\mathbf{C}^2 = \mathbf{I} - \mathbf{S}^2 \quad (3.43)$$

allowing us to relate the decompositions of \mathbf{Q} and \mathbf{P} with $\mathbf{Q} = \mathbf{R}\mathbf{C}\mathbf{T}^{\text{T}}$, with the understanding that the singular values c_j , which appear as diagonal elements of \mathbf{C} , are ordered with respect to increasing value, rather than decreasing as is conventional in the SVD. To make this relationship explicit we note that (3.43) allows us to write the diagonal elements of \mathbf{S} and \mathbf{C} —the singular values of \mathbf{P} and \mathbf{Q} —as $s_j = \sin \theta_j$ and $c_j = \cos \theta_j$, respectively, for $0 \leq \theta_{j+1} \leq \theta_j \leq \frac{\pi}{4}$. By combining these two decompositions we can now write down an SVD of the TLS weight matrix (3.40) as

$$\mathbf{W} = -\mathbf{L}\mathbf{A}\mathbf{R}^{\text{T}} \quad (3.44)$$

with singular values $\lambda_j = \tan \theta_j$ forming the diagonal elements of $\mathbf{A} = \mathbf{S}\mathbf{C}^{-1}$.

3.3.2 Detecting non-invertible maps with total least squares

We are now in a position to establish a definitive test for the existence (or not) of an invertible linear relationship in \mathbb{R}^{2p} . We base this test, naturally enough, on the directional errors $\epsilon_{\mathcal{A}}$ and $\epsilon_{\mathcal{B}}$, on the

grounds that (assuming ϵ_{\perp} is sufficiently small) a large $\epsilon_{\mathcal{A}}$ or $\epsilon_{\mathcal{B}}$ is indicative of a rank-deficiency in \mathbf{W} or \mathbf{W}^{-1} , respectively. To investigate this mechanism we make the additional partitions $\mathbf{U} = (\mathbf{U}_p; \overline{\mathbf{U}}_p)$ and $\Sigma = \begin{pmatrix} \Sigma_p & 0 \\ 0 & \overline{\Sigma}_p \end{pmatrix}$, so that $(\mathbf{A}; \mathbf{B})\overline{\mathbf{V}}_p = \mathbf{A}\mathbf{P} + \mathbf{B}\mathbf{Q} = \overline{\mathbf{U}}_p\overline{\Sigma}_p$, and write

$$\begin{aligned}\epsilon_{\mathcal{A}}^2 &= \sigma_{\mathcal{B}}^{-2} \|\mathbf{A}\mathbf{W} - \mathbf{B}\|_{\text{F}}^2 \\ &= \sigma_{\mathcal{B}}^{-2} \|(\mathbf{A}\mathbf{P} + \mathbf{B}\mathbf{Q})\mathbf{Q}^{-1}\|_{\text{F}}^2 \\ &= \sigma_{\mathcal{B}}^{-2} \|\overline{\mathbf{U}}_p\overline{\Sigma}_p\mathbf{T}\mathbf{C}^{-1}\mathbf{R}^{\text{T}}\|_{\text{F}}^2 \\ &= \sigma_{\mathcal{B}}^{-2} \|\overline{\Sigma}_p\mathbf{T}\mathbf{C}^{-1}\|_{\text{F}}^2\end{aligned}\tag{3.45}$$

and

$$\begin{aligned}\epsilon_{\mathcal{B}}^2 &= \sigma_{\mathcal{A}}^{-2} \|\mathbf{B}\mathbf{W}^{-1} - \mathbf{A}\|_{\text{F}}^2 \\ &= \sigma_{\mathcal{A}}^{-2} \|(\mathbf{A}\mathbf{P} + \mathbf{B}\mathbf{Q})\mathbf{P}^{-1}\|_{\text{F}}^2 \\ &= \sigma_{\mathcal{A}}^{-2} \|\overline{\mathbf{U}}_p\overline{\Sigma}_p\mathbf{T}\mathbf{S}^{-1}\mathbf{R}^{\text{T}}\|_{\text{F}}^2 \\ &= \sigma_{\mathcal{A}}^{-2} \|\overline{\Sigma}_p\mathbf{T}\mathbf{S}^{-1}\|_{\text{F}}^2\end{aligned}\tag{3.46}$$

If we denote by $\epsilon_j \in \mathbb{R}^p$ the columns of $\overline{\Sigma}_p\mathbf{T}$, and write $\epsilon_j = \|\epsilon_j\|$, then (3.45) and (3.46) become

$$\epsilon_{\mathcal{A}}^2 = \sigma_{\mathcal{B}}^{-2} \sum_{j=1}^p \frac{\epsilon_j^2}{\cos^2 \theta_j}, \quad \epsilon_{\mathcal{B}}^2 = \sigma_{\mathcal{A}}^{-2} \sum_{j=1}^p \frac{\epsilon_j^2}{\sin^2 \theta_j}\tag{3.47}$$

revealing their dependence on the reciprocals of the singular values $\sin \theta_j$ and $\cos \theta_j$. As an aside, it is interesting to note that we can also rewrite (3.38) in terms of ϵ_j to get

$$\epsilon_{\perp}^2 = \sum_{j=1}^p \epsilon_j^2\tag{3.48}$$

from which we see, in the (for our purposes) unlikely case that $p = 1$, that ϵ_{\perp} , $\epsilon_{\mathcal{A}}$ and $\epsilon_{\mathcal{B}}$ obey the simple reciprocal relationship

$$\frac{1}{\sigma^2 \epsilon_{\perp}^2} = \frac{1}{\sigma_{\mathcal{B}}^2 \epsilon_{\mathcal{A}}^2} + \frac{1}{\sigma_{\mathcal{A}}^2 \epsilon_{\mathcal{B}}^2}\tag{3.49}$$

A commonly used measure of the rank deficiency of a matrix is the condition number $\kappa(\cdot)$, defined as the ratio of the highest to lowest singular values of that matrix; for our purposes, an ill-conditioned matrix is one whose smallest singular value approaches the lower limit of numerical precision. The condition numbers of \mathbf{P} and \mathbf{Q} are

$$\kappa(\mathbf{P}) = \frac{\sin \theta_1}{\sin \theta_p}, \quad \kappa(\mathbf{Q}) = \frac{\cos \theta_p}{\cos \theta_1}\tag{3.50}$$

A large error $\epsilon_{\mathcal{A}}$ or $\epsilon_{\mathcal{B}}$ is symptomatic of an ill-conditioned submatrix \mathbf{Q} or \mathbf{P} , respectively, so we will examine both errors and condition numbers in the experiments to follow. It is interesting to note, from equation (3.44), that the condition number of \mathbf{W} is

$$\kappa(\mathbf{W}) = \frac{\tan \theta_1}{\tan \theta_p} = \kappa(\mathbf{P}) \kappa(\mathbf{Q}) \quad (3.51)$$

although this is not true for matrix products in general.

It turns out that we can relate these condition numbers directly to $\epsilon_{\mathcal{A}}$ and $\epsilon_{\mathcal{B}}$ if we consider what happens at small singular values, as θ_1 approaches $\frac{\pi}{4}$ or θ_p approaches zero: in the former case we get

$$\epsilon_{\mathcal{A}} \sim \frac{\epsilon_1}{\cos \theta_1} \sim \frac{\epsilon_1}{\cos \theta_p} \kappa(\mathbf{Q}) \quad (3.52)$$

and in the latter,

$$\epsilon_{\mathcal{B}} \sim \frac{\epsilon_p}{\sin \theta_p} \sim \frac{\epsilon_p}{\sin \theta_1} \kappa(\mathbf{P}) \quad (3.53)$$

so in other words we expect to see (in the limit) a linear relationship between $\epsilon_{\mathcal{A}}$ and $\kappa(\mathbf{Q})$, and between $\epsilon_{\mathcal{B}}$ and $\kappa(\mathbf{P})$. This prediction will be confirmed in the experiments which follow in chapter 4.

3.3.3 Numerical instability in the total least squares solution

Implicit in the TLS model is the assumption that the rank of $(\mathbf{A}; \mathbf{B})$ is exactly p —that is, the relationship between \mathbf{a} and $\mathbf{b} = \varphi_{\mathcal{N}} \circ \mathbf{f} \circ \varphi_{\mathcal{M}}^{-1}(\mathbf{a})$ is a linear one. If this were truly the case, then we would necessarily find a TLS error $\epsilon_{\perp} = 0$. In practice, however, with finite p we will invariably find that $\epsilon_{\perp} > 0$, and determining an appropriate value for $r = \text{rank}(\mathbf{A}; \mathbf{B})$ will become a somewhat qualitative endeavour. Nevertheless, provided that $r \geq p$, and ϵ_{\perp} is acceptably low, we might reasonably expect the method of TLS to be a valid technique. However, if $r < p$ a new problem arises: by insisting on a p -dimensional hyperplane \mathcal{H} we are forced to construct the projector $\overline{\mathbf{V}}_p$ from a degenerate subspace of dimension $2p - r > p$ in \mathbb{R}^{2p} . It is then possible for a small change in (say) a single element of the training set in \mathbb{R}^{2p} to result in a large change to the critical matrix $\overline{\mathbf{V}}_p$, simply because a different subset of the set of degenerate basis vectors happened to be selected on that occasion. This change, propagated to the submatrices \mathbf{P} and \mathbf{Q} , is likely to become extremely noticeable when either \mathbf{P} or \mathbf{Q} is ill-conditioned—which is precisely the case in which we are interested. To see how such a rank deficiency might come about, we define $s = \text{rank } \mathbf{A}$ and $t = \text{rank } \mathbf{B}$, and write

$$\max(s, t) \leq r \leq s + t \quad (3.54)$$

We can therefore expect to see this situation arise when both $s, t < p$ —in other words, when there are linear dependencies *within* both \mathcal{A} and \mathcal{B} . Experience has shown that this is indeed often the case; particularly when p must be sufficiently high for \mathcal{W} to be made as nearly linear as possible through the transformations $\varphi_{\mathcal{A}}$ and $\varphi_{\mathcal{B}}$. This is the case in the experiments to be described in the following chapters, in which we will find that small variations in the experimental parameter under observation do indeed give rise to large fluctuations in the condition numbers $\kappa(Q)$ and $\kappa(P)$, and hence in the errors $\epsilon_{\mathcal{A}}$ and $\epsilon_{\mathcal{B}}$.

Chapter 4

Maps on manifolds

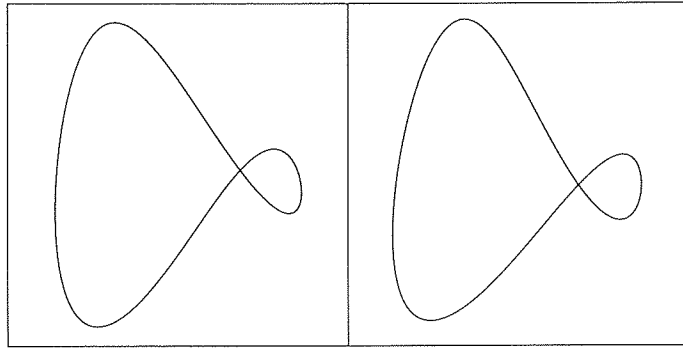
In this chapter we apply the fitting techniques described in the previous section to a search for diffeomorphisms between subsets of Euclidean space. We describe two experiments, each of which comprises a family of maps $\{f_\mu\}_{\mu \in \mathbb{R}}$, where μ is a control parameter determining whether or not f_μ is a diffeomorphism, to which we fit RBF maps \widehat{f}_μ and, where appropriate, \widehat{f}_μ^{-1} , using training and test sets generated by sampling f_μ . In the first experiment the maps in question will be defined on projections into the plane of a circle in \mathbb{R}^3 , and in the second they will be on 2-tori, also in \mathbb{R}^3 . Since neither the domain or range of these maps are delay embeddings of dynamical systems, they do not exhibit the shift property described in section 2.2, so we do not impose any additional constraints on their RBF approximations. We investigate both LS and TLS RBF maps, for purposes of comparison. In the LS case, we will see that in both the circle and torus experiments we are able to write down f_μ explicitly, and hence expect to be able to find an arbitrarily good LS approximation \widehat{f}_μ to f_μ , as discussed in the previous chapter. This expectation will be borne out in this chapter.

4.1 Embedding a circle in the plane

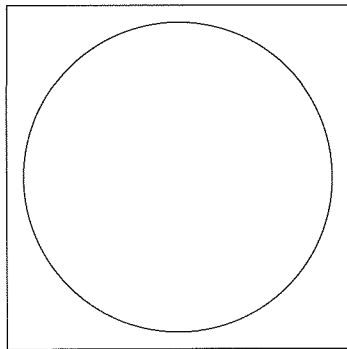
We start by describing a simple experiment in fitting maps between projections into the plane of the circle S^1 embedded in \mathbb{R}^3 . The data was generated by the map $\Phi: S^1 \rightarrow \mathbb{R}^3$ defined by

$$\Phi(\theta) = \begin{pmatrix} \sin \theta \\ \cos \theta \\ \sin 2\theta \end{pmatrix} \quad (4.1)$$

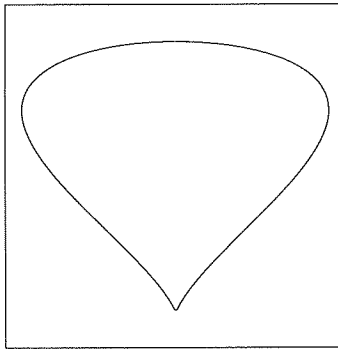
which has the effect of drawing a sinusoid around a cylinder in \mathbb{R}^3 , producing a manifold $\mathcal{S} = \{\Phi(\theta): 0 \leq \theta < 2\pi\}$ which is a topological circle. (We use the symbol ' Φ ' to denote this map to indicate that it is



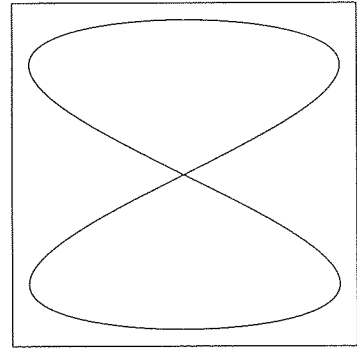
(a)



(b)



(c)



(d)

Figure 4.1 The circle embedded in \mathbb{R}^3 and its projections in the plane. Part (a) is a stereoscopic projection of $S \subset \mathbb{R}^3$; examples of its topologically distinct images, under projection into \mathbb{R}^2 , are (b) the circle S_0 and (d) the figure of eight S_{90} , and are separated by (c) a cusp which occurs at $\phi \approx 26.57$ degrees.

an embedding—although not a delay embedding—of \mathcal{S}^1 .) A stereoscopic projection of \mathcal{S} is illustrated in figure 4.1(a). This circle is then mapped under a one-parameter family of projections $\mathcal{F}_\phi: \mathcal{S} \subset \mathbb{R}^3 \rightarrow \mathbb{R}^2$ (the symbol ϕ is not to be confused here with its earlier usage as an RBF nonlinearity) realised by matrices of the form

$$\mathbf{F}_\phi = \begin{pmatrix} \cos \phi & 0 & \sin \phi \\ 0 & 1 & 0 \end{pmatrix} \quad (4.2)$$

The result is the composite map $\mathcal{G}_\phi = \mathcal{F}_\phi \circ \Phi$ written

$$\mathcal{G}_\phi(\theta) = \begin{pmatrix} \sin \theta \cos \phi + \sin 2\theta \sin \phi \\ \cos \theta \end{pmatrix} \quad (4.3)$$

We write $\mathcal{S}_\phi = \mathcal{F}_\phi \mathcal{S}$, and confine our interest to values of ϕ in the range $0 \leq \phi \leq \frac{\pi}{2}$. Values of ϕ outside of this range produce projections \mathcal{S}_ϕ which can be obtained by reflecting those obtained with ϕ inside that range about the horizontal axis in \mathbb{R}^2 . A sequence of projections was performed, with ϕ increasing in one degree increments. The form of \mathcal{F}_ϕ was chosen so that there is a critical parameter value ϕ^* such that \mathcal{F}_ϕ is an embedding on \mathcal{S} if $\phi < \phi^*$. At values of ϕ above this critical value the image \mathcal{S}_ϕ contains a self-intersection; we call such a set a figure-of-eight. This behaviour is illustrated in figure 4.1(b) through (d), which shows the sets \mathcal{S}_0 , \mathcal{S}_{25} and \mathcal{S}_{90} , where the subscripts are written in degrees for convenience.

At the critical value, the onset of self-intersection manifests itself as a cusp in \mathcal{F}_{ϕ^*} , as indicated in figure 4.1(c). Self-intersections occur when two points in \mathcal{S} are projected to the same point in \mathcal{S}_ϕ ; the occurrence of a cusp corresponds to a projection in the direction of the gradient $\nabla \Phi(\theta)$, given by

$$\nabla \Phi(\theta) = \begin{pmatrix} \cos \theta \\ -\sin \theta \\ 2 \cos 2\theta \end{pmatrix} \quad (4.4)$$

In other words, \mathcal{F}_{ϕ^*} is the projection which maps $\nabla \Phi(\theta)$ to the zero vector for some θ . We can thus determine ϕ^* analytically by solving the system of equations

$$\left. \begin{aligned} \cos \theta \cos \phi^* + 2 \cos 2\theta \sin \phi^* &= 0 \\ -\sin \theta &= 0 \end{aligned} \right\} \quad (4.5)$$

which has solutions $\phi^* = \tan^{-1} \pm \frac{1}{2}$. We disregard the negative solution to keep ϕ in the appropriate range, leaving $\phi^* \approx 0.46$, or approximately 26.57 degrees.

What we would like to do with this data is identify those sets \mathcal{S}_ϕ which contain self-intersections—that is, for which $\phi > \phi^*$ —by identifying those projections \mathcal{F}_ϕ which fail to be a diffeomorphism on \mathcal{S} . We could attempt to solve this problem by approximating the inverse, if it exists, to each \mathcal{F}_ϕ , acting on \mathcal{S}_ϕ , directly, but since the circle \mathcal{S}_0 is clearly topologically equivalent to \mathcal{S} we instead assert that \mathcal{F}_0 is a diffeomorphism, then simply approximate the map $f_\phi: \mathcal{S}_0 \rightarrow \mathcal{S}_\phi$, defined by $f_\phi = \mathcal{F}_\phi \circ \mathcal{F}_0^{-1}|_{\mathcal{S}_0}$, for

each ϕ to be considered. We have tackled this problem with both the LS and TLS RBF methods described in chapter 3. In order to keep the number of variables to a minimum, in all the experiments described in this section we will obtain the necessary error curves by constructing RBF maps with $p = 200$ (and, in the case of TLS, also $q = 200$) repulsive centers, selected from a training set of $N = 2000$ points in \mathcal{S}_ϕ , uniformly distributed in θ , and using cubic basis functions. Unless otherwise specified, in order to reduce the dependency of these results on the particular sets of centers so obtained we will follow the procedure adopted in the previous chapter in actually plotting the means (and, where appropriate, the standard deviations) of the errors arising from a collection of 500 distinct sets of centers, chosen by random selection (without replacement) of the repulsive seed.

4.1.1 The circle and least squares

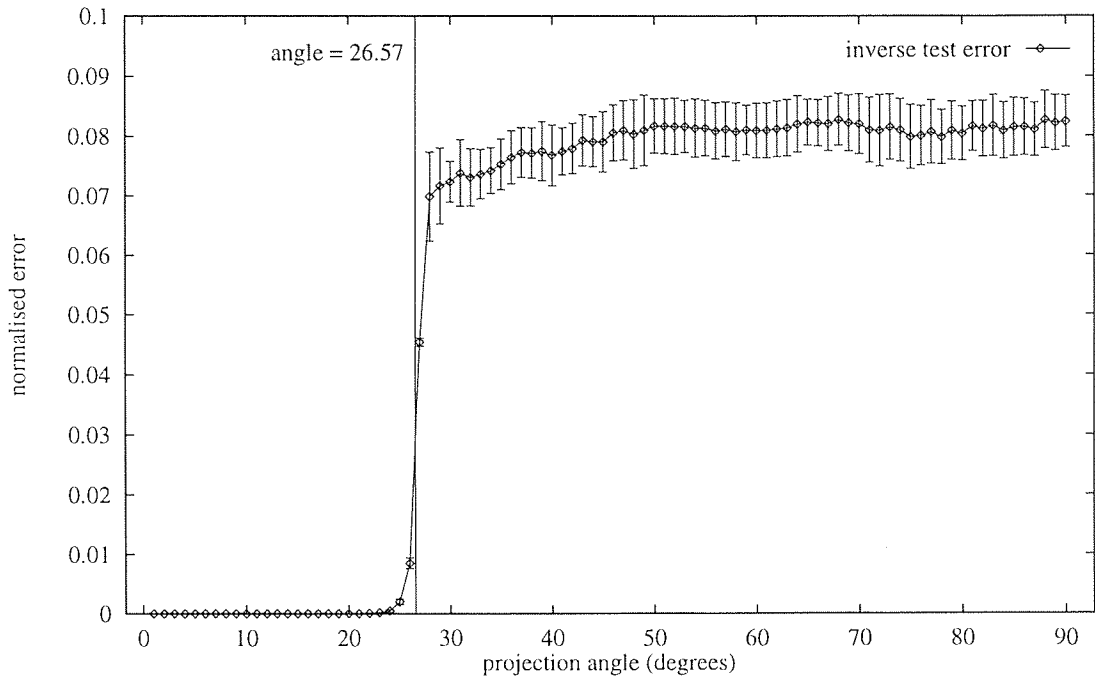
We will begin by applying the method of least squares to the problem, separately modelling the relationships $\mathcal{S}_0 \mapsto \mathcal{S}_\phi$ and $\mathcal{S}_\phi \mapsto \mathcal{S}_0$ with RBF approximations $\widehat{\mathbf{f}}_\phi: \mathcal{S}_0 \subset \mathbb{R}^2 \rightarrow \mathbb{R}^2$ and $\widehat{\mathbf{f}}_\phi^{-1}: \mathcal{S}_\phi \subset \mathbb{R}^2 \rightarrow \mathbb{R}^2$ (bearing in mind that for $\phi > \phi^*$ the label $\widehat{\mathbf{f}}_\phi^{-1}$ is merely a notational convenience). In order to determine whether or not the images $\widehat{\mathbf{f}}_\phi \mathcal{S}_0$ and $\widehat{\mathbf{f}}_\phi^{-1} \mathcal{S}_\phi$ are sufficiently close to their targets, respectively \mathcal{S}_ϕ and \mathcal{S}_0 , we adapt the notation introduced in chapter 3.2, defining a forward error function $\epsilon_0^{(\phi)} = \widehat{\mathbf{f}}_\phi - \mathbf{f}_\phi$, and a corresponding inverse error $\epsilon_\phi^{(0)} = \widehat{\mathbf{f}}_\phi^{-1} - \mathbf{f}_\phi^{-1}$, where the subscript on each function denotes its domain and the superscript denotes its co-domain. In analogy with equation (3.5), for $\mathbf{x} \in \mathcal{S}_0$ and $\mathbf{y} \in \mathcal{S}_\phi$ we write the corresponding normalised fitting (and, where appropriate, test) errors $\epsilon_0^{(\phi)}$ and $\epsilon_\phi^{(0)}$ as

$$\epsilon_0^{(\phi)2} = \sigma_\phi^{-2} \sum_{i=1}^N \|\epsilon_0^{(\phi)}(\mathbf{x}_i)\|^2, \quad \epsilon_\phi^{(0)2} = \sigma_0^{-2} \sum_{i=1}^N \|\epsilon_\phi^{(0)}(\mathbf{y}_i)\|^2 \quad (4.6)$$

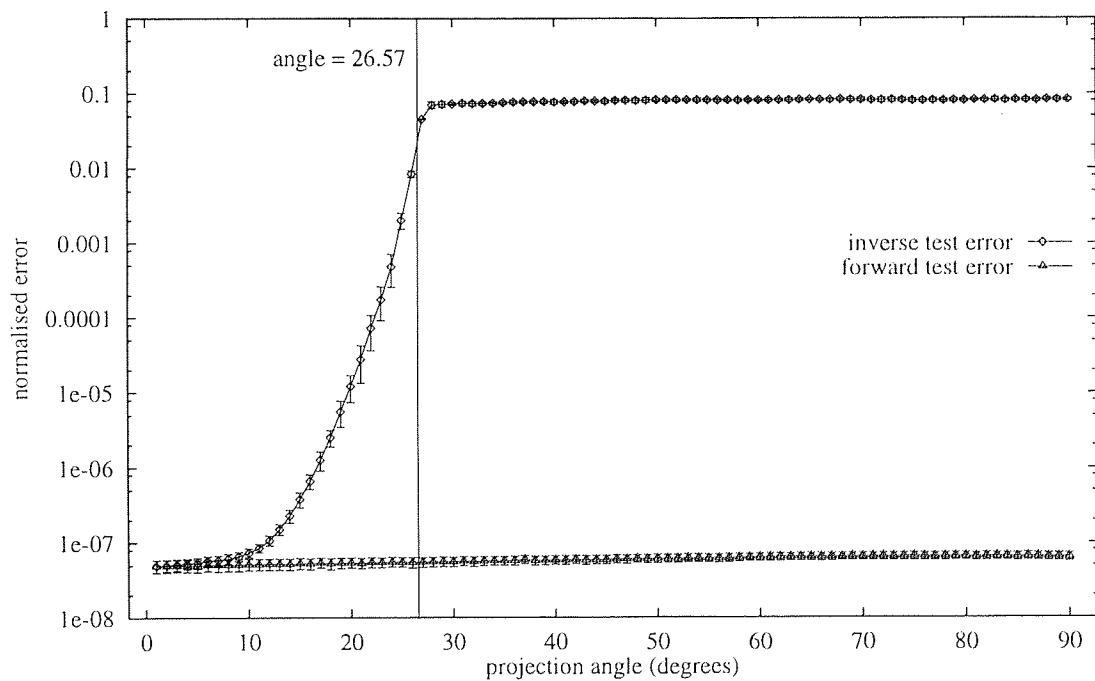
with normalisers σ_0 and σ_ϕ calculated on $\mathcal{S}_0 \subset \mathbb{R}^2$ and $\mathcal{S}_\phi \subset \mathbb{R}^2$ in the usual manner.

Since \mathbf{f}_ϕ is, by definition, single-valued for all ϕ we expect to find a good approximation $\widehat{\mathbf{f}}_\phi$, in terms of a small fitting/test error $\epsilon_0^{(\phi)}$, regardless of the value of ϕ . For values of $\phi < \phi^*$ we should also find a good approximation to its inverse, but beyond that critical value the presence of a self-intersecting set in \mathcal{S}_ϕ will force the RBF map to attempt to map points which are close together, in the region of the self-intersection, to points which are far apart in \mathcal{S}_0 . We therefore expect to see a large inverse error $\epsilon_\phi^{(0)}$ for $\phi \geq \phi^*$.

To verify these predictions we plot, in figure 4.2, the expected values of the normalised errors $\epsilon_0^{(\phi)}$ and $\epsilon_\phi^{(0)}$, calculated over the test set, as a function of ϕ in the range $0 \leq \phi \leq \frac{\pi}{2}$. (The ϕ -axis is calibrated in units of degrees for convenience.) The fitting errors are virtually indistinguishable from the test errors, even before averaging over RBF models, and are therefore not plotted. On a linear scale the forward error $\epsilon_0^{(\phi)}$ is undiscernible, being practically constant at $\epsilon_0^{(\phi)} \sim 10^{-7}$ for all choices of repulsive seed, so in part (a) we plot only the inverse error: it clearly bears out our expectations, rising sharply at around the



(a)



(b)

Figure 4.2 Comparing the mean, normalised test errors $\langle \epsilon_0^{(\phi)} \rangle$ and $\langle \epsilon_\phi^{(0)} \rangle$, versus ϕ , for the LS approximations \widehat{f}_ϕ and \widehat{f}_ϕ^{-1} on the projected circles \mathcal{S}_0 and \mathcal{S}_ϕ . (a) On a linear scale both $\epsilon_0^{(\phi)}$ and $\epsilon_\phi^{(0)}$ are negligible, on average, for $\phi < \phi^*$, but while the forward error (not plotted) remains so over the entire range of ϕ , the inverse error rises steeply at the critical value to saturate at $\langle \epsilon_\phi^{(0)} \rangle \approx 0.08$ for $\phi > \phi^*$; (b) a log-linear scale reveals a practically constant forward error, at $\epsilon_0^{(\phi)} \sim 10^{-7}$. On both scales the separation of fitting and test errors is virtually impossible to distinguish—even before averaging—so only the latter are plotted. Error bars denote one standard deviation in each direction.

critical value from the same floor of $\langle \epsilon_\phi^{(0)} \rangle \sim 10^{-7}$ to a ceiling of $\langle \epsilon_\phi^{(0)} \rangle \approx 0.08$ for $\phi \geq \phi^*$. The standard deviation of $\epsilon_\phi^{(0)}$ also grows significantly as ϕ increases beyond the critical value, but remains small in comparison to $\epsilon_\phi^{(0)}$ itself. In part (b) we adopt a log-linear scale, which also demonstrates the virtual independence on ϕ of the forward error $\epsilon_0^{(\phi)}$. That the inverse error increases gradually as it approaches ϕ^* from below, rather than making a discontinuous step at $\phi = \phi^*$, we ascribe to an inability of the LS RBF map to resolve separate points as they become increasingly closer together in \mathcal{S}_ϕ .

4.1.1.1 Analysis of the least squares solution

The results obtained from figure 4.2 appear to be satisfactory, in as much as they reveal—to a good approximation—the presence and extent of the interval $\phi > \phi^*$ within which f_ϕ fails to be a diffeomorphism. But this might not always be the case: as discussed in section 3.2.4, we feel that it could be dangerous to rely purely on the LS error measure $\epsilon_\phi^{(0)}$ to detect the presence of a self-intersection which may occur on a set of small, maybe vanishing measure in $\mathcal{S}_{\phi > \phi^*}$ and hence be effectively averaged out by the LS error minimisation algorithm. To better understand this process it will be instructive to examine the images of \mathcal{S}_0 and \mathcal{S}_ϕ under specific RBF maps \widehat{f}_ϕ and \widehat{f}_ϕ^{-1} directly, for values of ϕ on either side of the critical value. To this end we plot in figure 4.3 the approximating set $\widehat{f}_\phi^{-1}\mathcal{S}_\phi$, superimposed on the circle \mathcal{S}_0 , and also the set $\widehat{f}_\phi\mathcal{S}_0$, superimposed on the projected circle \mathcal{S}_ϕ , where \widehat{f}_ϕ and \widehat{f}_ϕ^{-1} are now obtained from a single set of repulsive centers, seeded so as to maximise $\|\mathbf{x}_i\|$ over the training set (as previously described in section 3.2.3). Parts (a) and (b) correspond to the case $\phi = 20$ degrees and parts (c) and (d) to $\phi = 40$ degrees.

The injectivity of f_ϕ is demonstrated in parts (b) and (d) of this figure by the close correspondence between the approximating set $\widehat{f}_\phi\mathcal{S}_0$ and the co-domain \mathcal{S}_ϕ of f_ϕ in each case. In part (a) we see an equally close correspondence between the set of $\widehat{f}_{20}^{-1}\mathcal{S}_{20}$ and its target \mathcal{S}_0 , confirming our conclusion that f_{20} is a diffeomorphism. In comparison, however, it is clear from part (c) that for $\phi = 40$ degrees we have been unable to find a good approximation \widehat{f}_{40}^{-1} to the one-to-many relationship $\mathcal{S}_{40} \mapsto \mathcal{S}_0$. As anticipated, although the source of this non-injectivity—a co-dimension-two self-intersecting subset of \mathcal{S}_{40} —is of vanishing measure, its effects are seen to be global in \mathcal{S}_0 , giving rise not only to a large oscillation in the vicinity of those points which map to the same point in \mathcal{S}_{40} , but also to a noticeably worse fit elsewhere in the figure, as the LS RBF map \widehat{f}_{40}^{-1} sacrifices the closeness of its fit over the rest of \mathcal{S}_0 in an attempt to reduce the errors in those affected regions.

In figure 4.4 we show the corresponding plots for the cases $\phi = 26$ and $\phi = 27$ degrees, closely straddling the critical value. Despite their proximity in ϕ these plots illustrate the effect, on the RBF map \widehat{f}_ϕ^{-1} , of the transition from the diffeomorphic relationship illustrated in parts (a) and (b) to the non-invertible one illustrated in parts (c) and (d). In particular, the breakdown in the injectivity of f_ϕ above its critical value results, in part (c), in an approximating set $\widehat{f}_{27}^{-1}\mathcal{S}_{27}$ which is clearly distinct from its target

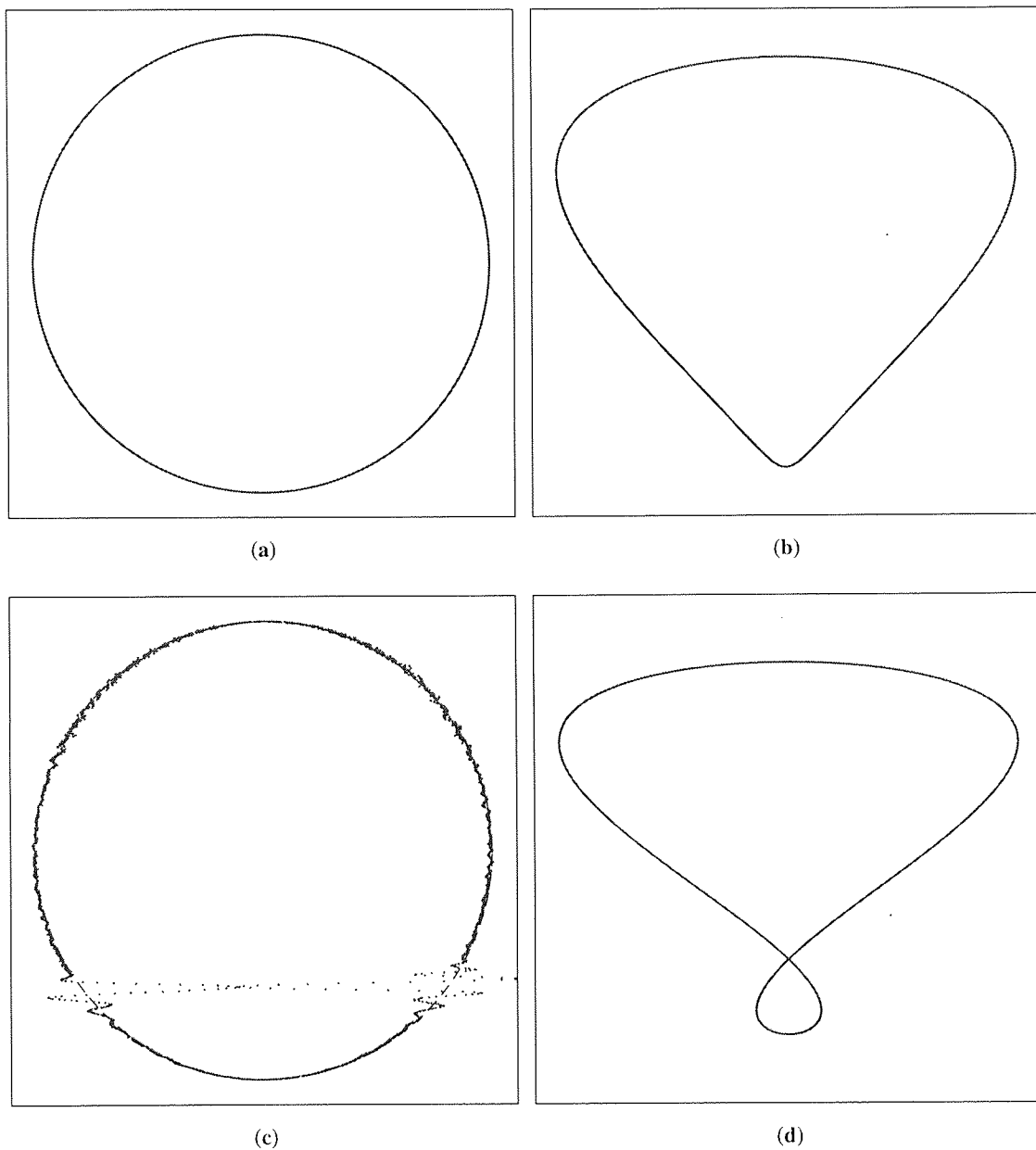


Figure 4.3 Approximating the map $f_\phi: S_0 \rightarrow S_\phi$ and its inverse, for $\phi = 20, 40$ degrees. Part (a) shows the image $\widehat{f_{20}^{-1}}S_{20}$, visually indistinguishable from S_0 (dashed), and in part (b), $\widehat{f_{20}}S_0$ is similarly close to S_{20} , in correspondence with our foreknowledge that f_{20} is a diffeomorphism. Conversely, in part (c) we see the image of S_{40} under $\widehat{f_{40}^{-1}}$ diverge markedly from S_0 in the region of those points which map together under f_{40} , while in (d) $\widehat{f_{40}}S_0$ sits neatly on top of S_{40} , in a vivid demonstration of the result of attempting to approximate the inverse of a many-to-one map.

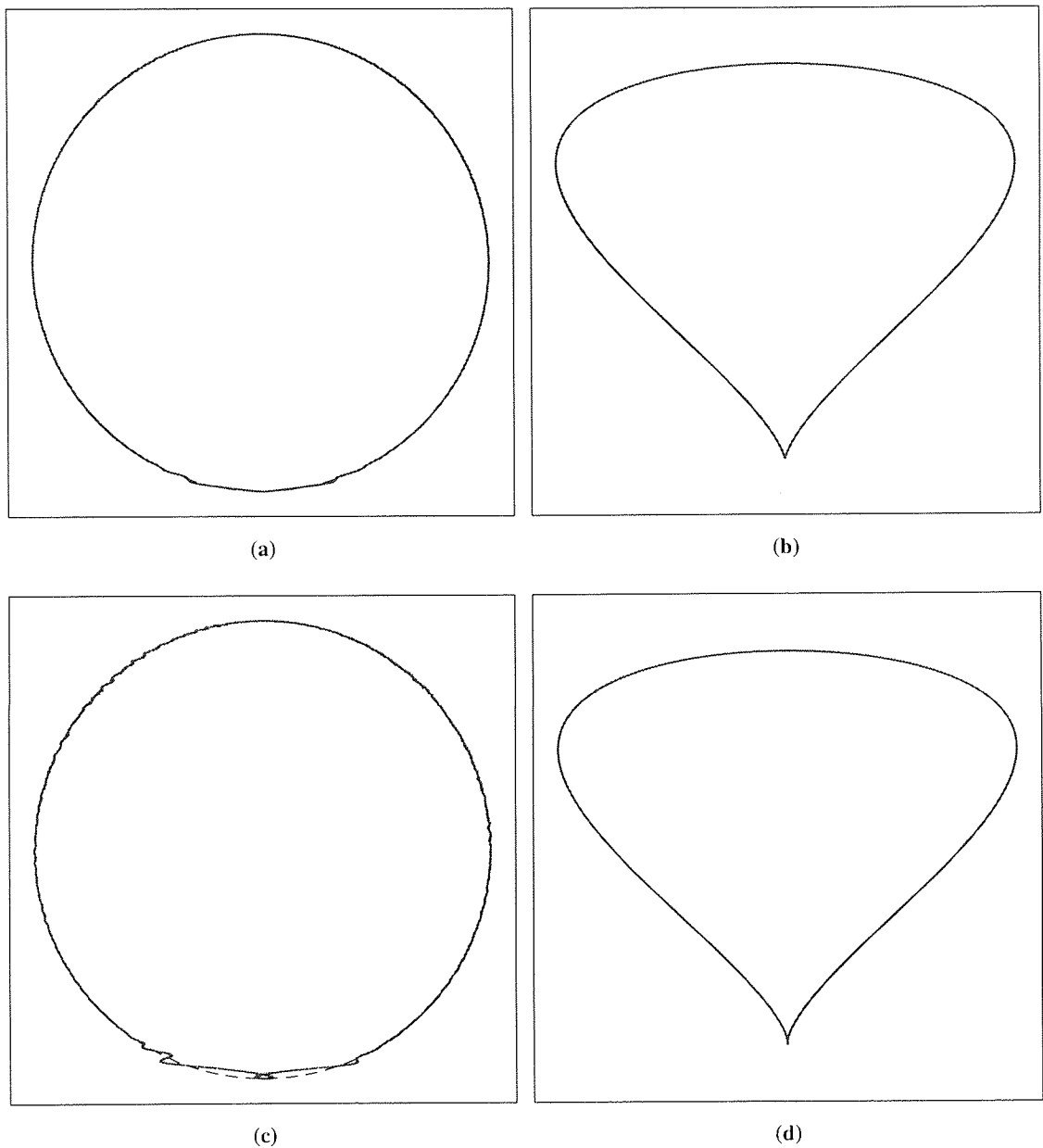


Figure 4.4 Approximating the map $f_\phi: \mathcal{S}_0 \rightarrow \mathcal{S}_\phi$ and its inverse, for $\phi = 26, 27$ degrees. Despite the proximity of ϕ to the critical value, part (a), showing the image $\widehat{f_{26}^{-1}} \mathcal{S}_{26}$ of \mathcal{S}_{26} , reveals an almost negligible difference from \mathcal{S}_0 (dashed) near the bottom of the plot, while in part (b) the image $\widehat{f_{26}} \mathcal{S}_0$ is indistinguishable from \mathcal{S}_{26} , as expected. In part (c), however, the image $\widehat{f_{27}^{-1}} \mathcal{S}_{27}$ is just beginning to diverge from \mathcal{S}_0 , corresponding to the onset of self-intersection in \mathcal{S}_{27} , which is again indistinguishable from $\widehat{f_{27}} \mathcal{S}_0$ in (d), as expected.

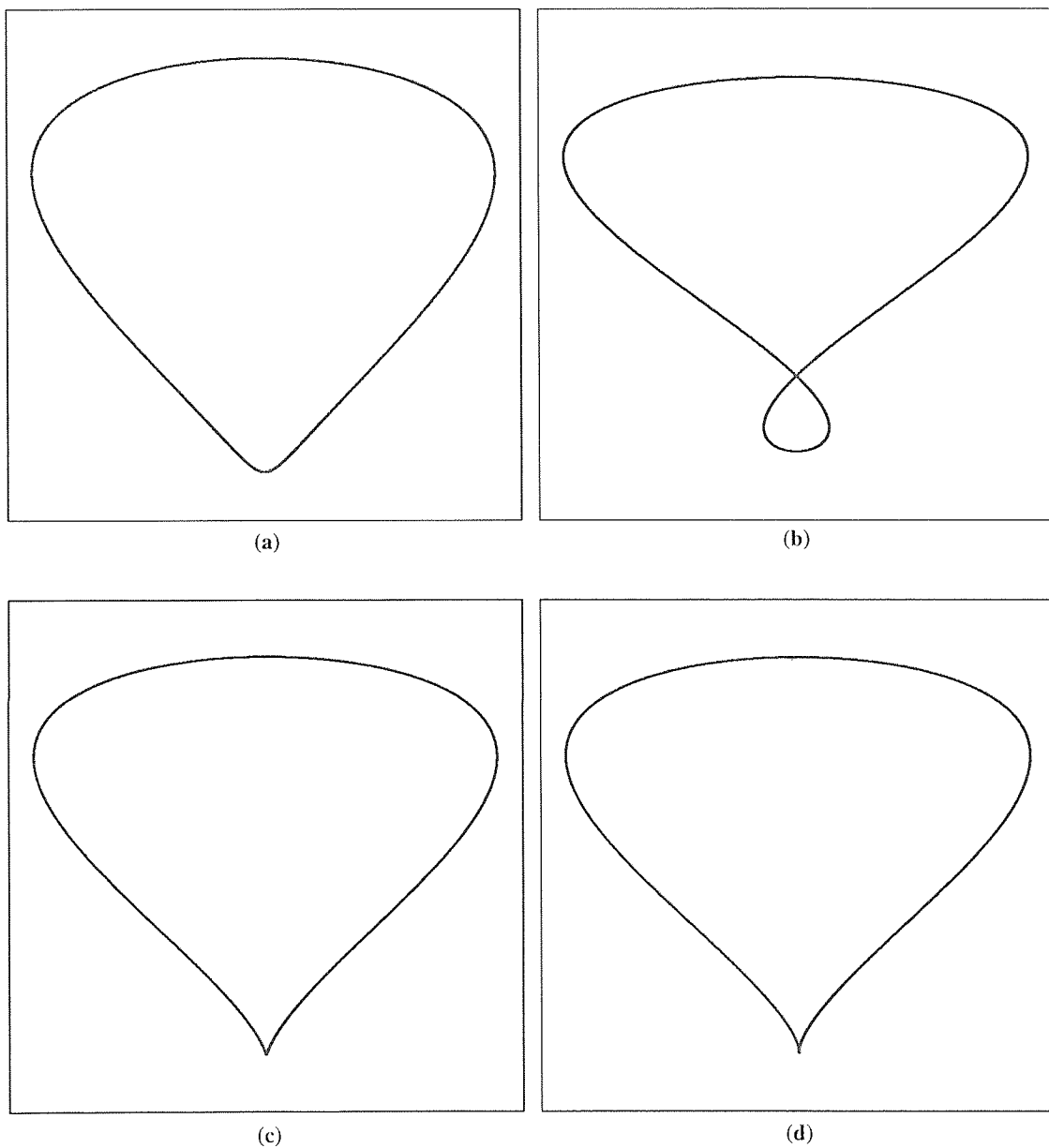


Figure 4.5 Colour-coding the projected circles \mathcal{S}_ϕ , for $\phi = 20, 26, 27, 40$ degrees, by the relative magnitudes of the per-point errors $\epsilon_\phi^{(0)}(\mathbf{y})$ to which they give rise under \widehat{f}_ϕ^{-1} . Part (a), with $\phi = 20$ degrees, does not exhibit a significant amount of deviation in colour-coding, which is to be expected below the critical value; part (b), on the other hand, with $\phi = 40$ degrees, reveals a strong, localised peak in error near the self-intersecting set, and parts (c) and (d), for $\phi = 26$ and 27 degrees, respectively, also show a significant concentration of large errors near the cusp.

\mathcal{S}_0 in the affected region. It is worth noting, in part (a), that the image of \mathcal{S}_{26} under $\widehat{f_{26}^{-1}}$ also exhibits a slight ‘ripple’ in this region, even though f_{26} is known to be a diffeomorphism. As previously remarked, we ascribe this behaviour to the finite dimensionality of the basis set on which these RBF maps have been defined.

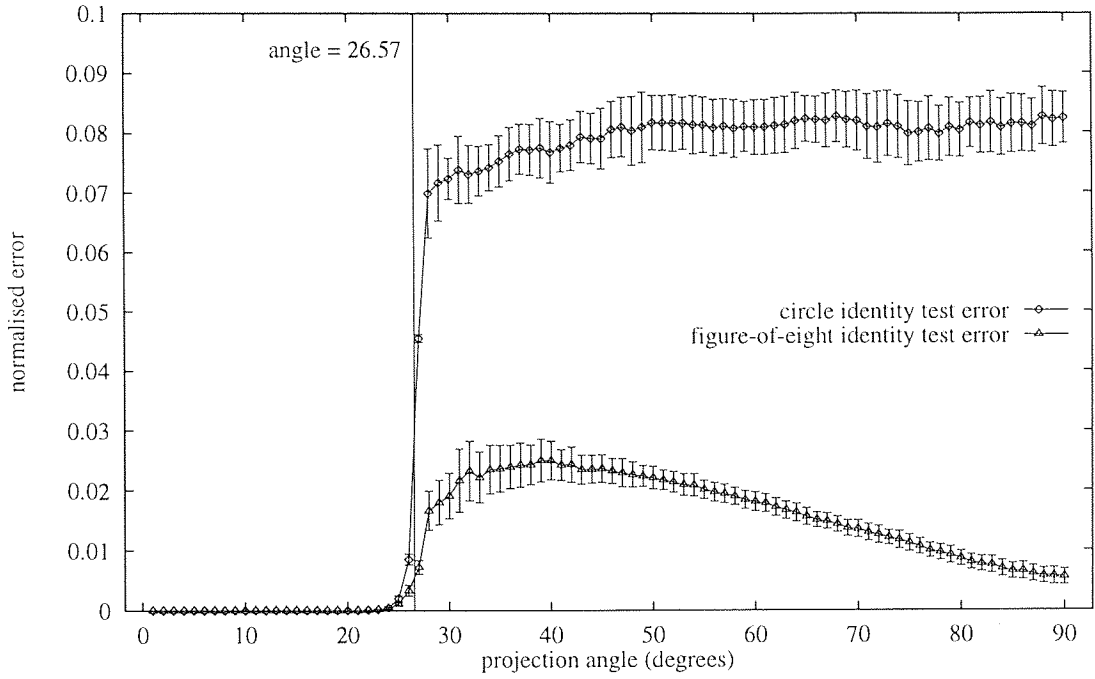
Another useful way in which to view the onset of non-injectivity in f_ϕ is illustrated in figure 4.5, in which we replot the projections \mathcal{S}_ϕ for $\phi = 20, 26, 27$ and 40 degrees. In this case, however, we colour-code the individual points $\mathbf{y} \in \mathcal{S}_\phi$ to reflect the per-point error magnitudes $\|\epsilon_\phi^{(0)}(\mathbf{y})\|$ to which they give rise under $\widehat{f_\phi^{-1}}$. These scalars index a linear colour map traversing a diagonal of the RGB colour cube, from blue to red, scaled so that points attaining a minimum error are plotted in blue and those attaining a maximum are plotted in red. (This scaling has been carried out individually for each RBF map, with the result that colour-codes cannot easily be compared between distinct plots; it should also be borne in mind that some points may be obscured by others, as they have been plotted according to the natural ordering of the data sets involved.) The colour-coding by $\epsilon_{20}^{(0)}$ of \mathcal{S}_{20} , in part (a) of this figure, is difficult to discern, but it is possible to make out a region near the bottom of the curve in which the error magnitudes are slightly larger than elsewhere. In part (b), however, which illustrates the errors arising from an attempt to fit the one-to-many relationship $\mathcal{S}_{40} \mapsto \mathcal{S}_0$ with the RBF map $\widehat{f_{40}^{-1}}$, we see a very definite, and strongly localised peak in colour-coding near the region of self-intersection in \mathcal{S}_{40} , confirming our analysis of the corresponding plot in figure 4.3(c). We see a similar effect in \mathcal{S}_{26} and \mathcal{S}_{27} , plotted in parts (c) and (d), respectively, with the peak in error magnitudes now concentrated in the vicinity of the cusp in each case.

4.1.1.2 Approximating the identity map

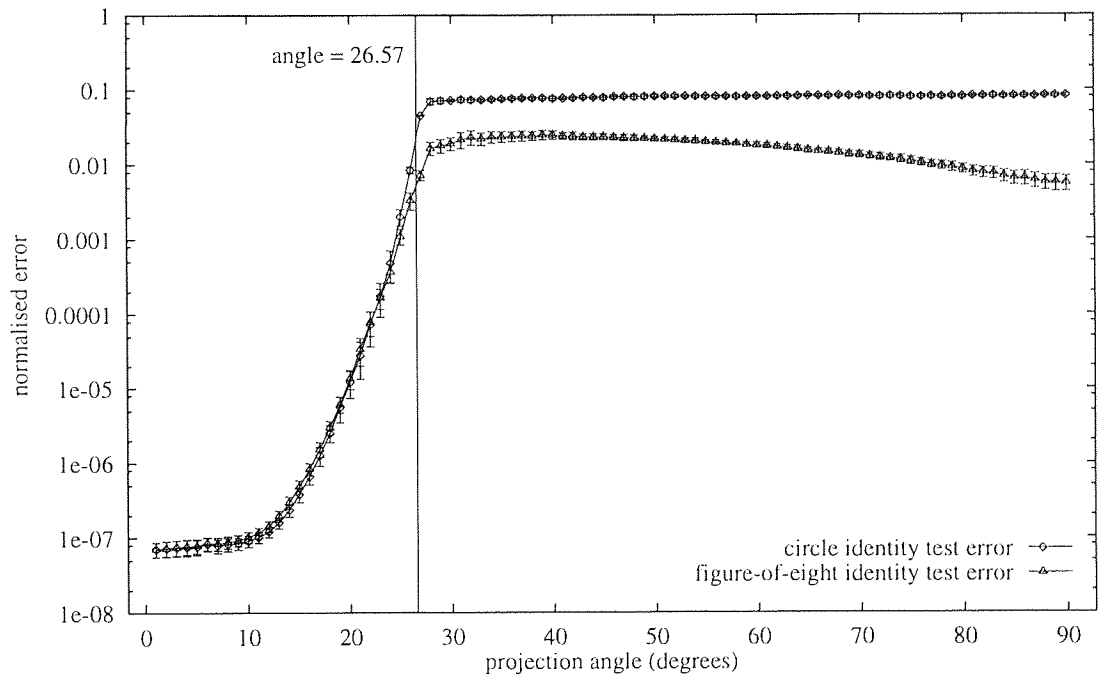
We would now like to perform the local error analysis, described in section 3.2.4, and attempt to establish the existence of Lipschitz constants U_ϕ and L_ϕ^{-1} for f_ϕ and its inverse (if it exists) by calculating experimental upper and lower bounds \widehat{U}_ϕ and \widehat{L}_ϕ for the growth of errors under $\widehat{f_\phi}$. In order to carry out this procedure we must first make the definitions $\widehat{I}_0^{(\phi)} = \widehat{f_\phi^{-1}} \circ \widehat{f_\phi}$ and $\widehat{I}_\phi^{(0)} = \widehat{f_\phi} \circ \widehat{f_\phi^{-1}}$, where the subscripts denote the domain as usual, but the superscripts now indicate the ‘route’ taken by the composition. These maps approximate the identity maps $I_\phi: \mathcal{S}_\phi \rightarrow \mathcal{S}_\phi$, and give rise to the identity error functions $\eta_0^{(\phi)} = \widehat{I}_0^{(\phi)} - I_0$ and $\eta_\phi^{(0)} = \widehat{I}_\phi^{(0)} - I_\phi$. We also define normalised identity errors $\eta_0^{(\phi)}$ and $\eta_\phi^{(0)}$, following (3.22) and (3.23), with

$$\eta_0^{(\phi)2} = \sigma_0^{-2} \sum_{i=1}^N \|\eta_0^{(\phi)}(\mathbf{x}_i)\|^2, \quad \eta_\phi^{(0)2} = \sigma_\phi^{-2} \sum_{i=1}^N \|\eta_\phi^{(0)}(\mathbf{y}_i)\|^2 \quad (4.7)$$

The expected values of these errors, calculated over the test set, are plotted in figure 4.6, on a linear scale in part (a) and a log-linear scale in part (b). The first thing we notice in this figure is that the mean identity error calculated on \mathcal{S}_0 is almost identical to the mean inverse error calculated on \mathcal{S}_ϕ , as plotted in figure



(a)



(b)

Figure 4.6 Comparing the mean, normalised test set errors $\langle \eta_0^{(\phi)} \rangle$ and $\langle \eta_\phi^{(0)} \rangle$, versus ϕ , arising from the LS identity approximations $\tilde{I}_0^{(\phi)}$ and $\tilde{I}_\phi^{(0)}$ on the projected circles \mathcal{S}_0 and \mathcal{S}_ϕ . (a) On a linear scale the identity error $\eta_0^{(\phi)}$, calculated on \mathcal{S}_0 , is virtually identical, on average, to the inverse error $\epsilon_\phi^{(0)}$ calculated on \mathcal{S}_ϕ , while the fact that $\langle \eta_\phi^{(0)} \rangle$ is consistently smaller than $\langle \eta_0^{(\phi)} \rangle$ indicates that \widehat{f}_ϕ is contractive in the region of $\widehat{f}_\phi^{-1} \mathcal{S}_\phi$; (b) the same plot, on a log-linear scale, is included for the sake of completeness. Error bars denote one standard deviation in each direction.

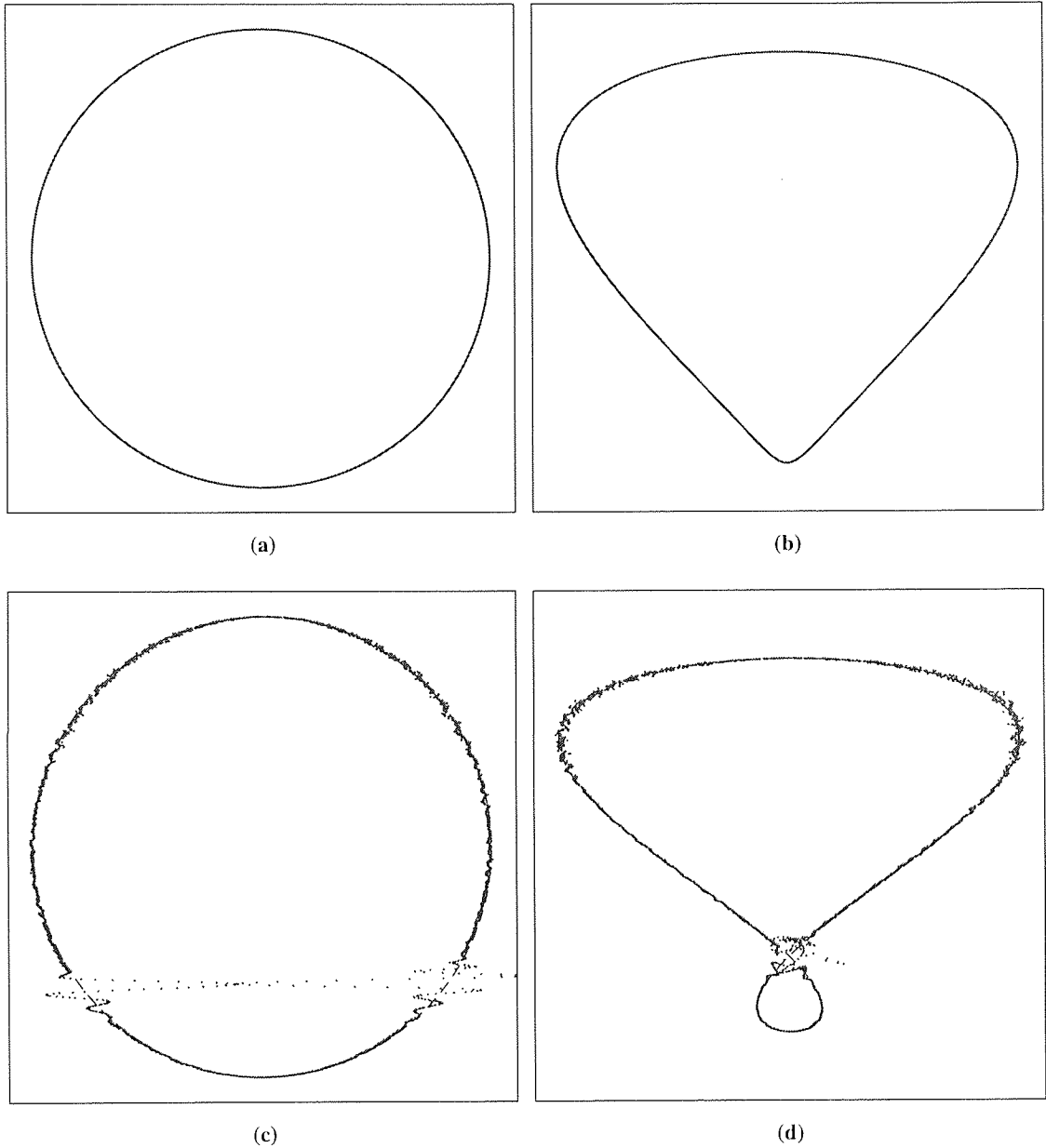


Figure 4.7 Approximating the identity maps between \mathcal{S}_0 and \mathcal{S}_ϕ , for $\phi = 20, 40$ degrees. Part (a) shows that $\tilde{\mathcal{I}}_0^{(20)}\mathcal{S}_0$ is a good fit to \mathcal{S}_0 (dashed), and part (b) shows that $\tilde{\mathcal{I}}_{20}^{(0)}\mathcal{S}_{20}$ is a similarly good fit to \mathcal{S}_{20} , justifying the claim that $f_{20}: \mathcal{S}_0 \rightarrow \mathcal{S}_{20}$ is a diffeomorphism. In contrast, in part (c), large errors in $\tilde{\mathcal{I}}_0^{(40)}\mathcal{S}_0$ are clearly visible with respect to \mathcal{S}_0 , in those regions which map through the self-intersection in \mathcal{S}_{40} and, in part (d), errors in $\tilde{\mathcal{I}}_{40}^{(0)}\mathcal{S}_{40}$ are also visible with respect to \mathcal{S}_{40} , leading us to the desired conclusion that $f_{40}: \mathcal{S}_0 \rightarrow \mathcal{S}_{40}$ is not a diffeomorphism. It is interesting to note that the relatively small errors in (d) indicate that f_{40} is contractive in the neighbourhood of those points in \mathcal{S}_0 which map together under f_{40} .

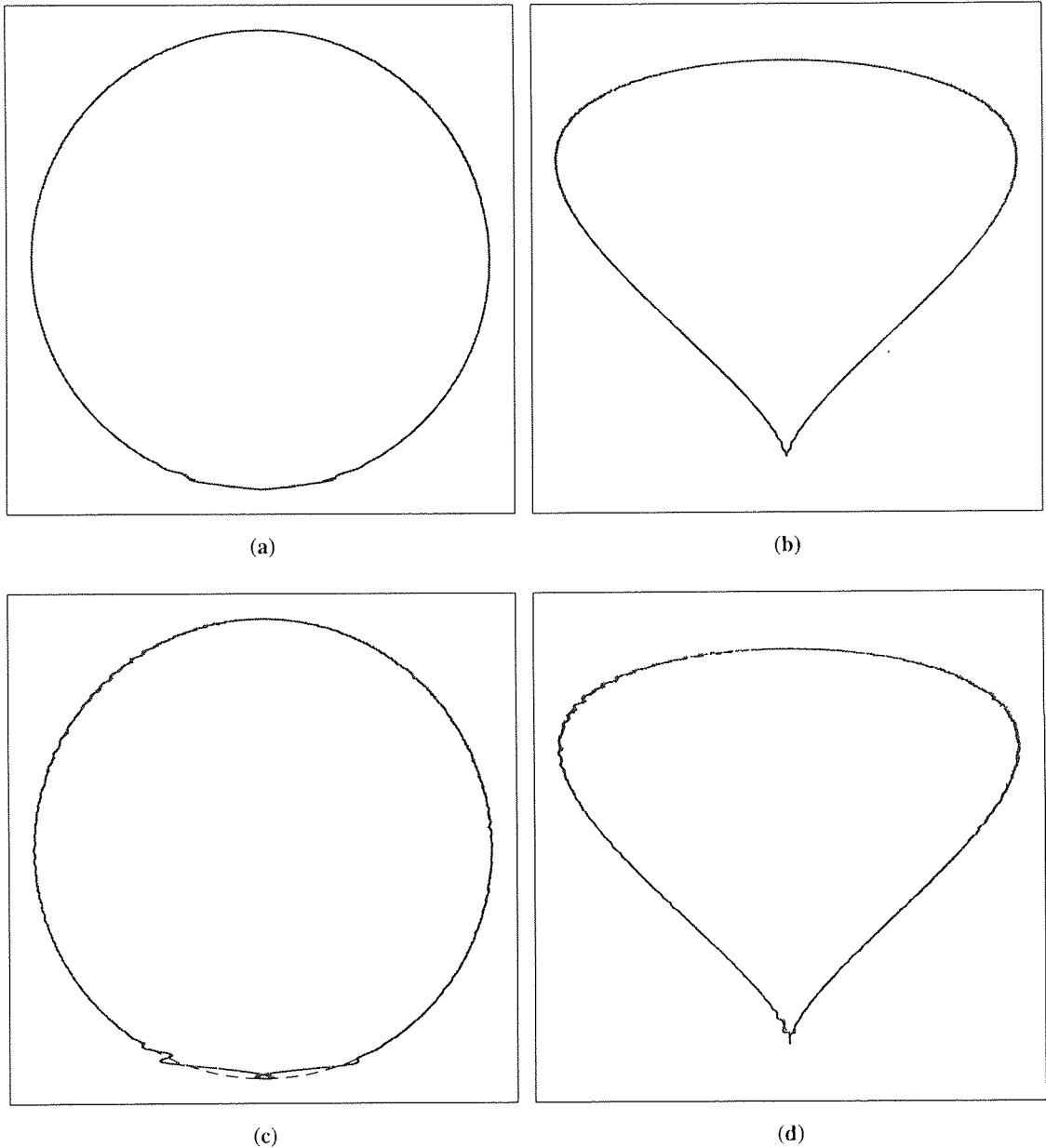


Figure 4.8 Approximating the identity maps between S_0 and S_ϕ , for $\phi = 26, 27$ degrees. Part (a) plots the image $\tilde{I}_0^{(26)} S_0$ of S_0 on S_0 itself (dashed), revealing negligible errors, and part (b) plots $\tilde{I}_{26}^{(0)} S_{26}$ over S_{26} , again with very small errors. As expected, this leads us to conclude that f_{26} is a diffeomorphism, although the RBF approximations in question are clearly beginning to break down due to the proximity of ϕ to ϕ^* . In parts (c), showing $\tilde{I}_0^{(27)} S_0$ plotted over S_0 , and (d), showing $\tilde{I}_{27}^{(0)} S_{27}$ over S_{27} , the errors are noticeably larger, as we expect, since f_{27} just fails to be a diffeomorphism.

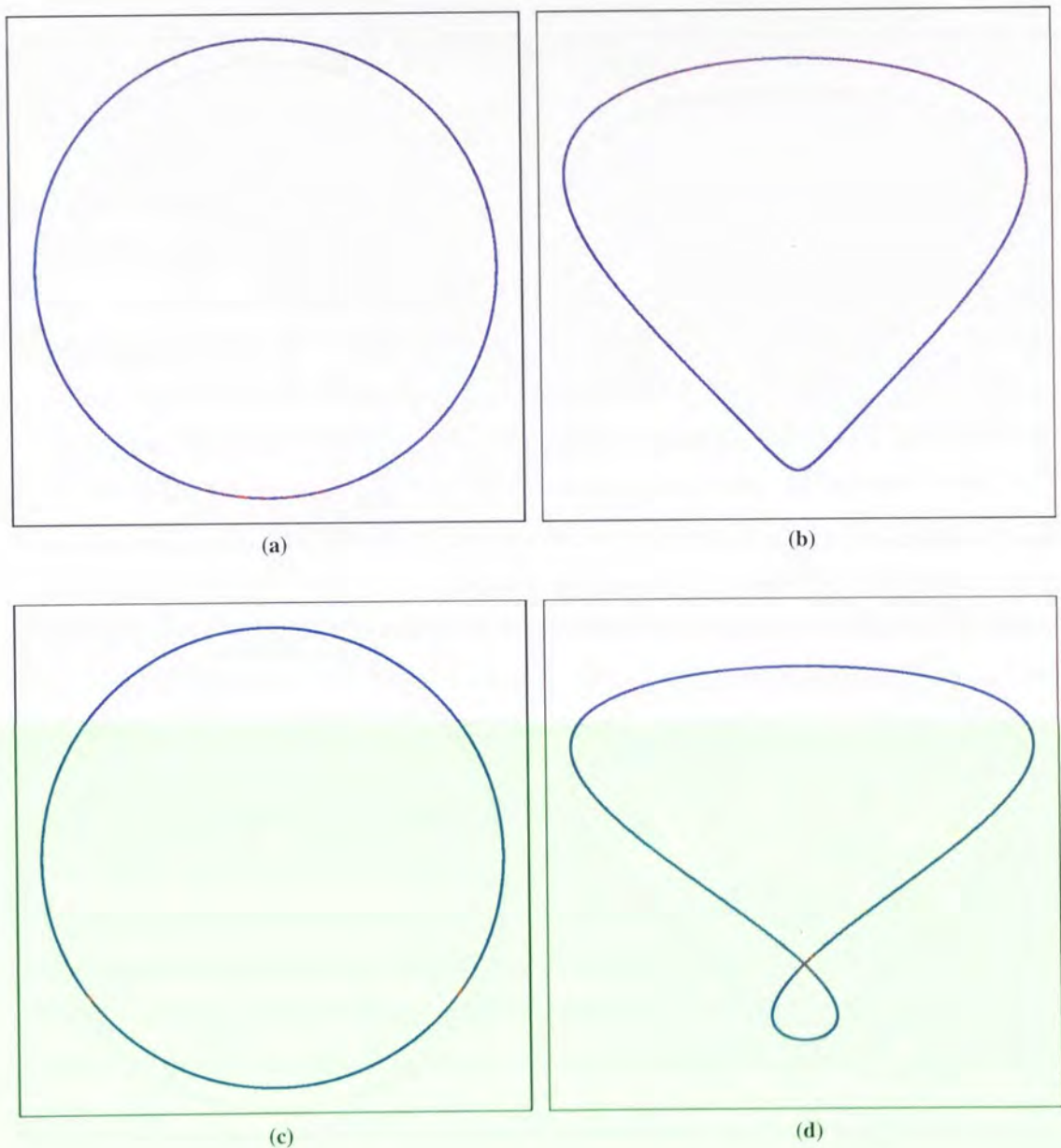


Figure 4.9 Colour-coding the identity maps between \mathcal{S}_0 and \mathcal{S}_ϕ , for $\phi = 20$ and 40 degrees, by the per-point error magnitudes $\|\boldsymbol{\eta}_0^{(\phi)}(\mathbf{x})\|$ and $\|\boldsymbol{\eta}_\phi^{(0)}(\mathbf{y})\|$. (a) At $\phi = 20$ degrees we see a significant localisation of errors near those points in \mathcal{S}_{20} which are mapped together under $\widehat{\mathbf{f}}_{20}$, resulting from the action of $\widehat{\mathbf{f}}_{20}^{-1}$; (b) in the other direction we see little evidence of localisation. (c) For $\phi = 40$ degrees there is a strong localisation of errors at those points in \mathcal{S}_0 which are mapped through the self-intersection in \mathcal{S}_{40} by $\widehat{\mathbf{I}}_0^{(40)}$; (d) in \mathcal{S}_{40} we see a somewhat broader peak, shifted slightly away from the self-intersection by the compensatory action of $\widehat{\mathbf{f}}_{40}$.

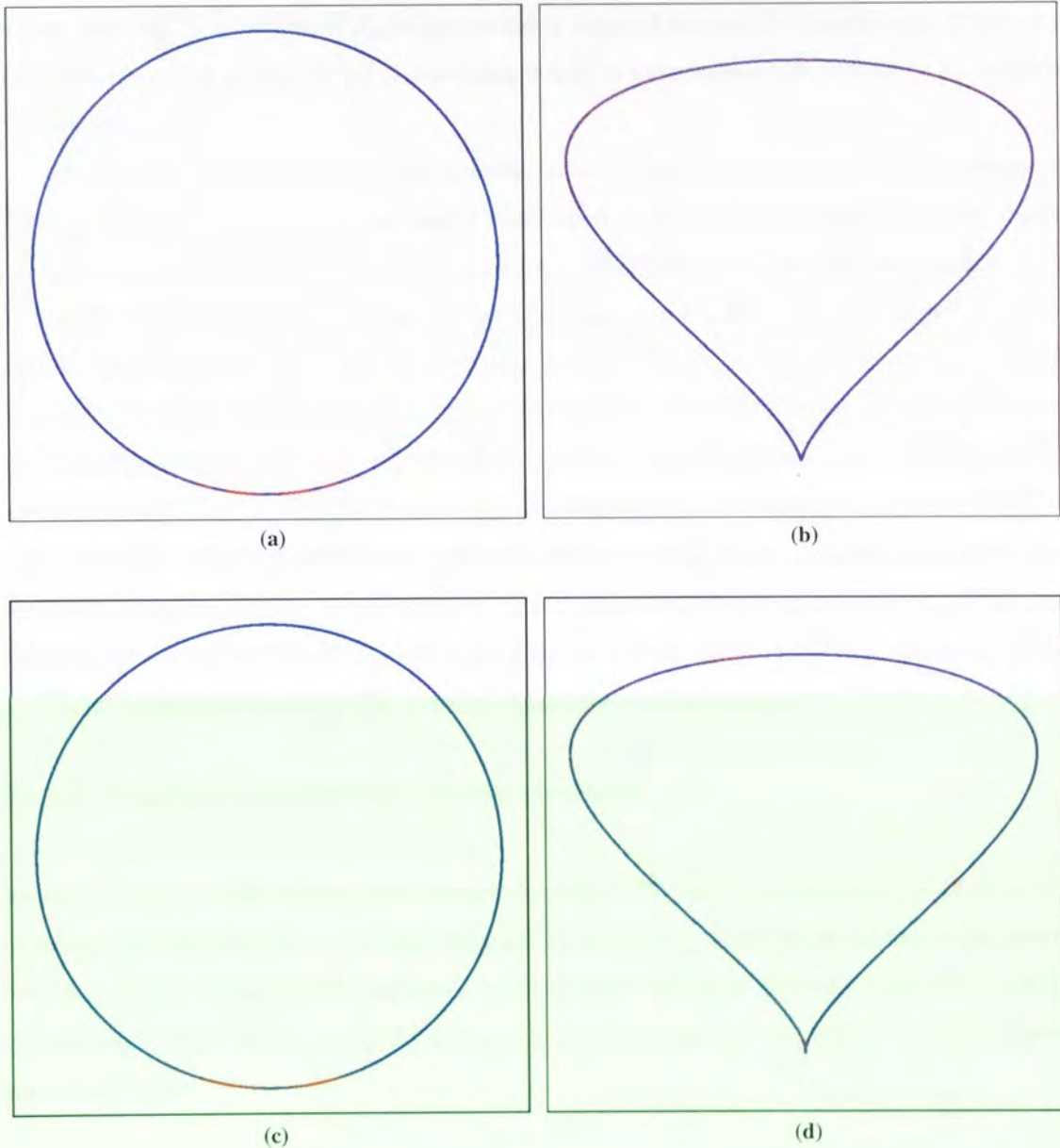


Figure 4.10 Colour-coding the identity maps between \mathcal{S}_0 and \mathcal{S}_ϕ , for $\phi = 26$ and 27 degrees, by the per-point error magnitudes $\|\boldsymbol{\eta}_0^{(\phi)}(\boldsymbol{x})\|$ and $\|\boldsymbol{\eta}_\phi^{(0)}(\boldsymbol{y})\|$. (a) At $\phi = 26$ degrees the errors are firmly concentrated near the bottom of \mathcal{S}_0 after being mapped through the cusp in \mathcal{S}_{26} ; (b) in \mathcal{S}_{26} itself, the error distribution has been spread out under the action of $\widehat{\boldsymbol{f}}_\phi$. (c) A similar effect is seen for $\phi = 27$ degrees, made a little more noticeable by transition through the critical value; (d) once again, there is no easily discernible concentration of errors in \mathcal{S}_{27} .

4.2; so similar, in fact, that we do not bother to replot the latter curve here for comparison. This observation is merely a restatement of the relationship $\eta_0^{(\phi)} \approx \epsilon_\phi^{(0)} \circ f_\phi$ which follows from equation (3.24), and holds even before the calculation of expectation values. It confirms that \widehat{f}_ϕ is an excellent approximation to f_ϕ and hence serves to validate the application of the Lipschitz analysis to this experiment. The mean identity error $\langle \eta_\phi^{(0)} \rangle$, calculated in \mathcal{S}_ϕ , is also plotted in figures 4.6(a) and (b). Interestingly, this error is consistently less than or equal to $\langle \eta_0^{(\phi)} \rangle$, indicating that \widehat{f}_ϕ is a *contractive* map, at least on the image of \mathcal{S}_ϕ under \widehat{f}_ϕ^{-1} .

For the sake of completeness we will examine some of these fits in more detail by superimposing $\widehat{I}_0^{(\phi)} \mathcal{S}_0$ on \mathcal{S}_0 and $\widehat{I}_\phi^{(0)} \mathcal{S}_\phi$ on \mathcal{S}_ϕ , again using a single set of repulsive centers. Figures 4.7 and 4.8 show examples of approximating these identity maps for $\phi = 20, 40, 26$ and 27 degrees, respectively. These figures clearly differ from their analogues of figures 4.3 and 4.4 in that for $\phi > \phi^*$ we can now see errors induced by both $\widehat{I}_0^{(\phi)}$ and $\widehat{I}_\phi^{(0)}$, arising from the poor fit of \widehat{f}_ϕ^{-1} to f_ϕ^{-1} in this parameter range. In figures 4.9 and 4.10 we colour-code each point in \mathcal{S}_0 and \mathcal{S}_ϕ by the magnitude of the identity error $\eta_0^{(\phi)}(\mathbf{x})$ and $\eta_\phi^{(0)}(\mathbf{y})$ to which it gives rise under $\widehat{I}_0^{(\phi)}$ and $\widehat{I}_\phi^{(0)}$, respectively, as in figure 4.5. Due to the action of \widehat{f}_ϕ^{-1} we see, in parts (a) and (c) of both figures, a significant concentration of large per-point errors $\eta_0^{(\phi)}(\mathbf{x})$ around those elements of \mathcal{S}_0 which are mapped close together under \widehat{f}_ϕ . In parts (b) and (d) of each figure, the distribution of errors appears to be more evenly spread across \mathcal{S}_ϕ than was the case in figure 4.5; this is due, for $\phi > \phi^*$, to the fact that \widehat{f}_ϕ , on composition with an approximation \widehat{f}_ϕ^{-1} to a one-to-many map f_ϕ^{-1} , is being applied to a subset of \mathbb{R}^2 sparsely represented in its training set.

4.1.1.3 Analytical calculation of Lipschitz constants

Before making use of these identity errors for the calculation of \widehat{U}_ϕ and \widehat{L}_ϕ , estimating U_ϕ and L_ϕ , we take advantage of the fact that we can—in this example—write down f_ϕ explicitly, which means that we can calculate U_ϕ and L_ϕ analytically. (In general, of course, this will not be the case.) Thus, if we consider the points $z_1, z_2 \in \mathcal{S}$, where $z_1 = \Phi(\theta_1)$ and $z_2 = \Phi(\theta_2)$ then, since $f_\phi = \mathcal{F}_\phi \circ \mathcal{F}_0^{-1}|_{\mathcal{S}_0}$, we can express equation (3.21) as

$$U_\phi = \max_{z_1, z_2 \in \mathcal{S}} \frac{\|\mathcal{F}_\phi(z_1) - \mathcal{F}_\phi(z_2)\|}{\|\mathcal{F}_0(z_1) - \mathcal{F}_0(z_2)\|}, \quad L_\phi = \min_{z_1, z_2 \in \mathcal{S}} \frac{\|\mathcal{F}_\phi(z_1) - \mathcal{F}_\phi(z_2)\|}{\|\mathcal{F}_0(z_1) - \mathcal{F}_0(z_2)\|} \quad (4.8)$$

If we now make the substitutions $\alpha = \frac{1}{2}(\theta_1 + \theta_2)$ and $\beta = \frac{1}{2}(\theta_1 - \theta_2)$, and the definition

$$R_\phi^2(\alpha, \beta) = (\cos \phi \cos \alpha + 2 \sin \phi \cos 2\alpha \cos \beta)^2 + \sin^2 \alpha \quad (4.9)$$

then it can be shown [2], via the standard trigonometric identities, that

$$U_\phi = \max_{\alpha, \beta} R_\phi(\alpha, \beta), \quad L_\phi = \min_{\alpha, \beta} R_\phi(\alpha, \beta) \quad (4.10)$$

Finally, with the further definitions

$$R_\phi^{(\pm)^2}(\alpha) = (\cos \phi \cos \alpha \pm 2 \sin \phi \cos 2\alpha)^2 + \sin^2 \alpha \quad (4.11)$$

and a little more algebraic manipulation we arrive, from geometrical considerations, at

$$U_\phi = \max_\alpha \begin{cases} R_\phi^{(-)}(\alpha), & \text{if } \cos 2\alpha \leq 0; \\ R_\phi^{(+)}(\alpha), & \text{if } \cos 2\alpha \geq 0. \end{cases} \quad (4.12)$$

and

$$L_\phi = \min_\alpha \begin{cases} R_\phi^{(+)}(\alpha), & \text{if } \cos 2\alpha \leq 0; \\ R_\phi^{(-)}(\alpha), & \text{if } \cos 2\alpha \geq 0. \end{cases} \quad (4.13)$$

Empirical estimates of U_ϕ and L_ϕ , obtained by varying α in a numerical simulation of (4.12) and (4.13) over the range $0 \leq \alpha < 2\pi$, in steps of $\frac{\pi}{500}$, are plotted in figure 4.11(a), for the full range of ϕ under investigation. This figure provides a graphic illustration of the breakdown in injectivity suffered by f_ϕ at the critical angle ϕ^* , as the lower bound L_ϕ drops almost linearly from $L_0 = 1$ towards $L_\phi = 0$ as ϕ approaches ϕ^* from above (note that L_ϕ is not defined for $\phi \geq \phi^*$). The fact that f_ϕ remains single-valued over the entire range is similarly illustrated by the U_ϕ curve, which peaks at the relatively low value of $U_\phi \approx 2$ at $\phi \approx \frac{\pi}{3}$; the apparent discontinuity shortly above this value is simply a consequence of equation (4.12). That both curves meet at $L_0 = U_0 = 1$ is, of course, attributable to f_0 being the identity map.

4.1.1.4 Approximation of the Lipschitz constants

We are finally ready to find upper and lower bounds \widehat{U}_ϕ and \widehat{L}_ϕ for the ratio of per-point error magnitudes $\|\epsilon_0^{(\phi)}(x)\|$ and $\|\epsilon_\phi^{(0)}(y)\|$ over all pairs x, y such that $y = f_\phi(x)$. We therefore rewrite equations (3.25) and (3.26), relating $\eta_0^{(\phi)}(x)$ and $\eta_\phi^{(0)}(y)$ with

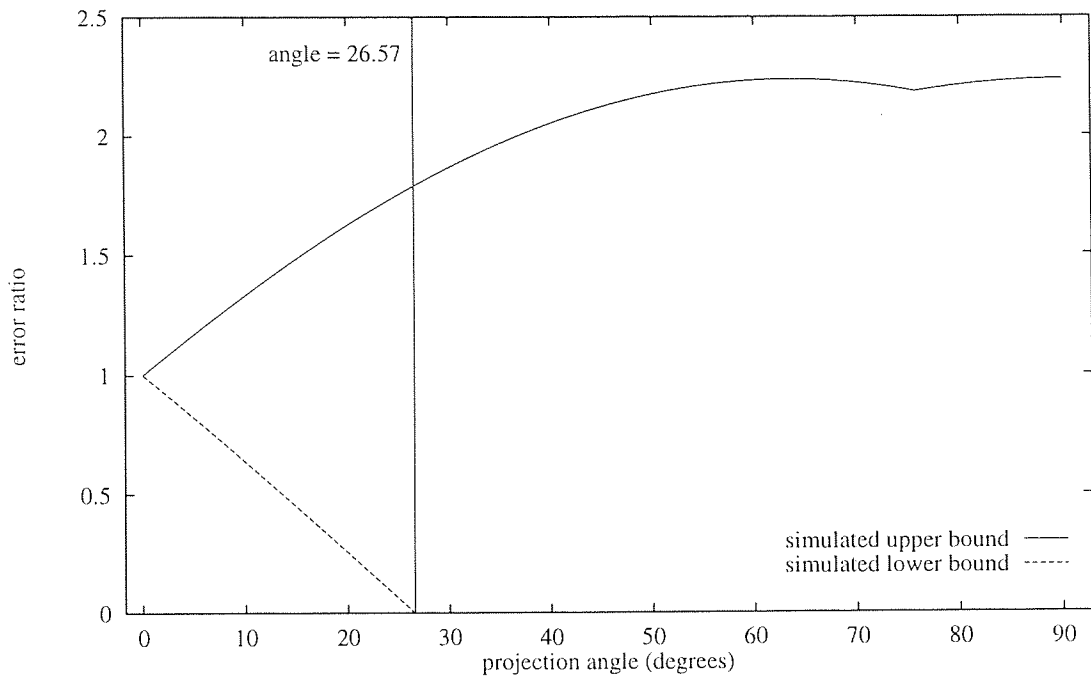
$$\eta_\phi^{(0)}(y) \approx \widehat{f}_\phi(x + \eta_0^{(\phi)}(x)) - \widehat{f}_\phi(x) \quad (4.14)$$

and, assuming that \widehat{f}_ϕ^{-1} exists,

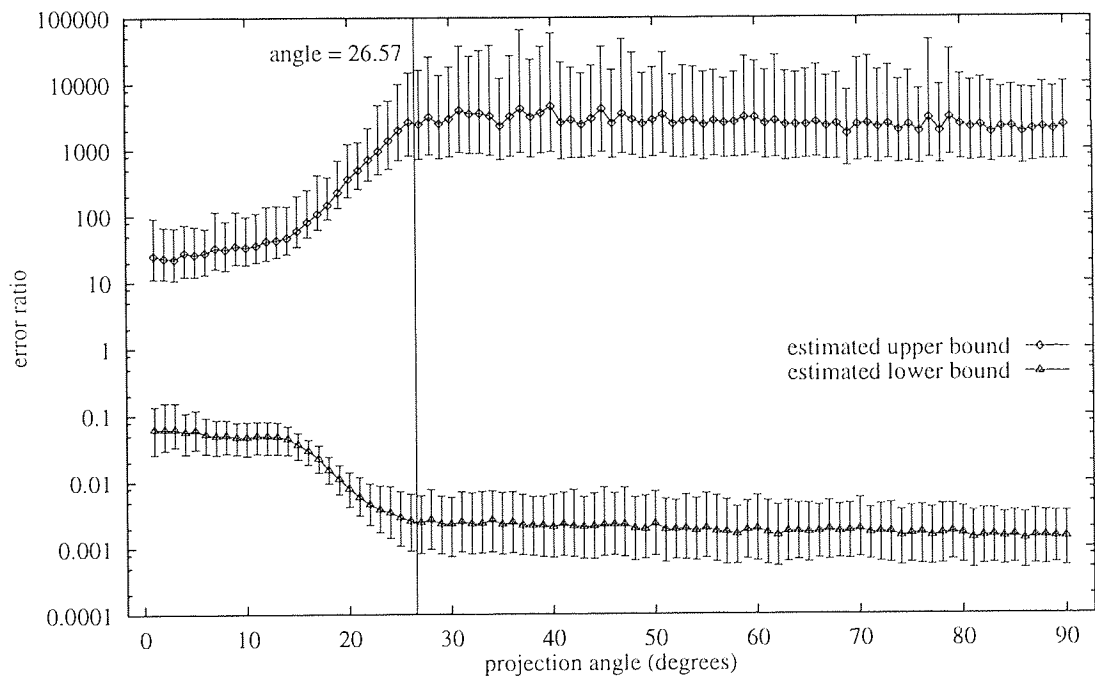
$$\eta_0^{(\phi)}(x) \approx \widehat{f}_\phi^{-1}(y + \eta_\phi^{(0)}(y)) - \widehat{f}_\phi^{-1}(y) \quad (4.15)$$

and calculate, following equation (3.29),

$$\widehat{U}_\phi = \max_{y=f_\phi(x)} \frac{\|\eta_\phi^{(0)}(y)\|}{\|\eta_0^{(\phi)}(x)\|}, \quad \widehat{L}_\phi = \min_{y=f_\phi(x)} \frac{\|\eta_\phi^{(0)}(y)\|}{\|\eta_0^{(\phi)}(x)\|} \quad (4.16)$$



(a)



(b)

Figure 4.11 Numerically simulated, and empirically estimated, upper and lower bounds for the growth of errors under $f_\phi: \mathcal{S}_0 \rightarrow \mathcal{S}_\phi$ and, where appropriate, its inverse. (a) The upper bound U_ϕ clearly indicates a Lipschitz f_ϕ ; the lower bound L_ϕ reaches a terminal value of $L_\phi = 0$ at the critical value $\phi = \phi^*$, since f_ϕ^{-1} is not defined for $\phi > \phi^*$. (b) On a log-linear scale, the mean upper and lower bounds $\langle \hat{U}_\phi \rangle$ and $\langle \hat{L}_\phi \rangle$ for error growth under \hat{f}_ϕ exhibit an unexpectedly reciprocal relationship. Error bars denote one standard deviation in each direction.

If a sufficiently small value of \widehat{U}_ϕ can be found, then it will be taken to approximate a Lipschitz constant U_ϕ of the Lipschitz map f_ϕ , and if a sufficiently large \widehat{L}_ϕ can be found then we will assume that f_ϕ^{-1} exists, and is Lipschitz, with Lipschitz constant L_ϕ^{-1} approximated by \widehat{L}_ϕ^{-1} . (The term ‘sufficiently’ will be interpreted in the context of the variation of \widehat{U}_ϕ and \widehat{L}_ϕ with ϕ , as in the preceding error analyses.) If one or both of these conditions can not be met for a given value of ϕ then we have ruled out the possibility that f_ϕ is a diffeomorphism; we can make no statement stronger than this because, as already stated, finding Lipschitz constants for both f_ϕ and its inverse is not a sufficient condition to show that f_ϕ is a diffeomorphism. In fact, since we already know that f_ϕ is Lipschitz throughout the entire considered range of ϕ , we expect to find evidence for the existence of a finite upper bound U_ϕ for all $0 \leq \phi \leq \frac{\pi}{2}$; we also know that f_ϕ does not have an inverse for $\phi > \phi^*$, so we expect \widehat{L}_ϕ to achieve a distinct minimum over that range.

We plot the expected values of these bounds, calculated over the test set, in figure 4.11(b), on a log-linear scale. Interestingly, the curves obtained in this manner vary somewhat from their (numerically simulated) analytic analogues of figure 4.11(a) in exhibiting a roughly reciprocal relationship: both $\langle \widehat{U}_\phi \rangle$ and $\langle \widehat{L}_\phi^{-1} \rangle$ rise swiftly from a value of $\langle \widehat{U}_0 \rangle, \langle \widehat{L}_0^{-1} \rangle \sim 10$ to a ceiling of $\langle \widehat{U}_\phi \rangle, \langle \widehat{L}_\phi^{-1} \rangle \sim 10,000$ at the critical value which would appear to indicate that neither \widehat{f}_ϕ nor \widehat{f}_ϕ^{-1} are Lipschitz for $\phi > \phi^*$. The significantly increased variability indicated by the error bars in figure 4.11(b)—compared to that of the error curves of figures 4.2 and 4.6—we assume to be a consequence of the fact that the assumption underlying our Lipschitz analysis—that $\widehat{f}_\phi = f_\phi$ for all ϕ —is not true in practice.

It is clear that the two sets of bounds—analytical and empirical—plotted in figure 4.11 demonstrate unexpectedly, and significantly different characteristics. The mechanism behind this somewhat disappointing disparity becomes a little clearer when we examine the behaviour of some specific RBF maps \widehat{f}_ϕ and \widehat{f}_ϕ^{-1} more closely. We do this by constructing scatter plots of the per-point errors $\|\eta_\phi \circ f_\phi(x)\|$ versus $\|\eta_0^{(\phi)}(x)\|$, averaged over 500 random seeds as usual. These plots are shown in figures 4.12 and 4.13. Where appropriate, on each of these graphs we have also plotted the analytically calculated Lipschitz constants, in the form of lines through the origin with gradients U_ϕ and L_ϕ . Naturally, an upper bound exists (with finite gradient U_ϕ) for $0 \leq \phi \leq \frac{\pi}{2}$. Within the parameter range $0 \leq \phi < \phi^*$ a lower bound also exists (with gradient $L_\phi > 0$), and we would like to see the experimental points stay as nearly between these two bounds as possible; for $\phi \geq \phi^*$, we expect only to see evidence for a finite upper bound.

The first of these figures, 4.12(a), shows the scatter plot corresponding to the RBF approximations to f_{20} and its inverse, chosen so that the value of $\phi = 20$ degrees is inside the interval in ϕ within which f_ϕ is a diffeomorphism. In this plot we can plainly see a concentration of points along the directions corresponding to both analytic bounds U_{20} and L_{20} ; indeed, there is virtually no data anywhere *but* along those two lines. In contrast, the case $\phi = 40$ degrees is illustrated in 4.12(b) and, in close-up, in 4.12(c). Once again, the existence, and approximate value of a finite upper bound is firmly indicated by a distinct clustering of points immediately below (and also slightly above) the line $\|\epsilon_{40}^{(0)} \circ f_{40}(x)\| = U_{40} \|\epsilon_0^{(40)}(x)\|$.

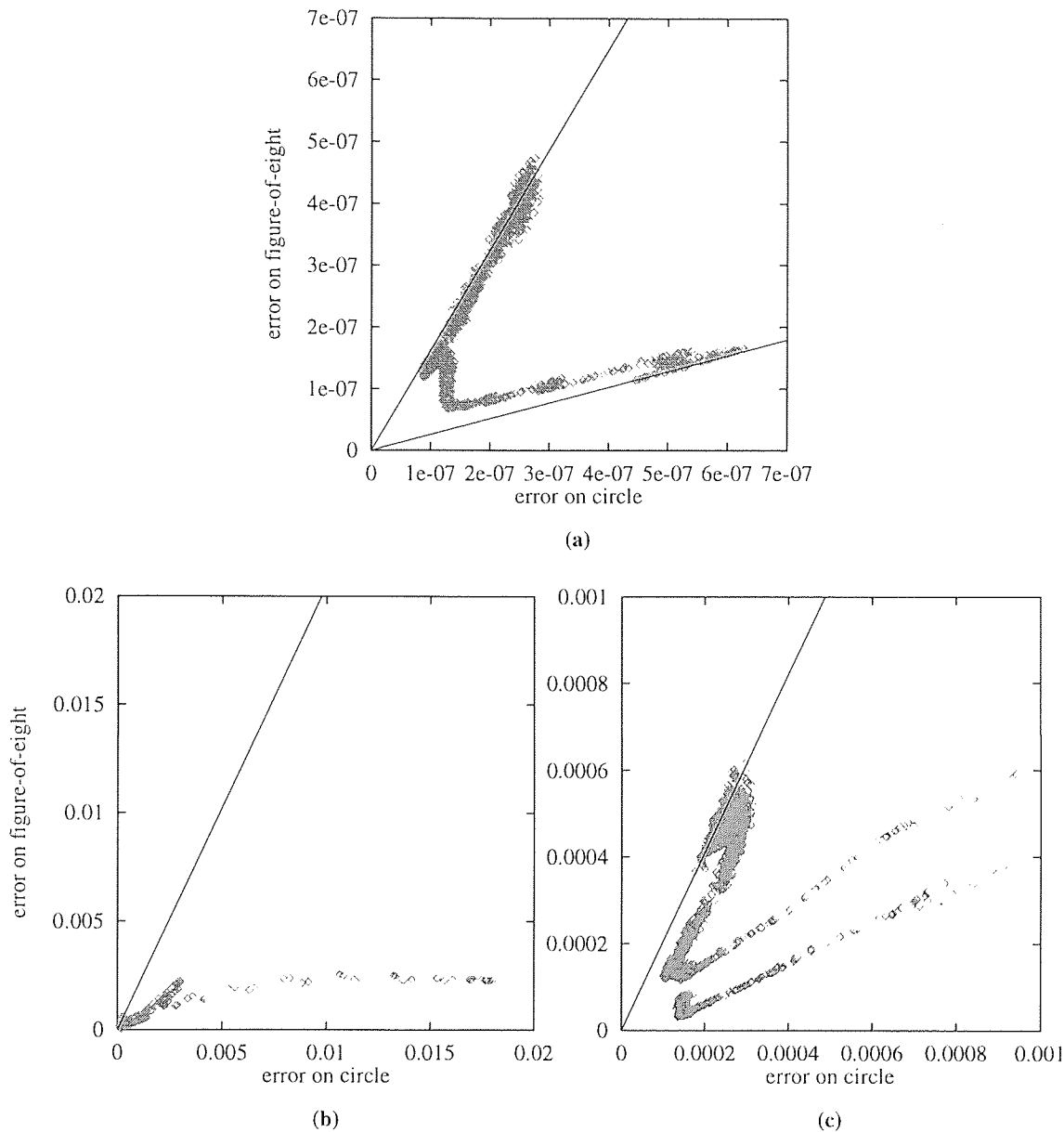


Figure 4.12 Scatter plot of per-point identity errors $\|\eta_\phi^{(0)}(\mathbf{y})\|$ versus $\|\eta_0^{(\phi)}(\mathbf{x})\|$ for $\phi = 20, 40$ degrees, averaged over 500 sets of randomly-seeded repulsive centers. Part (a) plots the errors arising from \mathbf{f}_{20} , in which both upper and lower bounds are clearly visible and agree closely with their analytical versions. Part (b) plots the errors arising from \mathbf{f}_{40} , for which we expect a lower bound of zero, but the data does not appear to directly support this conclusion. We investigate a restricted portion of this plot in part (c), which again gives clear empirical evidence for the analytical upper bound but provides little additional evidence for the expected lower bound other than the fact that the distribution of data points approaches the line $\|\epsilon_{40}^{(0)} \circ \mathbf{f}_{40}(\mathbf{x})\| = 0$ substantially more closely than it does the line $\|\epsilon_0^{(40)}(\mathbf{x})\| = 0$.

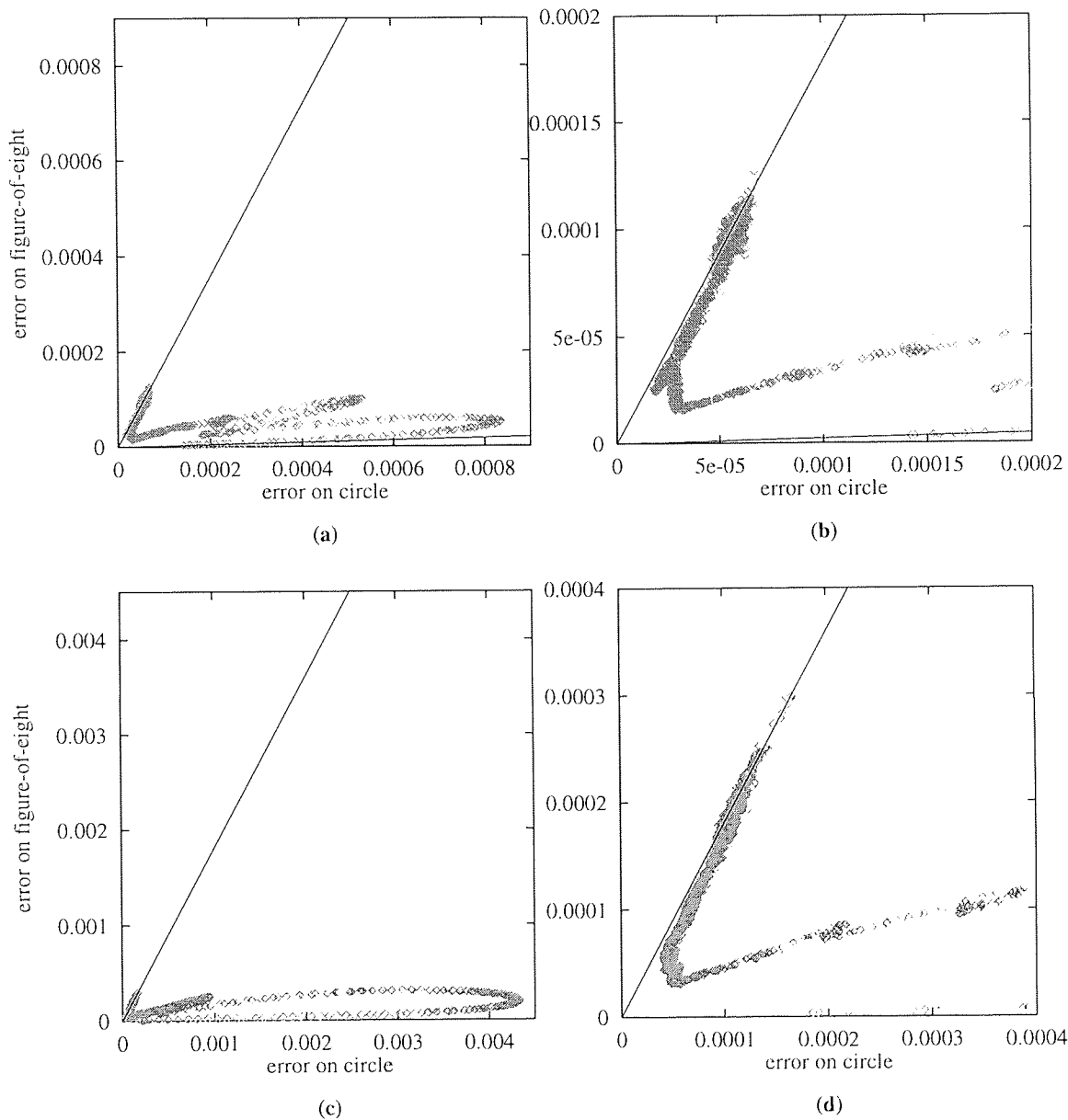


Figure 4.13 Scatter plot of per-point identity errors $\|\eta_\phi^{(0)}(y)\|$ versus $\|\eta_0^{(\phi)}(x)\|$ for $\phi = 26, 27$ degrees, averaged over 500 sets of randomly-seeded repulsive centers. Parts (a) and (b) plot the errors arising from f_{26} , ostensibly a diffeomorphism, the latter on an expanded scale. Again we can clearly see an empirical upper bound which agrees closely with its analytical equivalent. The existence of a small, but non-zero lower bound is trivially indicated in (b), but it is difficult to see any clear numerical distinction between this bound and the one indicated by parts (c) and (d), which show the corresponding plots for f_{27} .

Although we do not, in this case, expect to find a non-zero lower bound, we can clearly draw an empirical one, so it is not entirely clear that we have established the non-invertibility of f_ϕ^{-1} in this case.

In figure 4.13 we examine the behaviour of f_ϕ at either side of the critical value ϕ^* , with $\phi = 26$ degrees in parts (a) and (b) and $\phi = 27$ degrees in (c) and (d). In close-up, in both cases, we again see a characteristic concentration of data points almost directly along, and largely below, the lines corresponding to upper bounds U_{26} and U_{27} . The distinction between the existence of a non-zero lower bound U_{26} , in parts (a) and (b), and of $U_{27} = 0$, in parts (c) and (d), is rather more difficult to determine.

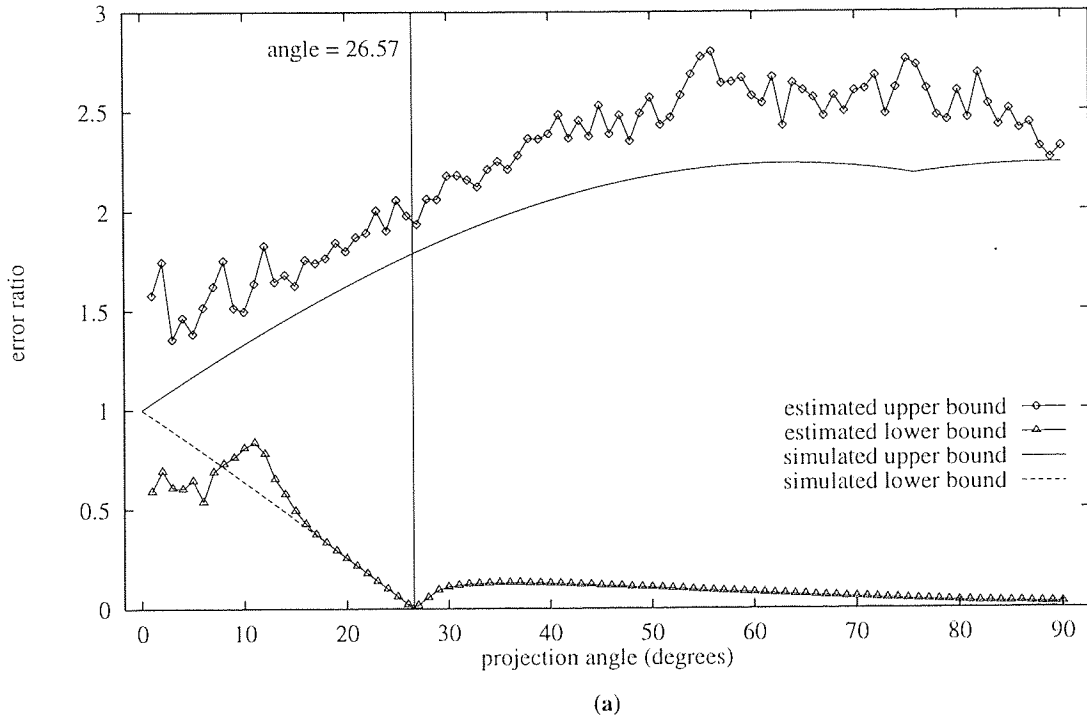


Figure 4.14 Empirically estimated upper and lower bounds \widehat{U}'_ϕ and \widehat{L}'_ϕ , obtained by moving the average over random seeds inside the calculation of extrema, superimposed on the numerically simulated analytical bounds U_ϕ and L_ϕ to which they correspond. These estimates are clearly substantially closer to their analytical analogues than were the estimates \widehat{U}_ϕ and \widehat{L}_ϕ previously calculated: the lower bound, in particular, is closely shadowed by \widehat{L}'_ϕ over a region shortly before the critical value ϕ^* , which also exhibits an interesting increase immediately above this value, presumably due to the non-transversal nature of the intersection in \mathcal{S}_ϕ .

As already mentioned, it is fairly immediate that a significant disparity exists between the distributions plotted in figures 4.12 and 4.13 and the estimated upper and lower bounds plotted in figure 4.11(b): the values of U_ϕ and L_ϕ estimated by the latter are (respectively) substantially greater, and substantially smaller, than the estimates which could clearly be obtained from the scatter plots themselves; the scalar measures $\langle \widehat{U}_\phi \rangle$ and $\langle \widehat{L}_\phi \rangle$ appear insufficient to adequately convey the information contained in the the distributions of mean per-point error ratios. This is almost certainly due to the fact that in estimating U_ϕ and L_ϕ according to equation (4.16) we are calculating minima and maxima of error ratios which have not yet been smoothed by averaging over the model. It seems likely, therefore, that more reliable estimates of U_ϕ and L_ϕ might be obtained by moving the average over random seeds *inside* the calculation of both

minima and maxima, so as to replace equation (4.16) with the corresponding formulae

$$\widehat{U}'_\phi = \max_{\mathbf{y}=\mathbf{f}_\phi(\mathbf{x})} \left\langle \frac{\|\boldsymbol{\eta}_\phi^{(0)}(\mathbf{y})\|}{\|\boldsymbol{\eta}_0^{(\phi)}(\mathbf{x})\|} \right\rangle, \quad \widehat{L}'_\phi = \min_{\mathbf{y}=\mathbf{f}_\phi(\mathbf{x})} \left\langle \frac{\|\boldsymbol{\eta}_\phi^{(0)}(\mathbf{y})\|}{\|\boldsymbol{\eta}_0^{(\phi)}(\mathbf{x})\|} \right\rangle \quad (4.17)$$

These new estimates are plotted in figure 4.14. As hoped, \widehat{U}'_ϕ and \widehat{L}'_ϕ are substantially closer to U_ϕ and L_ϕ than are the estimates plotted in figure 4.11(b). In particular, \widehat{L}'_ϕ appears to coincide almost exactly with the analytical lower bound over the (degree) range $15 \lesssim \phi \lesssim \phi^*$, although it exhibits an interestingly localised increase immediately after the critical value, followed by a decay towards zero as $\phi \rightarrow \frac{\pi}{2}$, which would appear to indicate that the composition of RBF maps $\widehat{\mathbf{f}}_\phi^{-1}$ and $\widehat{\mathbf{f}}_\phi$ is actually somewhat more sensitive to the non-differentiability of \mathbf{f}_{ϕ^*} —which is injective, but not immersive—than it is to the non-injectivity of $\mathbf{f}_{\phi > \phi^*}$. This is almost certainly due to the fact that the cusp in \mathcal{S}_{ϕ^*} actually results in more points which are close together in \mathcal{S}_{ϕ^*} , but whose pre-images in \mathcal{S}_0 are comparatively far apart, than do the transversal intersections which occur in \mathcal{S}_ϕ as ϕ increases beyond ϕ^* .

4.1.2 The circle and total least squares

Now that we have demonstrated that we can successfully use the method of LS to find the critical value ϕ^* it is time to apply the method of TLS to the same problem. We will now be approximating \mathbf{f}_ϕ with the symmetrical RBF map $\widehat{\mathbf{f}}_\phi = \varphi_\phi^{-1} \circ \mathcal{W}_\phi \circ \varphi_0$, in analogy with equation (3.30), where $\varphi_0: \mathcal{S}_0 \subset \mathbb{R}^2 \rightarrow \mathbb{R}^{200}$ and $\varphi_\phi: \mathcal{S}_\phi \subset \mathbb{R}^2 \rightarrow \mathbb{R}^{200}$ have the usual form, with repulsive centers and cubic nonlinearities, as already established.

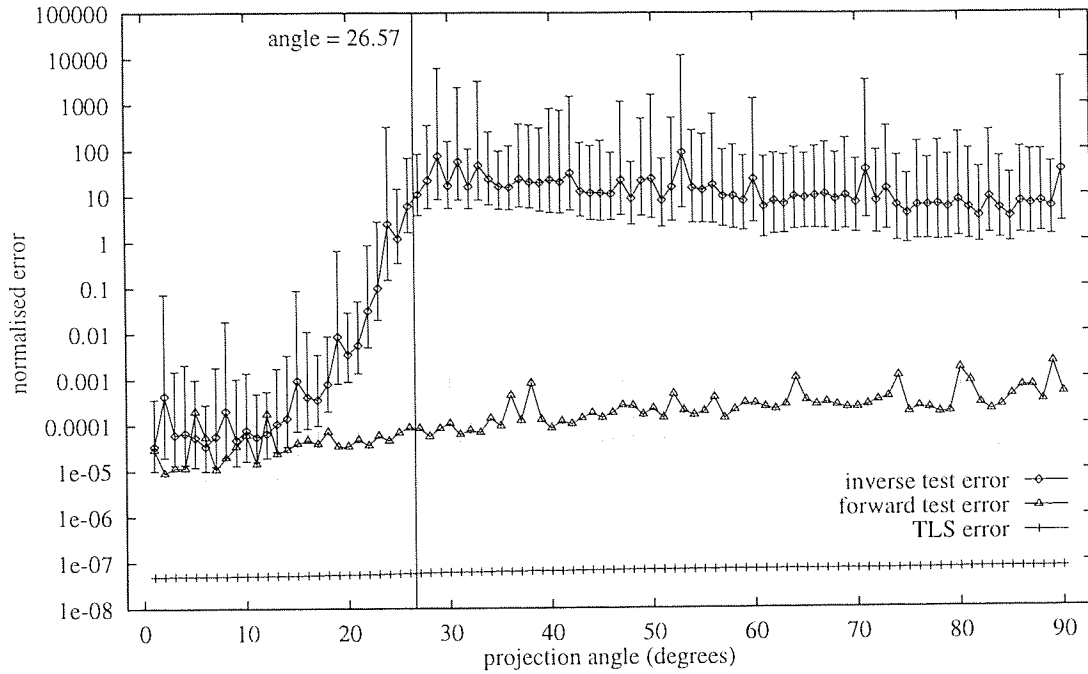
As discussed in section 3.3, having fixed the nonlinear maps φ_0 and φ_ϕ , for a given ϕ , we now concentrate our attention on the linear map $\mathcal{W}_\phi: \mathbb{R}^{200} \rightarrow \mathbb{R}^{200}$, defined by $\mathcal{W}_\phi(\boldsymbol{\varphi}) = \overline{\boldsymbol{\varphi}}_\phi + \mathbf{W}_\phi^T(\boldsymbol{\varphi} - \overline{\boldsymbol{\varphi}}_0)$. We adapt the forward and inverse errors $\epsilon_0^{(\phi)}$ and $\epsilon_\phi^{(0)}$, calculated in the co-domain \mathbb{R}^{200} of φ_ϕ and φ_0 , respectively, from equation (3.32) to get

$$\epsilon_0^{(\phi)2} = \sigma_\phi^{-2} \sum_{i=1}^N \|\mathcal{W}_\phi(\mathbf{a}_i) - \mathbf{b}_i\|^2, \quad \epsilon_\phi^{(0)2} = \sigma_0^{-2} \sum_{i=1}^N \|\mathcal{W}_\phi^{-1}(\mathbf{b}_i) - \mathbf{a}_i\|^2 \quad (4.18)$$

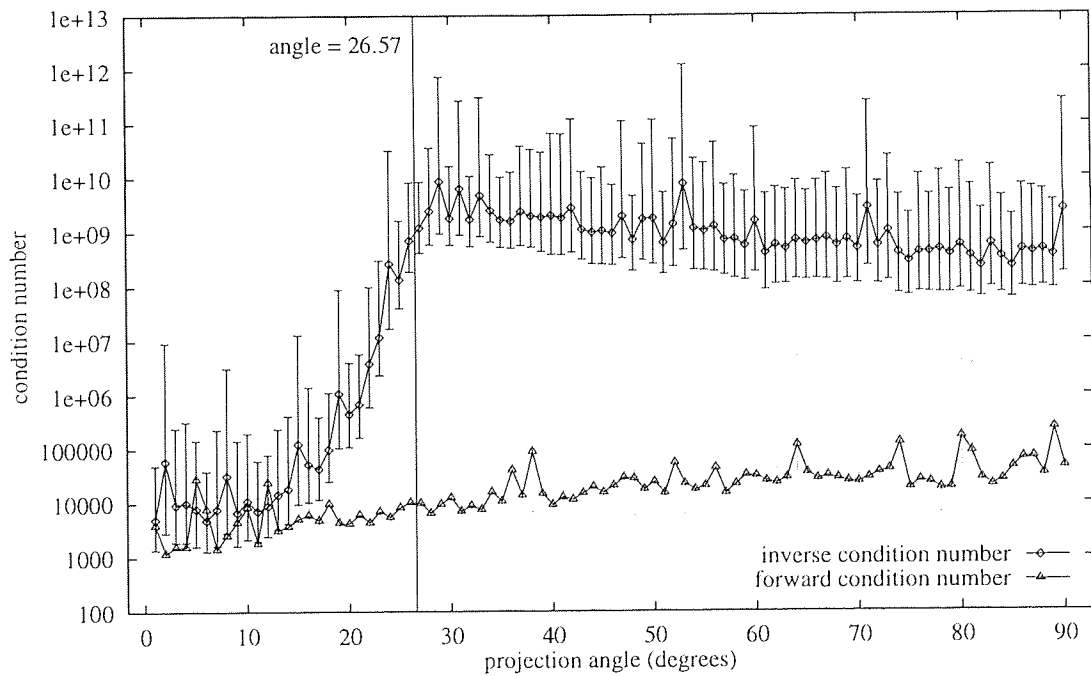
where $\mathbf{a} = \varphi_0(\mathbf{x}) - \overline{\boldsymbol{\varphi}}_0$ and $\mathbf{b} = \varphi_\phi(\mathbf{y}) - \overline{\boldsymbol{\varphi}}_\phi$; the normalising constants σ_0 and σ_ϕ (not to be confused with the σ_0 and σ_ϕ , calculated on \mathcal{S}_0 and \mathcal{S}_ϕ , of the previous section) are calculated on $\varphi_0\mathcal{S}_0 \subset \mathbb{R}^{200}$ and $\varphi_\phi\mathcal{S}_\phi \subset \mathbb{R}^{200}$, respectively. The TLS error $\epsilon_\perp^{(\phi)}$, calculated in the product space \mathbb{R}^{400} , follows similarly from (3.36) as

$$\epsilon_\perp^{(\phi)} = \sigma_\perp^{(\phi)-2} \sum_{i=1}^N \|\mathbf{P}_\phi^T \mathbf{a}_i + \mathbf{Q}_\phi^T \mathbf{b}_i\|^2 \quad (4.19)$$

where $\sigma_\perp^{(\phi)} = \sigma_0 + \sigma_\phi$ and the 200 by 200 matrices \mathbf{P}_ϕ and \mathbf{Q}_ϕ satisfy $\mathbf{W}_\phi = -\mathbf{P}_\phi \mathbf{Q}_\phi^{-1}$ as usual, in analogy with (3.40).



(a)



(b)

Figure 4.15 Comparing the mean, normalised forward and inverse errors and condition numbers for the TLS approximation to $f_\phi: \mathcal{S}_0 \rightarrow \mathcal{S}_\phi$, calculated over the test set. (a) On a log-linear scale the TLS error, at $\langle \epsilon_\perp^{(\phi)} \rangle \approx 10^{-7}$, indicates a uniformly good fit: both it and the forward error $\langle \epsilon_\phi^{(\phi)} \rangle$ are negligible compared to the inverse error $\langle \epsilon_\phi^{(0)} \rangle$. However, the distinction between $\epsilon_0^{(\phi)}$ and $\langle \epsilon_\phi^{(0)} \rangle$ is less clear-cut than in the LS case due to the extreme sensitivity characteristic of the TLS algorithm. The fitting and test errors are all but indistinguishable, even before averaging, and hence only the latter are plotted. (b) The condition numbers $\langle \kappa(Q_\phi) \rangle$ and $\langle \kappa(P_\phi) \rangle$ exhibit the expected correspondence to $\langle \epsilon_\phi^{(\phi)} \rangle$ and $\langle \epsilon_\phi^{(0)} \rangle$, before and after averaging over the random seeds. Error bars denote a one-sided standard deviation in both directions; those on the forward test error and condition number have been lightened for clarity.

We plot the expected values of these errors, calculated over the test set, in figure 4.15(a), for ϕ in the range $0 \leq \phi \leq \frac{\pi}{2}$ as usual; owing to the presence of significantly large outliers associated with many of the random repulsive seeds, each of the error bars plotted in this figure actually corresponds to two one-sided standard deviations (obtained by restricting the variance summation to points respectively above, or below, the mean). The fitting errors are virtually indistinguishable from their test set analogues, and therefore not plotted. On a log-linear scale we see an appreciable amount of variation in the mean inverse error $\langle \epsilon_\phi^{(0)} \rangle$, against which the forward error $\langle \epsilon_0^{(\phi)} \rangle$ is negligibly small, while the product space error $\langle \epsilon_\perp^{(\phi)} \rangle$ is almost constant at $\langle \epsilon_\perp^{(\phi)} \rangle \sim 10^{-7}$, independent of ϕ . Recalling the definition (3.38) of the TLS error, we conclude that the relationship in \mathbb{R}^{400} between $\varphi_0(\mathbf{x})$ and $\varphi_\phi(\mathbf{y})$ in the training set is extremely close to a linear one. Despite the large standard deviations associated with these curves, it is clear that while the forward error remains (relatively) low over the entire range plotted, the inverse error rises, for $\phi \approx \phi^*$, to a substantially higher level. This leads us to conclude that the linear map \mathcal{W}_ϕ , and hence also \mathbf{f}_ϕ , is not invertible for ϕ above the critical value, and hence—once again—that $\mathbf{f}_{\phi > \phi^*}$ is no not a diffeomorphism. The characteristic instability referred to in section 3.3.3 is clearly visible in this plot: although there is an easily discernible separation between forward and inverse error for $\phi \gtrsim \phi^*$, it is significantly less clear-cut than in the corresponding LS case. In part (b) we plot the expected values of the condition numbers $\kappa(\mathbf{Q}_\phi)$ and $\kappa(\mathbf{P}_\phi)$, corresponding closely, on a log-linear scale, to their equivalent errors $\epsilon_0^{(\phi)}$ and $\epsilon_\phi^{(0)}$, respectively. (In fact this correspondence holds for individual RBF maps, before the averaging process has taken place.) This experimental result confirms the analysis in 3.3.2, and demonstrates that (at least in this case) we can use condition numbers and fitting errors interchangeably for the purpose of determining whether or not a given map \mathbf{f}_ϕ is a diffeomorphism.

Finally, in figure 4.16, we illustrate the result of mapping \mathcal{S}_0 and \mathcal{S}_ϕ through the symmetrical RBF maps $\widehat{\mathbf{f}}_\phi$ and $\widehat{\mathbf{f}}_\phi^{-1}$, for $\phi = 20$ degrees, using in each case a single repulsive center seed chosen by maximising $\|\mathbf{x}_i\|$ and $\|\mathbf{y}_i\|$ over the training set in \mathcal{S}_0 and \mathcal{S}_ϕ , respectively. Although we are specifically interested in the images of these sets in the 200-dimensional co-domain of φ_0 and φ_ϕ we actually plot, in each case, the resulting of inverting the appropriate nonlinear map, so as to obtain a two-dimensional image. We must therefore be aware that any visible errors are likely to be due, in part, to the inevitable approximation made in attempting to invert a given φ_0 or φ_ϕ for data not in the image of a compact subset of \mathbb{R}^2 , as discussed in appendix A. In the case of the RBF approximation to \mathbf{f}_{20} this additional source of error is necessarily small: we already know from figure 4.15(b) that \mathbf{P}_{20} is a (relatively) well-conditioned matrix on average, and hence the result—shown in part (a)—of transforming the non self-intersecting set \mathcal{S}_{20} through \mathcal{W}_{20}^{-1} is a good approximation to the circle \mathcal{S}_0 . The condition number $\langle \kappa(\mathbf{Q}_{20}) \rangle$ is even lower, as can be seen from the image of \mathcal{S}_0 , which coincides exactly with \mathcal{S}_{20} in part (b). Corresponding plots for \mathbf{f}_{40} are not shown, because although the image of \mathcal{S}_0 under \mathcal{W}_{40} is equally close to \mathcal{S}_{40} , the ill-conditioning of \mathbf{P}_{40} results in a transformed set $\widehat{\mathbf{f}}_{40}^{-1}$ which appears hopelessly ‘tangled’, and in no way resembles the circle \mathcal{S}_0 which it attempts to approximate. The results of approximating the maps \mathbf{f}_{26} and

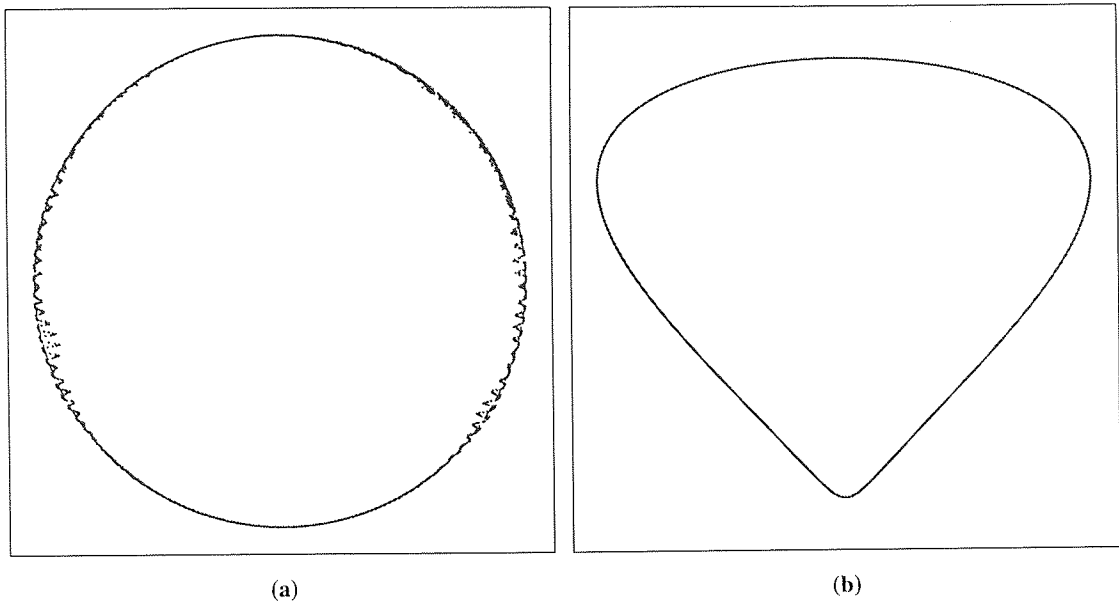


Figure 4.16 Approximating the map $f_\phi: S_0 \rightarrow S_\phi$ and its inverse, for $\phi = 20$, with a symmetrical RBF map. Part (a) shows the image $\widehat{f_{20}^{-1}} S_{20}$ superimposed over S_0 (dashed); the errors visible in this plot are due to the approximation involved in inverting the nonlinearity φ_0 . The image of S_0 under $\widehat{f_{20}}$ in (b) is indistinguishable from S_{20} , in agreement with the low inverse error observed for this value.

f_{27} are similarly messy, and we also do not bother to plot them.

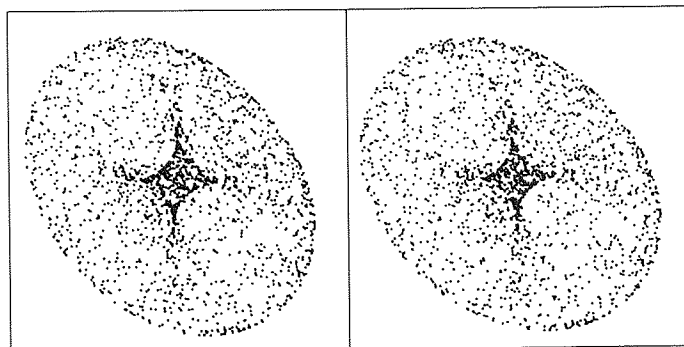
4.2 Embedding a torus in three dimensions

As a brief, further test of this technique, in a somewhat more complicated domain, we have also examined the behaviour of a family of 2-tori in \mathbb{R}^4 under a fixed projection into \mathbb{R}^3 . The 2-torus \mathcal{T}^2 is easily embedded in \mathbb{R}^4 with the map $\Phi_r: \mathcal{T}^2 \rightarrow \mathbb{R}^4$, defined by

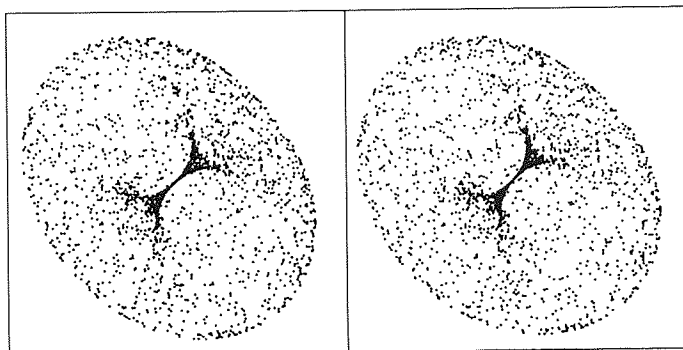
$$\Phi_r(\theta, \phi) = \begin{pmatrix} \cos \theta \\ \sin \theta \\ r \sin \phi \\ r \cos \phi \end{pmatrix} \quad (4.20)$$

where $r \in \mathbb{R}^+$ is a control parameter; we write $\mathcal{T}_{4,r} = \{\Phi_r(\theta, \phi) : 0 \leq \theta, \phi < 2\pi\}$. Φ_r is an embedding provided that $r > 0$, but at the critical value $r^* = 0$ the image of \mathcal{T}^2 under Φ_r is a topological circle. It is a little trickier to embed a 2-torus in \mathbb{R}^3 ; for instance, the result of projecting out one component of $\mathcal{T}_{4,r}$ is a 2-cylinder in \mathbb{R}^3 . To achieve the desired effect we use the nonlinear transformation $\mathcal{F}: \mathbb{R}^4 \rightarrow \mathbb{R}^3$, defined for some $z = \Phi_r(\theta, \phi)$ by

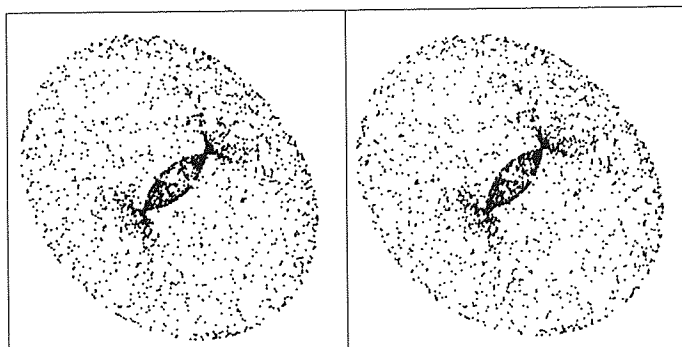
$$\mathcal{F}(z) = \begin{pmatrix} (1 + z_4)z_1 \\ (1 + z_4)z_2 \\ z_3 \end{pmatrix} \quad (4.21)$$



(a)



(b)



(c)

Figure 4.17 Illustrating the 2-tori $\mathcal{T}_{3,r}$ for $r = 0.8, 1$ and 1.2 . The sets in parts (a), (b) and (c) were generated by mapping a joint random distribution in θ and ϕ under the maps $\mathcal{G}_{0.8}$, \mathcal{G}_1 and $\mathcal{G}_{1.2}$, respectively. In part (a) we see that $\mathcal{T}_{3,0.8}$ is embedded in \mathbb{R}^3 , but in part (b) the set $\mathcal{T}_{3,1}$ is not embedded, owing to a co-dimension two self-intersection at the origin, and the set $\mathcal{T}_{3,1.2}$ in part (c) fails to embed at a co-dimension one self-intersecting set centered at the origin. Plotted in stereo.

and resulting in a composite map $\mathcal{G}_r = \mathcal{F} \circ \Phi_r$ with

$$\mathcal{G}_r(\theta, \phi) = \begin{pmatrix} (1 + r \cos \phi) \cos \theta \\ (1 + r \cos \phi) \sin \theta \\ r \sin \phi \end{pmatrix} \quad (4.22)$$

We write $\mathcal{T}_{3,r} = \{x_r(\theta, \phi) : 0 \leq \theta, \phi < 2\pi\}$. This set is embedded in \mathbb{R}^3 provided that $0 < r < 1$. At $r = 0$ it is again a topological circle, and for $r \geq 1$ it contains a self-intersecting subset near the origin. Figure 4.17 depicts the images $\mathcal{T}_{3,0.8}$, $\mathcal{T}_{3,1}$ and $\mathcal{T}_{3,1.2}$ of \mathcal{T}^2 , showing how the self-intersection arises at $r = 1$ when the circle of points of minimal radius about the origin in $\mathcal{T}_{3,r}$ collapses to a point at the origin.

The experiment which we wish to perform is to detect the onset of this self-intersection in $\mathcal{T}_{3,r}$. This task is subtly different from that of the previous section in that \mathcal{F} is a fixed function, independent of the control parameter r , and it is its restriction $\mathcal{F}|_{\mathcal{T}_{4,r}}: \mathcal{T}_{4,r} \rightarrow \mathcal{T}_{3,r}$ and, if it exists, $\mathcal{F}^{-1}|_{\mathcal{T}_{3,r}}: \mathcal{T}_{3,r} \rightarrow \mathcal{T}_{4,r}$ to the particular domain $\mathcal{T}_{4,r}$ and range $\mathcal{T}_{3,r}$ which we must investigate for a given r . For simplicity, we will write $f_r = \mathcal{F}|_{\mathcal{T}_{4,r}}$ and $f_r^{-1} = \mathcal{F}^{-1}|_{\mathcal{T}_{3,r}}$. Clearly, f_r is a function for all values of r to be considered, but f_r^{-1} only exists for $r < 1$. Once again, we will attempt to confirm the existence of this critical value $r^* = 1$ by fitting LS and TLS RBF approximations to f_r . As in the previous experiment, these approximations will be constructed by selecting 500 sets of randomly-seeded repulsive centers (with p and/or $q = 200$), from training sets of $N = 2000$ points in each of $\mathcal{T}_{3,r}$ and $\mathcal{T}_{4,r}$, and plotting means and standard deviations of the errors so obtained. We will obtain these data sets by mapping a joint random distribution in θ and ϕ through Φ_r and \mathcal{G}_r , respectively, with r in the range $0 \leq r \leq 2$, incremented in steps of 0.02.

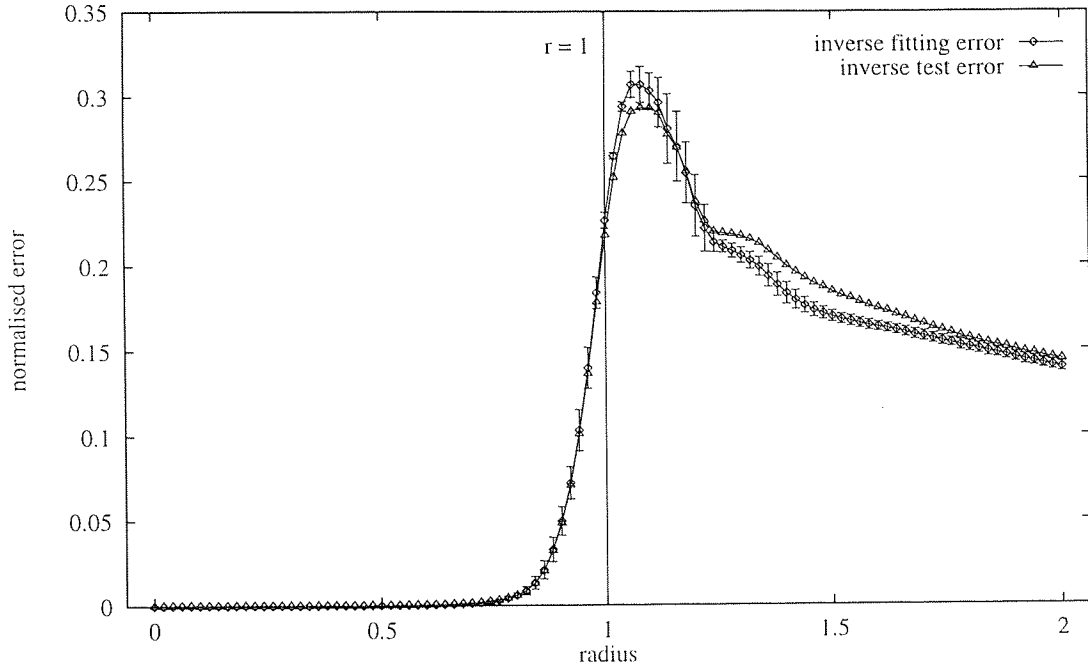
4.2.1 The torus and least squares

We begin again by constructing LS RBF approximations $\widehat{f}_r: \mathcal{T}_{4,r} \subset \mathbb{R}^4 \rightarrow \mathbb{R}^3$ and $\widehat{f}_r^{-1}: \mathcal{T}_{3,r} \subset \mathbb{R}^3 \rightarrow \mathbb{R}^4$ to the maps $f_r: \mathcal{T}_{4,r} \rightarrow \mathcal{T}_{3,r}$ and $f_r^{-1}: \mathcal{T}_{3,r} \rightarrow \mathcal{T}_{4,r}$. We define error functions $\epsilon_{4,r} = \widehat{f}_r - f_r$ and $\epsilon_{3,r} = \widehat{f}_r^{-1} - f_r^{-1}$, giving rise to normalised fitting and/or test errors $\epsilon_{4,r}$ and $\epsilon_{3,r}$, following equation (3.5), by

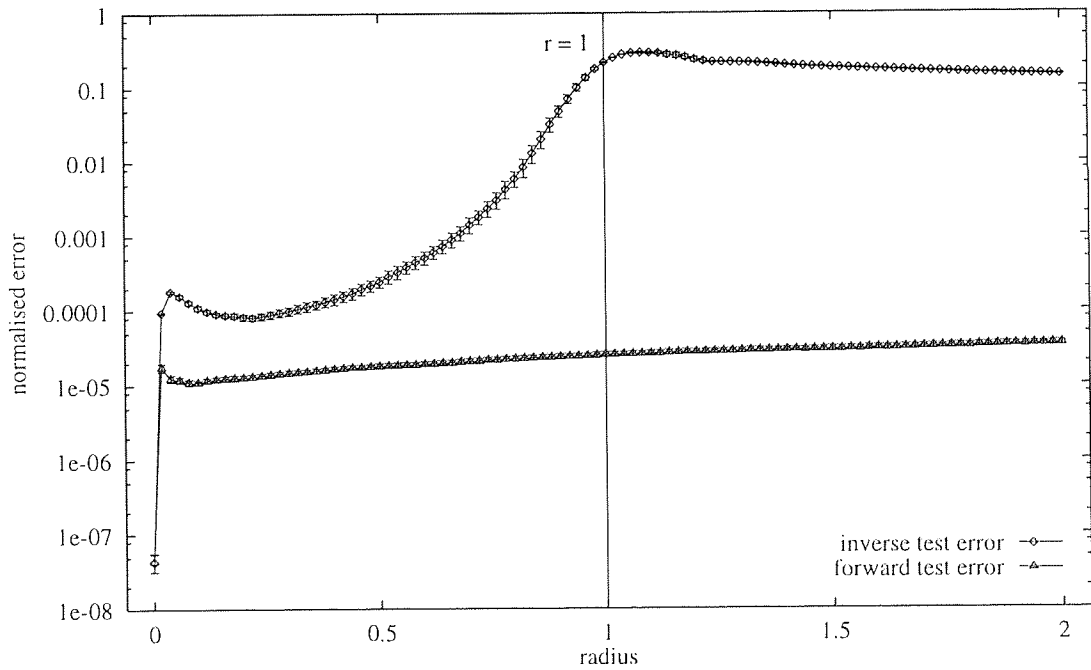
$$\epsilon_{4,r}^2 = \sigma_{3,r}^{-2} \sum_{i=1}^N \|\epsilon_{4,r}(x_i)\|^2, \quad \epsilon_{3,r}^2 = \sigma_{4,r}^{-2} \sum_{i=1}^N \|\epsilon_{3,r}(y_i)\|^2 \quad (4.23)$$

where $x \in \mathcal{T}_{4,r}$, $y \in \mathcal{T}_{3,r}$. Normalisers $\sigma_{4,r}$ and $\sigma_{3,r}$ are calculated on $\mathcal{T}_{4,r} \subset \mathbb{R}^4$ and $\mathcal{T}_{3,r} \subset \mathbb{R}^3$, as usual. As in the previous example (section 4.1.1), we expect f_r to be readily amenable to LS approximation, and hence $\epsilon_{4,r}$ to be uniformly small, for all values of r to be considered. Owing to the presence of a self-intersecting set in $\mathcal{T}_{3,r}$ for $r \geq r^*$, however, we only expect to be able to find a good approximation to the relationship $\mathcal{T}_{3,r} \mapsto \mathcal{T}_{4,r}$ for $r < r^*$.

We plot the expected values of $\epsilon_{4,r}$ and $\epsilon_{3,r}$, calculated over both training and test sets, in figure 4.18. In part (a), on a linear scale, the forward error $\epsilon_{4,r}$ is negligible for both sets; for this reason we



(a)



(b)

Figure 4.18 Comparing the mean, normalised errors $\langle \epsilon_{4,r} \rangle$ and $\langle \epsilon_{3,r} \rangle$, versus r , for the LS approximations \widehat{f}_r and \widehat{f}_r^{-1} on the tori $\mathcal{T}_{4,r}$ and $\mathcal{T}_{3,r}$, calculated over training and test sets. (a) On a linear scale both $\langle \epsilon_{4,r} \rangle$ and $\langle \epsilon_{3,r} \rangle$ are negligible for $r < r^*$ but while the forward error (not plotted) remains so over the entire range of r , the mean inverse error rises steeply at the critical value to saturate at $\langle \epsilon_{3,r} \rangle \sim 0.8$ for $r > r^*$; on this scale the fitting and test errors diverge slightly beyond this critical value. (b) A log-linear scale reveals a practically constant forward error, at $\langle \epsilon_{4,r} \rangle \sim 10^{-7}$; on this scale the separation of fitting and test errors is virtually impossible to distinguish, so only the test error is plotted. Error bars denote one standard deviation in each direction; those on the inverse test error have been lightened for clarity.

do not bother to plot it. The inverse error is also small for $r \lesssim 0.8$, but grows to a relatively large peak value of $\langle \epsilon_{3,r} \rangle \approx 0.3$ at $r \approx 1$. This behaviour once again confirms our conclusion that, although f_r is a diffeomorphism for $r < 1$, for $r \geq 1$ there is no functional relationship taking $\mathcal{T}_{3,r}$ to $\mathcal{T}_{4,r}$. There is a noticeable separation between fitting and test errors for $r > r^*$ on this scale, indicating that a small degree of over-fitting is taking place. In part (b) we plot both $\langle \epsilon_{4,r} \rangle$ and $\langle \epsilon_{3,r} \rangle$ on a log-linear scale (on this scale it is virtually impossible to distinguish between fitting and test errors, so only the latter is plotted); we now see a nearly uniform forward error $\langle \epsilon_{4,r} \rangle \sim 10^{-5}$, as expected. Also revealed by this plot, a dip in both forward and inverse errors at $r = 0$ to $\langle \epsilon_{4,0} \rangle \sim \langle \epsilon_{3,0} \rangle \sim 10^{-8}$ occurs when both tori are collapsed to their respective origins in \mathbb{R}^4 and \mathbb{R}^3 , respectively.

Curiously, having reached its peak at $r \approx 1$, the mean inverse error then actually decreases by several standard deviations from this value as r increases further. To understand this behaviour it is necessary to refer back to the tori $\mathcal{T}_{3,1}$ and $\mathcal{T}_{3,1.2}$ illustrated in parts (b) and (c) of figure 4.17 and consider the nature of the self-intersection which occurs in each case. For $r > 1$, as demonstrated by part (c), this intersection takes place transversally, on a two-point set in \mathbb{R}^3 . In contrast, for $r = 1$, part (b) shows that the self-intersecting set contains a single point in \mathbb{R}^3 and, in particular, is non-transversal, with a single tangent vector aligned with the axis of the torus. It therefore seems reasonable to assume that the region within which the self-intersection in $\mathcal{T}_{3,1}$ has a significant effect on the approximating map $\widehat{f_r^{-1}}$ is greater in extent than in the $r > 1$ case; we would expect this effect to be most noticeable for $r = 1$ and become steadily less significant as r increases beyond 1, as seen in figure 4.18.

4.2.1.1 Analysis of the least squares solution

The image of $\mathcal{T}_{4,r}$ under f_r turns out to be indistinguishable from $\mathcal{T}_{3,r}$ for all r , as expected. Furthermore, as already noted, any projection of $\mathcal{T}_{4,r}$ onto three of its basis vectors gives rise to a cylinder in \mathbb{R}^3 . This, coupled with the random nature of the training and test sets, makes a visual comparison of $\mathcal{T}_{3,r}$ with its image under f_r^{-1} of limited use. Suffice it to say that for $r < 1$ there is no discernible difference between the two, while for $r \geq 1$ the difference is considerable. More appealing, from a visual point of view, is a plot of $\mathcal{T}_{3,r}$, colour-coded by the per-point error magnitudes $\|\epsilon_{3,r}\|$ to which it gives rise under f_r^{-1} , as in figures 4.9 and 4.10. We therefore plot these in figure 4.19, for $r = 0.8, 1$ and 1.2 , this time constructing $\widehat{f_r^{-1}}$ from a single set of repulsive centers in $\mathcal{T}_{3,r} \subset \mathbb{R}^3$, whose seed maximises $\|\mathbf{y}_i\|$ over the training set. (Again, we neglect to make the corresponding plots for f_r applied to $\mathcal{T}_{4,r}$, since the errors in each such case are negligible and the manifolds in question are poorly represented in three dimensions.) Since the area of interest is localised to the ‘hub’ of each of these tori we actually plot a close-up, cross-sectional view of each object (enabling a particularly vivid illustration of the distinction between transversal and non-transversal intersections discussed above). When studying this figure it should again be borne in mind

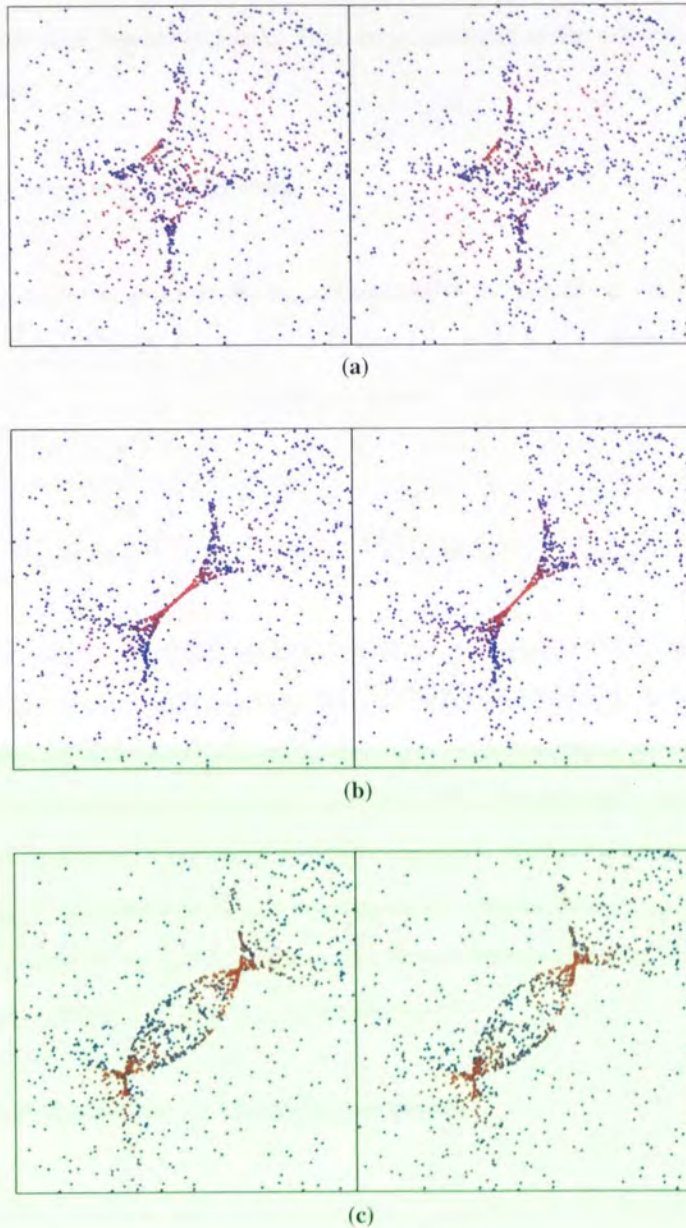


Figure 4.19 Illustrating the errors arising from the approximation $\widehat{f}_r^{-1}: \mathcal{T}_{3,r} \subset \mathbb{R}^3 \rightarrow \mathbb{R}^4$ by colour-coding the elements of $\mathcal{T}_{3,r}$, for $r = 0.8, 1$ and 1.2 . (a) With $r = 0.8$ the error magnitudes are negligibly small and distributed throughout the manifold $\mathcal{T}_{3,0.8}$, as $f_{0.8}$ is a diffeomorphism. (b) At $r = 1$, however, we see a concentration of large error magnitudes in the vicinity of the self-intersection in $\mathcal{T}_{3,1}$, and (c) at $r = 1.2$ we see two distinct regions of large error, concentrated at the two one-dimensional intersections in $\mathcal{T}_{3,1.2}$. Plotted in stereo.

that the colour-coding is normalised for each individual transformation, making a comparison of colour-codes between different plots unhelpful. In particular, the errors in part (a), for $r = 0.8$, are negligible; this plot is included for the purposes of geometrical comparison with parts (b) and (c), whose errors can be seen to be localised to the region(s) of self-intersection in each case, as expected. In particular, in 4.19(c), we actually see two distinct regions of locally large error, centered at the self-intersecting point-sets in $\mathcal{T}_{3,1,2}$, also as expected.

4.2.1.2 Approximating the identity map

Once again it is instructive to examine the approximations to the identity maps $\mathbf{I}_{4,r}: \mathcal{T}_{4,r} \rightarrow \mathcal{T}_{4,r}$ and $\mathbf{I}_{3,r}: \mathcal{T}_{3,r} \rightarrow \mathcal{T}_{3,r}$ constructed with $\widehat{\mathbf{I}}_{4,r} = \widehat{\mathbf{f}}_r^{-1} \circ \widehat{\mathbf{f}}_r$ and $\widehat{\mathbf{I}}_{3,r} = \widehat{\mathbf{f}}_r \circ \widehat{\mathbf{f}}_r^{-1}$. With identity error functions $\eta_{4,r} = \widehat{\mathbf{I}}_{4,r} - \mathbf{I}_{4,r}$ and $\eta_{3,r} = \widehat{\mathbf{I}}_{3,r} - \mathbf{I}_{3,r}$ the normalised identity errors $\eta_{4,r}, \eta_{3,r} \in \mathbb{R}$ follow from equations (3.22) and (3.23) as

$$\eta_{4,r}^2 = \sigma_{4,r}^{-2} \sum_{i=1}^N \|\eta_{4,r}(\mathbf{x}_i)\|^2, \quad \eta_{3,r}^2 = \sigma_{3,r}^{-2} \sum_{i=1}^N \|\eta_{3,r}(\mathbf{y}_i)\|^2 \quad (4.24)$$

We plot the expected values of these errors, calculated over the test set as usual, in figure 4.20. We immediately notice—in analogy with equation (3.24)—the relationship $\eta_{4,r} \approx \epsilon_{3,r} \circ \mathbf{f}_r$, which again holds well enough that the difference between the two curves is impossible to perceive (and therefore not shown) even on the linear scale of 4.20(a), and even before the averaging takes place. Application of the Lipschitz analysis of section 3.2.4 is thereby justified once more. Unlike the previous example, however, the identity error $\langle \eta_{3,r} \rangle$, calculated on $\mathcal{T}_{3,r}$, is not consistently smaller than $\langle \eta_{4,r} \rangle$, so we cannot claim in this case that $\widehat{\mathbf{f}}_r$ is contractive on $\widehat{\mathbf{f}}_r^{-1} \mathcal{T}_{3,r}$. Due to the reasons discussed above, we do not plot the actual images of $\mathcal{T}_{4,r}$ and $\mathcal{T}_{3,r}$, under $\widehat{\mathbf{I}}_{4,r}$ and $\widehat{\mathbf{I}}_{3,r}$, in this instance.

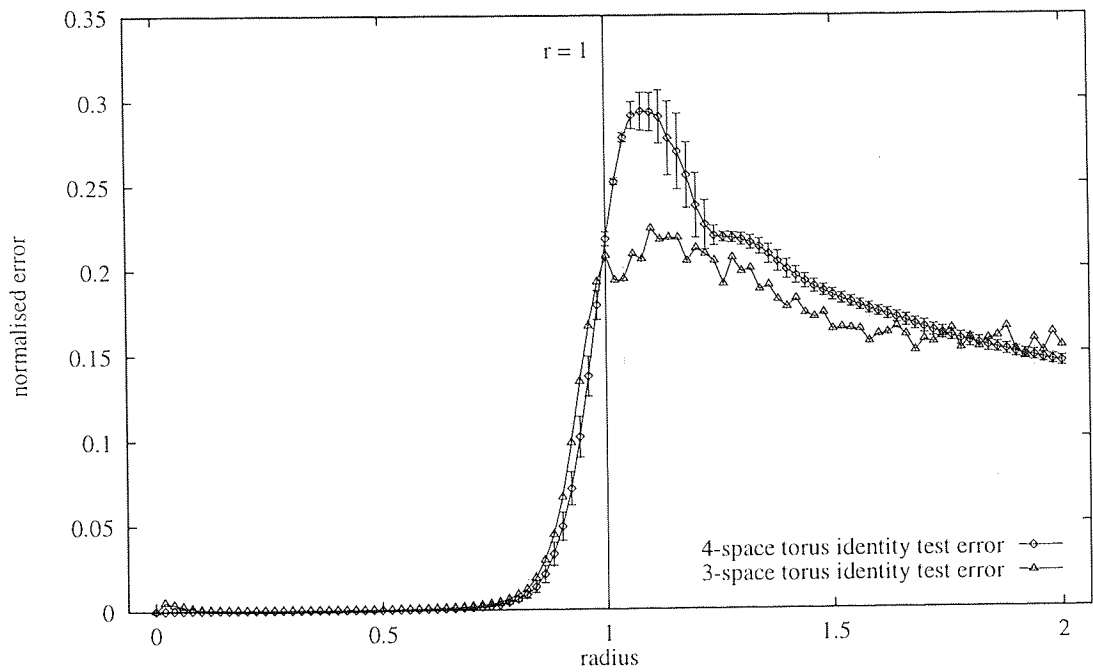
4.2.1.3 Analytical calculation of Lipschitz constants

In this, as in the previous example, we can derive U_r and L_r analytically. We therefore consider the points $\mathbf{x}_1, \mathbf{x}_2 \in \mathcal{T}_{4,r}$ and $\mathbf{y}_1, \mathbf{y}_2 \in \mathcal{T}_{3,r}$ and write (3.21) as

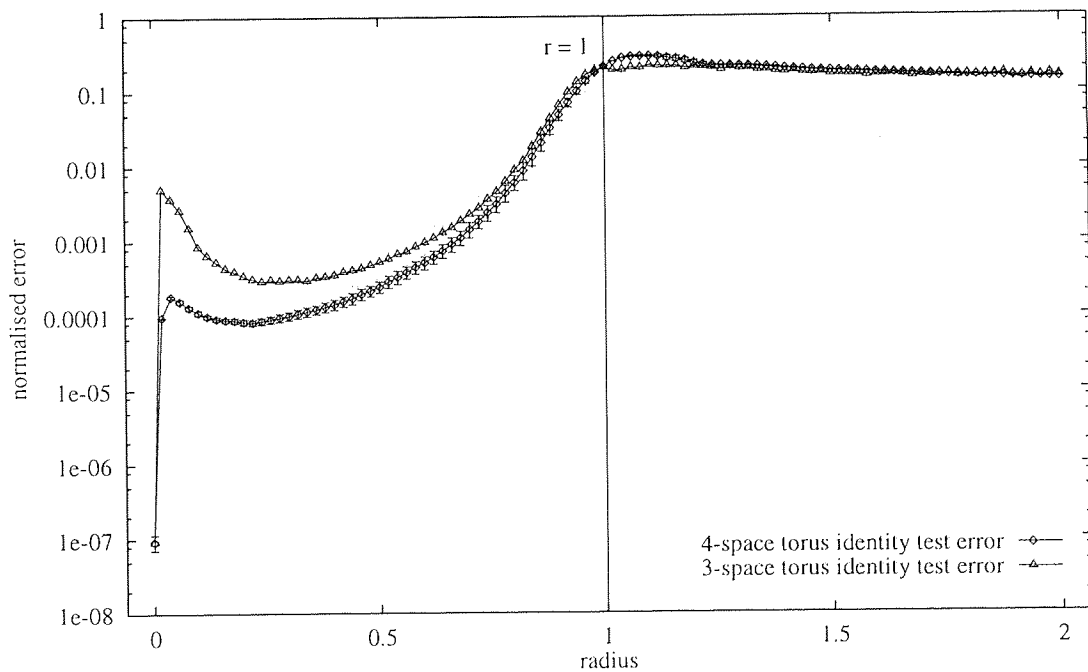
$$U_r = \max_{\mathbf{y}=\mathbf{f}_r(\mathbf{x})} \frac{\|\mathbf{y}_2 - \mathbf{y}_1\|}{\|\mathbf{x}_2 - \mathbf{x}_1\|}, \quad L_r = \min_{\mathbf{y}=\mathbf{f}_r(\mathbf{x})} \frac{\|\mathbf{y}_2 - \mathbf{y}_1\|}{\|\mathbf{x}_2 - \mathbf{x}_1\|} \quad (4.25)$$

Writing \mathbf{x} and \mathbf{y} directly in terms of θ and ϕ , following equation (4.22), and making the substitutions $\alpha = \frac{1}{2}(\phi_2 + \phi_1)$, $\beta = \frac{1}{2}(\phi_2 - \phi_1)$ and $\gamma = \frac{1}{2}(\theta_2 - \theta_1)$ this becomes

$$U_r = \max_{\alpha, \beta, \gamma} R_r(\alpha, \beta, \gamma), \quad L_r = \min_{\alpha, \beta, \gamma} R_r(\alpha, \beta, \gamma) \quad (4.26)$$



(a)



(b)

Figure 4.20 Comparing the mean, normalised test set errors $\langle \eta_{4,r} \rangle$ and $\langle \eta_{3,r} \rangle$ versus r for the LS identity approximations $I_{4,r}$ and $I_{3,r}$ on $\mathcal{T}_{4,r}$ and $\mathcal{T}_{3,r}$. (a) On a linear scale we once again see an identity error $\langle \eta_{4,r} \rangle$ nearly identical to $\langle \epsilon_{3,\phi} \rangle$, while $\langle \eta_{3,r} \rangle$ exhibits a sensitive dependence on the choice of repulsive seed. (b) A log-linear scale is included for completeness. Error bars denote one standard deviation in each direction; those on $I_{3,r}$ have been lightened for clarity.

where

$$R_r^2(\alpha, \beta, \gamma) = \frac{r^2 \sin^2 \beta + [1 + 2r \cos \alpha \cos \beta + r^2(\cos^2 \alpha + \cos^2 \beta - 1)] \sin^2 \gamma}{r^2 \sin^2 \beta + \sin^2 \gamma} \quad (4.27)$$

With a fair amount of further analysis it can now be shown [2] that (4.27) is maximised, independently of γ , when $\alpha = \beta = 0$, with the result that $U_r = 1 + r$. A similar analysis results in a minimum of $L_r = 1 - r$, with the understanding that L_r is only defined for $r < 1$. These values have been verified by numerical simulation, as in the previous example, but for reasons of scale, they are not plotted until figure 4.23.

4.2.1.4 Approximation of the Lipschitz constants

The Lipschitz constants of \widehat{f}_r and, if it exists, \widehat{f}_r^{-1} follow from equation (3.29) with

$$\widehat{U}_r = \max_{\mathbf{y}=\widehat{f}_r(\mathbf{x})} \frac{\|\eta_{3,r}(\mathbf{y})\|}{\|\eta_{4,r}(\mathbf{x})\|}, \quad \widehat{L}_r = \min_{\mathbf{y}=\widehat{f}_r(\mathbf{x})} \frac{\|\eta_{3,r}(\mathbf{y})\|}{\|\eta_{4,r}(\mathbf{x})\|} \quad (4.28)$$

since, in analogy with (3.25) and (3.26),

$$\eta_{3,r}(\mathbf{y}) \approx \widehat{f}_r(\mathbf{x} + \eta_{4,r}(\mathbf{x})) - \widehat{f}_r(\mathbf{x}) \quad (4.29)$$

and

$$\eta_{4,r}(\mathbf{x}) \approx \widehat{f}_r^{-1}(\mathbf{y} + \eta_{3,r}(\mathbf{y})) - \widehat{f}_r^{-1}(\mathbf{y}) \quad (4.30)$$

In this case, we expect to find that \widehat{U}_r is finite, irrespective of r , but \widehat{L}_r^{-1} should become effectively infinite (up to the limit of numerical precision) as $r \rightarrow 1$. The expected values of \widehat{U}_r and \widehat{L}_r , calculated over the test set as usual, are plotted in figure 4.21, and are clearly substantially larger than the analytical values which they are intended to estimate; they also exhibit a relationship similar to the somewhat reciprocal one we saw in figure 4.11(b).

In figure 4.22 we plot the analytic bounds U_r and, where appropriate, L_r , superimposed on scatter plots of $\|\eta_{3,r} \circ \widehat{f}_r(\mathbf{x})\|$ versus $\|\eta_{4,r}(\mathbf{x})\|$, averaged over 500 random repulsive seeds as usual, for $r = 0.8, 1$ and 1.2 . In all three of these plots the data appears to agree quite closely with the analytic upper bound (although the upper bound $U_{0.8}$ is noticeably underestimated), so it seems likely that it might again be useful to bring the average over models inside the extremum calculation, as in the previous experiment. We therefore define the quantities \widehat{U}'_r and \widehat{L}'_r , in analogy with equation (4.17), by

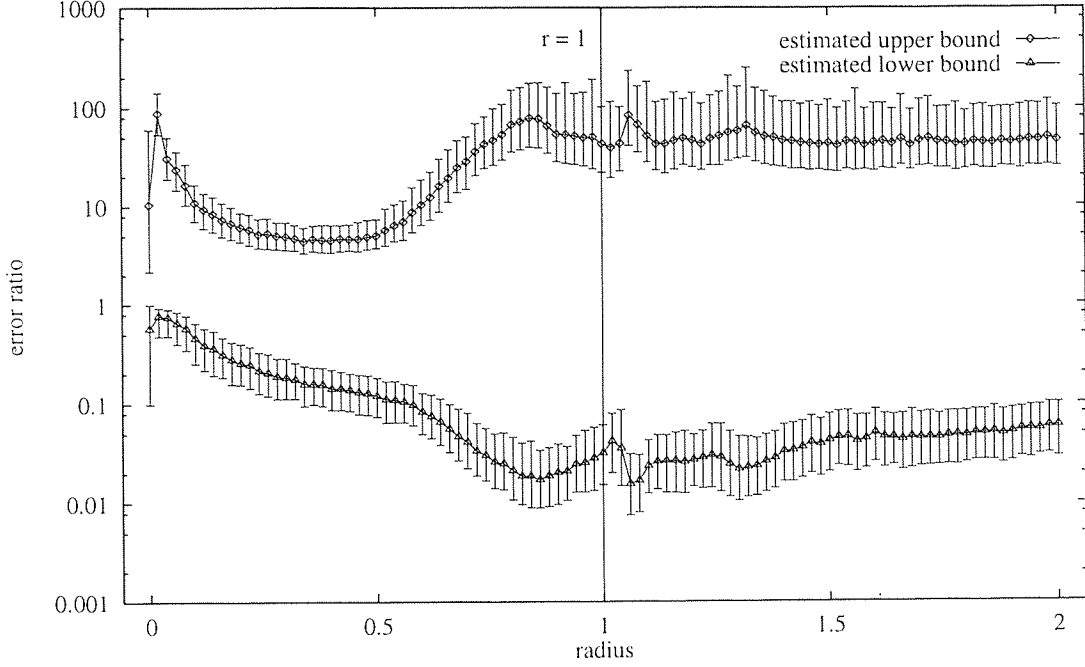


Figure 4.21 Empirically estimated upper and lower bounds for the growth of errors under $f_r: \mathcal{T}_{4,r} \rightarrow \mathcal{T}_{3,r}$. On a log-linear scale we see the same reciprocal relationship between the mean upper bound (\widehat{U}_r) and the mean lower bound (\widehat{L}_r) that we saw in the previous example, with a mean upper bound substantially greater than anticipated. Error bars denote one standard deviation in each direction.

$$\widehat{U}'_r = \max_{\mathbf{y}=f_r(\mathbf{x})} \left\langle \frac{\|\boldsymbol{\eta}_{3,r}(\mathbf{y})\|}{\|\boldsymbol{\eta}_{4,r}(\mathbf{x})\|} \right\rangle, \quad \widehat{L}'_r = \min_{\mathbf{y}=f_r(\mathbf{x})} \left\langle \frac{\|\boldsymbol{\eta}_{3,r}(\mathbf{y})\|}{\|\boldsymbol{\eta}_{4,r}(\mathbf{x})\|} \right\rangle \quad (4.31)$$

The result of calculating these alternative estimates is plotted in figure 4.23. At $r = 0$, both \widehat{U}'_r and \widehat{L}'_r appear to coincide almost exactly with $U_0 = L_0 = 1$, as might be expected. However, both estimates then surprisingly jump in value, the upper bound to $\widehat{U}'_r \sim 100$, and the lower bound to $\widehat{L}'_r \sim 10$, as soon as r increases past zero (or in other words, as soon as the tori in \mathbb{R}^4 and \mathbb{R}^3 acquire a non-zero width), before asymptoting back towards their analytical analogues as $r \rightarrow r^*$. This behaviour is presumably symptomatic of the resolving power of the RBF algorithm. Beyond the critical value $r^* = 1$, the estimated upper bound follows U_r very closely, while the lower bound (which in this region should theoretically approach zero) appears to reach a floor substantially higher than that observed in the previous experiment.

4.2.2 The torus and total least squares

A similar analysis applies to the TLS case as it did in the previous section. We write (3.30) in this case as $\widehat{f}_r = \varphi_{3,r}^{-1} \circ \mathcal{W}_r \circ \varphi_{4,r}$, with $\varphi_{4,r}: \mathcal{T}_{4,r} \subset \mathbb{R}^4 \rightarrow \mathbb{R}^{200}$ and $\varphi_{3,r}: \mathcal{T}_{3,r} \subset \mathbb{R}^3 \rightarrow \mathbb{R}^{200}$ defined as usual and $\mathcal{W}: \mathbb{R}^{200} \rightarrow \mathbb{R}^{200}$ defined by $\mathcal{W}_r(\varphi) = \overline{\varphi_{3,r}} + \mathcal{W}_r^T(\varphi - \overline{\varphi_{4,r}})$. Forward and inverse errors follow from equation (3.32) as

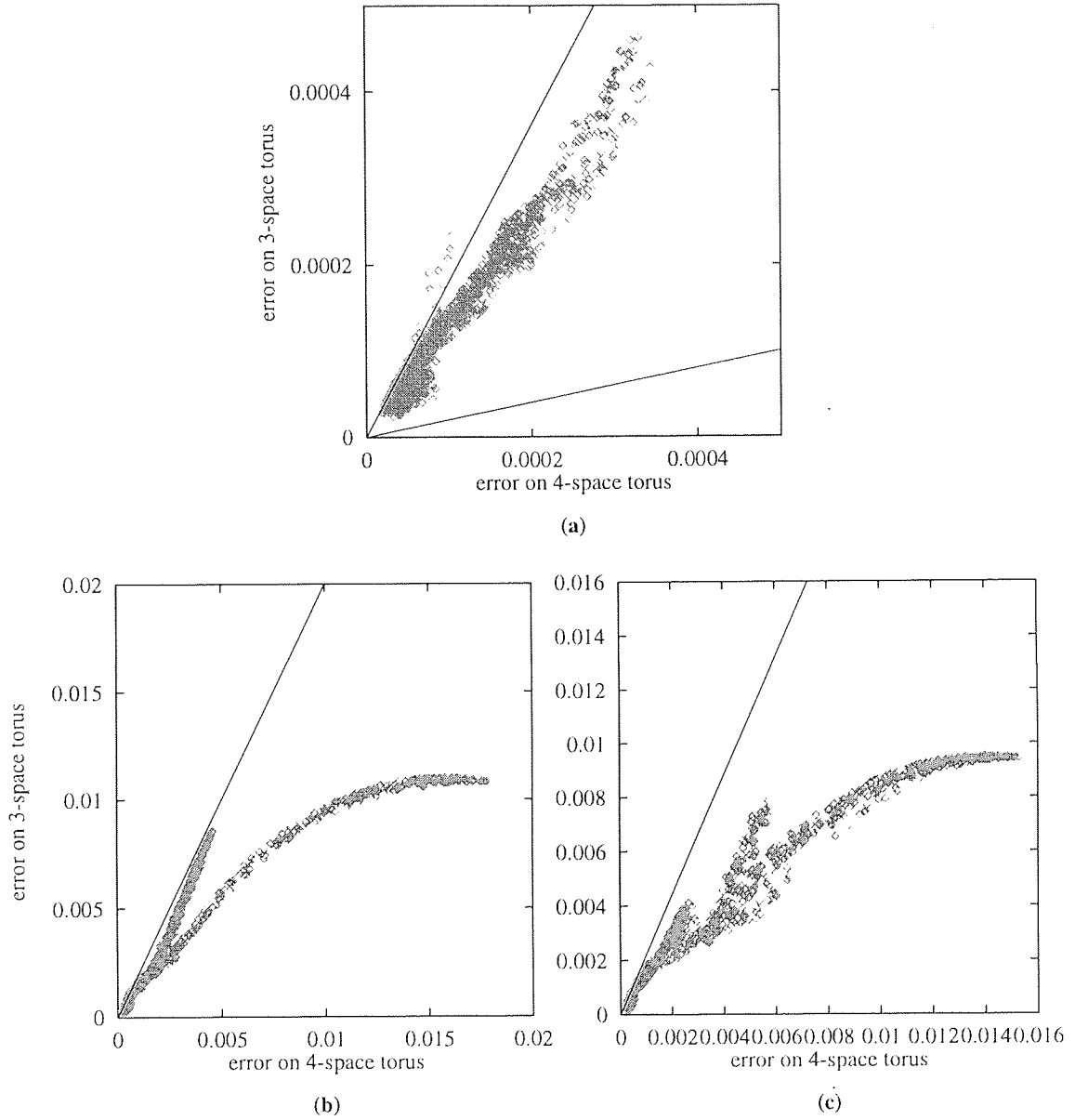


Figure 4.22 Scatter plot of per-point identity errors $\|\eta_{3,r}(y)\|$ versus $\|\eta_{4,r}(x)\|$ for $r = 0.8, 1$ and 1.2 , averaged over 500 sets of randomly-seeded repulsive centers. (a) For $r = 0.8$, the upper bound $U_{0.8} = 1.8$ and lower bound $L_{0.8} = 0.2$ seem to bracket the data fairly well, with the exception of a little ‘leakage’ above $U_{0.8}$; (b) for $r = 1$, the upper bound $U_1 = 2$ and (c) the upper bound $U_{1.2} = 2.2$ also appear to be a good fit to the data.

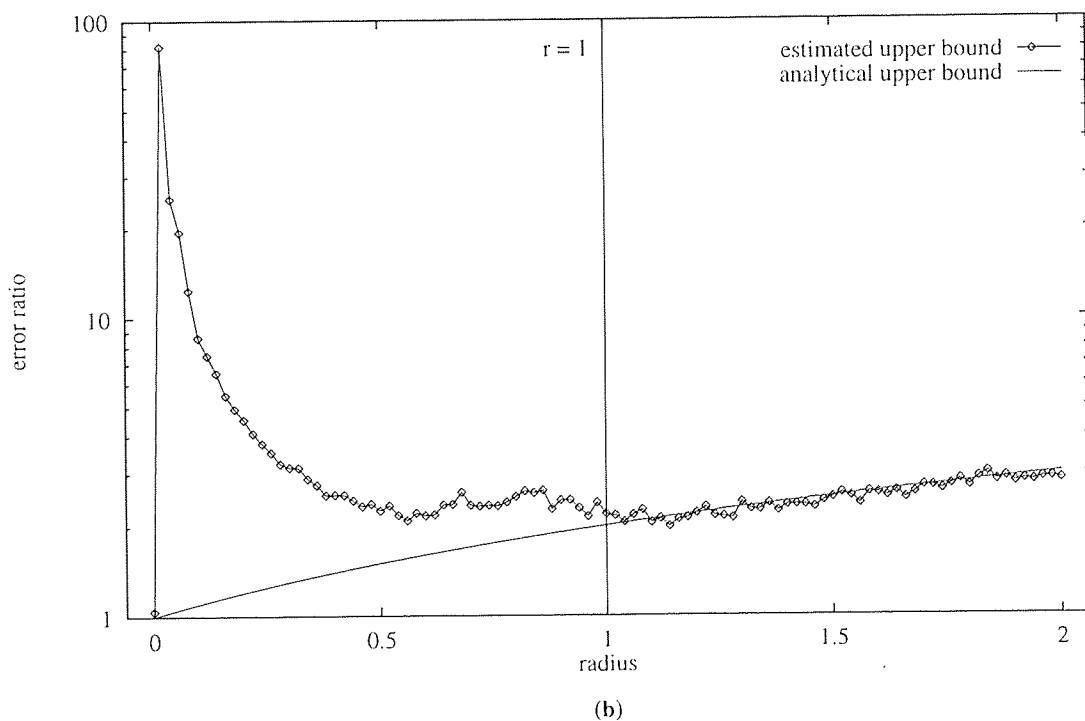
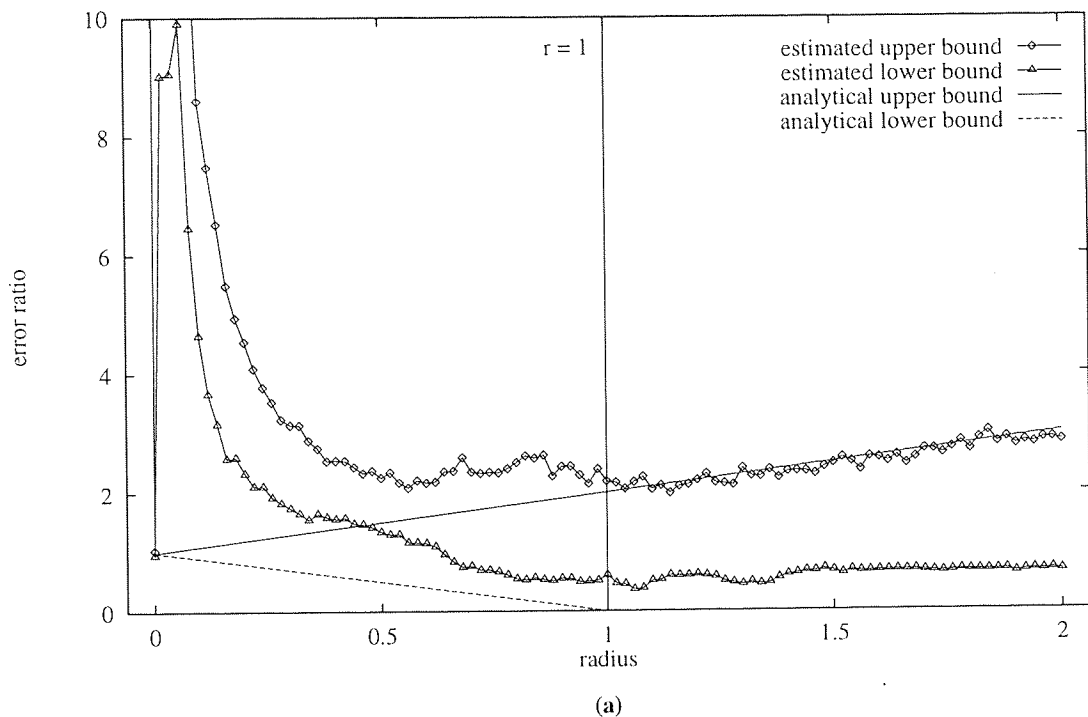
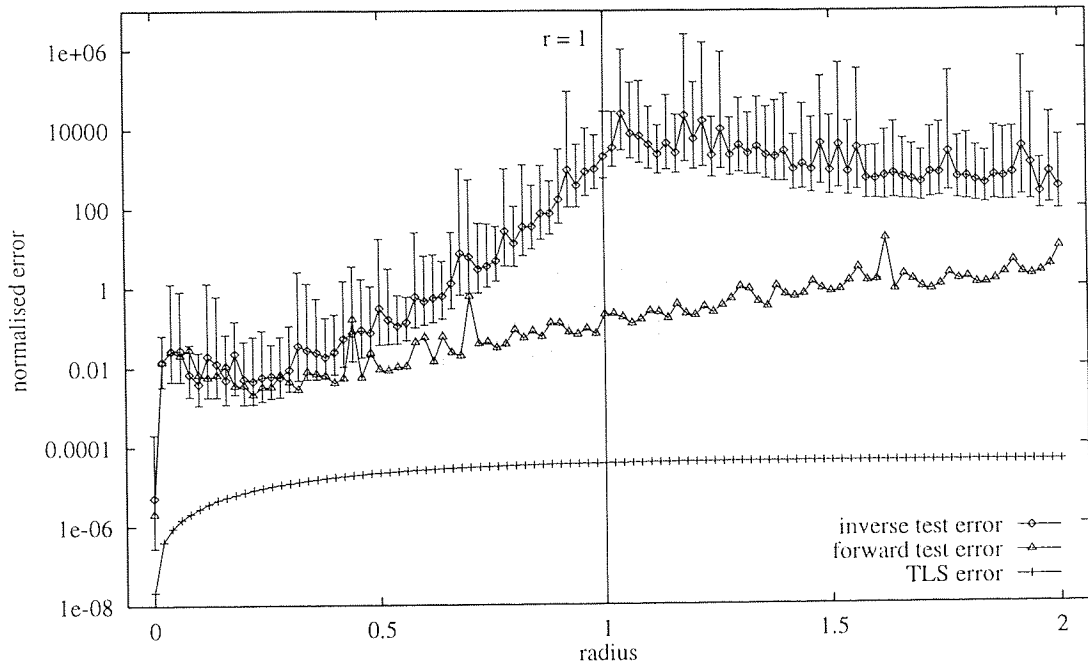
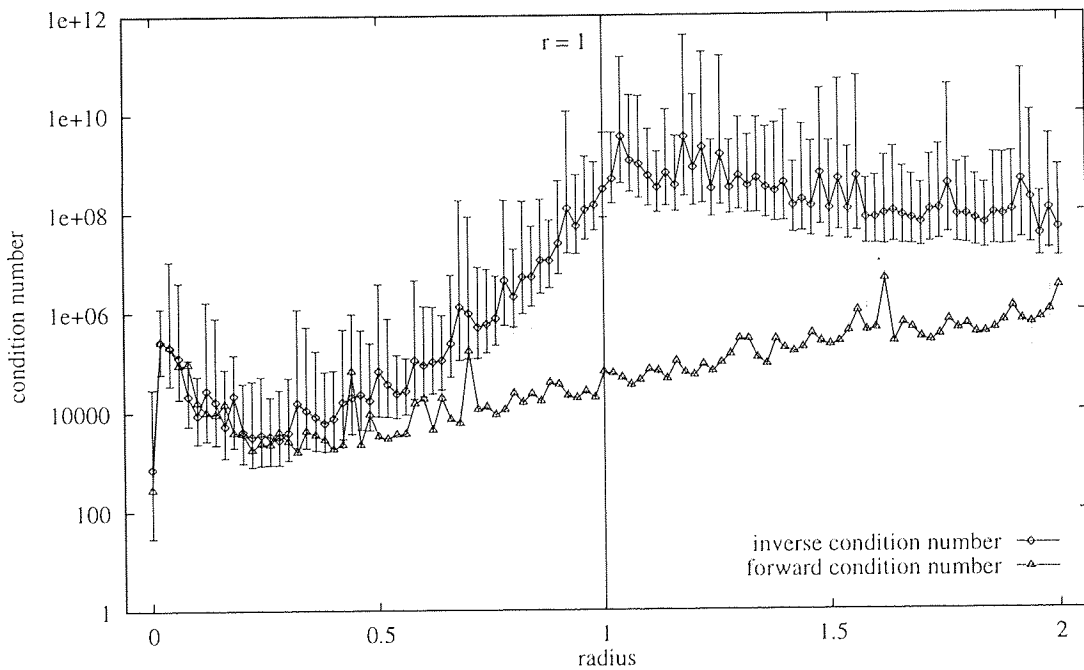


Figure 4.23 Empirically estimated upper and lower bounds \widehat{U}'_r and \widehat{L}'_r , obtained by moving the average over random seeds inside the calculation of extrema, superimposed on the analytical bounds U_r and L_r to which they correspond. (a) Both \widehat{U}'_r and \widehat{L}'_r jump from the expected values of $\widehat{U}'_0 \approx \widehat{L}'_0 \approx 1$ to initial peaks of $\widehat{U}'_r \sim 100$ (not visible on this truncated linear scale) and $\widehat{L}'_r \sim 10$, before asymptoting down towards U_r and L_r , respectively. For $r > 1$, the estimated upper bound is very close to its analytical value, while the estimated lower bound reaches a floor of $\widehat{L}'_r \sim 1$. (b) The peak value of \widehat{U}'_r is revealed by a log-linear scale.



(a)



(b)

Figure 4.24 Comparing the mean, normalised forward and inverse errors and condition numbers for the TLS approximation to $f_r: \mathcal{T}_{4,r} \rightarrow \mathcal{T}_{3,r}$, calculated over the test set. **(a)** On a log-linear scale, the inverse error ($\epsilon_{3,r}$) is small throughout the expected interval but becomes large for $r \gtrsim r^*$, lending some evidence for our theoretical knowledge that $\mathcal{T}_{4,r}$ is only diffeomorphic to $\mathcal{T}_{3,r}$ for $r < r^*$. In comparison, the product space and forward errors ($\epsilon_{\perp}^{(r)}$) and ($\epsilon_{4,r}$) are negligibly small. Both forward and inverse errors exhibit the characteristic TLS instability. The fitting and test errors are all but indistinguishable, even on a log-linear scale, and hence only the latter are plotted. **(b)** The condition numbers, respectively $\langle \kappa(Q_r) \rangle$ and $\langle \kappa(P_r) \rangle$, correspond closely to $\langle \epsilon_{4,r} \rangle$ and $\langle \epsilon_{3,r} \rangle$ as expected. Error bars denote a one-sided standard deviation in both directions; those on the forward test error and condition number have been lightened for clarity.

$$\epsilon_{4,r}^2 = \sigma_{3,r}^{-2} \sum_{i=1}^N \|\mathbf{W}_r(\mathbf{a}_i) - \mathbf{b}_i\|^2, \quad \epsilon_{3,\phi}^2 = \sigma_{4,r}^{-2} \sum_{i=1}^N \|\mathbf{W}_r^{-1}(\mathbf{b}_i) - \mathbf{a}_i\|^2 \quad (4.32)$$

where $\mathbf{a} = \varphi_{4,r}(\mathbf{x}) - \overline{\varphi_{4,r}}$ and $\mathbf{b} = \varphi_{3,r}(\mathbf{y}) - \overline{\varphi_{3,r}}$, with normalisation constants $\sigma_{4,r}$ and $\sigma_{3,r}$ calculated on $\varphi_{4,r}\mathcal{T}_{4,r} \subset \mathbb{R}^{200}$ and $\varphi_{3,r}\mathcal{T}_{3,\phi} \subset \mathbb{R}^{200}$, respectively (again, these should not be confused with their LS analogues). Finally, the TLS error of (3.36) is defined as

$$\epsilon_{\perp}^{(r)} = \sigma_r^{-2} \sum_{i=1}^N \|\mathbf{P}_r^{\top} \mathbf{a}_i + \mathbf{Q}_r^{\top} \mathbf{b}_i\|^2 \quad (4.33)$$

where $\sigma_r = \sigma_{4,r} + \sigma_{3,r}$ and, from (3.40), $\mathbf{W}_r = -\mathbf{P}_r \mathbf{Q}_r^{-1}$.

The expected values of these errors are plotted in figure 4.24, for $0 \leq r \leq 2$; once again, the error bars each correspond to two one-sided standard deviations (as described in section 4.1.2). The errors measured on training and test sets are once again indistinguishable so we only plot the latter here. On a log-linear scale in part (a) we again see $\langle \epsilon_{3,r} \rangle$ oscillating wildly, although on this scale $\langle \epsilon_{3,r < r^*} \rangle$ is negligible, as is $\langle \epsilon_{4,r} \rangle$ for all r . However, $\langle \epsilon_{3,r} \rangle$ does indeed rise significantly at about the critical value r^* , while $\langle \epsilon_{4,r} \rangle$ remains relatively small throughout, indicating—as expected—that \mathbf{f}_r is a diffeomorphism only for $r < r^*$. Both, however, suffer from the characteristic instability discussed in section 3.3.3. The product error $\langle \epsilon_{\perp}^{(r)} \rangle$, on the other hand, is small and smooth, rising from $\langle \epsilon_{\perp}^{(r)} \rangle \sim 10^{-8}$ to an asymptotic value of $\langle \epsilon_{\perp}^{(r)} \rangle \sim 10^{-4}$. Not surprisingly, all three curves meet at $r = 0$, where both $\mathcal{T}_{4,0}$ and $\mathcal{T}_{3,0}$ become diffeomorphic to the circle \mathcal{S}^1 . The expected values of the condition numbers $\kappa(\mathbf{Q}_r)$ and $\kappa(\mathbf{P}_r)$ are illustrated in figure 4.24(b); once again, before and after averaging, these closely mirror the corresponding errors, respectively $\epsilon_{4,r}$ and $\epsilon_{3,r}$, as predicted by the analysis in section 3.3.2.

Chapter 5

Maps on dynamical systems

We now consider the application of RBF maps to the detection of diffeomorphisms between embedded dynamical systems. We will be interested in maps of the form $f: \mathcal{M}_m \rightarrow \mathcal{M}_n$, where the compact subset $\mathcal{M}_m \subset \mathbb{R}^m$ is the image of a delay embedding $\Phi_{v,m}$, obtained from a measurement function $v: \mathcal{M} \rightarrow \mathbb{R}$ on the dynamical system (\mathcal{M}, ψ) , and f is the restriction to \mathcal{M}_m of the linear transformation $\mathcal{F}: \mathbb{R}^m \rightarrow \mathbb{R}^n$, as described in section 2.3, with $\mathcal{M}_n = \mathcal{F}\mathcal{M}_m$. Our task will be to determine whether or not f is a diffeomorphism—that is, whether or not $\mathcal{F} \circ \Phi_{v,m}$ is a ‘filtered’ embedding on (\mathcal{M}, ψ) , as also defined in section 2.3.

Since \mathcal{F} is a linear map we already know that f is a function, so we expect to find—as demonstrated in the previous chapter—an arbitrarily good LS RBF approximation \hat{f} to f . This being the case, we will not usually attempt to fit f explicitly but concentrate instead on the properties of an RBF approximation to its inverse. As a consequence of the delay structure induced by $\Phi_{v,m}$ in the domain of f , it will not be necessary to construct an explicit approximator $\widehat{f^{-1}}$ in the examples to follow. Indeed, such a procedure may not even be physically appropriate, as it will *not* generally exhibit the shift property itself. Instead, in the interest of computational efficiency, we will fit one or more individual components $(f^{-1})_j$ of f^{-1} , any of which could be used, in principle, to construct an inverse for f with delay structure intact. We will therefore rely solely on the LS RBF map in the experiments to follow. Our use of prior knowledge about the structure of f and its inverse to restrict our attention to individual components of f^{-1} rules out the use of the Lipschitz analysis of section 3.2.4, and also of the somewhat less reliable TLS method, in the experiments described in this chapter.

In the remainder of this chapter we describe four related applications of this constrained form of diffeomorphism detection. The first involves the determination of a minimum embedding dimension for a dynamical system of unknown topological dimension, and requires the construction of a time series

predictor. The second deals with projections of an embedded system into a singular subspace, and attempts to similarly identify the minimum basis set required for such a projection to be a filtered embedding. The final two incorporate the more general form of filtered embedding in which a FIR filter is applied to the measurement function before embedding the system under investigation. Two applications of this form are studied: the first in detecting the existence of periodic orbits in the system's state space and the second in applying the results of this procedure to the separation of encoded messages from additive chaos.

5.1 Determination of a minimum embedding dimension

In section 2.2 we defined the method of delays in terms of the topological dimension d of the manifold \mathcal{M} on which the state of the dynamical system (\mathcal{M}, ψ) evolves under $\psi: \mathcal{M} \rightarrow \mathcal{M}$, namely that $\Phi_{v,m}: \mathcal{M} \rightarrow \mathbb{R}^m$ is generically an embedding if $m > 2d$. In practice, however, d will generally be an unknown quantity, so how can we be sure that we have chosen m large enough to embed \mathcal{M} in the first place? We clearly need an experimental method, involving no prior knowledge of the topology of \mathcal{M} , with which to determine the minimum value of m —say $m = m^*$ —necessary for $\Phi_{v,m}$ to be an embedding on \mathcal{M} .

Several solutions to this problem have been suggested in the literature. One such approach, by Broomhead, Jones and King [5], relies on a 'local' analysis of the manifold $\mathcal{M}_m = \Phi_{v,m}\mathcal{M}$ to estimate d directly (\mathcal{M}_m is assumed to have been reconstructed in a sufficiently high-dimensional space $\mathbb{R}^{m \gg d}$ for $\Phi_{v,m}$ to be an embedding on \mathcal{M}). In this approach, the singular spectrum of a set of points contained within an open ball \mathcal{B}_r in \mathcal{M}_m is calculated (via SVD) as the radius r of \mathcal{B}_r is increased from zero. For r small enough, the authors find that a subset of the singular spectrum scales linearly with r , while the remaining singular values stay roughly constant within the noise floor identified by a global SVD of \mathcal{M}_m . As r increases beyond a local, critical value these latter singular values begin to grow as well, as the curvature of \mathcal{M}_m makes its effect felt within \mathcal{B}_r . The number of singular values which scale linearly below that critical value of r is taken to be an estimate of d , and a suitable value for m^* follows directly from Takens' theorem. (Strictly speaking, of course, the value of m^* so calculated is a *sufficient*, rather than a *minimum* embedding dimension, and it may actually be possible to embed (\mathcal{M}, ψ) with $d \leq m \leq m^*$, as previously noted.)

In another approach, related to that of Pecora et al which we described in chapter 1, Kennel, Brown and Abarbanel [22] apply the method of 'false nearest neighbours', in which they attempt to successively 'unfold' the reconstructed attractor $\mathcal{M}_m \in \mathbb{R}^m$ until it is determined to be embedded in \mathbb{R}^m for a large enough m . In order to make this determination, the distance between each point \mathbf{x} and its nearest neighbour in \mathcal{M}_m is calculated for increasing values of m . This process is continued for as long as one or more of these nearest neighbour distances continues to grow, by a sufficiently large factor, with the addition of each successive dimension, yielding a direct, empirical estimate of m^* . A similar method, espoused by

Kaplan [21], consists of calculating an equivalent distance measure between neighbouring points *solely* in the additional dimension: in other words, between the scalar elements of the time series itself.

These approaches are based on a direct analysis (via delay map) of the time series itself. The empirical solution which we explore in this section relies instead on the analysis of successive LS RBF approximations to maps between the images $\mathcal{M}_m \subset \mathbb{R}^m$ and $\mathcal{M}_{m+1} \subset \mathbb{R}^{m+1}$ of the delay maps $\Phi_{v,m}$ and $\Phi_{v,m+1}$, for increasing m . We make use of the fact that if $\Phi_{v,m}$ and $\Phi_{v,m+1}$ both embed \mathcal{M} then the composite map $f_m: \mathcal{M}_{m+1} \rightarrow \mathcal{M}_m$, defined by $f_m = \Phi_{v,m} \circ \psi^{-1} \circ \Phi_{v,m+1}^{-1}|_{\mathcal{M}_{m+1}}$, is necessarily a diffeomorphism. (Our motive for incorporating an inverse iterate of ψ will shortly be made clear.) Conversely, if we can find some m^* such that f_m is a diffeomorphism for all $m \geq m^*$, but not for $m < m^*$, then we can infer that m^* is the minimum reconstruction dimension necessary to embed \mathcal{M} , using the time series in question. We must, in principle, consider all values of $m \geq m^*$ because it is conceivable that a particular f_m , so defined, might be a diffeomorphism even though neither $\Phi_{v,m}$ nor $\Phi_{v,m+1}$ actually embed \mathcal{M} . Consider, for example, the experiment described in section 4.2, in which we examined the effect of the parameter r on whether or not a 2-torus was embedded in three or four dimensions: for $r = 0$ the image of this torus was a circle in both \mathbb{R}^3 and \mathbb{R}^4 , neither of which were diffeomorphic to the torus, but both of which were unarguably diffeomorphic to each other.

For implementational purposes, we can view f_m as the restriction to \mathcal{M}_{m+1} of the linear transformation $\mathcal{F}_m: \mathbb{R}^{m+1} \rightarrow \mathbb{R}^m$ defined, following section 2.3.1, by the m by $(m+1)$ matrix \mathbf{F}_m formed by removing the first row from the $(m+1)$ by $(m+1)$ identity matrix \mathbf{I}_{m+1} . We certainly do not need to find a LS RBF approximation to f_m in order to determine that it is a one-to-one map. Its inverse, on the other hand—if it exists—is a map taking each $\mathbf{y}_i \in \mathcal{M}_m$, with components $\mathbf{y}_i = (v_i, \dots, v_{i-(m-1)})^T$, to an $\mathbf{x}_{i+1} \in \mathcal{M}_{m+1}$ defined by $\mathbf{x}_{i+1} = (v_{i+1}, v_i, \dots, v_{i-(m-1)})^T$. All but one component of this map is present in its domain, so it is unnecessary to construct an explicit RBF approximation to f_m^{-1} . Instead, we concentrate on the first component $w_m \equiv (\mathbf{f}_m^{-1})_1$, and investigate the question of whether or not f_m is a diffeomorphism by constructing its RBF approximation $\widehat{w}_m: \mathcal{M}_m \subset \mathbb{R}^m \rightarrow \mathbb{R}$, which is a one-step predictor for the time series itself, with $\widehat{v}_{i+1} = \widehat{w}_m(\mathbf{y}_i)$. In other words, we are asking whether or not we can predict the next element of the time series from a delay reconstruction using the previous m elements. (If f_m had been defined without an inverse iterate of ψ then we would be fitting a one-step-behind predictor instead.)

Although we could, in principle, treat w_m as an induced measurement function on \mathcal{M}_m , using the delay structure in \mathcal{M}_{m+1} to construct the approximation $\widehat{f}_m^{-1} \equiv \mathbf{F}_m^T + \mathbf{e}_1 \widehat{w}_m$ to f_m^{-1} , where $\mathbf{e}_1 = (1, 0, \dots, 0)^T$ is an $(m+1)$ -dimensional unit vector, we do not actually need to take this final step in order to determine whether or not f_m is a diffeomorphism. We do, however, take this approach in Potts and Broomhead [32], in which we use the first m components of \widehat{f}_m^{-1} , so defined, to construct an approximator $\widehat{\psi}_m: \mathbb{R}^m \rightarrow \mathbb{R}^m$ for ψ_m and then investigate the stability of $\widehat{\psi}_m$ under iteration by calculating its characteristic exponents [12, 30].

In the following sections we attempt to calculate minimum embedding dimensions for the the Ikeda and Hénon systems, the Lorenz system (integrated with step sizes of 0.1 and 0.01) and the laser system. In analogy with equation (3.5) we write the normalised fitting (or test) error ϵ_m as

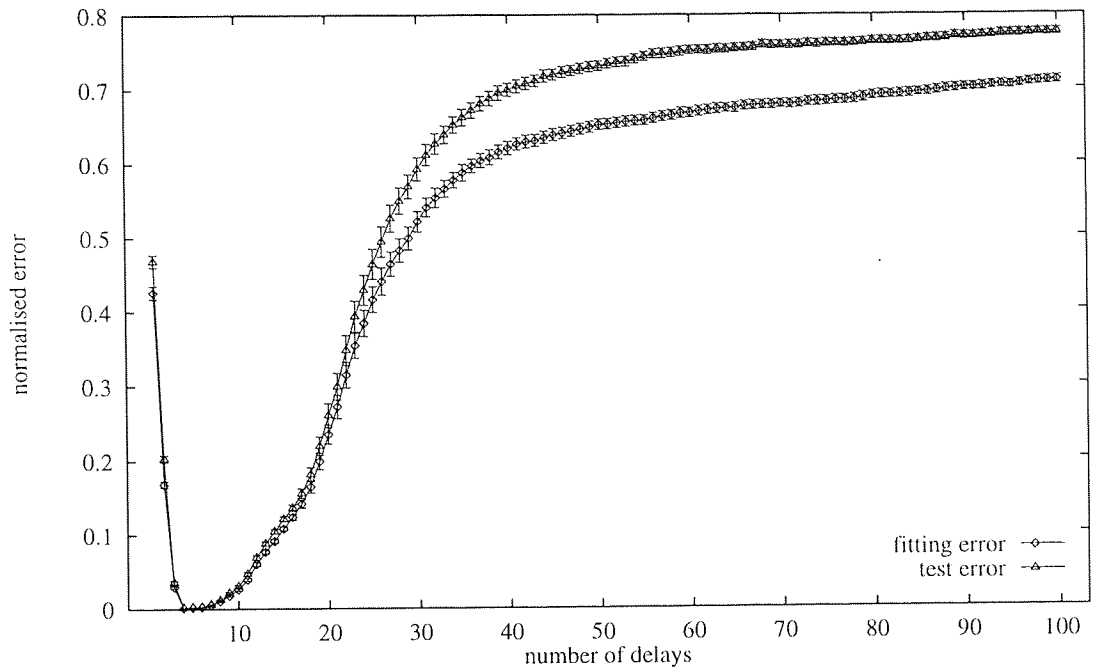
$$\epsilon_m^2 = \sigma_v^{-2} \sum_{i=1}^N \|\widehat{w}_m(\mathbf{y}_i) - v_{i+1}\|^2 \quad (5.1)$$

where the normalisation constant σ_v^2 is N times the variance of the time series itself, calculated over the training set. As usual, in order to eliminate as many sources of uncertainty as possible, in each experiment we will be plotting the expected value $\langle \epsilon_m \rangle$ of the error obtained from 500 sets of $p = 200$ repulsive centers, with cubic basis functions, and with repulsive seeds selected at random, without replacement, from a training set of $N = 2000$ points in \mathcal{M}_m .

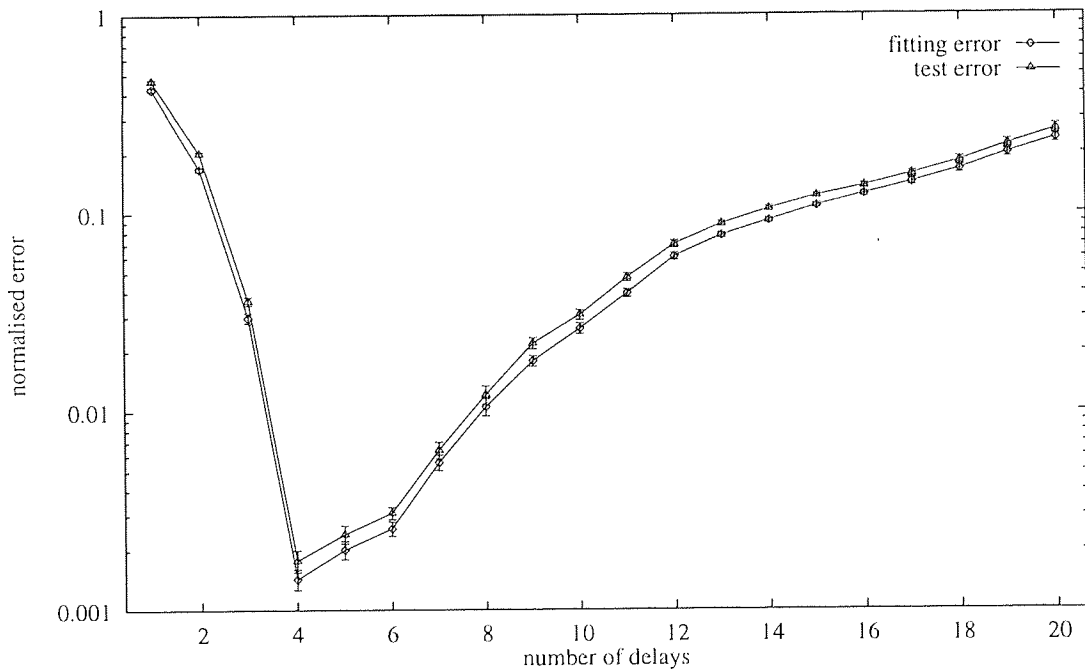
In analysing the results of these experiments we must keep in mind the fact that there are other contributions to the LS error than whether or not the underlying manifold has been embedded in \mathbb{R}^m . An obvious example is that, as the dimensionality m of the domain of w_m increases, so the ‘coverage’ of that space by a fixed number of centers decreases, leading us to expect both fitting and test errors to asymptote to 1 in the limit of $m \rightarrow \infty$. More generally still, we have no a priori reason even to believe that—neglecting this asymptotic effect—the error will vary monotonically for a given choice of manifold, or training and test sets (although the dependence on a specific RBF model, at least, has hopefully been taken into account by the average over random repulsive seeds). For instance, the presence of periodicities in $\{v_i\}$ may adversely affect the fitting error for some values of m more than others. In certain case, therefore, the precise identification of a suitable value for m^* may require an element of judgement.

5.1.1 Embedding the Ikeda system

We begin with the Ikeda map (2.1), for which Takens predicts an embedding dimension of $m^* = 5$. In figure 5.1 we plot, versus m , the mean error $\langle \epsilon_m \rangle$ calculated over both training and test sets, for a delay reconstruction from the time series plotted in figure 2.3(a). In part (a) of this figure, for $1 \leq m \leq 100$, we illustrate the asymptotic behaviour of these errors, both of which rise steadily with increasing m and appear almost saturated at $\langle \epsilon_{100} \rangle \approx 0.7$ and 0.8, respectively. We attribute this upwards trend to the result of constructing an RBF map with a fixed number of centers in an increasingly high-dimension domain, as discussed above. We plot the first twenty values of $\langle \epsilon_m \rangle$ on a log-linear scale in part (b), in which we see strong empirical evidence for a minimum embedding dimension of $m^* = 4$. The (relatively) large error at $m = 3$ is probably due to the presence of a self-intersecting point-set, visible toward the lower half of figure 2.6(a), in the reconstructed attractor; this self-intersection is eliminated by the inclusion of a fourth delay. Despite the fact that the errors under analysis were obtained from a full-rank LS RBF approximation, no significant over-fitting is revealed in figure 5.1. That the minimum occurs at $m = 4$,



(a)



(b)

Figure 5.1 Establishing a minimum embedding dimension for the Ikeda system by plotting the mean time series prediction error (ϵ_m), versus m , calculated over both training and test sets. (a) With $1 \leq m \leq 100$ both curves achieve an early minimum before presumably asymptoting to 1; (b) on a log-linear scale, a close up of the interval $1 \leq m \leq 20$ provides an estimated minimum embedding dimension of $m^* = 4$. A comparison of fitting and test errors reveals no significant over-fitting. Error bars denote one standard deviation in each direction.

and not $m = 5$, as we might expect from Takens, is of no great consequence: as already noted, Takens' theorem applies generically, not specifically.

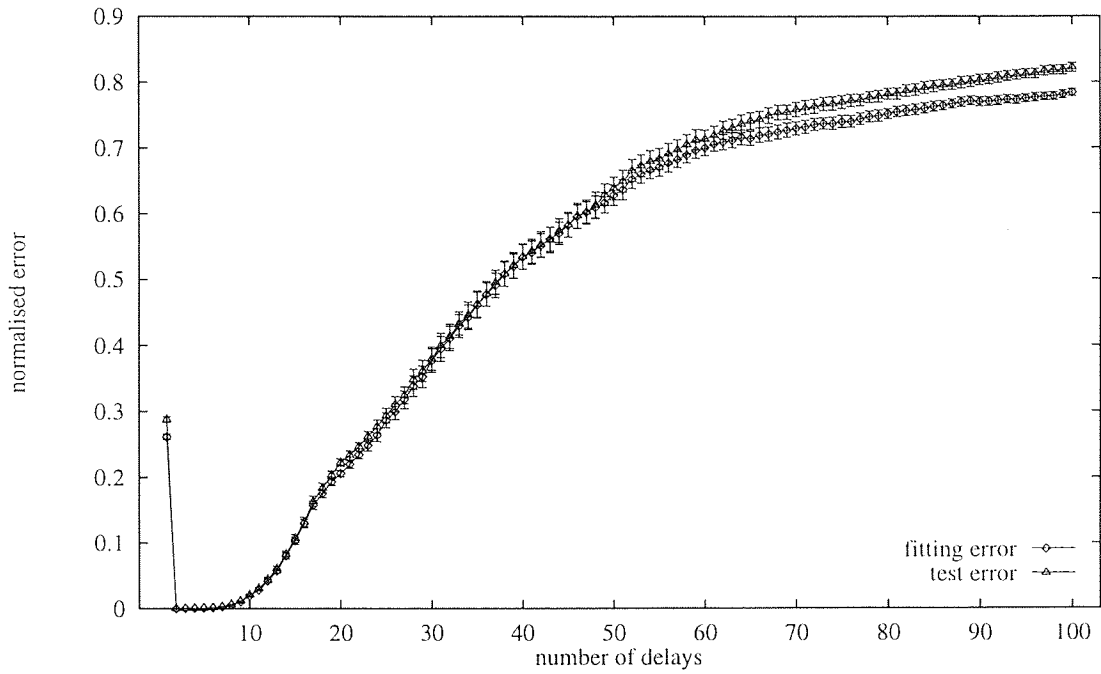
5.1.2 Embedding the Hénon system

The same approach, applied to the Hénon map (2.2), via the time series illustrated in figure 2.3(b), yields the error curves plotted in figure 5.2. Once again, in part (a), we see the expected asymptotic behaviour, but in part (b), showing a close-up of the interval $1 \leq m \leq 20$ on a log-linear scale, we now see a more clearly indicated minimum embedding dimension of $m^* = 2$. The attractor embedded in \mathbb{R}^2 has already been plotted, in figure 2.6(b). This result supports the analytical discussion of section 2.2.1, in which we showed, through equation (2.9), that the evolution of the Hénon map is exactly determined by two delays. Once again, no significant over-fitting is observed.

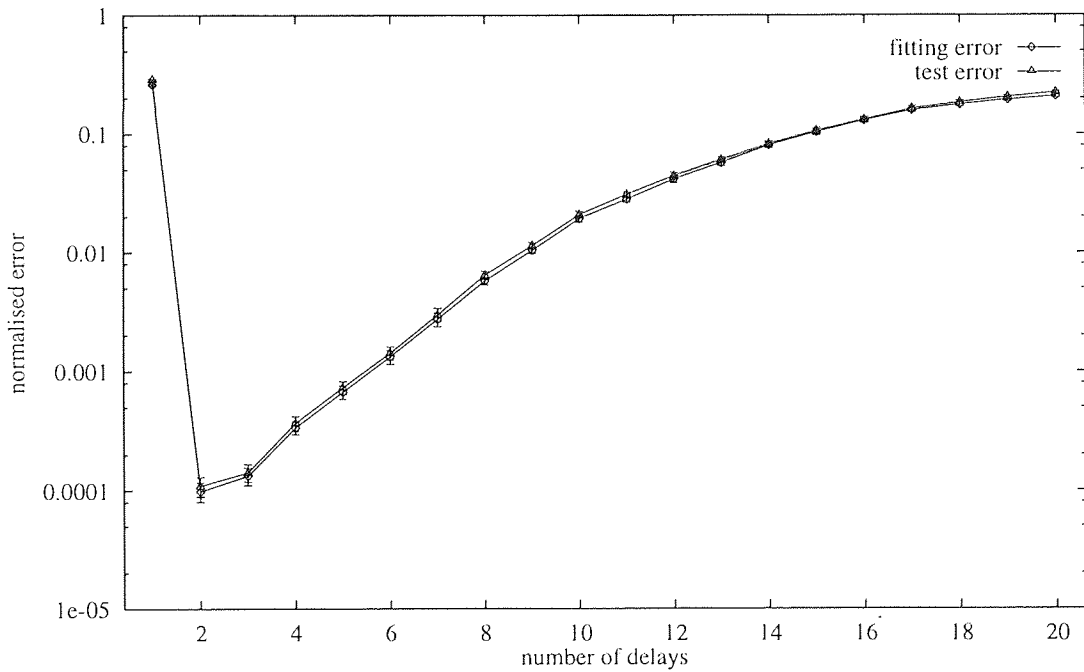
5.1.3 Embedding the Lorenz system

The errors arising from the 0.01-step Lorenz system, reconstructed from the time series of figure 2.4(a), are plotted in figure 5.3. In part (a) of this figure both errors remain relatively small, compared to those calculated on the Ikeda and Hénon systems, as a result of the extremely short integration step (ie. sampling interval) $\tau = 0.01$ with which the time series in question was obtained: the function $w_m: \mathcal{M}_m \rightarrow \mathbb{R}$ is nearly linear at this scale. It is, however, over-fitting the time series in question to a correspondingly larger degree. Because the predictor in this case is such a simple one, the minimum embedding dimension of $m^* = 2$ indicated by part (b) is significantly smaller than the value of $m \geq 5$ suggested by Takens (based on the fractal dimension of 2.06 estimated by Lorenz [24]): the 'true' minimum is effectively hidden by local fluctuations in error between individual RBF fits. Interestingly, these fluctuations appear to be driven, to a limited extent, by one or more unstable periodic orbits in the Lorenz attractor.

That we can so easily embed the 0.01-step Lorenz system is largely due to the fact that the time series in question is noise-free (neglecting quantisation errors). This will not usually be true in practice, so we have also examined the time series obtained by 'corrupting' the original with an additive stochastic component, normally distributed with a standard deviation approximately one tenth that of the uncorrupted time series. Referring back to figure 2.7(a), it is clear that a noise component of this size must make a unit-lag embedding impossible, due to the tendency of the reconstructed attractor to 'hug' the diagonal in \mathbb{R}^m . We plot the errors in figure 5.4. In part (a) of this figure we immediately notice that both are substantially larger than their noise-free equivalents of figure 5.3; part (b) is similarly inconclusive. It is interesting to note that the degree of over-fitting exhibited by figure 5.3 is substantially larger than in the noise-free case; this is because the LS RBF map, given enough degrees of freedom, will model variations due to noise in the training set and hence fail to generalise well on the test set; this time series would be a

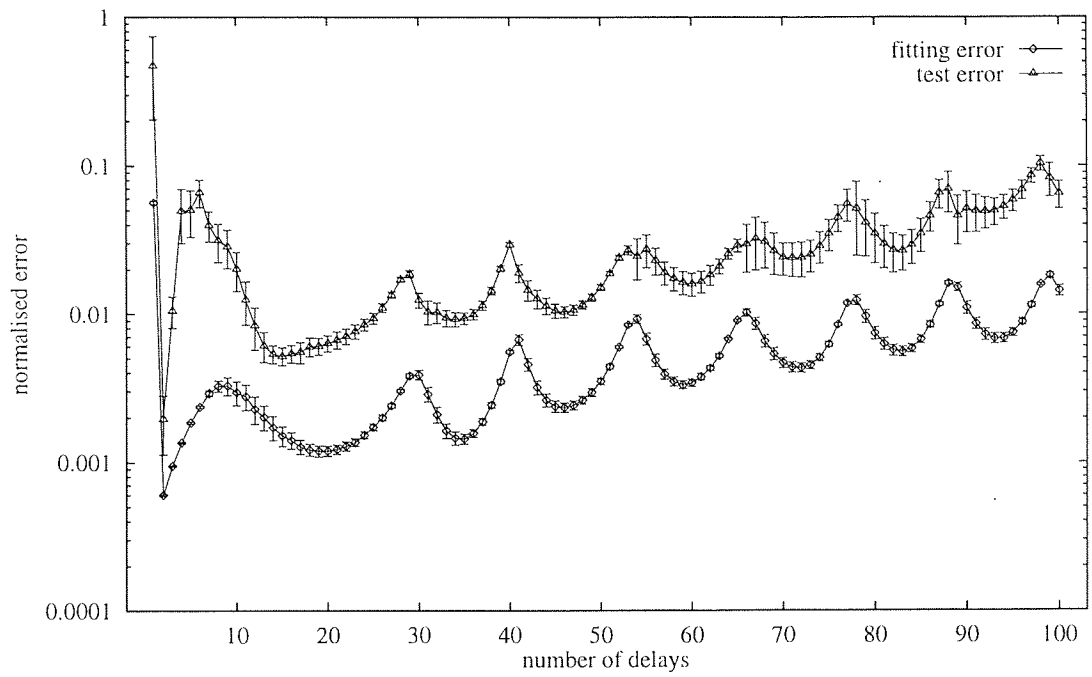


(a)

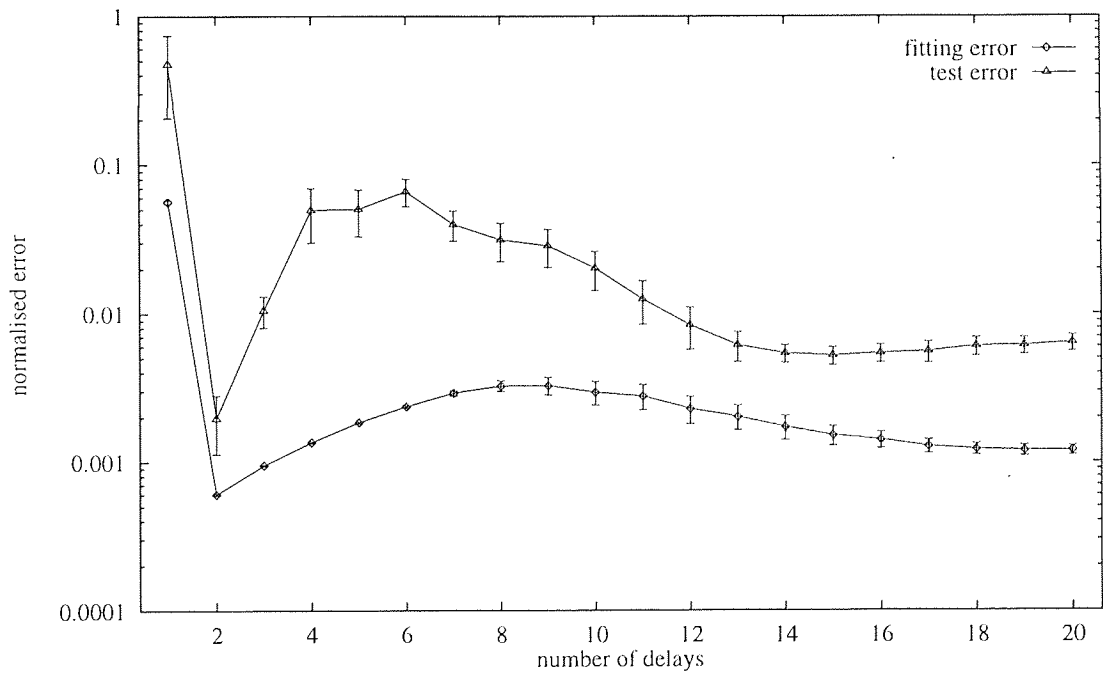


(b)

Figure 5.2 Establishing a minimum embedding dimension for the Hénon system by plotting the mean, normalised prediction error (ϵ_m), versus m , over training and test sets. (a) Variation of the reconstruction dimension in the interval $1 \leq m \leq 100$ reveals properties similar to those of the Ikeda system, such as the asymptotic convergence to $\langle \epsilon_m \rangle \sim 1$; (b) a close up of $1 \leq m \leq 20$, on a log-linear scale, clearly indicates a minimum embedding dimension of $m^* = 2$. No significant over-fitting is observed; error bars denote one standard deviation in each direction.

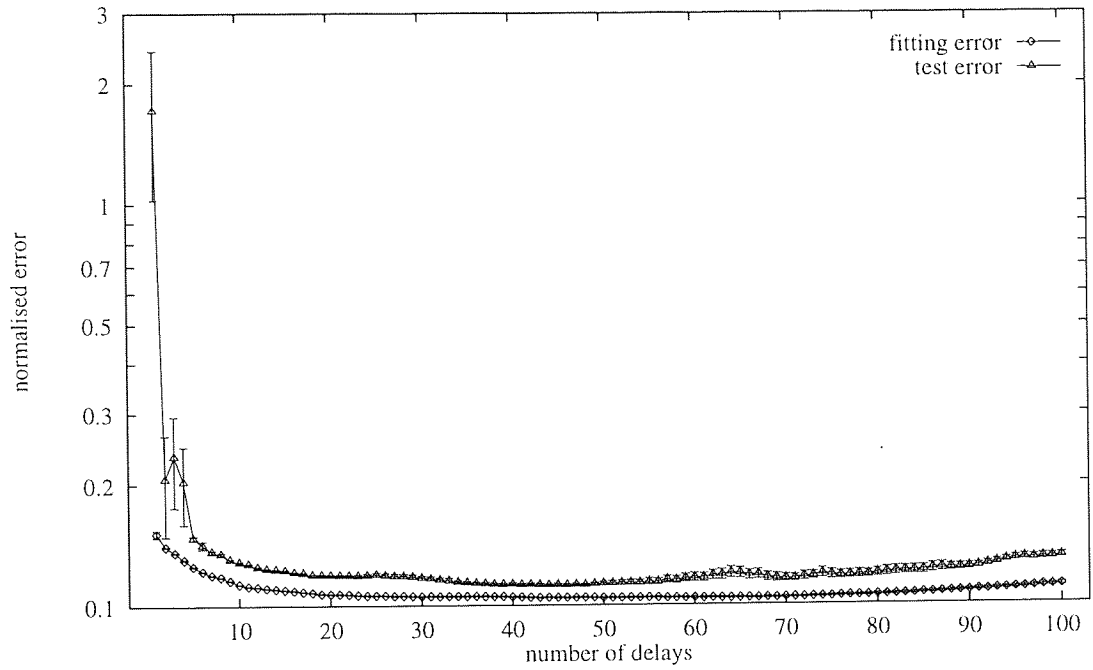


(a)

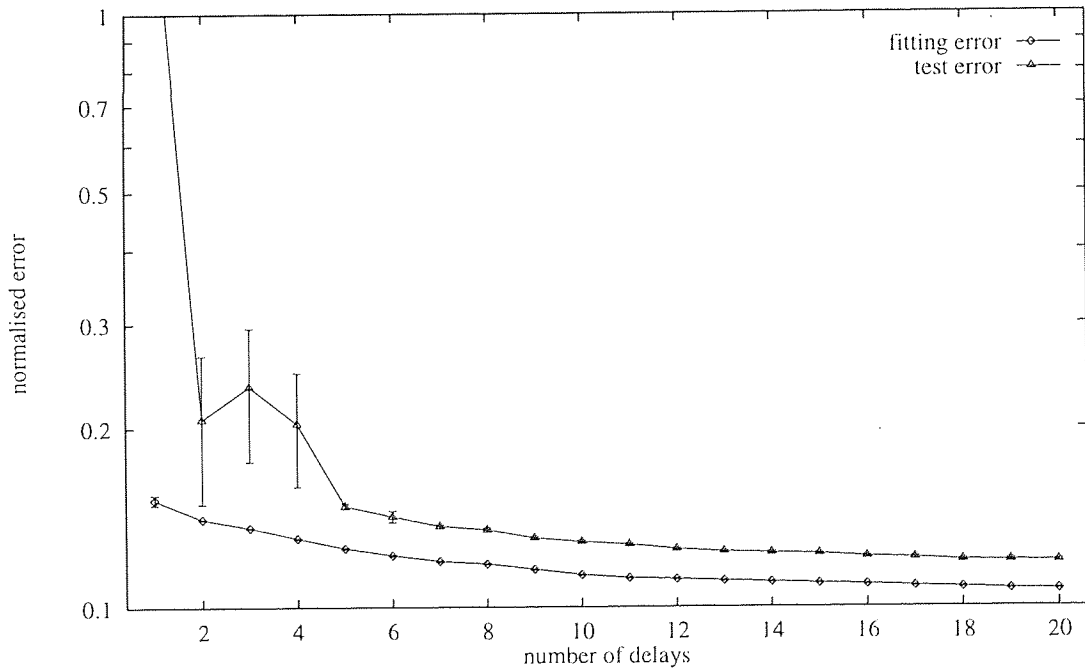


(b)

Figure 5.3 Establishing a minimum embedding dimension for the 0.01-step Lorenz system by plotting the mean, normalised prediction error (ϵ_m), versus m , over training and test sets. (a) At $m = 100$, on a log-linear scale, neither fitting nor test error has begun to saturate, but both reveal an interesting periodicity of extremely small amplitude; (b) in close-up an embedding dimension of $m^* = 2$ is clearly indicated by both errors. Error bars denote one standard deviation in each direction.



(a)



(b)

Figure 5.4 Establishing a minimum embedding dimension for the 0.01-step Lorenz system, generated from a noisy time series, by plotting the mean, normalised prediction error (ϵ_m), versus m , over training and test sets. (a) Neither train nor test error appears to reach a sufficient minimum, and it is not clear that the system is embedded for any m ; (b) a close-up merely confirms this impression. Error bars denote one standard deviation in each direction.

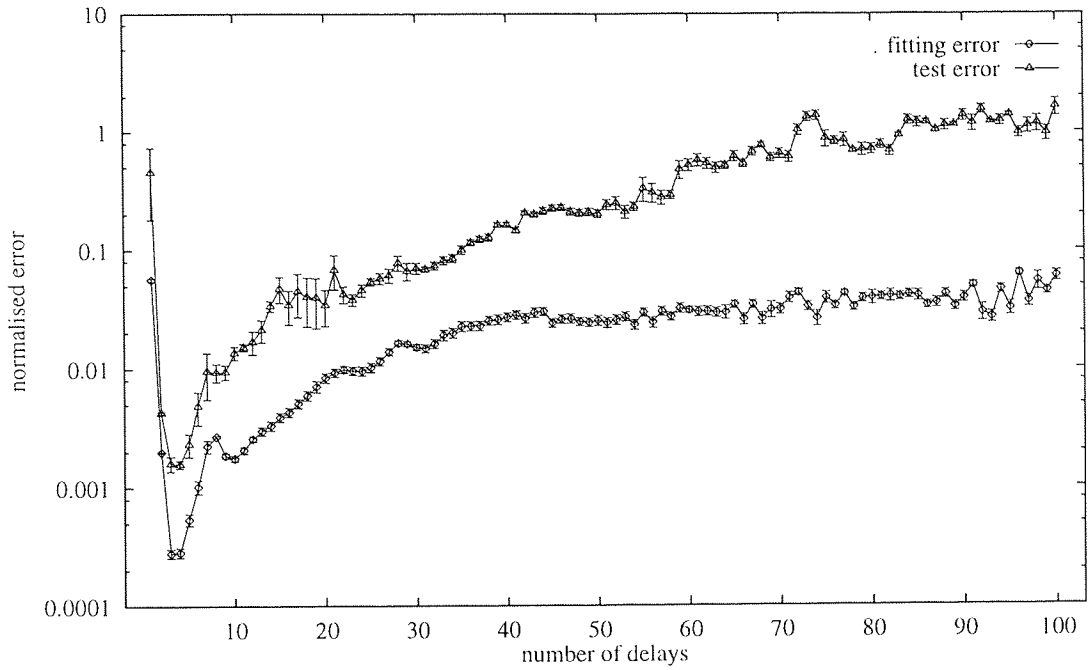
good candidate for the rank-reduction methods described in section 3.2.3.

One way in which to avoid the situation in which the variance of the reconstructed attractor is concentrated along the diagonal is to sub-sample the time series in question with the lagged delay map $\Phi_{v,m,\tau}: \mathcal{M} \rightarrow \mathbb{R}$ defined in equation (2.16); we have already illustrated, in figure 2.8(a), the result of embedding the 0.01-step Lorenz system in three dimensions with a lag of $\tau = 10$ samples. In figure 5.5 we plot the errors arising from an RBF approximation to the Lorenz time series obtained in this manner. In part (a) of this figure we once again see the test error asymptote to $\langle \epsilon_m \rangle \sim 1$, although the fitting error remains relatively small over the interval plotted, indicating that some over-fitting is going on. A close-up, in part (b), strongly indicates a minimum embedding dimension of $m^* = 4$, although we might argue that $m^* = 3$ is sufficient. Again, we have come up with an estimate substantially smaller than the value of $m^* = 5$ suggested by Takens' theorem; this is merely an indication that the first component of the Lorenz map is a relatively benign measurement function—for the purposes of delay reconstruction—preserving intact the two unstable fixed points and their nearby orbits, as demonstrated by a comparison of figure 2.8(a) with the original attractor in 2.2.

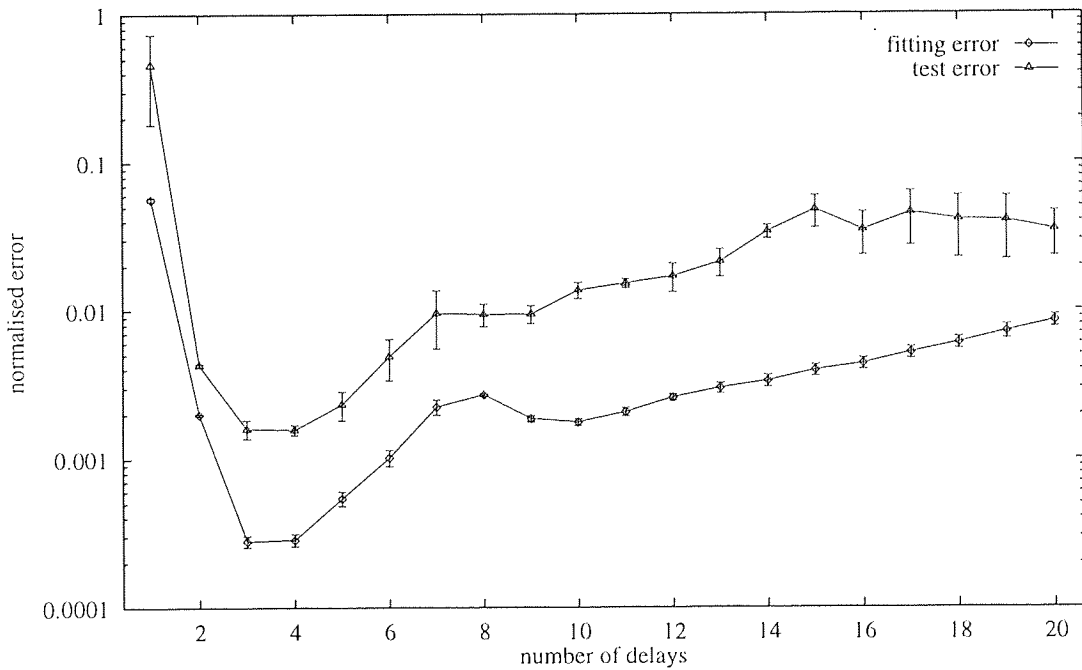
On duplicating this experiment with the noise-corrupted time series already described we obtained the errors plotted in figure 5.6. We are once more unable to identify a minimum embedding dimension, as the fitting and test errors are nowhere sufficiently small to indicate a good predictor for the time series; the over-fitting is severe. Although we might expect the trajectory obtained from a lagged delay embedding to be less susceptible to the presence of stochastic noise than a lagged embedding, because it occupies a larger volume in pseudo-phase space, the errors obtained in this experiment are of the same order as those obtained with a sample lag of $\tau = 1$. We ascribe this result to a trade-off between that effect and the extreme predictability of the finely-sampled time series.

5.1.4 Embedding the laser system

Finally, in figure 5.7, we consider the experimental time series of figure 2.5(b). As discussed in section 2.2, this measurement function is thought to correspond to the square of the first component of the Lorenz map, and therefore results in a reconstructed attractor whose unstable fixed points are superimposed, and we therefore do not expect to be able to embed this system. Nevertheless, the time series does possess a significant degree of predictability, as indicated by the relatively small test error at $m = 9$. Compared to the 0.1-step Lorenz system, however—as illustrated in figure 5.5—we see that both fitting and test errors are some two orders of magnitude larger than might otherwise be expected were the delay map truly an embedding. This result illustrates the potential pitfalls of the predictive approach to establishing the existence of a minimum delay embedding, in that it is not entirely clear how well we must approximate the induced measurement function w_m in order to justify the claim that f_m is a diffeomorphism.

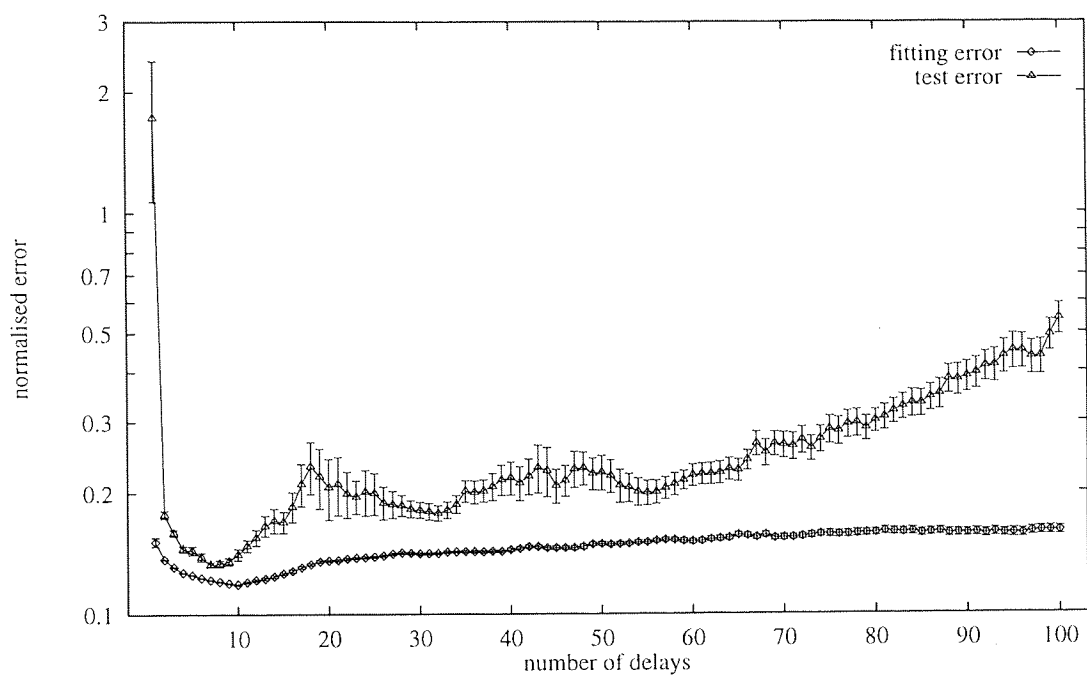


(a)

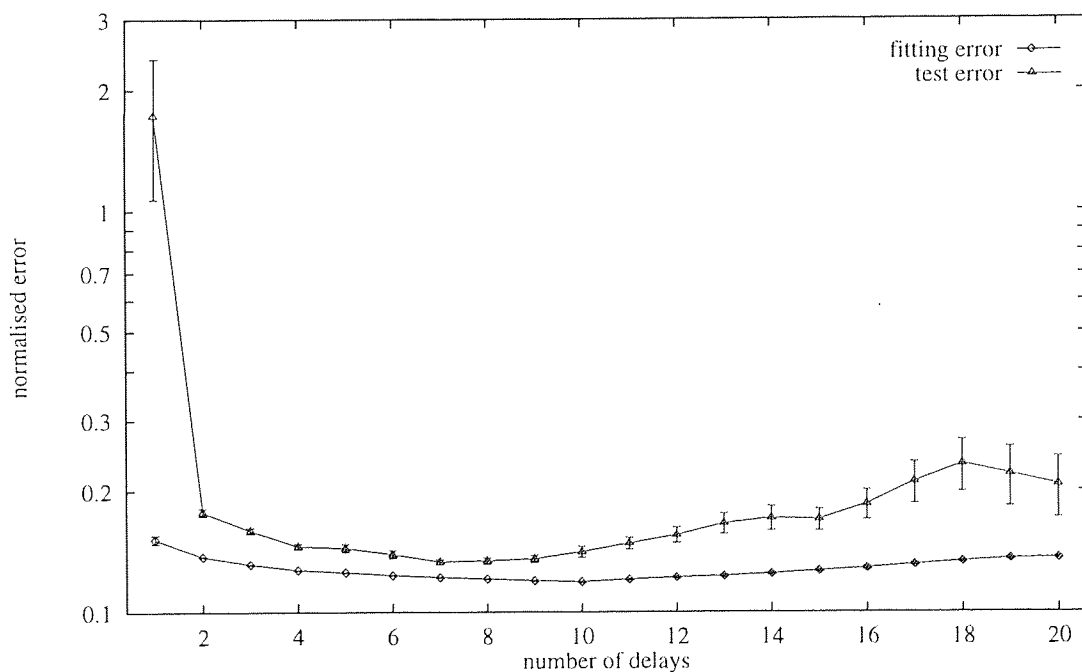


(b)

Figure 5.5 Establishing a minimum embedding dimension for the 0.01-step Lorenz system, reconstructed with a lag of $\tau = 10$, by plotting the mean, normalised prediction error $\langle \epsilon_m \rangle$, versus m , over training and test sets. (a) On a log-linear scale there is a marked separation of fitting and test errors as $m \rightarrow 100$; (b) in a close-up of this plot the test error appears to indicate an embedding dimension of $m^* = 4$, although the decrease in error between $m = 3$ and $m = 4$ is arguably small. Error bars denote one standard deviation in each direction.

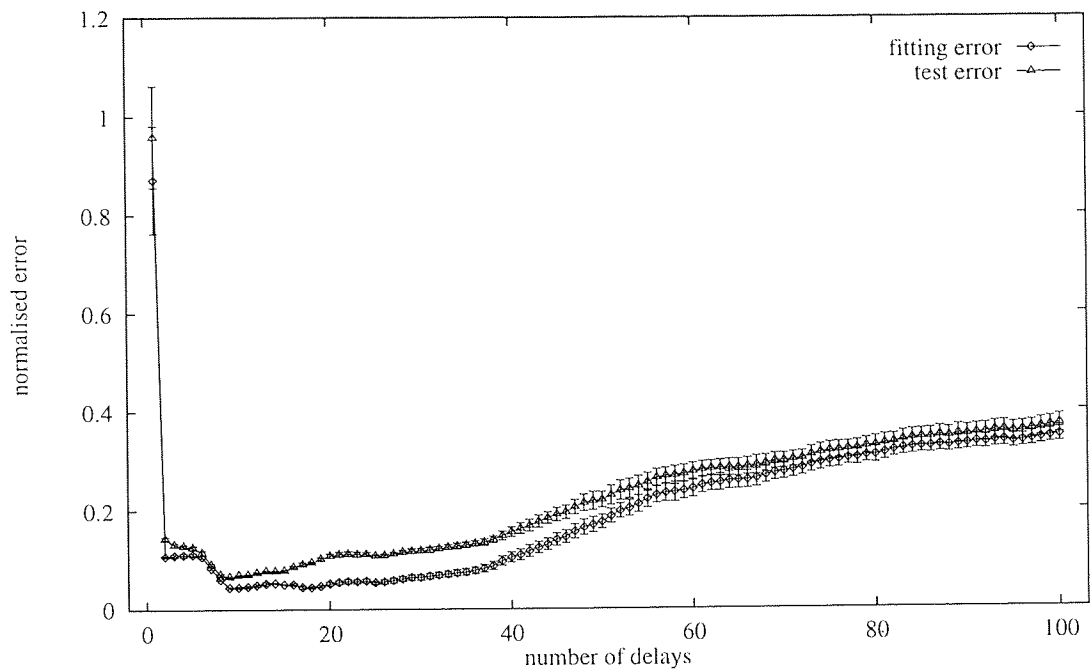


(a)

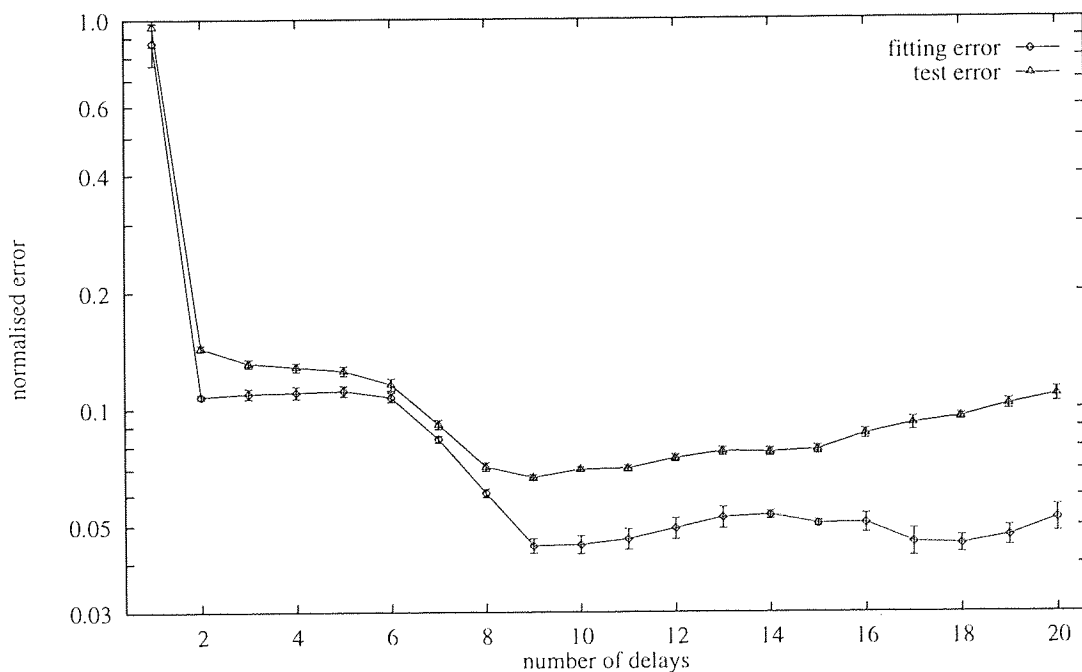


(b)

Figure 5.6 Establishing a minimum embedding dimension for the 0.01-step Lorenz system, reconstructed with a lag of $\tau = 10$ from a noisy time series, by plotting the mean, normalised prediction error (ϵ_m), versus m , over training and test sets. (a) In this case we again see an increasing separation of fitting and test errors with $m \gtrsim 8$, with serious over-fitting clearly visible; (b) in close-up we might set m^* anywhere in the interval $2 \leq m^* \leq 8$. Error bars denote one standard deviation in each direction.



(a)



(b)

Figure 5.7 Establishing a minimum embedding dimension for the laser system by plotting the mean, normalised prediction error $\langle \epsilon_m \rangle$, versus m , over training and test sets. (a) The fitting and test errors are just beginning to saturate at $m = 100$, with no significant over-fitting in evidence; (b) a close-up appears to indicate that $m^* = 9$ is a minimum dimension for successful reconstruction of this system, although in actual fact we already know that the measurement function in question is not generic, in the sense of Takens' theorem. Error bars denote one standard deviation in each direction.

5.2 Constructing a singular subspace

In attempting to predict the Lorenz time series in the previous section we found that the presence of stochastic noise affected not only our ability to embed the system in the first place, but also the generalisation ability of the predictor. In section 2.3.2 we discussed the use of a linear transformation $\mathcal{F}_{m,n}: \mathbb{R}^m \rightarrow \mathbb{R}^n$, in the form of an n by m matrix $\mathbf{F}_{m,n} = (\mathbf{v}_1, \dots, \mathbf{v}_n)^\top$ of singular vectors $\mathbf{v}_k \in \mathbb{R}^m$, principal components of the distribution in \mathbb{R}^m , in reducing the impact of stochastic noise on a delay embedding $\Phi_{v,m}: \mathcal{M} \rightarrow \mathbb{R}^m$. If $\mathcal{M}_m = \Phi_{v,m} \mathcal{M}$ is diffeomorphic to \mathcal{M} then provided $n > 2d$, and for generic choices of $\mathcal{F}_{m,n}$, we also expect $\mathcal{M}_{m,n} = \mathcal{F}_{m,n} \circ \Phi_{v,m} \mathcal{M}$ to be diffeomorphic to \mathcal{M} .

We will now attempt to test the viability of this method by constructing the map $f_{m,n}: \mathcal{M}_m \rightarrow \mathcal{M}_{m,n}$ defined by $f_{m,n} = \mathbf{F}_{m,n} \circ \Phi_{v,m} \circ \psi^{-1} \circ \Phi_{v,m}^{-1}|_{\mathcal{M}_m}$. If we can show that there exists a minimum singular subspace dimension $n^* \leq m$ such that $f_{m,n}$ is a diffeomorphism for $n \geq n^*$, but not for $n < n^*$, then we will consider the test to have succeeded. As before, since $\mathcal{F}_{m,n}$ is a linear map we already know that $f_{m,n}$ is a function, so we can restrict our analysis to its inverse. Although in principle—thanks to the delay structure in \mathcal{M}_m —we could construct an approximation to $f_{m,n}^{-1}$ from an RBF model of any one of its components, in this case such a procedure would not be sufficient to ensure that $f_{m,n}^{-1}$ is itself a function: the projection $\mathcal{F}_{m,n}$ might preserve one or more components of its domain more or less intact in \mathbb{R}^n , making the task of estimating that component an artificially simple one (as, trivially, applies to all but the first component of f_m^{-1} in the previous section). Rather than build an explicit model of $f_{m,n}^{-1}$, however, we have once more incorporated a single inverse iterate of ψ into the definition of $f_{m,n}$ so that its first component $w_{m,n} \equiv (f_{m,n}^{-1})_1$ is again the one-step time series generator $v_{i+1} = w_{m,n}(\mathbf{y}_i)$, where $\mathbf{y}_i = \mathcal{F}_{m,n}(\mathbf{x}_i)$. Since v_{i+1} is not an element of \mathbf{x}_i we can rely on its RBF approximation to determine whether or not $f_{m,n}$ is a diffeomorphism.

We model $w_{m,n}$ with the LS RBF approximation $\widehat{w}_{m,n}: \mathcal{M}_{m,n} \subset \mathbb{R}^n \rightarrow \mathbb{R}$, constructed once again from $p = 200$ repulsive centers selected from a training set of $N = 2000$ points in $\mathcal{M}_{m,n}$, and using cubic basis functions. We evaluate the success of this approximation via the expected values of the normalised fitting and test errors

$$\epsilon_{m,n}^2 = \sigma_v^{-2} \sum_{i=1}^N \|\widehat{w}_{m,n}(\mathbf{y}_i) - v_{i+1}\|^2 \quad (5.2)$$

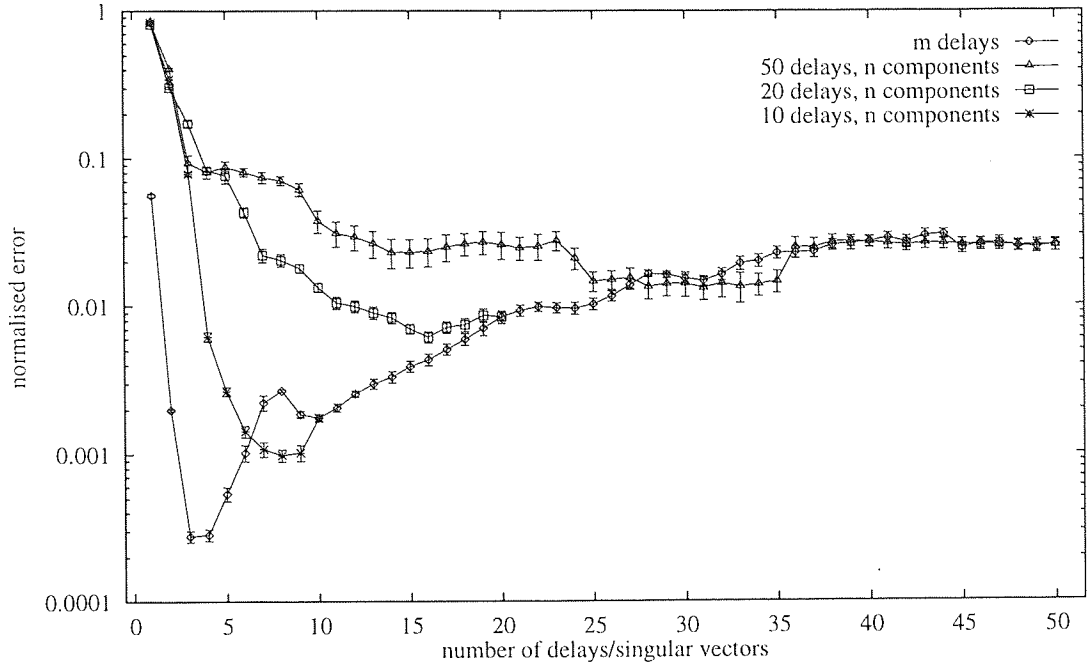
calculated over the standard 500 random choices of repulsive seed; the normalisation constant σ_v^2 is unchanged from the previous section. As before, due to the vagaries of fitting individual RBF maps (albeit mediated by the average over random repulsive seeds), and the expected increase in fitting error with increasing dimensionality n for a constant number of centers, we do not necessarily expect a monotonic decrease in $\epsilon_{m,n}$ as n increases towards the critical value n^* . What we are looking for is a region $n^* \leq n \leq m$ within which $\langle \epsilon_{m,n} \rangle \approx \langle \epsilon_m \rangle$, where ϵ_m is the prediction error on \mathcal{M}_m defined in (5.1),

indicating that those singular vectors $v_{j>n^*}$ which are missing from \mathcal{F}_{m,n^*} contribute negligibly to the fit. For this reason, in each of the following experiments, we plot the mean m -delay prediction error $\langle \epsilon_m \rangle$, over the range $1 \leq m \leq 50$, then superimpose the mean singular subspace prediction errors $\langle \epsilon_{m,n} \rangle$, obtained by fixing m at 10, 20 and 50 delays, respectively, and varying $1 \leq n \leq m$ in each case. Naturally, we must find $\epsilon_{m,m} = \epsilon_m$. (Due to our choice of random seeds, the same need not be strictly true of $\langle \epsilon_{m,n} \rangle$ and $\langle \epsilon_m \rangle$, although the differences will turn out to be negligible.) If the noise is sufficiently pathological we might even hope to find $\epsilon_{m,n} < \epsilon_m$ for some values of $n \geq n^*$, as components containing noise-induced self-intersections are eliminated from the image of $\mathcal{F}_{m,n}$.

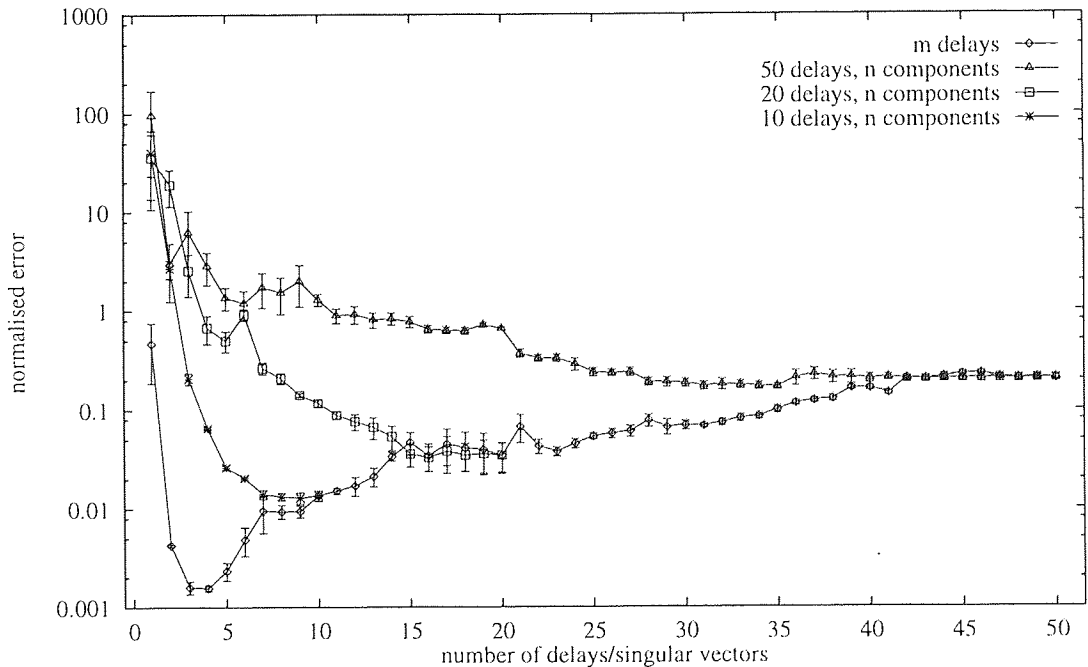
5.2.1 Singular subspaces of the embedded Lorenz system

For our initial investigation we again choose to embed the 0.01-step Lorenz system by sub-sampling the time series of figure 2.4(a) with a lag of $\tau = 10$, to get the relatively smooth trajectory pictured in figure 2.8(a). In figure 5.8(a) we replot, versus m , the first 50 values of $\langle \epsilon_m \rangle$ from figure 5.5. Superimposed on this curve we plot the expected values of the fitting errors $\epsilon_{50,n}$, $\epsilon_{20,n}$ and $\epsilon_{10,n}$, obtained via projection into singular subspaces of \mathbb{R}^{50} , \mathbb{R}^{20} and \mathbb{R}^{10} , respectively, versus the appropriate interval in n ; the test errors are plotted separately in figure 5.8(b). (Of course, if we were really concerned with obtaining the best possible predictor for a delay window of length 50τ , and we had sufficient processing power available, we could eliminate τ completely, and simply apply a 500-delay window directly to the un-lagged time series.) As the time series in this example was generated by numerical simulation the reconstructed trajectories in \mathcal{M}_m are effectively noise-free, so for a given choice of centers we do not expect $\epsilon_{m,n < m}$ to display any significant improvement over its baseline value ϵ_m , due to the elimination of noise-dominated dimensions in \mathcal{M}_m . We might hope, however, to find an $n^* \leq m$ such that $\epsilon_{m,n} \approx \epsilon_m$, for $n^* < n < m$. Figure 5.8 confirms these tentative expectations: for example, the mean fitting error $\langle \epsilon_{50,n} \rangle$ remains within an order of magnitude of $\langle \epsilon_{50} \rangle$ over the interval $3 \leq n \leq 50$, even dropping below that value for certain n . This is presumably due in part to the beneficial effect on the RBF map resulting from a reduction in the dimensionality of its domain. For $m = 20$ and $m = 10$ the advantage of singular subspace projection is less obvious, although it is clear that a few dimensions may be safely ‘shaved off’ the pseudo-phase space even for m on the order of the minimum embedding dimension, $m^* = 4$, of the unfiltered system. Not surprisingly, however, given the exceptionally low noise floor associated with this time series, the best prediction for $\{v_i\}$ is still that obtained from $\widehat{w}_4: \mathcal{M}_4 \rightarrow \mathbb{R}$, as defined in the previous section.

A slightly more appropriate problem domain is the noisy version of this time series. We plot the usual errors, $\langle \epsilon_{m,n} \rangle$ and $\langle \epsilon_m \rangle$, for filtered and unfiltered delay reconstructions, in figure 5.9. Although, as shown in the previous section, the Lorenz system cannot actually be embedded from this time series—the noise component is too large—this example still demonstrates the effectiveness of $\mathcal{F}_{m,n}$ at eliminating noise-dominated dimensions in \mathcal{M}_m , with all three singular subspace prediction errors, calculated over

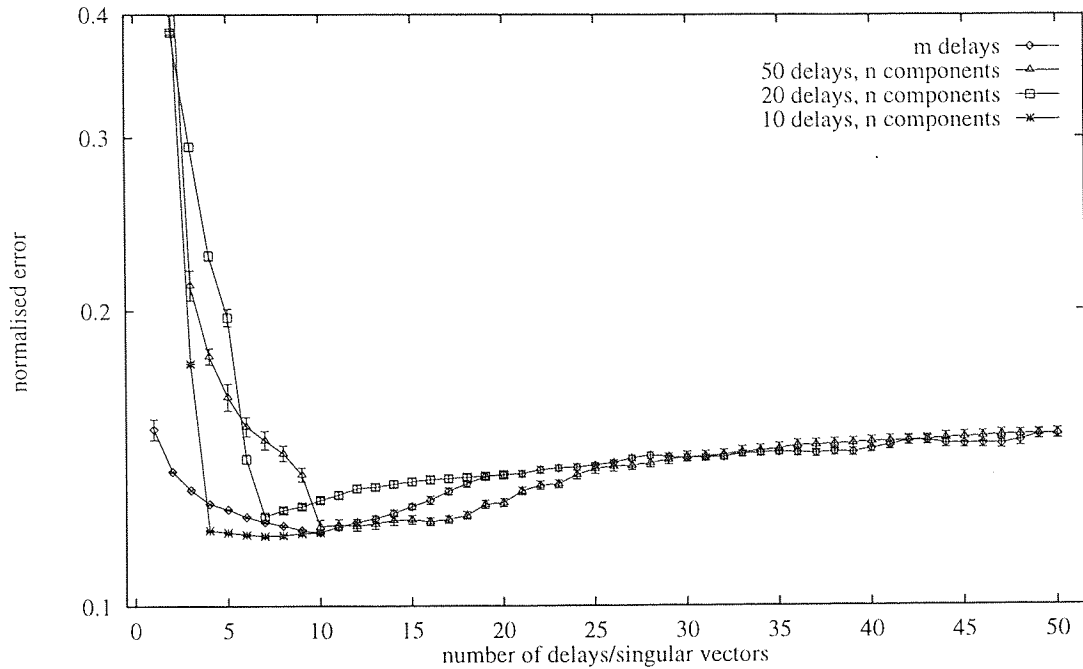


(a)

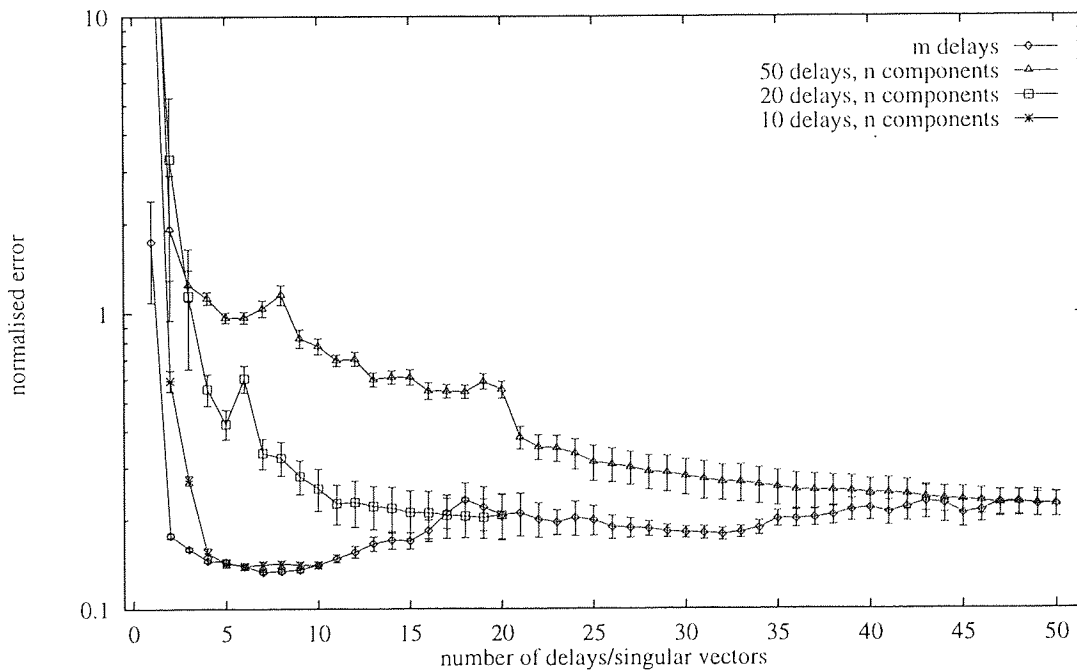


(b)

Figure 5.8 Comparing the mean, normalised prediction errors $\langle \epsilon_{m,n} \rangle$ for the 0.01-step Lorenz system, reconstructed in singular subspaces of \mathbb{R}^{50} , \mathbb{R}^{20} and \mathbb{R}^{10} with a lag-10 delay map. (a) Mean values of the fitting errors $\epsilon_{50,n}$, $\epsilon_{20,n}$ and $\epsilon_{10,n}$ are plotted, versus the appropriate interval of n , on top of a replot of the corresponding error $\langle \epsilon_m \rangle$ obtained from the unfiltered embedding, for $1 \leq m \leq 50$. Although a given level of error can clearly be maintained (up to an order of magnitude) through projection into an appropriately lower-dimensional subspace (particularly in the case of $m = 50$), we still find that $\langle \epsilon_{m,n} \rangle > \langle \epsilon_n \rangle$ almost everywhere. (b) The corresponding test errors tell a similar story. Error bars denote one standard deviation in each direction.



(a)



(b)

Figure 5.9 Comparing the mean, normalised prediction errors $\langle \epsilon_{m,n} \rangle$, obtained with $m = 10, 20$ and 50 , for the 0.01 -step Lorenz system, generated with a lag- 10 delay map from a noisy time series. **(a)** The fitting errors in this case all achieve a minimum of approximately the same value as that of the unfiltered reconstruction error $\langle \epsilon_m \rangle$, which occurs at $m = 10$. In particular, at $n = 4$, the minimum of $\langle \epsilon_{10,n} \rangle$ indicates that projection of the image \mathcal{M}_{10} of $\Phi_{v,10}$ into a singular subspace of only four dimensions has a negligible effect on the resulting predictor, whilst reducing the complexity of the RBF map in question by a factor of more than two. **(b)** This conclusion is confirmed by the corresponding test error (also hard-limited) although a value of $n = 5$ is more firmly indicated. These results serve as a proof of concept only, as we have already established that the Lorenz system cannot be embedded from this particular time series. Error bars denote one standard deviation in each direction.

the training set, reaching minimum values in figure 5.6(a) of $\langle \epsilon_{50,10} \rangle, \langle \epsilon_{20,7} \rangle, \langle \epsilon_{10,4} \rangle \approx \langle \epsilon_{10} \rangle$. As before, we see a significant degree of over-fitting in this experiment.

5.2.2 Singular subspaces of the embedded laser system

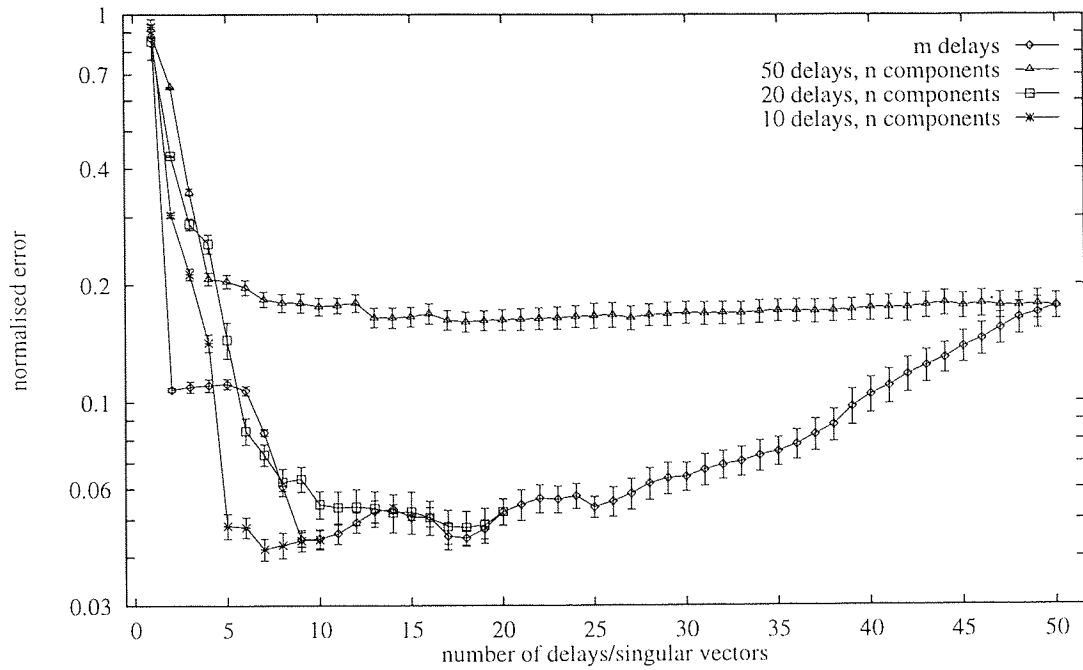
We see a similar effect in our final example, in which we apply the singular subspace technique to the laser time series of figure 2.5(b), plotting the errors in figure 5.10. We have already suggested, both by analogy with the Lorenz map (section 2.2) and by experiment (section 5.1.4) that this time series does not represent a suitable measurement function for delay reconstruction of \mathcal{M} . Nevertheless, in part (a) of figure 5.10, for both $m = 10$ and $m = 20$ we can (approximately) achieve the minimum fitting error, arising from an unfiltered 10-delay reconstruction, with $\langle \epsilon_{10,5} \rangle, \langle \epsilon_{20,14} \rangle \approx \langle \epsilon_{10} \rangle$; the corresponding test errors in part (b) are once more in rough agreement. With 50 delays, on the other hand, we can merely maintain the relatively large error of $\langle \epsilon_{50,n} \rangle \approx \langle \epsilon_{50} \rangle$, for $4 \leq n \leq 50$, on training and test sets alike. The filtered attractor obtained with the projection $\mathcal{F}_{10,3}$ has already been plotted, in figure 2.8(b).

Performing rank-reducing projections into singular subspaces in this manner is, of course, highly reminiscent of the RBF generalisation techniques described in section 3.2.3. It is, in fact, exactly equivalent to the blind truncation criterion, where the order in which individual principal components are incorporated into the model takes no account of the ensuing errors. Taking this analogy a little further, we could clearly attempt to use some other, less ad hoc criterion, such as the targeted truncation method of section 3.2.3. This is, however, impractical for the purposes to which we would put it, as it would necessarily involve some form of nonlinear optimisation strategy—such as error gradient descent—owing to the presence of the nonlinear transformation $\varphi: \mathbb{R}^n \rightarrow \mathbb{R}^p$ in $\widehat{w_{m,n}}$, so we do not pursue that idea here.

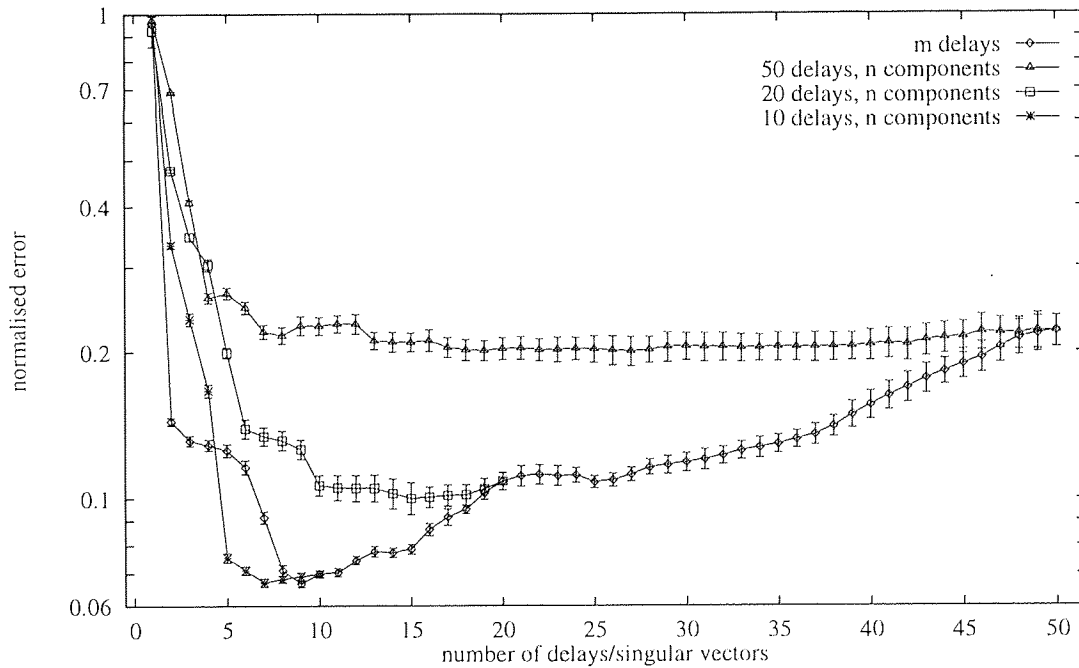
5.3 Detecting unstable periodic orbits

We will now consider the class of linear transformations, introduced by Broomhead, Huke and Muldoon [3] and described in section 2.3.3, which relates the images of \mathcal{M} under the delay maps $\Phi_{v,m}: \mathcal{M} \rightarrow \mathbb{R}^m$ and $\Phi_{u,n}: \mathcal{M} \rightarrow \mathbb{R}^n$, where the time series $\{u_i\}$ is obtained from $\{v_i\}$ through the c -coefficient FIR filter defined in equation (2.17), and $n = m - c + 1$. As discussed in section 2.3.3, the transformation $\Phi_{v,n} \circ \Phi_{u,m}^{-1}$ is generically a diffeomorphism on \mathcal{M}_m provided that $n > 2d$, so provided that $\{v_i\}$ is a generic measurement function (in the sense of Takens) on \mathcal{M} then Takens' theorem applies equally well—or equally badly—to the time series obtained from $\{v_i\}$ through almost any FIR filter as it does to $\{v_i\}$ itself: either may be used to construct a diffeomorphic copy of \mathcal{M} in a sufficiently high-dimensional Euclidean space.

Non-generic filters, on the other hand, give rise to delay maps $\Phi_{u,n}$ which do *not* embed \mathcal{M} , no matter



(a)



(b)

Figure 5.10 Comparing the mean, normalised prediction errors $\langle \epsilon_{m,n} \rangle$, obtained with $m = 10, 20$ and 50 , for the laser system. (a) In this case, we find that a delay reconstruction with both $m = 10$ and $m = 20$ provides the opportunity for a singular subspace projection to achieve the minimum fitting error of $\langle \epsilon_m \rangle$ with, respectively, $n = 5$ and $n = 14$, although at $m = 50$ we can maintain, but not reduce, the value of $\langle \epsilon_{50,n} \rangle \approx \langle \epsilon_{50} \rangle$ for $4 \leq n \leq 50$. (b) The corresponding test errors allow a similar analysis. Error bars denote one standard deviation in each direction.

how large n is chosen to be. We will now exploit this non-genericity, in the form of the single-parameter family of FIR filters with coefficients $(1, -2 \cos 2\pi\nu, 1)$; it can be shown [26] that a three-coefficient filter of this form has a zero response at the frequency ν in $\{v_i\}$. For fixed m , such that $n = m - 2$ is sufficiently large for $\Phi_{v,n}$ to embed \mathcal{M} , and writing $\mathcal{M}_m = \Phi_{v,m}\mathcal{M}$ and $\mathcal{M}_{n,\nu} = \Phi_{u,n}\mathcal{M}$ (at the risk of some slight confusion, we will use the subscripts ν and u interchangeably to identify a particular FIR filter), a useful effect of the corresponding linear transformation $\mathcal{F}_{n,\nu}: \mathcal{M}_m \rightarrow \mathcal{M}_{n,\nu}$ is to ‘collapse’ any orbits of period $\frac{1}{\nu}$ samples in \mathcal{M} onto a fixed point in $\mathcal{M}_{n,\nu}$. (Such orbits will be a fundamental characteristic of (\mathcal{M}, ψ) , so a filter of this form will persist in its non-genericity independently of the choice of measurement function v .) Thus, if we can establish that $f_{n,\nu}: \mathcal{M}_m \rightarrow \mathcal{M}_{n,\nu}$, defined by $f_{n,\nu} = \Phi_{u,n} \circ \Phi_{v,m}^{-1}|_{\mathcal{M}_m}$, is a diffeomorphism for generic values of ν , but *not* for a particular value ν^* , then we can conclude that $\Phi_{n,\nu}$ collapses an orbit of period $\frac{1}{\nu^*}$ in \mathcal{M} .

As usual, we know that $f_{n,\nu}$ is a function since it is the restriction to \mathcal{M}_m of the linear map $\mathcal{F}_{n,\nu}$, so we base our decision on whether or not $f_{n,\nu}$ is a diffeomorphism on an RBF analysis of its inverse. (We assume that the orbit in question is visited often enough for its presence to show up in a LS RBF fit.) Although the delay structure in \mathcal{M}_m enables us, in principle, to construct an approximation to $f_{n,\nu}^{-1}$ from any one of its components, we must ensure that every component is amenable to RBF approximation before we can claim that $f_{n,\nu}$ is a diffeomorphism. In this case, however, rather than incorporate a forward or inverse iterate of ψ into $f_{n,\nu}$ we explicitly fit an RBF map $\widehat{w_{n,\nu}^{(j)}}: \mathcal{M}_{n,\nu} \rightarrow \mathbb{R}$ to each component $w_{n,\nu}^{(j)} \equiv (f_{n,\nu}^{-1})_j$ of $f_{n,\nu}^{-1}$ mapping $\mathbf{y}_i = (u_i, \dots, u_{i-n+1})^T$ to $v_{i-j+1} = w_{n,\nu}^{(j)}(\mathbf{y}_i)$. We define the resulting errors $\epsilon_{n,\nu}^{(j)}$ by

$$\epsilon_{n,\nu}^{(j)2} = \sigma_v^{(j)-2} \sum_{i=1}^N \|\widehat{w_{n,\nu}^{(j)}}(\mathbf{y}_i) - v_{i-j+1}\|^2 \quad (5.3)$$

where the normalisers $\sigma_v^{(j)}$ are calculated over the appropriate interval in $\{v_i\}$ (neglecting end effects, these constants are effectively identical). In each of the figures to follow we will superimpose these errors, for $1 \leq j \leq m$, on a single set of axes, but for practicality we will not visually distinguish them from each other; (we will also not make use of the average over models in this section, nor in the next one). Since the RBF map is linear in its basis functions the error which would have arisen on fitting all m components of $f_{n,\nu}^{-1}$ simultaneously is easily obtained as

$$\epsilon_{n,\nu} = \frac{1}{m} \sum_{j=1}^m \epsilon_{n,\nu}^{(j)} \quad (5.4)$$

Each estimator $\widehat{w_{n,\nu}^{(j)}}$ represents a nonlinear inverse to the effect on \mathcal{M}_m of the FIR filter through which $\mathcal{M}_{n,\nu}$ has been constructed. That we are able to construct such an inverse is a consequence of the deterministic nature of $\{v_i\}$. In particular, chaotic time series generally exhibit broad-band power spectra, as a result of the broad spread of unstable periodic orbits embedded in the attractors of their generating

systems. For this reason, even though we have removed all of the power at the target frequency ν , we expect to retain sufficient power at other frequencies with which to reconstruct (up to a certain point) a diffeomorphism of the original attractor; we will see examples of this behaviour in section 5.4.

5.3.1 Filtering the Ikeda attractor

We begin with the Ikeda system, constructing filtered delay maps $\Phi_{u,5}: \mathcal{M} \rightarrow \mathbb{R}^5$ from the time series plotted in figure 2.3(a) by sweeping the middle coefficient $a_1 = -2 \cos 2\pi\nu$ through the interval $-1 \leq a_1 \leq 1$. In figure 5.11 we plot the seven error curves $\epsilon_{5,\nu}^{(j)}$ obtained (over the test set) from a LS RBF fit to $w_{5,\nu}^{(j)}: \mathcal{M}_{5,\nu} \rightarrow \mathbb{R}$, for $j = 1, \dots, 7$. For convenience, we use a_1 , rather than ν , to index the x -axis in this, and subsequent plots. Although all seven errors are small over most of the interval, they all peak sharply at $a_1^* = 0$. The FIR filter $(1, 0, 1)$ has a zero response at the frequency $\nu^* = \frac{1}{4}$, leading us to deduce the existence of a period-4 orbit in the Ikeda attractor; careful examination of figure 2.3(a) confirms this conclusion. Before attaching any significance to this result, however, we must ensure that $\mathcal{M}_{5,\nu}$ is an embedding of \mathcal{M} for generic ν . In fact, the mean error level $\epsilon_{5,\nu}$ turns out to be on the same order of magnitude as the prediction error ϵ_4 which we measured in section 5.1.1.

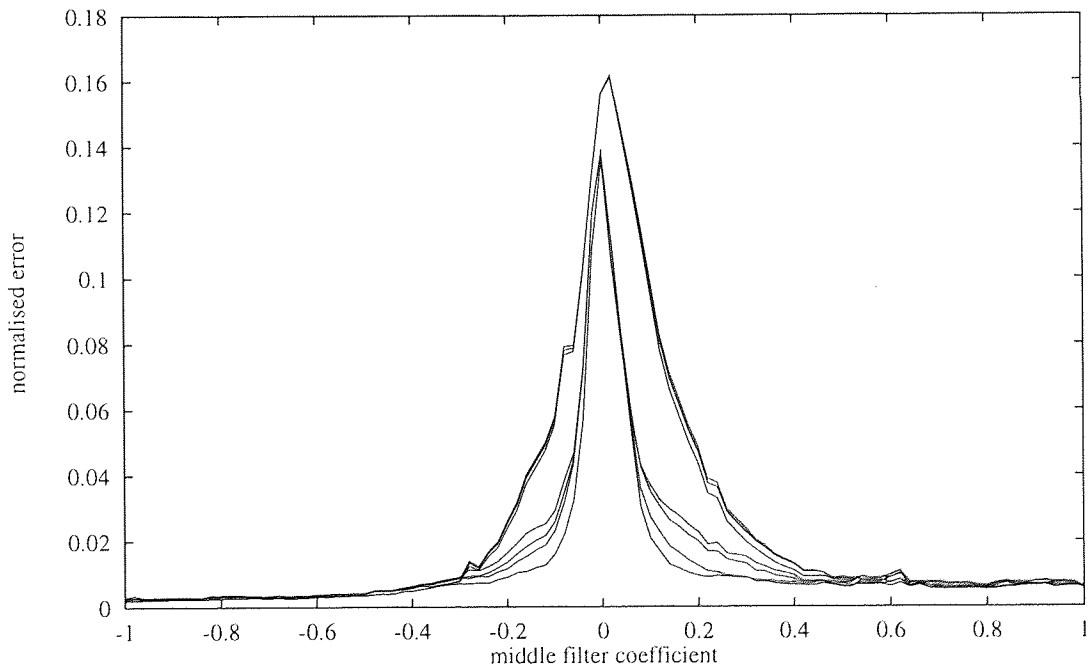


Figure 5.11 Detecting periodic orbits in the Ikeda attractor by plotting the test error obtained by approximating the inverse to a family of FIR filtered delay embeddings into \mathbb{R}^5 . A non-diffeomorphic relationship is clearly indicated at $a_1 = 0$, which corresponds to a period of 4 samples.

In figure 5.12 we illustrate the images of the filtered delay maps $\Phi_{u,5}$, for specific values of filter coefficient a_1 , colour-coding each point by the relative magnitude of the per-point error to which it gives rise under the zero-offset predictor $\widehat{w_{5,\nu}^{(1)}}$, as in the previous chapter. For the purposes of comparison, in

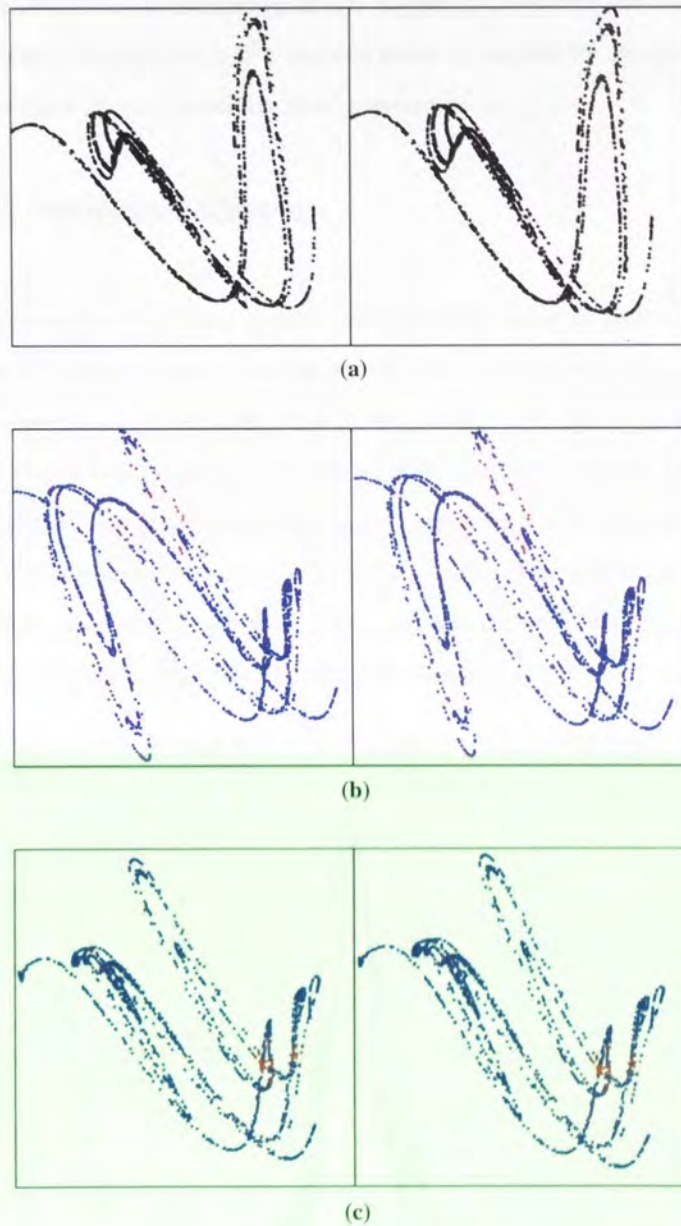


Figure 5.12 Colour-coding the first three components of the Ikeda attractor, reconstructed with filter coefficients $a_1 = -1$ and $a_2 = 0$, by the errors to which they give rise in predicting the unfiltered Ikeda time series. (a) For comparison, the attractor obtained from a direct delay reconstruction is actually embedded in \mathbb{R}^4 ; (b) the generic filter $(1, -1, 1)$ provides an embedding of \mathcal{M} in \mathbb{R}^5 , with correspondingly small and evenly distributed per-point errors; (c) the Ikeda system cannot be embedded after the non-generic filter $(1, 0, 1)$ has been applied to $\{v_i\}$, and the resulting errors are localised in the resulting set of self-intersections. Plotted in stereo.

part (a) of this figure we plot the first three components of the attractor in $\mathcal{M}_7 = \Phi_{\nu,7}\mathcal{M}$, which we estimated in section 5.1.1 to be embedded in \mathbb{R}^4 . In part (b) we show the attractor obtained from the filtered embedding corresponding to a coefficient of $a_1 = -1$, which appears more convoluted than that of \mathcal{M}_7 but is also known to be an embedding of \mathcal{M} . In part (c), however, the filtered object \mathcal{M}_{5,ν^*} is *not* embedded in \mathbb{R}^5 (and certainly not in \mathbb{R}^3), and as a result we see that the per-point errors are strongly localised to a self-intersecting set of relatively small measure in \mathcal{M}_{5,ν^*} .

5.3.2 Filtering the Hénon attractor

We now move on to examine the Hénon system, using the time series of figure 2.3(b) to construct an embedding of \mathcal{M} in \mathbb{R}^6 which is then projected into \mathbb{R}^4 by the linear map $\mathcal{F}_{4,\nu}$. The test errors $\epsilon_{4,\nu}^{(j)}$ resulting from the approximations $\widehat{w_{4,\nu}^{(j)}}: \mathcal{M}_{4,\nu} \rightarrow \mathbb{R}$, for $j = 1, \dots, 6$ and $-1 \leq a_1 \leq 1$, are plotted in figure 5.13, and tell a remarkably similar story to those in the previous example. Once again, although it appears that both FIR filter and reconstruction dimension are suitable for a successful embedding of the Hénon system provided that $|a_1| > 0$, at the critical value $a_1^* = 0$, corresponding to $\nu^* = \frac{1}{4}$, we are again unable to find a diffeomorphism $f_{4,\nu^*}: \mathcal{M}_6 \rightarrow \mathcal{M}_{4,\nu^*}$. A period-4 orbit is visible in the unfiltered time series, although not as frequently occurring as in the previous example.

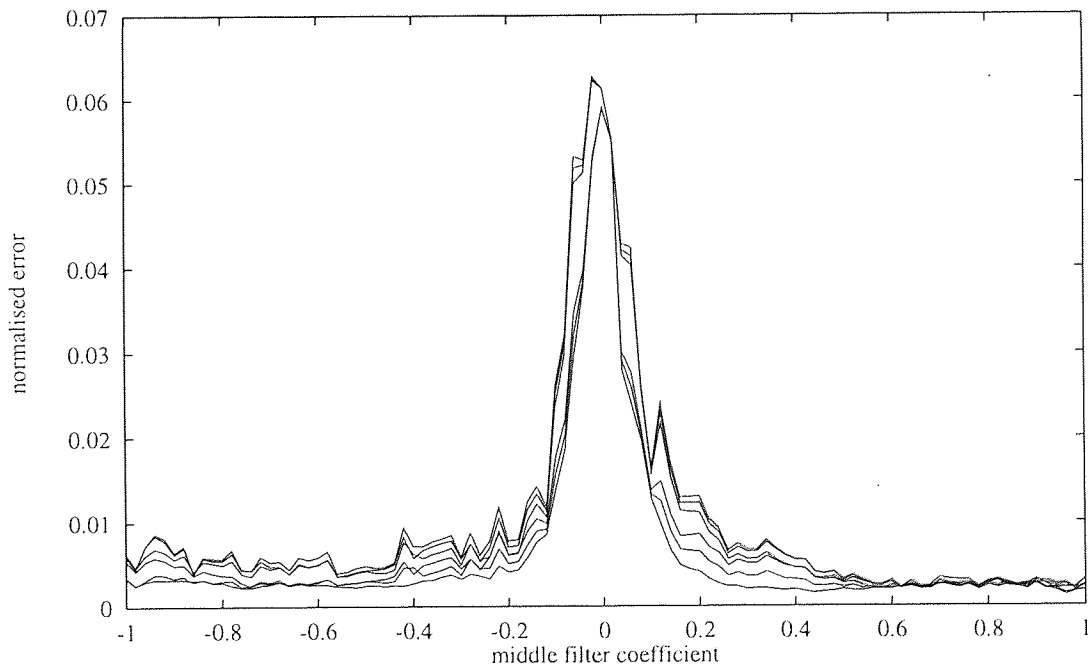


Figure 5.13 Prediction error for a FIR filtered delay embedding of the Hénon attractor into \mathbb{R}^4 . Once again, a period-4 orbit is identified by a peak at $a_1 = 0$.

We illustrate the attractors obtained from these filtered delay maps in figure 5.14, plotting only the first three components of each, as usual. In part (a), for comparison purposes, we show a direct, three-delay

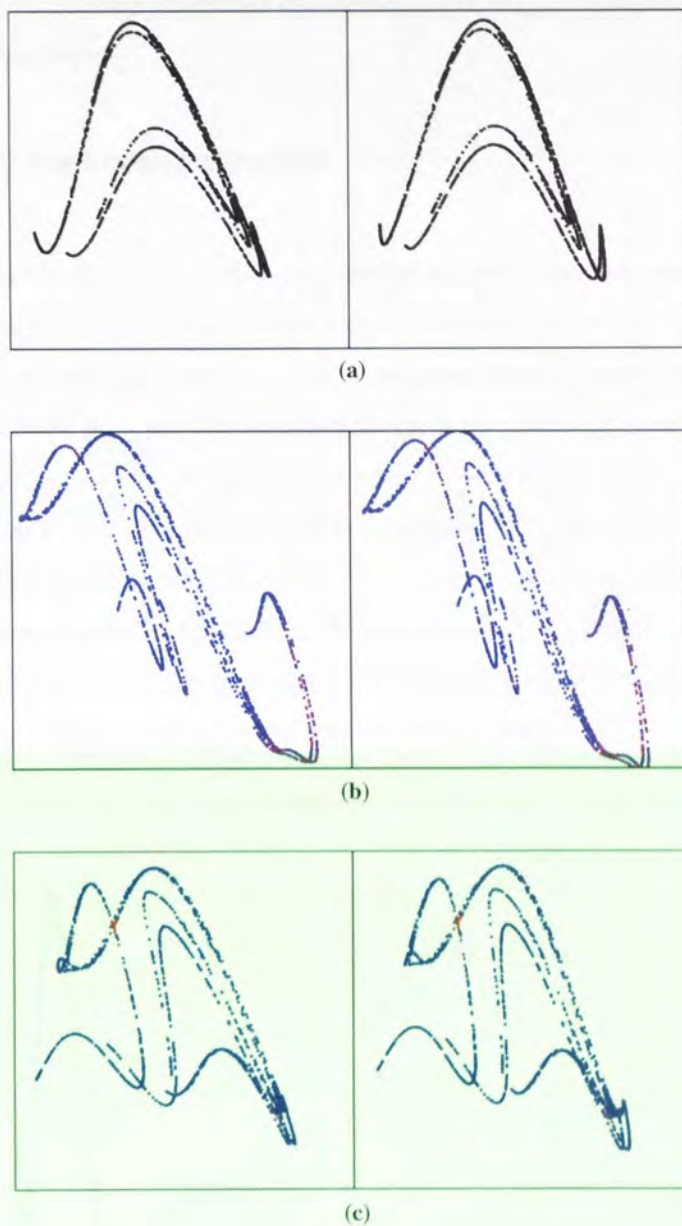


Figure 5.14 Colour-coding the first three components of the Hénon attractor, reconstructed with filter coefficients $a_1 = -1$ and $a_2 = 0$, by the errors to which they give rise in predicting the unfiltered Hénon time series. (a) For comparison, the direct delay reconstruction is again illustrated in \mathbb{R}^3 , although the Hénon attractor actually embeds in \mathbb{R}^2 ; (b) the generic filter $(1, -1, 1)$ embeds \mathcal{M} in \mathbb{R}^4 , with small and fairly evenly distributed prediction errors; (c) the non-generic filter $(1, 0, 1)$ prevents an embedding of \mathcal{M} in any number of dimensions, however, and the self-intersection responsible is clearly visible in this case. Plotted in stereo.

embedding of the Hénon system, although the time series $\{v_i\}$ actually enables us to embed this attractor in \mathbb{R}^2 , as we already know. The filtered embedding obtained from the FIR filter $(1, -1, 1)$ is shown in part (b), colour-coded by its fairly uniformly small per-point errors. In part (c), however, the image of the period-4 orbit in \mathcal{M} is clearly visible as a self-intersection in \mathcal{M}_{4,ν^*} , highlighted by an appropriately strong peak in per-point error.

5.3.3 Filtering the Lorenz attractor

Finally, we come to the Lorenz system, which we embed using the 0.1-step time series illustrated in figure 2.4(b). Because this time series was generated with an integration step of 0.1, rather than 0.01 as in the corresponding time series of figure 2.4(a), we do not need to use a lagged delay map in this case. With a filtered embedding $\Phi_{9,u}$ into \mathbb{R}^9 , we obtain the prediction errors $e_{9,\nu}^{(j)}$ plotted in figure 5.15, for $j = 1, \dots, 11$. In this case we vary $0 \leq \nu \leq \frac{1}{2}$ to get a filter coefficient $-2 \leq a_1 \leq 2$, and see two distinct peaks, at $a_1 = -2$ and -1.2 . The former of these corresponds to a filter with coefficients $(1, -2, 1)$, which maps fixed points in \mathbb{R}^m to the origin in \mathbb{R}^{m-2} . It is no surprise that this filter is non-generic with respect to the time series under investigation, as we have already seen that the Lorenz attractor has two such unstable fixed points. The latter picks out orbits of frequency approximately 0.17, or period 6.8, which is roughly the orbital period close to each of the two fixed points.

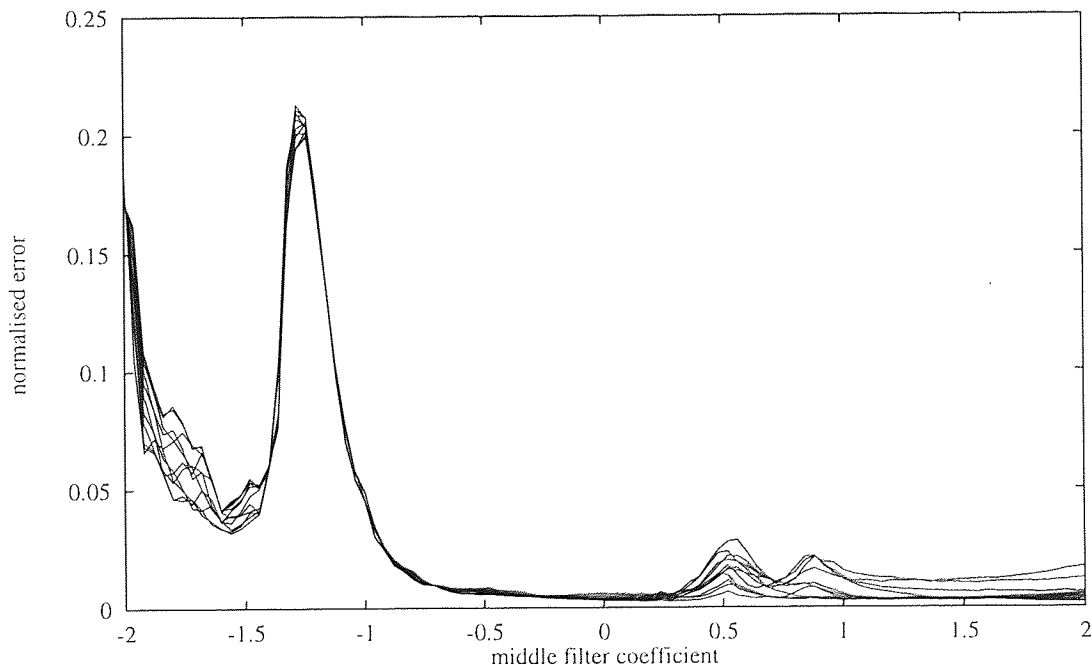
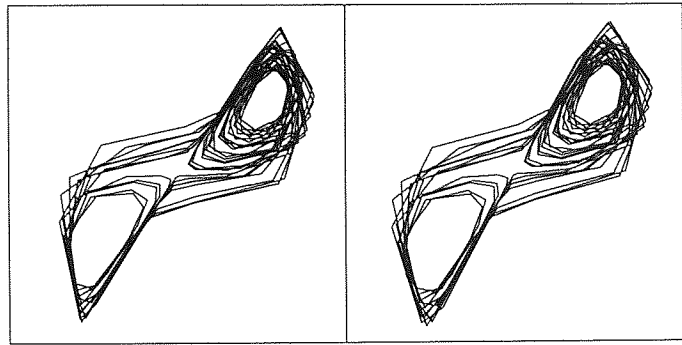
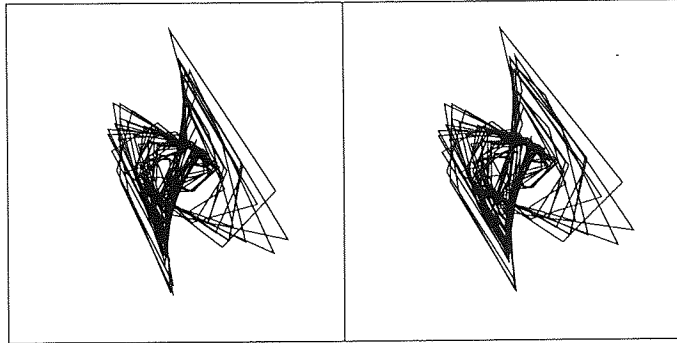


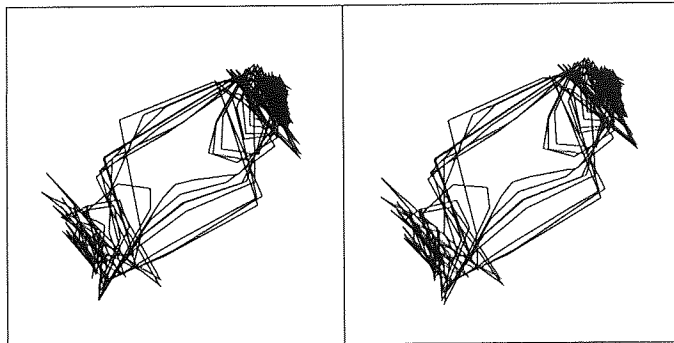
Figure 5.15 Prediction error obtained from a 9-delay filtered embedding of the 0.1-step Lorenz attractor. This plot reveals two peaks: one, at $a_1 = -2$, corresponds to a filter which maps the two fixed points onto the origin; the other, at $a_1 \approx -1.2$, collapses periodic orbits near those fixed points.



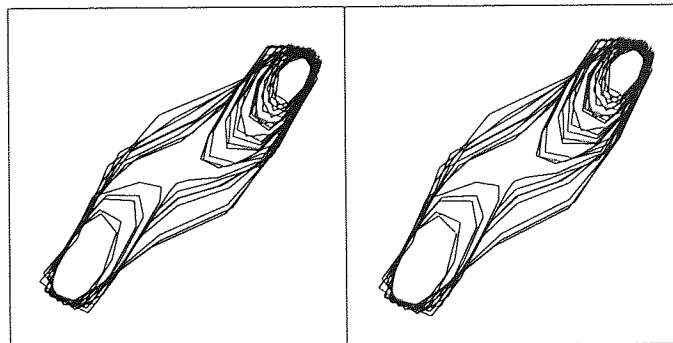
(a)



(b)

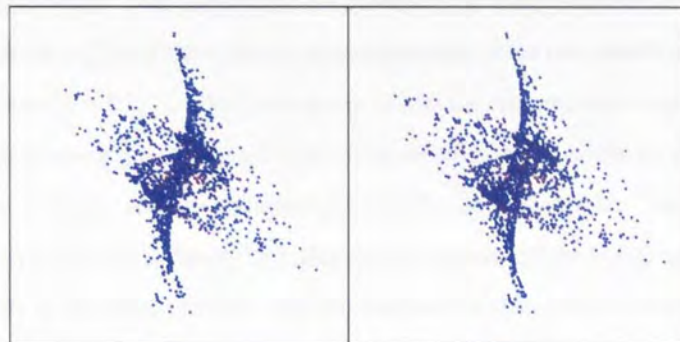


(c)

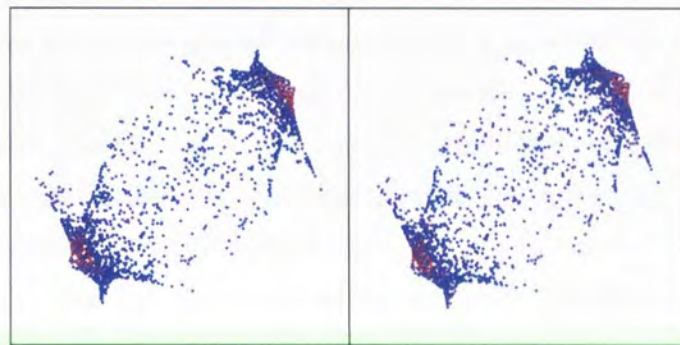


(d)

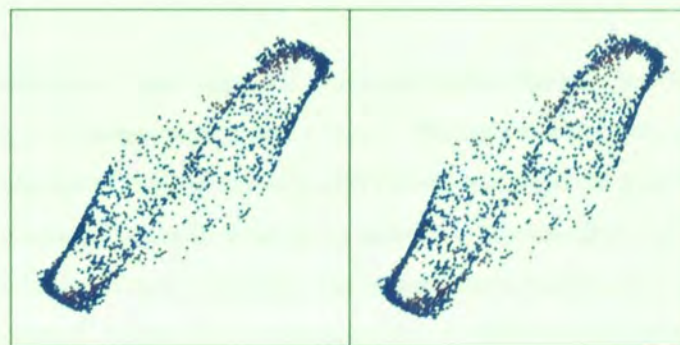
Figure 5.16 Delay reconstructions of the 0.1-step Lorenz system, obtained with various FIR filter coefficients. (a) A direct embedding of \mathcal{M} into \mathbb{R}^3 is compared with the reconstructed attractors obtained with (b) the fixed-point collapsing filter $(1, -2, 1)$, (c) the equally non-generic filter $(1, -1.2, 1)$ and (d) the generic filter $(1, 0, 1)$. Plotted in stereo.



(a)



(b)



(c)

Figure 5.17 Colour-coding the reconstructions of the 0.1-step Lorenz system, obtained with non-generic and generic FIR filters, by the corresponding prediction errors. (a) With a middle coefficient of $a_1 = -2$, the unstable fixed points in \mathcal{M} are superimposed in \mathbb{R}^9 , affecting the reconstructed dynamics to such an extent that the prediction errors are large but not particularly localised. (b) The filter $(1, -1.2, 1)$ should collapse the orbits near each fixed point; the strongly localised error confirm this prediction. (c) With $a_1 = 0$ the filter is now generic, and the errors are again small and uniformly distributed. Plotted in stereo.

In figure 5.16 we illustrate the attractors obtained from these filtered embeddings. To achieve the greatest possible clarity, we use lines to join consecutive points in these plots, and postpone the colour-coding exercise until the next figure. In part (a) we once again plot the attractor embedded directly into \mathbb{R}^3 using the unfiltered time series. Then, in part (b), we show the first three components of the reconstruction obtained with the filter $(1, -2, 1)$, illustrating the superimposition of the two unstable fixed points. Part (c) corresponds to the filter $(1, -1.2, 1)$, which adversely affects the reconstructed trajectory in the vicinity of each fixed point in the anticipated manner. The attractor of part (d), however, which corresponds to the now *generic* filter $(1, 0, 1)$, appears diffeomorphic to \mathcal{M}_3 , as it should be. The attractors of figures 5.16(b), (c) and (d) are replotted in figure 5.17, this time as points, colour-coded by the per-point errors to which they give rise in the usual manner. Part (a), with fixed points mapped onto the origin, does not actually reveal any strongly localised errors: this is most probably due to the fact that the collapse of the fixed points has resulted in a self-intersecting set which comprises a large proportion of the reconstructed object (bearing in mind that the trajectory does not actually pass *through* these fixed points, only nearby). Part (b), on the other hand, graphically illustrates the way in which the orbits of period approximately 6.8 are collapsed by the FIR filter $(1, -1.2, 1)$, with a region of large error clearly visible at each end of the reconstructed object. In part (c), however, the errors are more uniformly spread along the trajectory in \mathbb{R}^3 , reflecting the genericity of the filter $(1, 0, 1)$ in this case.

5.4 Signal separation

In section 5.2 we examined a linear approach to controlling the effect of stochastic noise on a delay embedding (\mathcal{M}_m, ψ_m) of the dynamical system (\mathcal{M}, ψ) . This statistical approach—projection of $\mathcal{M}_m \in \mathbb{R}^m$ onto a singular subspace of reduced dimensionality defined with respect to a noise floor in the variance in \mathbb{R}^m —is necessary because stochastic noise is, by definition, unpredictable. But what if a significant component of the ‘noise’ in question *only appears* to be statistically random, but is actually the output of a chaotic dynamical system: might we not then exploit this determinism and model the ‘noise process’ explicitly? Unfortunately, this will not be possible unless we can somehow separate the dynamics of the noise process from those of the system under observation.

To put this problem into a more convenient context we consider the situation in which a signal (the ‘message’), which may or may not be deterministic, becomes corrupted—for instance, during transmission down some channel—by additive ‘noise’ of a deterministic origin (the ‘chaos’), and it is our task to separate these two components. We will write $v_i = v_i^{(msg)} + v_i^{(chs)}$, where $v_i^{(msg)}$ represents the message and $v_i^{(chs)}$ is obtained from a measurement function $v_{chs}: \mathcal{M} \rightarrow \mathbb{R}$ on the noise process (\mathcal{M}, ψ) . We now make one further assumption: the message must be band-limited to a known frequency interval. This need not be too onerous an assumption, as it merely requires the message to be encoded in an appropriate fashion

before transmission. We can now design a c -coefficient FIR filter, using the techniques described in section 5.3, to eliminate as much of the power in this interval as possible, ideally zeroing the message component of the time series entirely. This operation results in a time series $\{u_i\}$ which is, to an approximation dependent on how successfully the message was removed, a filtered copy of $\{v_i^{(chs)}\}$ only. In other words, if we write $u_i = u_i^{(msg)} + u_i^{(chs)}$, where

$$u_i^{(msg)} = \sum_{k=0}^{c-1} a_k v_{i-k}^{(msg)}, \quad u_i^{(chs)} = \sum_{k=0}^{c-1} a_k v_{i-k}^{(chs)} \quad (5.5)$$

then the coefficients a_0 through a_{c-1} should be chosen so that $u_i^{(msg)}$ is everywhere as small as possible.

Assuming that this filter does a sufficiently good job of removing the message from $\{u_i\}$, we should now be able to find a value n sufficiently large that the delay map $\Phi_{u,n}: \mathcal{M} \rightarrow \mathbb{R}^n$ is an embedding on \mathcal{M} (provided that the filter turns out to be a generic one, in the sense of Broomhead, Huke and Muldoon), writing $\mathcal{M}_{n,\nu} = \Phi_{u,n}\mathcal{M}$. (We will assume that the message has been encoded so that it can be removed with a filter with a spectral null at the frequency ν .) What we would like to do now is follow the procedure described in the previous section and approximate the inverse $f_{n,\nu}^{-1}$ of the diffeomorphism $f_{n,\nu}: \mathcal{M}_m \rightarrow \mathcal{M}_{n,\nu}$, where \mathcal{M}_m is the image of \mathcal{M} under the delay embedding $\Phi_{v,m}: \mathcal{M} \rightarrow \mathbb{R}^{m=n+c-1}$ and $f_{n,\nu}$ is the restriction to \mathcal{M}_m of the linear transformation $\mathcal{F}_{n,\nu}$ defined as in equation (2.18). Since we have already assumed that $f_{n,\nu}$ is a diffeomorphism, we need only construct an RBF approximation to a single component of $f_{n,\nu}^{-1}$. In the examples to follow we choose the zero-step time series generator $w_{n,\nu} \equiv (f_{n,\nu}^{-1})_1$, which maps $\mathbf{y}_i = (u_i, \dots, u_{i-n+1})^T$ to $v_i^{(chs)} = w_{n,\nu}(\mathbf{y}_i)$, although any component should suffice.

We quantify the degree of success with which the zero-step predictor $\widehat{w_{n,\nu}}$ approximates $w_{n,\nu}$ with the error

$$\epsilon_{n,\nu}^2 = \sigma_{chs}^{-2} \sum_{i=1}^N \|\widehat{w_{n,\nu}}(\mathbf{y}_i) - v_i^{(chs)}\|^2 \quad (5.6)$$

where σ_{chs}^2 is calculated over the appropriate interval in $\{v_i^{(chs)}\}$ as usual. As we have defined the signal separation problem, however, we do not actually have access to the chaos $v_i^{(chs)}$ alone, but only as a component of the composite signal $v_i = v_i^{(msg)} + v_i^{(chs)}$. We therefore have no means by which to construct $\Phi_{v,m}$, and hence a direct approximation of $w_{n,\nu}$, by minimising equation (5.6), is not possible. To overcome this limitation we note that, to the extent that $v_i^{(msg)}$ and $v_i^{(chs)}$ are uncorrelated with each other, the effect of the former component on an RBF approximation to the relationship $\mathbf{y}_i \mapsto v_i$ will tend to average out over a sufficiently large training set, resulting in a misleadingly large fitting error which does not fully reflect the success with which $w_{n,\nu}$ may actually have been approximated. We call this type of fitting ‘blind’ prediction, and when we describe its application to the experiments below we will quote the errors defined by equation (5.6), rather than the error which was actually minimised.

The other potential source of error in this scenario arises when the message is not completely eliminated by the filter: the presence of residual message components in $\mathcal{M}_{n,\nu}$ is likely to enable the RBF map $\widehat{w}_{n,\nu}$ to model the unfiltered message components in $\{v_i\}$ to a greater degree than might otherwise be the case. This situation could be made to work to our advantage, in that if we did not know, a priori, the frequency at which the message was encoded, we might be able to locate it experimentally by sweeping ν through an appropriate interval, as in the previous section, and looking for a particularly large peak in the resulting error curve(s) as evidence that the message component had been successfully eliminated in $\{u_i\}$ for a given frequency ν .

In actual fact, it is not entirely out of the question that we might gain access to the chaos alone, most likely by arranging for a constant message to be transmitted for a given length of time. Access to a ‘reference signal’ of this nature would enable us to construct a predictor for $u_i^{(chs)}$ by minimising (5.6) directly; this could be used more successfully to predict the chaos component in later messages. We therefore call this method ‘targeted’ prediction. Although the presence of a non-zero residual message in u_i will still cause some degradation in the resulting RBF fit, it is unlikely that this effect will be as noticeable as in the case of fitting a blind predictor. Of course, in order to apply an targeted predictor to an incoming signal we must have previously trained that predictor on a reference signal, so all estimates $\widehat{v}_i^{(chs)}$ thus obtained will be out-of-sample estimates by definition. However, since $w_{n,\nu}$ does not incorporate an iterate of ψ it would be entirely possible, should we choose to model $v_i^{(chs)}$ with a blind predictor, to retrain that predictor on each successive portion of signal as it arrives; the only cost incurred by this procedure would be the requirement that prediction be performed offline, unless a sufficiently fast method of optimising the RBF map was available. (One such candidate is the systolic array proposed by McWhirter, Broomhead and Shepherd [27].)

We have now effectively separated message from chaos, with a degree of success dependent on the accuracy of the RBF fit: our estimate of the chaos is directly obtained as $\widehat{v}_i^{(chs)} = \widehat{w}_{n,\nu}(\mathbf{y}_i)$, and the message follows as $\widehat{v}_i^{(msg)} = v_i - \widehat{v}_i^{(chs)}$. This approach to signal separation has been documented separately in Broomhead, Huke and Potts [4]. In the subsection immediately below we will describe the result of applying this technique, using both blind and targeted predictive methods, to the extraction of a sinusoidal message from a chaotic time series generated from the Ikeda map. We will then explore a slightly more complicated example, in which the message is generated from a binary sequence by phase modulation and the chaos is generated by the Lorenz system. Although we will be making extensive use of Fourier power spectra, to illustrate the effects of the FIR filters and their nonlinear inverses in the figures to follow, it should be borne in mind that spectral analysis is of limited use in characterising nonlinear systems, as the decomposition of a signal into a superposition of independent spectral modes is a purely linear concept. On a more practical note, the frequency scale in all such plots will be labelled in units of the sampling rate, so that a normalised frequency of 0.5 corresponds to the Nyquist rate.

5.4.1 Isolating a sinusoidal message from Ikeda chaos

For this first example we adopt the Ikeda system as the noise process, using the time series plotted in figure 2.3(a) to form $\{v_i^{(chs)}\}$. The message $\{v_i^{(msg)}\}$ is a sinusoid with an amplitude approximately equal to the standard deviation of the Ikeda time series, at 0.19, and a frequency equal to the deliberately non-commensurate fraction $\nu^* = \frac{9}{32}$ of the Ikeda sampling interval (or a period of approximately 3.56 samples). This frequency is stopped by the FIR filter $(1, 0.3902, 1)$, and was chosen so as to be some distance (in terms of the middle coefficient a_1) from the non-generic filter $(1, 0, 1)$ revealed by figure 5.11.

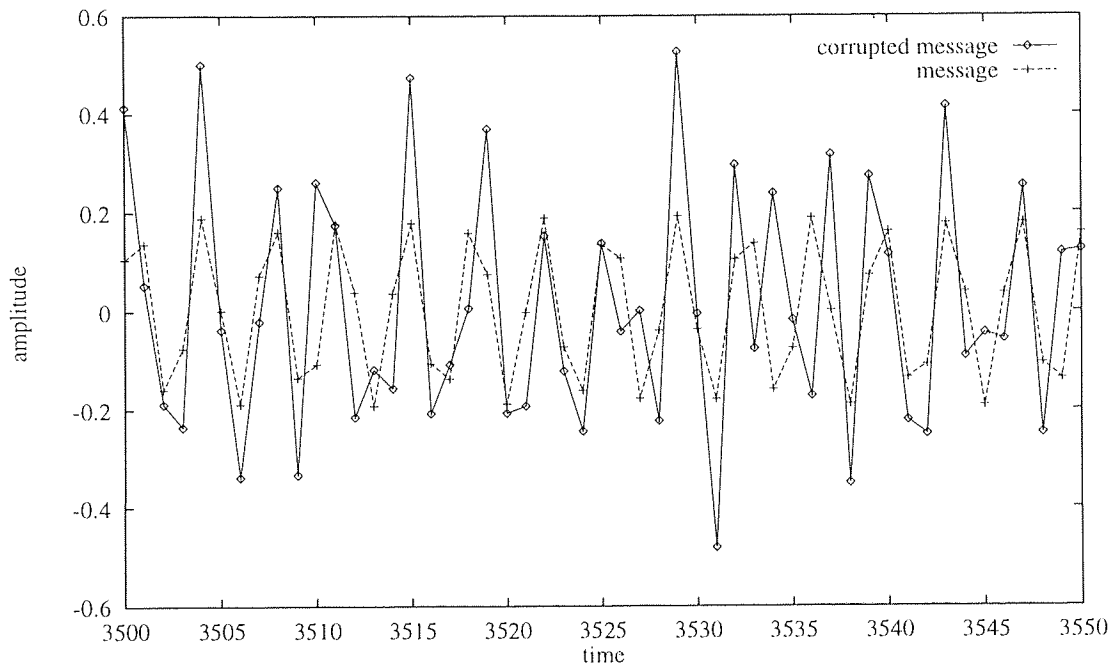
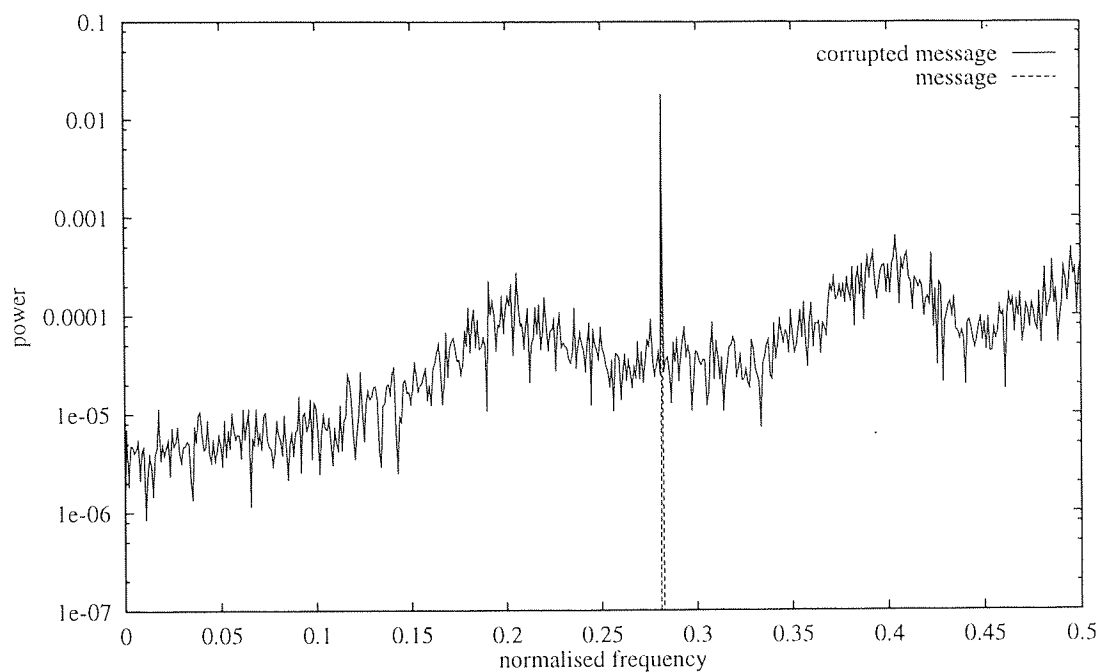


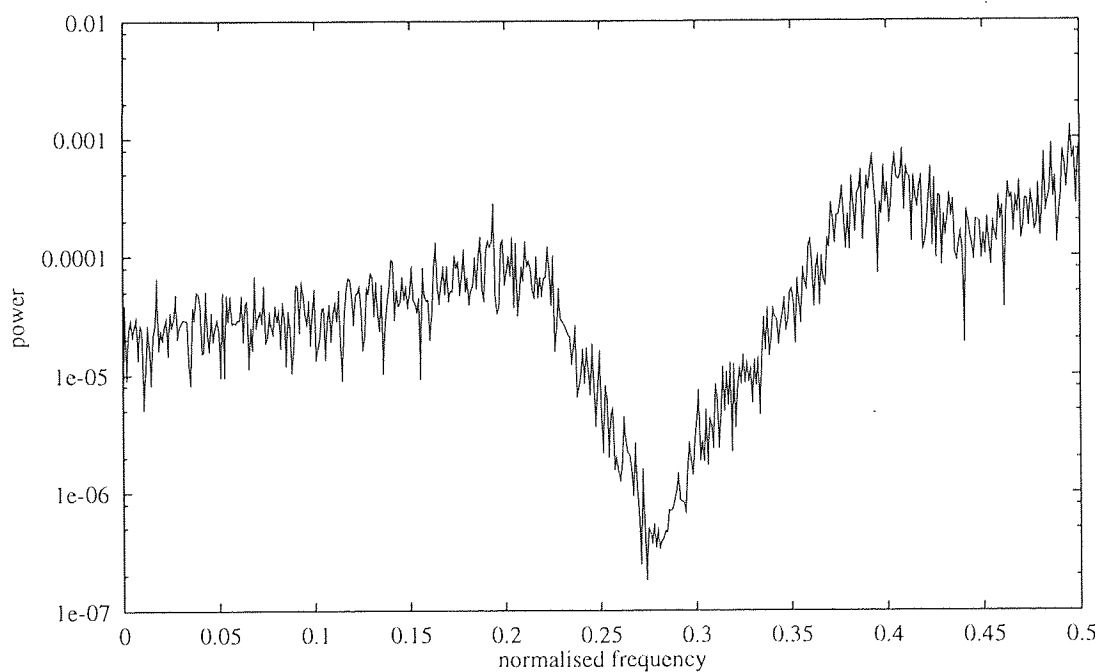
Figure 5.18 Coarsely sampled sinusoid with additive deterministic noise generated from the Ikeda map, superimposed on the sinusoid itself. The addition of the chaos, at a slightly higher standard deviation than the sinusoid, effectively hides the message.

A portion of this message is plotted in figure 5.18, along with the composite time series $\{v_i\}$. It is clear from this figure that the message component is significantly degraded by the presence of the chaos. In figure 5.19(a) we plot the corresponding power spectra, obtained by Fourier decomposition, of both message and composite time series. Not surprisingly, the power in the message component is concentrated entirely in a narrow peak, centered at the expected frequency. When added to the Ikeda time series, this peak emerges from what is otherwise a typically broadband chaotic power spectrum. In figure 5.19(b) we plot the power spectrum obtained from the filtered time series $\{u_i\}$, from which it is immediately evident that the message component has been almost completely removed. A comparison with figure 5.19(a) also illustrates the inevitable fact that the effect of the FIR filter is by no means limited to the vicinity of its target frequency, but extends throughout the entire spectrum.

We now construct the filtered embedding $\Phi_{u,5}: \mathcal{M} \rightarrow \mathbb{R}^5$. In the first instance, we use the composite



(a)



(b)

Figure 5.19 Power spectra for a sinusoid with additive Ikeda chaos, before and after filtering. (a) Superimposing the corrupted sinusoid on the sinusoid itself, the combined power spectrum is unchanged, other than at the frequency $\nu = \frac{9}{32}$, and exhibits the broadband nature typical of chaotic time series; (b) after filtering to remove the sinusoid, the spectrum has a large hole at the expected frequency, and is somewhat altered from its unfiltered form throughout the frequency range.

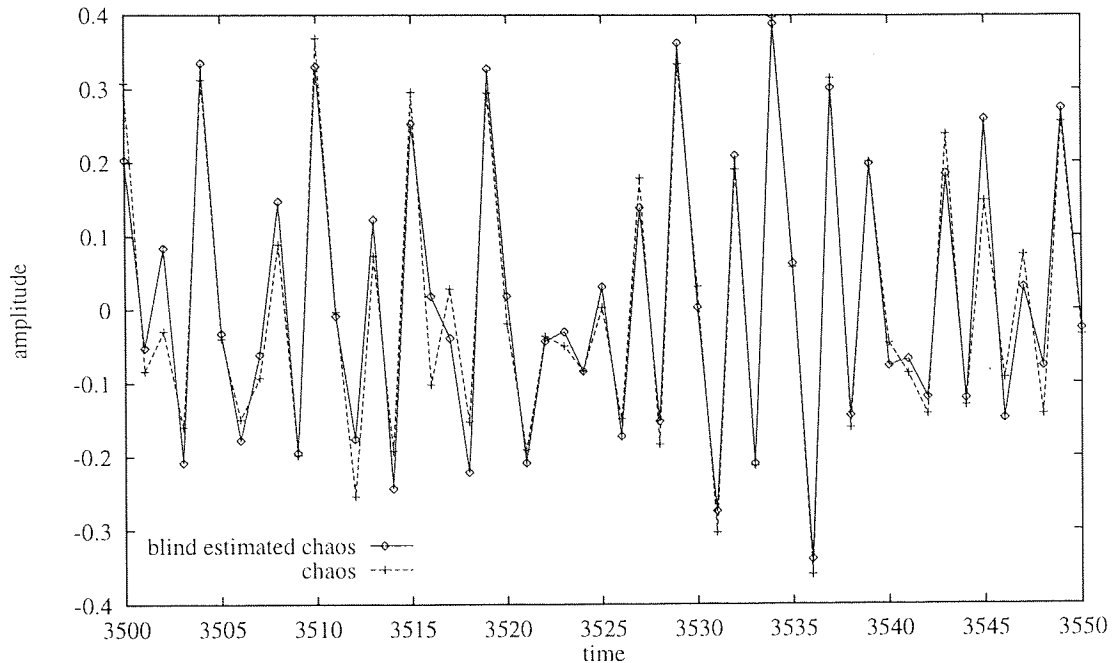
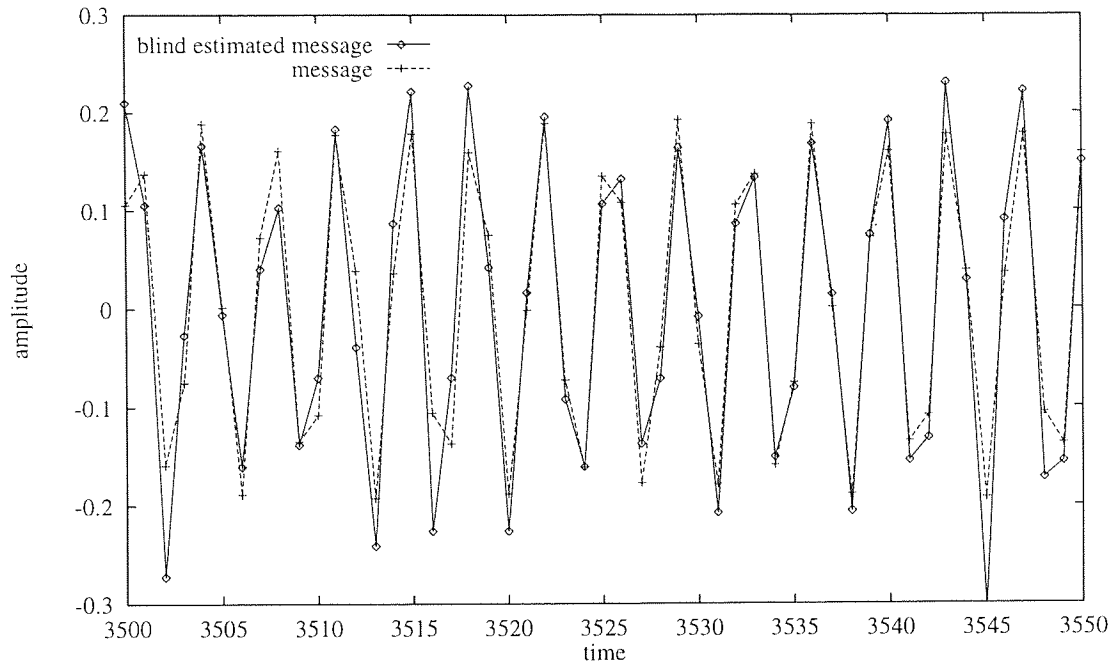


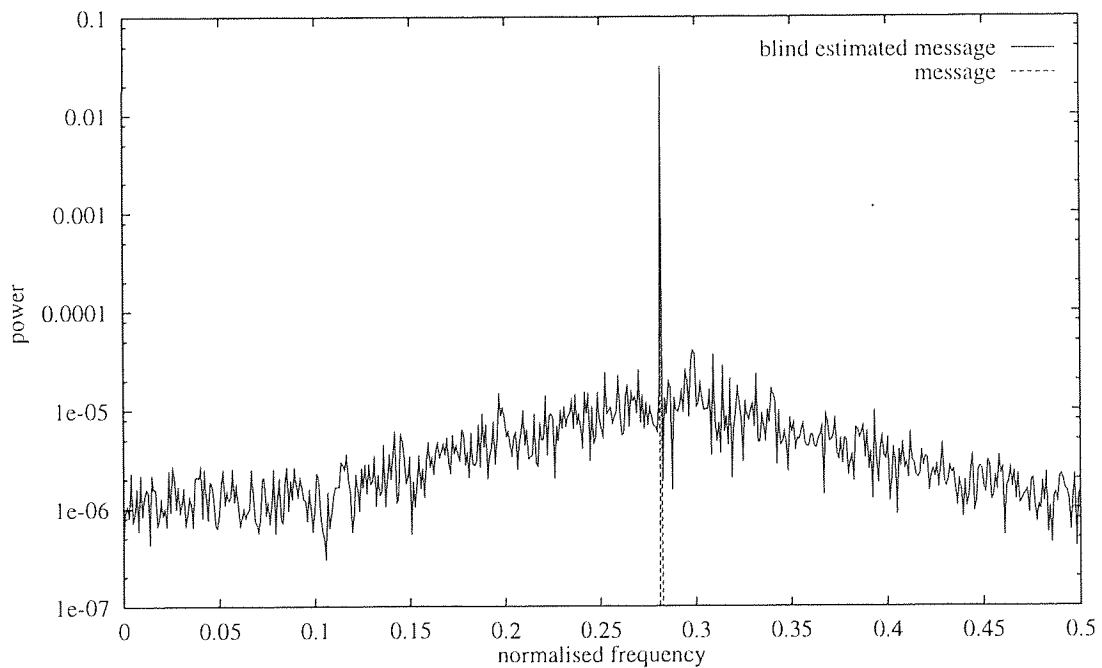
Figure 5.20 Isolated chaotic signal obtained through blind prediction from the filtered time series. The estimated time series is a reasonably good approximation to the chaos.

time series $\{v_i\}$ to construct a blind predictor \widehat{w}_{5,ν^*} for the chaotic component $v_i^{(chs)} = w_{5,\nu^*}(\mathbf{y}_i)$. The errors, calculated from equation (5.6) on training and test sets, of $\epsilon_{5,\nu^*} \approx 0.23$ and 0.25, respectively, are reflected in the reconstructed values $\widehat{v}_i^{(chs)} = \widehat{w}_{5,\nu^*}(\mathbf{y}_i)$, plotted in figure 5.20 on top of their target values over part of the test set. Compared to the original, corrupted time series of figure 5.18, this estimate is relatively close to its target value of $v_i^{(chs)}$ almost everywhere, although certain regions in \mathcal{M}_{5,ν^*} apparently give rise to significantly larger errors than others. Nevertheless, this is an encouraging result, particularly as it represents the *blind* approach to signal separation, which relies on the averaging aspect of the LS RBF algorithm to counteract the presence of the $v_i^{(msg)}$ component in v_i , and is illustrated by an out-of-sample data set.

The corresponding estimates $\widehat{v}_i^{(msg)} = v_i - \widehat{v}_i^{(chs)}$ of $v_i^{(msg)}$ are plotted in figure 5.21(a), superimposed on the message itself: blind prediction has apparently succeeded in isolating the message to a fair degree of accuracy. This is because the *unnormalised* prediction error $\sigma_{chs}\epsilon_{n,\nu^*}$, in this experiment, is small compared to σ_{msg} , where σ_{msg}^2 is N times the variance of the message. If these terms had been on the same order of magnitude then we would have been unable to extract a useful estimate of $v_i^{(msg)}$ from v_i with the predictor in question. In other words, in any experiment of this nature, the ratio of $\sigma_{chs}\epsilon_{n,\nu}$ to σ_{msg} represents the limiting factor in our ability to reconstruct the message by modelling the chaos. On the other hand, for the purposes of predicting the chaos itself, the value of σ_{chs} is immaterial, provided that we are able consistently to achieve a filtered value of $u_i^{(msg)} \ll u_i^{(chs)}$. We also plot, in figure 5.21(b), the power spectra of original and predicted message. A comparison of this plot with that of the original and corrupted messages, in figure 5.19(a), reveals a decrease in power of one or two orders of magnitude



(a)



(b)

Figure 5.21 Isolating a sinusoidal message by subtracting the blind predicted chaos from the corrupted signal, plotted in both time and frequency domains. (a) the recovered message is comfortably close to the original; (b) its power spectrum now has a 'noise floor' an order of magnitude below the level of the chaotic power spectrum obtained from the unfiltered signal.

everywhere other than at $\nu = \nu^*$. Blind prediction has effectively lowered the 'noise floor' by removing a substantial part of the chaos from the composite signal.

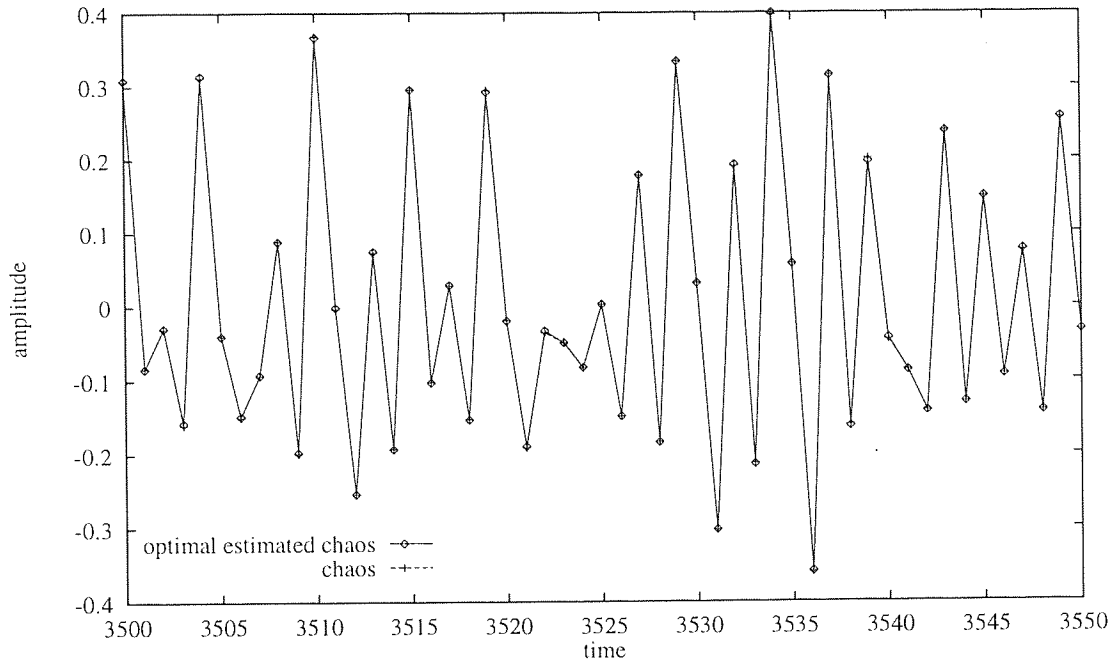
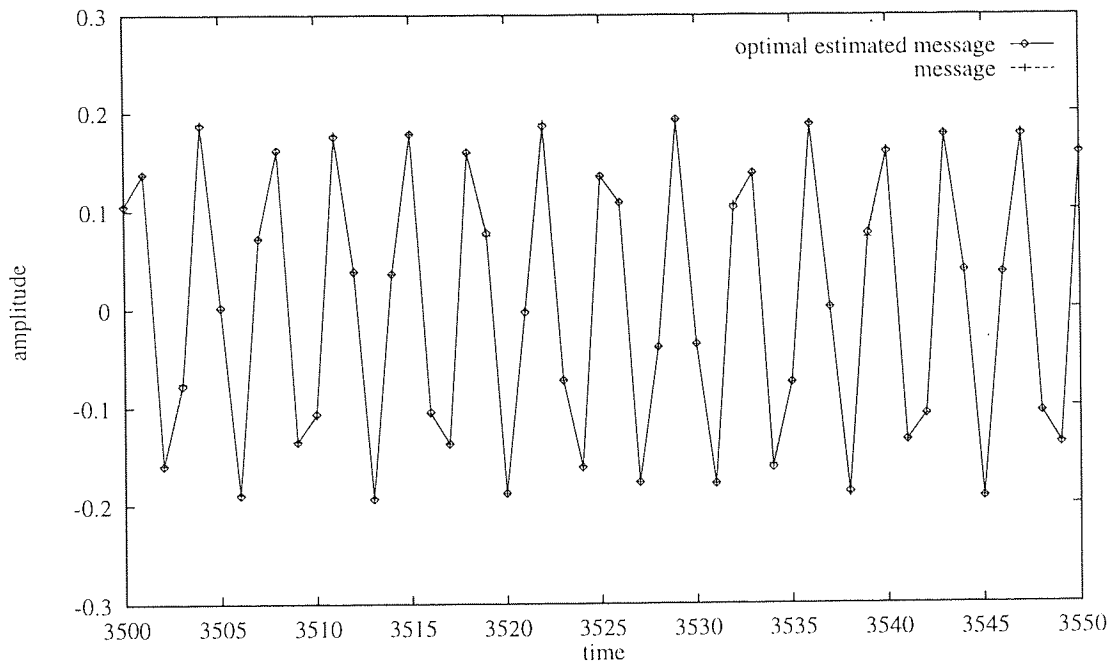


Figure 5.22 Isolated chaotic signal obtained through targeted prediction from a filtered time series. The estimated time series is now an excellent approximation to the chaos.

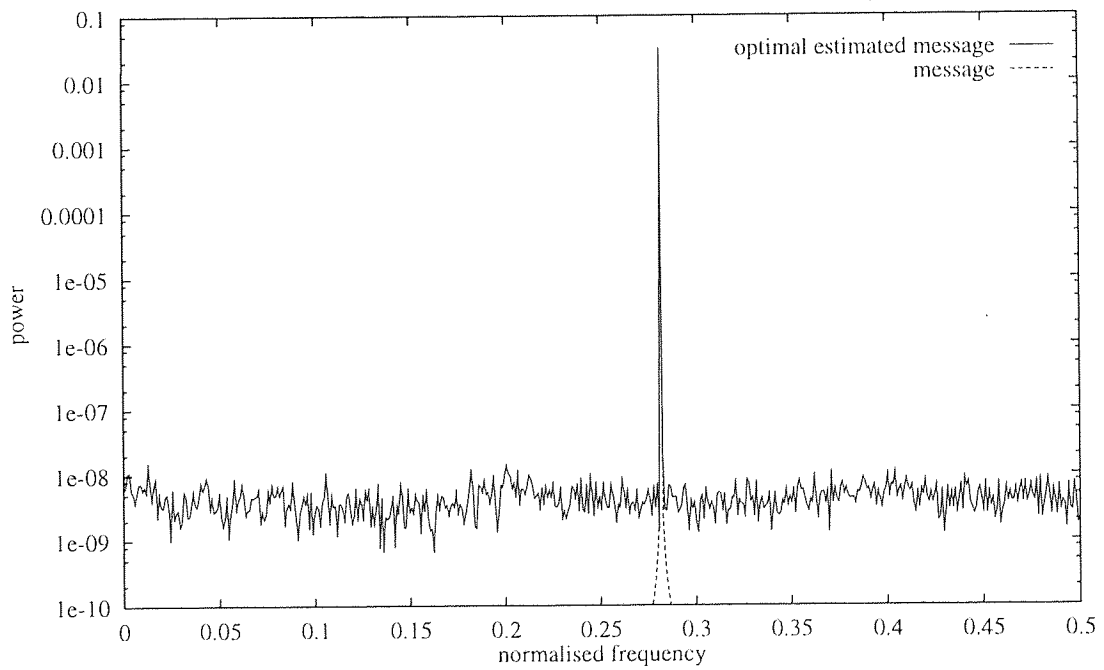
For the sake of completeness, we have also applied the targeted prediction technique to this same data set, replacing v_i with $v_i^{(chs)}$ in the training set for the \widehat{w}_{5,ν^*} . The resulting errors, now obtained by direct minimisation of equation (5.6), are $\epsilon_{5,\nu^*} \approx 0.0066$ and 0.0074 , on training and test set respectively. Given the small size of these errors, it is not surprising that the image $\widehat{v}_i^{(chs)}$ of $\mathbf{y}_i \in \mathcal{M}_{5,\nu^*}$, plotted in figure 5.22, is indistinguishable from the the chaos itself, confirming that the effect of the FIR filter $(1, 0.3902, 1)$ on the delay embedded Ikeda system is trivial to undo with an RBF inverse to \mathbf{f}_{5,ν^*} . Once again, in figures 5.23(a) and (b) we plot the corresponding actual and estimated message values and power spectra, respectively; naturally, these exhibit the same degree of accuracy as does the previous figure. In part (b), in particular, we notice that the effective (chaotic) noise floor has been substantially reduced from its original level in $\{v_i\}$ by this method.

5.4.2 Isolating a phase modulated message from Lorenz chaos

We explore this approach further by attempting to recover a phase modulated message from the chaotic time series $\{v_i^{(chs)}\}$, plotted in figure 2.4(b), obtained by integrating the Lorenz system with a step size of 0.1. Our first step is to construct the time series $\{s_i = \pm 1\}$, where the sign of s_i is allowed to change randomly once every 30 samples. This sequence is used to modulate a sinusoidal carrier wave,



(a)



(b)

Figure 5.23 Isolating a sinusoidal message by subtracting the targeted prediction of the chaos from the corrupted signal, plotted in both time and frequency domains. (a) the recovered message is indistinguishable from the original; (b) its power spectrum has a true noise floor at the level of numerical precision.

with frequency $\nu^* \approx 0.17$, to give $v_i^{(msg)} = s_i \cos 2\pi\nu^*i$, which is the message we will be attempting to reconstruct in this example. This method of encoding is known as phase shift keying (PSK), as each change of sign in the binary sequence $\{s_i\}$ gives rise to a corresponding phase change of π radians in the modulated waveform $\{v_i^{(msg)}\}$, and is commonly used in the transmission of binary information. The demodulation process consists of a second multiplication of v_i by the carrier wave, followed by a low-pass filter with $\nu \leq \nu^*$. It should be noted that a waveform encoded in this manner will necessarily lose some of its high frequency components as a result of the demodulation process. However, as it is the binary sequence itself in which we are interested, we can clearly accept any degree of smoothing provided the recovered signal exhibits the necessary zero-crossings. The PSK process is described in more detail in Taub and Schilling [40].

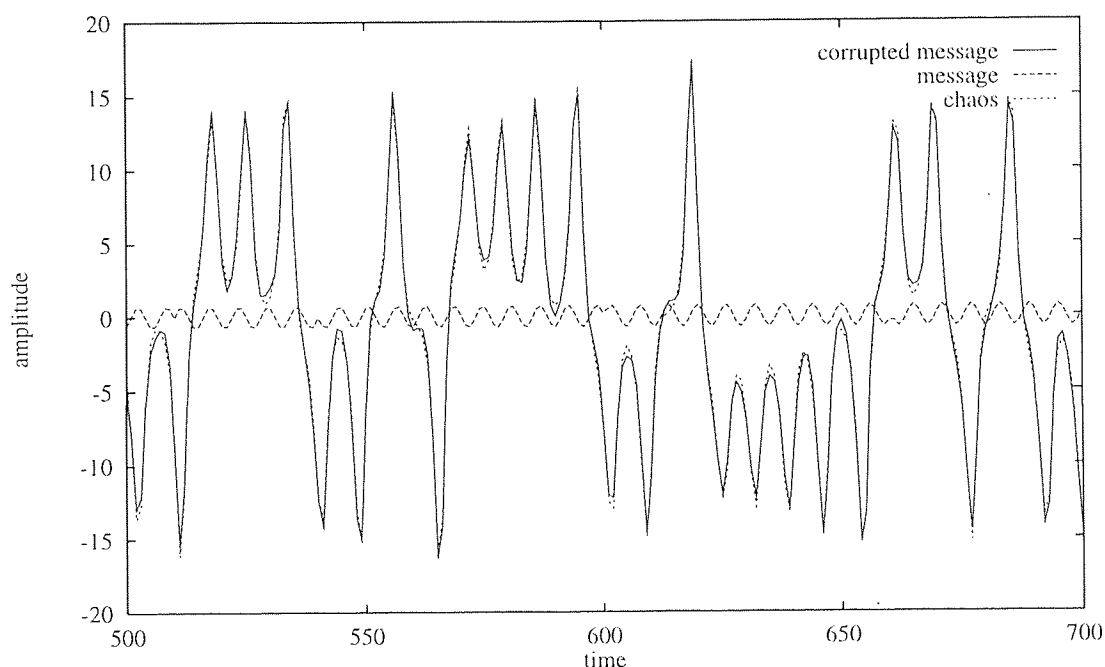
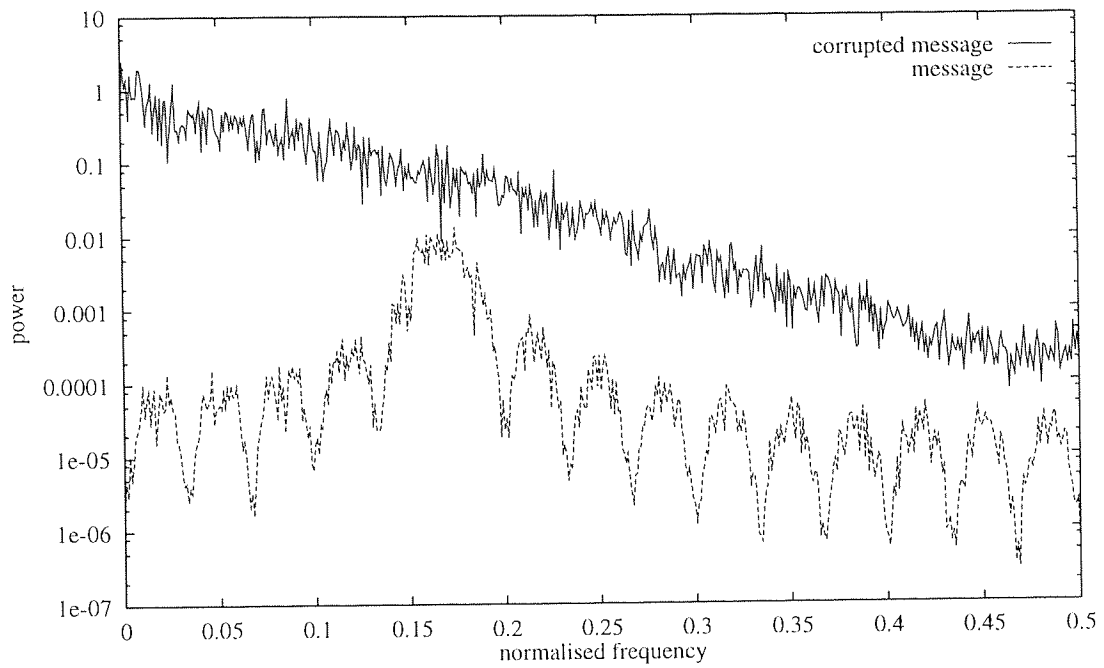
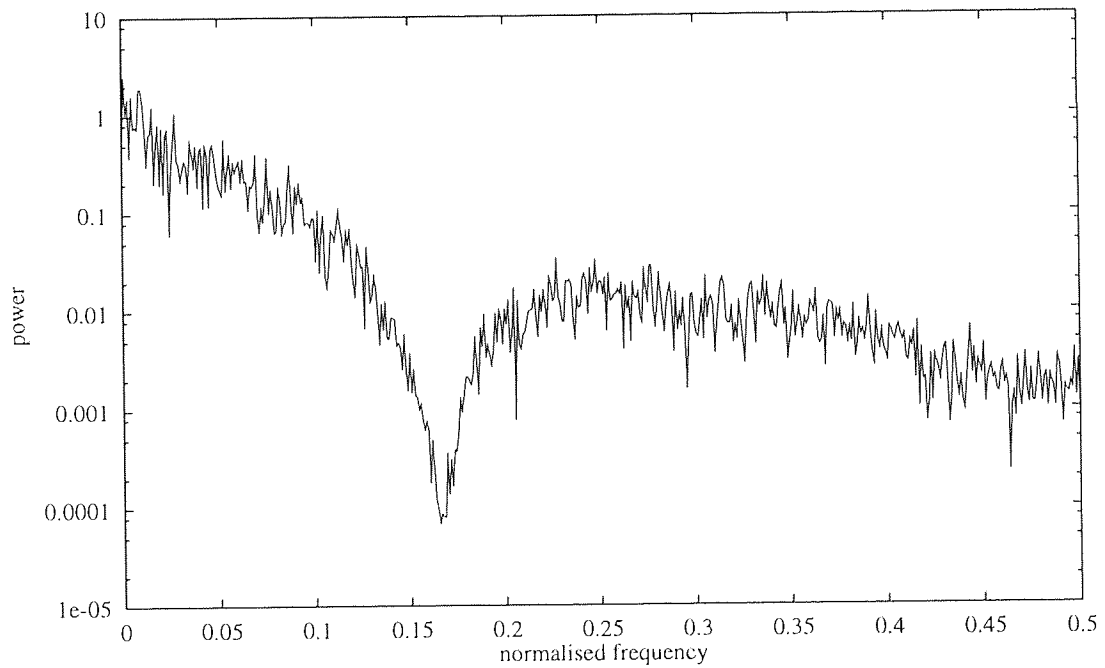


Figure 5.24 A phase modulated message corrupted by additive Lorenz chaos is superimposed on both its individual message and chaos components. Owing to the large amplitude difference, the corrupted signal is nearly indistinguishable from the chaos in this case.

The carrier frequency ν^* was chosen so that the FIR filter $(1, -1.007, 1)$, with a spectral null at 0.17, also removes a significant amount of power at $\nu \approx 0.15$, whose corresponding filter has already been shown to be non-generic with respect to the Lorenz system. This should make the task of inverting its effect on a delay embedding of the Lorenz system a little harder than might otherwise be the case. A further complication is due to the artificially high frequencies generated by the phase changes in $\{v_i^{(msg)}\}$, which will *not* be zeroed by the chosen FIR filter. The amplitude of $v_i^{(msg)}$ was chosen to be one tenth the standard deviation 7.92 of the chaos, which should also have a negative effect on our ability to recover it by predicting the chaos, as previously noted. We plot the composite time series elements v_i , superimposed on both the message $v_i^{(msg)}$ and the chaos $v_i^{(chs)}$ in figure 5.24.



(a)



(b)

Figure 5.25 Power spectra for a phase modulated message with additive Lorenz chaos, before and after filtering. (a) Superimposing the corrupted message on the message itself, we see that the message is at least an order of magnitude below the chaos, whose spectrum is virtually unchanged by the addition of the message; (b) the filter leaves a hole at the expected frequency, removing the carrier along with a substantial proportion of the overall power from the signal.

In figure 5.25(a) we plot the power spectra obtained from both original (phase modulated) and corrupted messages. In contrast to the corresponding spectra calculated in the previous experiment (figure 5.19(a)), the contribution of the message component is hardly noticeable in the combined spectrum, which exhibits a roughly exponential fall-off with increasing frequency. Even at its peak, the power in the message is an order of magnitude lower than that in the chaos. The power spectrum of the filtered time series $\{u_i\}$ is shown in figure 5.25(b). As expected, in addition to completely cancelling the carrier frequency, this filter also removes a substantial proportion of the power at $\nu = 0.15$.

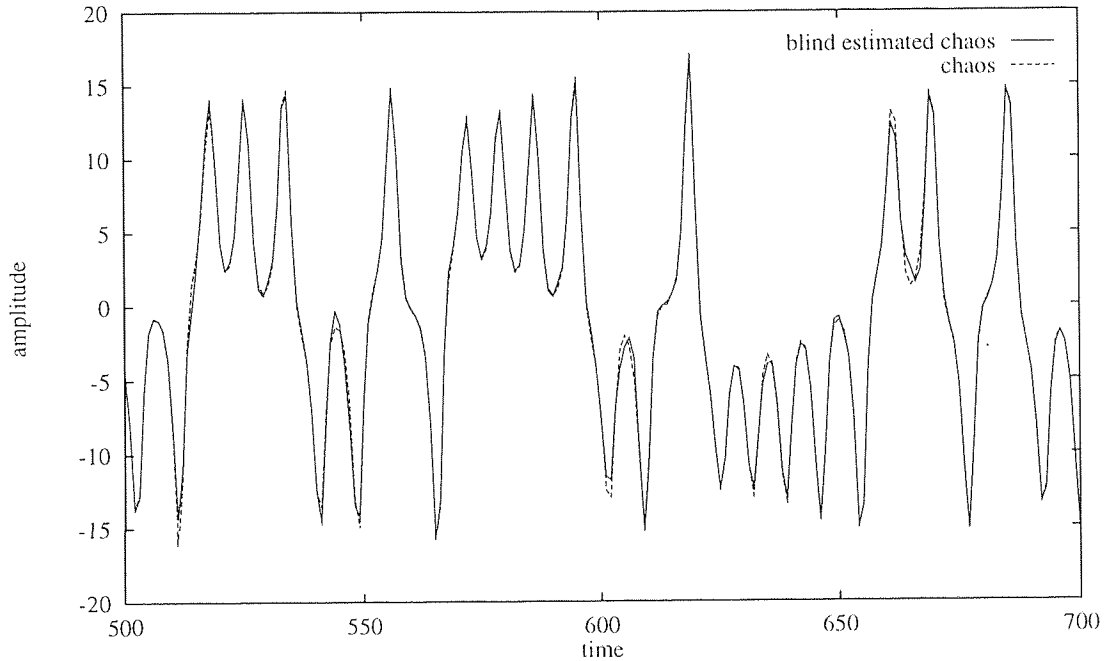
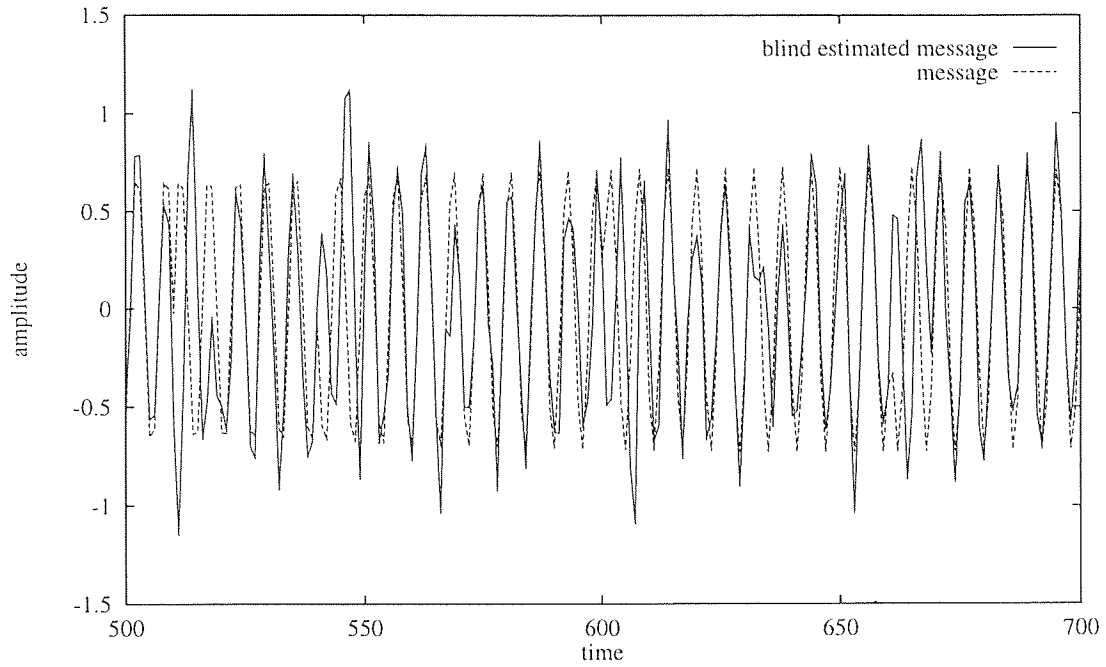


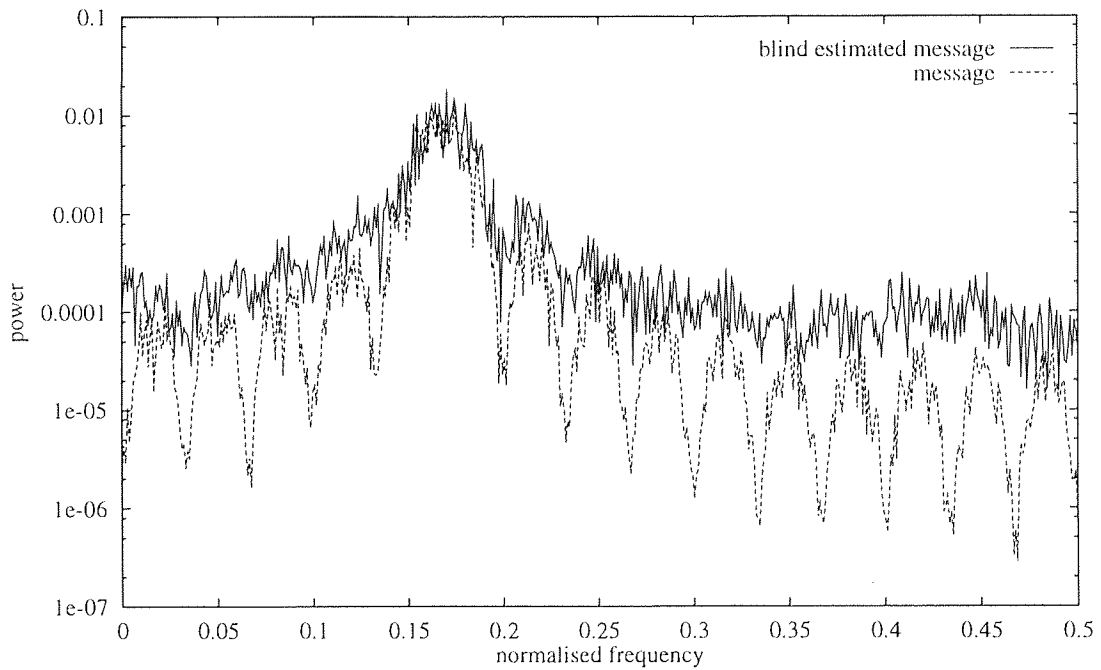
Figure 5.26 Chaotic time series estimated by blind prediction, superimposed on the original. The predicted curve is observably close to its target.

We now construct a filtered embedding of \mathcal{M} in \mathbb{R}^9 and, using the time series $\{v_i\}$ as target, train the blind predictor \widehat{w}_{g,ν^*} to reproduce the chaotic component $v_i^{(chs)} = w_{g,\nu^*}(y_i)$, relying on the LS algorithm to take care of the message component of v_i , assuming that it is uncorrelated with the chaos. The resulting errors, calculated over training and test sets from equation (5.6), of $\epsilon_{g,\nu^*} \approx 0.058$ and 0.060 , respectively, are reflected in figure 5.26, which superimposes $\widehat{v}_i^{(chs)}$ on $v_i^{(chs)}$ over a portion of the training set (there is no over-fitting, to speak of, so a similar picture is obtained from the test set).

As this figure, and the small error values demonstrate, we have succeeded in fitting the chaos to a reasonable approximation, despite having deliberately chosen a FIR filter close, in frequency space, to the non-generic filter $(1, -1.2, 1)$. However, when the errors visible in this plot are compared with those due to the presence of the message component in 5.24, it is not obvious that this blind estimate of $v_i^{(chs)}$ is sufficiently accurate to provide as good a fit $\widehat{v}_i^{(msg)} = v_i - \widehat{v}_i^{(chs)}$ to $v_i^{(msg)}$ as was obtained in the previous example. We plot these time series in part (a) of figure 5.27, which reveals a substantially worse



(a)



(b)

Figure 5.27 Reconstructing a message by subtracting the chaos estimated by blind prediction, in both time and frequency domains. (a) The recovered message is a reasonably good approximation to the original in some places, bad in others; (b) nevertheless, on comparing power spectra we see that the carrier frequency has been pulled out of the background chaos, to a certain extent, with only the harmonics still buried.

fit than was evident in the Ikeda experiment. Message and blind estimate are nevertheless recognizably close, over most of the figure. Furthermore, the power spectrum of $\{\hat{v}_i^{(msg)}\}$, plotted in part (b), is similar to that of $\{v_i^{(msg)}\}$, confirming that \widehat{w}_{9,ν^*} has succeeded in bringing the chaotic background well below the level of the carrier frequency ν^* , although the harmonics are effectively buried.

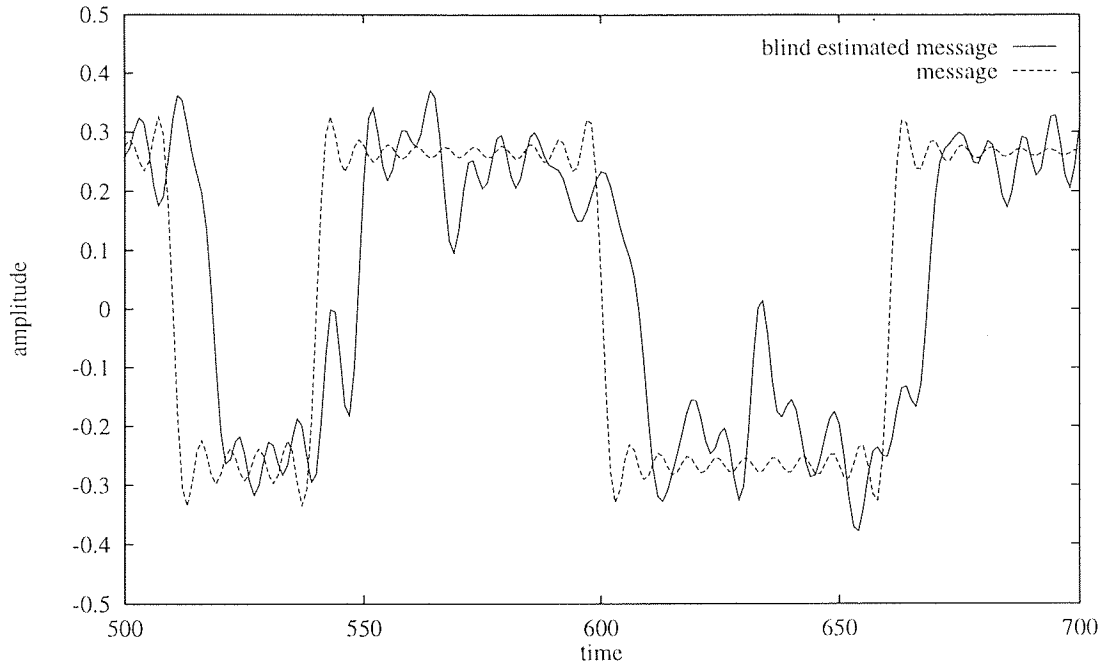


Figure 5.28 Result of demodulating the original and recovered messages. Demodulation of both messages causes a loss of high-frequency information, but despite the large errors incurred in the blind prediction process we have succeeded in recovering the binary sequence to a good approximation.

To complete our analysis of the blind predictor we plot, in figure 5.28, the result of demodulating both the original message and its estimate. This figure provides a useful illustration of the low-pass filtering which is an inevitable consequence of the demodulation process. The demodulated estimate largely preserves the gross structure of the binary sequence it is intended to approximate, although it is subject to rather more high-frequency oscillation. If we were to merely concern ourselves with the *sign* of the reconstructed signal, however—tracking only the zero-crossings, in other words—we would find that it was a relatively faithful copy of the original sequence. (Assuming that we had access to some common time frame of reference for the changes in sign we might go even further, in this particular application, using our prior knowledge of the 30-sample delay between sign changes to eliminate a percentage of any spurious zero-crossings which might arise.)

In the case of the targeted predictor, trained by direct minimisation of equation (5.6), we obtain only a marginally better estimate $\hat{v}_i^{(chs)}$ of $v_i^{(chs)}$ than in the blind case, with normalised fitting and test errors of $\epsilon_{9,\nu^*} \approx 0.055$ and 0.059 , respectively; any improvement in the estimated time series is virtually impossible to discern. It is, of course, no surprise that as the message component becomes smaller (with respect to the chaos) the advantages of having access to the chaotic signal on its own decrease in proportion; in any

case, it is the blind predictive method on which we would typically be forced to rely in practice, and we feel that its usefulness has been successfully demonstrated by this, and the previous example.

Chapter 6

Conclusions

We set out in this thesis to investigate the application of radial basis function maps to the detection of diffeomorphisms between delay embedded dynamical systems. To this end, we defined the diffeomorphism as an invertible, differentiable map whose inverse is also differentiable, using it in turn to define the d -dimensional (differentiable) manifold as a topological space *locally* diffeomorphic to \mathbb{R}^d . We then defined the dynamical system (\mathcal{M}, ψ) , consisting of a manifold \mathcal{M} and a diffeomorphism $\psi: \mathcal{M} \rightarrow \mathcal{M}$, and explained how, through the use of a time series $\{v_i\}$ obtained from almost any measurement function $v: \mathcal{M} \rightarrow \mathbb{R}$, Takens' theorem allows us to construct a differentiably equivalent copy (\mathcal{M}_m, ψ_m) of (\mathcal{M}, ψ) in \mathbb{R}^m with the delay embedding $\Phi_{v,m}: \mathcal{M} \rightarrow \mathbb{R}^m$, provided that $m > 2d$. With this framework in place, we went on to discuss the construction of filtered embeddings, where a linear transformation $\mathcal{F}: \mathbb{R}^m \rightarrow \mathbb{R}^n$ is applied to (\mathcal{M}_m, ψ_m) to obtain a further copy (\mathcal{M}_n, ψ_n) of (\mathcal{M}, ψ) . If \mathcal{F} embeds \mathcal{M}_m then its restriction $f: \mathcal{M}_m \rightarrow \mathcal{M}_n$ to \mathcal{M}_m is necessarily a diffeomorphism. Conversely, therefore, if we could show that f , so defined, is a diffeomorphism then we would have established the existence of a differentiable equivalence between (\mathcal{M}_m, ψ_m) and (\mathcal{M}_n, ψ_n) . (Because we are specifically interested in dynamical systems, we insisted that the correspondence between $x_i \in \mathcal{M}_m$ and $y_i \in \mathcal{M}_n$ be determined by a common time index i .) Of course, establishing that \mathcal{F} embeds \mathcal{M}_m does not necessarily tell us anything about $\Phi_{v,m}$, but we overcame this issue with appropriately designed experiments.

In order to develop a tool with which to determine whether or not a given f is a diffeomorphism we introduced the RBF map, a powerful nonlinear model, linear in its basis functions, which we use to construct approximations $\hat{f}: \mathbb{R}^m \rightarrow \mathbb{R}^n$ and $\hat{f}^{-1}: \mathbb{R}^n \rightarrow \mathbb{R}^m$ to f and its inverse. We showed that in its classical form \hat{f} is generically an embedding of compact sets, provided that $p > m$, which means that even if f is not a diffeomorphism, its RBF model, optimised by LS error minimisation, almost certainly will be. We therefore formulated a test for diffeomorphism which consists of examining the errors arising

in fitting both \widehat{f} and \widehat{f}^{-1} , suitably optimised by the method of LS, to determine whether or not either approximates a map which is *not* an injective immersion. We also considered the possibility that the LS error might prove insufficiently sensitive to the self-intersections in \mathbb{R}^m or \mathbb{R}^n responsible for a loss of diffeomorphism. To detect such cases, we described an analysis of the distribution of per-point errors, through which estimates of the Lipschitz constants of f and its inverse—if they exist—might be obtained.

Before proceeding further, we analysed the method of LS in some detail, deriving its solution in terms of a fixed set of centers and also through Chen’s adaptive algorithm, known as forward selection. For the former case, we introduced a compelling ad hoc center selection method, which we called repulsive selection, and evaluated its performance on a particular function approximation problem by comparing the errors resulting from a straightforward linear LS minimisation, using distributions of both randomly-seeded repulsive and purely random centers, with those resulting from forward selection of an equivalent number of centers. The results clearly showed that the repulsive method outperforms random selection (on average) by a substantial margin, although the forward selection method achieved a similar improvement over repulsive selection. The implementation of forward selection is, however, substantially more time-consuming than that of linear LS, so we decided to concentrate on repulsive selection for the purposes of this thesis.

We then investigated the phenomenon of over-fitting in LS maps, in which the model fits relationships in the training set but not the test set, typically due to the presence of noise in the former. The hallmark of over-fitting is a systematic divergence of training and test errors with increasing numbers of centers (degrees of freedom), and a traditional approach to its elimination is to restrict the rank of the linear part of the RBF map through projection out of singular subspaces whose variance is below some predetermined noise floor. We described an enhancement to this technique which consists of reordering the basis vectors by their contribution to the overall fitting error before eliminating those with the smallest contribution. We compared the result of applying both methods to a particular map, and found that the enhanced version consistently produced the smallest error, for a given rank, on both training and test sets. Despite this success, however, we did not attempt to constrain the rank of the RBF models used in later experiments, as we would be more concerned with the variation in error with some experimental parameter, and rank truncation would introduce an unnecessary level of complication at that stage.

As an alternative to the potential shortcomings of the LS error as indicator of a non-diffeomorphic relationship, we also considered a symmetrical form of the RBF map, trained by minimising Van Huffel’s TLS error. In this form, we showed how the question of whether or not f is a diffeomorphism can be reduced to an analysis of the degree of ill-conditioning exhibited by two square matrices which, in composition, form the linear part of the symmetrical RBF map. We noted, however, that the TLS method was likely to suffer from numerical instability brought on by rank deficiencies in either its domain or range.

Having discussed the LS and TLS approaches to RBF approximation we went on to compare their usefulness in detecting diffeomorphisms between manifolds *not* obtained by delay embedding. These

experiments consisted of fitting RBF approximations to maps f_μ between projections of a circle into \mathbb{R}^2 , for different angles of projection, and between projections of a 2-torus into \mathbb{R}^3 and \mathbb{R}^4 , for various ratios of radii. In both cases, the problem domain was designed around a critical parameter value μ^* so that the map in question was a diffeomorphism only for $\mu < \mu^*$, and when the diffeomorphism broke down it was through loss of immersivity (at $\mu = \mu^*$) or loss of injectivity (for $\mu \geq \mu^*$) in the forward direction. We therefore hoped to find good approximations to both f_μ and its inverse for $\mu < \mu^*$, but for $\mu \geq \mu^*$ we hoped to find $\widehat{f_\mu^{-1}}$ observably unable to find a good fit to its training set. In both of these experiments we found that the LS error, calculated in both forward and inverse directions, was actually an extremely good indicator of diffeomorphism: the forward error was consistently small, while the inverse made the expected, clearly defined transition from small to large error as μ increased through its critical value. Our expectation that a LS RBF map should be able to approximate a map f_μ which, although not necessarily a diffeomorphism itself, is nevertheless a function, was confirmed by the uniformly small forward error. That the inverse error should be so successful at detecting the non-invertibility of $f_{\mu \geq \mu^*}$, however, was less expected, but is clearly the result of attempting to approximate a one-to-many relationship with the RBF map $\widehat{f_\mu^{-1}}$. The success of this approach—using a scalar LS error to detect non-diffeomorphic maps of this kind—is very encouraging, and stems from our decision to construct *independent* approximations $\widehat{f_\mu}$ and $\widehat{f_\mu^{-1}}$ to the relationship in question. It is also a little surprising, as we had been concerned that the averaging process inherent in the LS error calculation might make it insensitive to large per-point errors occurring on a self-intersecting set of limited extent. That this turns out not to be the case (at least in these two examples) is apparently due to the fact that a small LS error can *not* generically be found by constructing an RBF ‘approximation’ (in the limited sense of chapter 4) to a non-injective, or even merely non-immersive f_μ .

On the assumption that insensitivity of the LS error measure to localised per-point errors might become an issue in future experiments, we applied the Lipschitz analysis to the maps between both circles and tori. The upper and lower bounds, calculated in this manner from individual RBF approximations $\widehat{f_\mu}$ and $\widehat{f_\mu^{-1}}$, were unfortunately extremely noisy when plotted versus μ . As a result, in neither experiment could we clearly discern the behaviour predicted by the corresponding analytical calculations (although to a very rough approximation the curves exhibited similar trends). This is, however, to be expected on moving from an average (the LS error) to an upper or lower bound; however, when we examined the error *distributions*, averaged over several RBF models, we saw fairly strong—if subjective—evidence for the predicted bounds; this led us to move the average over random repulsive center seeds *inside* the extremum calculation, to obtain new empirical bounds more closely resembling their analytical analogues.

The symmetrical RBF map, trained by minimising the TLS error and applied to both circles and tori, exhibited a similarly large variance with respect to the choice of centers, presumed to be due to the sensitivity to rank deficiency noted earlier. As predicted, the condition numbers mirrored the errors closely, and could clearly be used in their place. The overall form of these curves, however, was roughly

as expected in both experiments, so the method certainly deserves further investigation.

Having found that the LS error, calculated in both forward and inverse directions, actually produced the best quantitative results in both circle and torus problem domains, we finally applied this technique to detecting diffeomorphisms between delay embedded dynamical systems. Using both numerically simulated and experimentally generated time series $\{v_i\}$, we looked at four families of maps f , each the restriction to $\mathcal{M}_m = \Phi_{v,m}\mathcal{M}$ of a linear transformation $\mathcal{F}: \mathbb{R}^m \rightarrow \mathbb{R}^n$. By virtue of the linearity of \mathcal{F} , these experiments were once again designed so that a loss of diffeomorphism between \mathcal{M}_m and $\mathcal{F}\mathcal{M}_m$ would occur only through a loss of injectivity in f , so we were able to use this prior knowledge to avoid constructing an explicit RBF approximation to f . And, as a consequence of the delay structure in \mathcal{M}_m , we were able further to restrict our task to fitting one or more individual components of f^{-1} .

In the first such experiment we attempted to determine a minimum embedding dimension m^* for the system (\mathcal{M}, ψ) by fitting a predictor for the time series generator $w_m: \mathcal{M}_m \rightarrow \mathbb{R}$. This experiment was quite successful: in the case of the noise-free, numerically simulated Ikeda and Hénon systems, and the 0.1-step Lorenz system, we were able to find—in each case—a clearly defined estimate of m^* which agreed closely with that predicted by Takens (bearing in mind that Takens essentially establishes an upper bound for the minimum embedding dimension). In the case of the 0.01-step Lorenz system, however, we were unable to determine a suitable value for m^* due to the limited size of the delay window used. In the presence of noise, we found that the Lorenz systems integrated with both 0.01 and 0.1 step sizes gave rise to an error substantially larger than in the noise-free cases, leading us to pose the open question of just how small such an error should be before we ascribe it to the approximation of a diffeomorphism between embedded dynamical systems.

When we applied the same procedure to the projection of \mathcal{M}_m into singular subspaces of \mathbb{R}^m we found that for m large enough we could indeed determine a minimum subspace dimension n^* such that the predictor $\widehat{w}_{m,n}: \mathbb{R}^n \rightarrow \mathbb{R}$ achieved approximately the same error for $n \geq n^*$ as did the corresponding predictor $\widehat{w}_m: \mathbb{R}^m \rightarrow \mathbb{R}$ on the unfiltered manifold. However, we did *not* see $\widehat{w}_{m,n} < \widehat{w}_m$ for any of the three systems examined, indicating that dealing with the variance in the orthogonal complement of \mathbb{R}^n did not—at least in these examples—improve the predictive ability of the resulting map, although it undeniably resulted in a simpler map (in the sense of fewer degrees of freedom). The smallest error was always obtained from the unfiltered delay reconstruction, but in only one of the experiments performed was \mathcal{M}_m actually embedded in \mathbb{R}^m to begin with (the noisy Lorenz time series and laser intensity time series have already been shown to be non-generic observations, in the sense of Takens), so it is not unlikely that further experiments might reveal a more positive result for the singular subspace projective method.

The third experiment was a little more complex, involving a linear transformation which implemented a delay reconstruction from a FIR filtered copy $\{u_i\}$ of the original time series $\{v_i\}$. Our aim was to detect periodic orbits in (\mathcal{M}, ψ) by tuning the FIR filter to kill a particular frequency in $\{v_i\}$, then looking for a significantly large error in predicting the original time series as a function of the filtered delay vectors in the

image of $\Phi_{u,n}$. The rationale was that a FIR filter generically preserves the dynamical information in $\{v_i\}$ necessary for $\Phi_{u,n}$ to be an embedding on \mathcal{M} , so if we find that $f = \Phi_{u,n} \circ \Phi_{v,m}$ is *not* a diffeomorphism only for a filter designed to kill the frequency ν^* then we can assume that there is an orbit of period $\frac{1}{\nu^*}$ in \mathcal{M} visited often enough to influence the RBF approximation under investigation. This approach proved successful in the analysis of the Ikeda, Hénon and Lorenz systems, identifying orbital periods which were manifestly present in their respective attractors.

Having established that we could generically ‘undo’ the effect of a FIR filter on a delay embedding by constructing an approximator for $\{v_i\}$ in this manner, we went on to examine an interesting application of the technique, in the form of the separation of deterministic noise, in the form of a time series $\{v_i^{(chs)}\}$, from an encoded message $\{v_i^{(msg)}\}$. Relying solely on the requirement that the message be encoded in such a way that it could be largely eliminated by a FIR filter tuned to stop the frequency ν^* , we found that we were able to successfully ‘invert’ the filter, and hence reconstruct the chaotic component, with the RBF approximation \widehat{w}_{n,ν^*} to the time series generated by w_{n,ν^*} . The degree of accuracy to which we were thereby able to extract the message itself, by subtracting the estimated chaos, was seen to depend on the ratio of its standard deviation to that of the residual $\{\widehat{v}_i^{(chs)} - v_i^{(chs)}\}$. In constructing \widehat{w}_{n,ν^*} we relied on the incoherence of chaos and message to take account of the presence of the latter in the RBF training set, but we noted also that by arranging for periodic interruptions in the message broadcast we could reduce the residual further by approximating the chaos directly. This approach turned out to be very successful, in terms of predicting the chaotic component, for both a sine wave corrupted with Ikeda noise and for a phase-modulated signal corrupted with Lorenz noise. In the former example, the extracted message was also extremely close to the original, using both blind and targeted predictive methods, and in the latter it was somewhat less successful, due to the relatively small message amplitude, but still recognizably intact.

In conclusion, then, we have shown that the radial basis function map may successfully be used to detect the existence, or not, of diffeomorphic relationships between delay embedded dynamical systems, despite having also shown that it is itself generically an embedding of compact sets, and therefore at first sight unsuited for such a purpose. We overcame this apparent limitation by fitting LS RBF maps in both forward and inverse directions, yielding an error measure which we have demonstrated to be sensitive to those self-intersections, possibly of small measure, which may occur as the result of (for instance) a failed delay embedding. We have investigated several applications of this method to maps between delay embedded manifolds, and achieved encouraging results in all cases. These experiments are of more than academic interest, as they deal with issues commonly raised in the course of experimental time series analysis, such as the determination of a minimum embedding dimension and the choice of a singular subspace for removal of stochastic noise. The final experiment in deterministic noise cancellation, combining the results of Takens on delay embedding with those of Broomhead, Huke and Muldoon on FIR filtering, is of particular interest, as it has potential application to the more general field of signal processing.

References

- 1 R. Badii, G. Broggi, B. Derighetti and M. Ravani, "Dimension Increase in Filtered Chaotic Signals", *Phys. Rev. Lett.* **60** 11 (1988) 979–982.
- 2 D. S. Broomhead, *Private communication* (1995).
- 3 D. S. Broomhead, J. P. Huke and M. R. Muldoon, "Linear filters and non-linear systems", *J. R. Statist. Soc. B* **54** 2 (1992) 373–382.
- 4 D. S. Broomhead, J. P. Huke and M. A. S. Potts, "Controlling deterministic noise by constructing nonlinear inverses to linear filters", *Physica D* **89** (1996) 439–458.
- 5 D. S. Broomhead, R. Jones and G. P. King, "Topological dimension and local coordinates from time series data", *J. Phys. A: Math. Gen.* **20** (1987) L563–L569.
- 6 D. S. Broomhead and D. Lowe, "Multivariable functional interpolation and adaptive networks", *Complex Systems* **2** (1988) 321–355.
- 7 D. S. Broomhead and G. P. King, "Extracting qualitative dynamics from experimental data", *Physica D* **20** (1986) 217–236.
- 8 M. Casdagli, "Nonlinear prediction of chaotic time series", *Physica D* **35** (1989) 335–356.
- 9 S. Chen, C. F. N. Cowan and P. M Grant, "Orthogonal Least Squares Learning Algorithm for Radial Basis Function Networks", *IEE Trans. on Neural Networks* **2** 2 (1991) 302–309.
- 10 D. R. J. Chillingworth, "Differential topology with a view to applications", *Research Notes in Mathematics* **9**, Pitman Publishing (1976).
- 11 E. S. Chng, B. Mulgrew and S. Chen, "Backtracking Orthogonal Least Squares Algorithm for Model Selection", *IEE Digest 1994/034 Mathematical Aspects of Digital Signal Processing* (1994) 10/1–10/6.
- 12 J.-P. Eckmann, "Ergodic theory of chaos and strange attractors", *Rev. Mod. Phys.* **57** 3 (1985)

- 617–656.
- 13 J. D. Farmer and J. J. Sidorowich, “Exploiting chaos to predict the future and reduce noise”, in *Evolution, Learning and Cognition*, Y. C. Lee, Editor, World Scientific Press (1988) 277.
 - 14 J. F. Gibson, J. D. Farmer, M. Casdagli and S. Eubank, “An analytic approach to practical state space reconstruction”, *Physica D* **57** (1992) 1–30.
 - 15 G. H. Golub and C. F. Van Loan, “An analysis of the total least squares problem”, *SIAM J. Numer. Anal.* **17** 6 (1980) 883–893.
 - 16 J. Guckenheimer and P. Holmes, “Nonlinear oscillations, dynamical systems, and bifurcations of vector fields”, *Applied Mathematical Sciences* **42**, Springer-Verlag (1983).
 - 17 M. Hénon, “A two-dimensional-mapping with a strange attractor”, *Communications in Mathematical Physics* **50** (1976) 69–77.
 - 18 M. W. Hirsch, “Differential Topology”, Springer-Verlag (1976).
 - 19 U. Hubner, C-O. Weiss, N. B. Abraham and D. Tang, “Lorenz-like chaos in NH_3 -FIR lasers”, *Proc. SPIE XV Time Series Prediction: Forecasting the Future and Understanding the Past*, Addison-Wesley (1994) 73–104.
 - 20 K. Ikeda, “Multiple-valued stationary state and its instability of the transmitted light by a ring cavity system”, *Opt. Commun.* **30** (1979) 257.
 - 21 D. T. Kaplan, “A model-independent technique for determining the embedding dimension”, *Proc. SPIE 2038 Chaos in Communications* (1983) 236–240.
 - 22 M. B. Kennel, R. Brown, H. D. I. Abarbanel, “Determining embedding dimension for phase-space reconstruction using a geometrical construction”, *Phys. Rev. A* **45** No. 6 (1992) 3403–3411.
 - 23 C. L. Lawson and R. J. Hanson, “Solving Least Squares Problems”, Prentice-Hall (1974).
 - 24 E. N. Lorenz, “Deterministic nonperiodic flow”, *J. Atmos. Sciences* **20** (1963) 130–141.
 - 25 D. Lowe and A. R. Webb, “Time series prediction by adaptive networks: a dynamical systems perspective”, *Proc. IEEE* **138** 1 (1991) 17–24.
 - 26 P. A. Lynn, “An Introduction to the Analysis and Processing of Signals”, MacMillan Press (1982).
 - 27 J. G. McWhirter, D. S. Broomhead and T. J. Shepherd, “A Systolic Array for Nonlinear Adaptive Filtering and Pattern Recognition”, *J. VLSI Signal Processing* **3** (1991) 69–75.
 - 28 J. W. Milnor, “Topology from the Differentiable Viewpoint”, The University Press of Virginia (1990).
 - 29 M. R. Muldoon, R. S. MacKay, J. P. Huke and D. S. Broomhead, “Topology from time series”, *Physica D* **65** (1993) 1–16.
 - 30 Oseledec, V. I., “A multiplicative ergodic theorem. Ljapunov characteristic numbers for dynamical systems”, *Trans. Moscow Math. Soc.* **19** (1968) 197–231.
 - 31 L. M. Pecora, T. L. Carroll and J. F. Heagy, “Statistics for mathematical properties of maps between time series embeddings”, to appear in *Phys. Rev. E*.

- 32 M. A. S. Potts and D. S. Broomhead, “Time series prediction with a Radial Basis Function neural network”, in *Adaptive Signal Processing*, Simon Haykin, Editor, *Proc. SPIE* **1565** (1991) 255–266.
- 33 M. J. D. Powell, “The theory of radial basis function approximation in 1990”, *University of Cambridge Numerical Analysis Reports DAMTP 1990/NA11* (1990).
- 34 W. H. Press, B. P. Flannery, S. A. Teukolsky, and W. T. Vetterling, “Numerical Recipes in C: The Art of Scientific Computing”, Cambridge University Press (1988).
- 35 D. Ruelle, “Chaotic Evolution and Strange Attractors”, Cambridge University Press (1990).
- 36 T. Sauer, J. A. Yorke and M. Casdagli, “Embedology”, *J. Stat. Phys.* **65** 3/4 (1991) 579–616.
- 37 L. A. Smith, “Identification and prediction of low dimensional dynamics”, *Physica D* **58** (1992) 50–76.
- 38 G. W. Stewart, “Introduction to Matrix Computations”, Academic Press (1973).
- 39 F. Takens, “Detecting strange attractors in turbulence”, *Dynamical Systems and Turbulence, Lecture Notes in Mathematics* **898**, Springer, Berlin (1981) 366–381.
- 40 H. Taub and D. L. Schilling, “Principles of Communication Systems”, McGraw-Hill (1971).
- 41 S. Van Huffel and J. Vandewalle, “The Total Least Squares Problem: Computational Aspects and Analysis”, *Frontiers in Applied Mathematics* **9**, SIAM (1991).
- 42 A. Weigend and N. Gershenfeld, “Time Series Prediction: Forecasting the Future and Understanding the Past”, *Proc. SFI* **XV**, Addison-Wesley (1994).

Appendix A

Inverting radial basis functions

In this appendix we prove that the nonlinear transformation $\varphi: \mathbb{R}^m \rightarrow \mathbb{R}^p$, introduced in chapter 3, is an injective immersion for $p > m$ and hence, by the result stated in chapter 2, an embedding of compact subsets $\mathcal{M} \subset \mathbb{R}^m$. We then go on to discuss the implementational aspects of calculating the inverse of $\varphi|_{\mathcal{M}}$ by the method of least squares.

A.1 Proof that φ is an injective immersion

To show that φ is immersive we must show that it is differentiable and that its derivative $D\varphi: \mathbb{R}^m \rightarrow \mathbb{R}^p$ is injective. The differentiability is easy: each component $\varphi_j(\mathbf{x}) = \phi(\|\mathbf{x} - \mathbf{c}_j\|)$ is differentiable provided that ϕ itself is differentiable on \mathbb{R}^+ and $\phi'(0) = 0$. For $p \geq m$ we can show that $D\varphi$ is injective by showing that it has full rank. We therefore write, for some $\mathbf{x} \in \mathbb{R}^m$,

$$\frac{\partial \varphi_j}{\partial x_k} = \phi'(r_j) \frac{(\mathbf{r}_j)_k}{r_j} \quad (\text{A.1})$$

where $\mathbf{r}_j = \mathbf{x} - \mathbf{c}_j$ and $r_j = \|\mathbf{r}_j\|$. We can thus decompose $D\varphi$ as the matrix product

$$D\varphi = \begin{pmatrix} \phi'(r_1) & 0 & \cdots & 0 \\ 0 & \phi'(r_2) & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \phi'(r_p) \end{pmatrix} \begin{pmatrix} \hat{\mathbf{r}}_1 \\ \hat{\mathbf{r}}_2 \\ \vdots \\ \hat{\mathbf{r}}_p \end{pmatrix} \quad (\text{A.2})$$

where $\hat{\mathbf{r}}_j \equiv \mathbf{r}_j/r_j$ denotes a unit vector in the direction of \mathbf{r}_j . We now require the p by m matrix of unit vectors $\hat{\mathbf{r}}_j$ to have rank m , which means that $p > m$ and at least m of the $\hat{\mathbf{r}}_j$ must be linearly independent; the condition on p is a strict inequality to take into account the situation when $\hat{\mathbf{r}}_j = \mathbf{0}$, that is $\mathbf{x} = \mathbf{c}_j^{(c)}$

for some \mathbf{c}_j . Now consider the diagonal matrix of derivatives $\phi'(r_j)$. The case when $\phi'(0) = 0$ does not concern us since it only occurs, by definition, when $\hat{\mathbf{r}}_j = \mathbf{0}$. However, we can clearly not let $\phi'(r_j) = 0$ anywhere else because that would zero a non-zero unit vector and might reduce the rank of $\mathbf{D}\varphi$. So, a sufficient condition for immersive φ is that ϕ be strictly monotonic on \mathbb{R}^+ .

It only remains to show that φ is injective. To do this we must show that if $\mathbf{a} = \varphi(\mathbf{x})$ then \mathbf{x} is unique. We therefore write the components of \mathbf{a} as

$$\begin{aligned} a_j &= \phi(\|\mathbf{x} - \mathbf{c}_j\|) \\ \Rightarrow \phi^{-2}(a_j) &= \|\mathbf{x} - \mathbf{c}_j\|^2 \\ &= \|\mathbf{x}\|^2 - 2\mathbf{c}_j \cdot \mathbf{x} + \|\mathbf{c}_j\|^2 \\ \Rightarrow h_j &= \mathbf{c}_j \cdot \mathbf{x} - \frac{1}{2}\|\mathbf{x}\|^2 \end{aligned} \tag{A.3}$$

where

$$h_j \equiv \frac{1}{2}[\|\mathbf{c}_j\|^2 - \phi^{-2}(a_j)] \tag{A.4}$$

The solution to this set of equations ($j = 1, \dots, p$) is the intersection of p distinct $(m - 1)$ -spheres in \mathbb{R}^m , with widths $\phi^{-2}(a_j)$ and centered at the \mathbf{c}_j . However, the intersection of two $(m - 1)$ -spheres in \mathbb{R}^m lies in an $(m - 1)$ -plane, so we if subtract any one of these equations from the rest, notionally ‘fixing’ a center, \mathbf{c}_k , we get $p - 1$ linear equations of the form

$$(\mathbf{c}_j - \mathbf{c}_k) \cdot \mathbf{x} = h_j - h_k \tag{A.5}$$

for $j = 1, \dots, p$, where $j \neq k$, the solution to which is the intersection of $p - 1$ distinct $(m - 1)$ -planes in \mathbb{R}^m . Then, provided that the $\mathbf{c}_j - \mathbf{c}_k$ span \mathbb{R}^m , \mathbf{x} is uniquely determined if $p - 1 \geq m$, that is, if $p > m$.

A.2 Approximating the inverse of φ

It is clear that the solution for $\mathbf{x} = \varphi^{-1}(\mathbf{a})$ given above only exists if $\mathbf{a} \in \varphi\mathbb{R}^m$. In practice this may often not be the case: for instance, \mathbf{a} may be the image of a linear transformation in a symmetric RBF map. If this is the case then a fairly arbitrary, but intuitively appealing course of action is to make the smallest possible adjustment to \mathbf{a} such that it lies in the image of \mathbb{R}^m , and then find that $\mathbf{x} \in \mathbb{R}^m$ which is the inverse image under φ of that adjusted point. We will denote this point by $\hat{\mathbf{a}} \in \mathbb{R}^p$. We therefore wish to find $\hat{\mathbf{a}}$ so that $\|\hat{\mathbf{a}} - \mathbf{a}\|$ is a minimum, and hence find $\mathbf{x} = \varphi^{-1}(\hat{\mathbf{a}})$.

This is a nonlinear optimisation problem, so we cast it in the form of equation (A.5) to make it linear. We therefore define the p by m matrix \mathbf{C} , whose rows are the transposed centers \mathbf{c}_j and the p -vector \mathbf{h} , whose elements are given by equation (A.4), enabling us to recast (A.5) as

$$C_k \mathbf{x} = \mathbf{h}_k \quad (\text{A.6})$$

where $C_k = C - \mathbf{1}c_k^T$ and $\mathbf{h}_k = \mathbf{h} - \mathbf{1}h_k$, with $\mathbf{1}$ standing for a p -vector of ones. We now define the p -vector $\widehat{\mathbf{h}}$, whose elements \widehat{h}_j are given by

$$\widehat{h}_j \equiv \frac{1}{2} [\|\mathbf{c}_j\|^2 - \phi^{-2}(\widehat{a}_j)] \quad (\text{A.7})$$

as the LS approximation to \mathbf{h} , written

$$C_k \mathbf{x} = \widehat{\mathbf{h}}_k \quad (\text{A.8})$$

and found by minimising $\|\widehat{\mathbf{h}} - \mathbf{h}\|^2$. This equation holds for any value of k , so rather than choose a single value, we write

$$\begin{aligned} \sum_{k=1}^p C_k \mathbf{x} &= \sum_{k=1}^p \widehat{\mathbf{h}}_k \\ \Rightarrow (pC - \mathbf{1} \sum_{k=1}^p c_k^T) \mathbf{x} &= p\widehat{\mathbf{h}} - \mathbf{1} \sum_{k=1}^p \widehat{h}_k \\ \Rightarrow (C - \mathbf{1}c_0^T) \mathbf{x} &= (\widehat{\mathbf{h}} - \mathbf{1}\widehat{h}_0) \end{aligned} \quad (\text{A.9})$$

where c_0 and \widehat{h}_0 represent the means of the distributions of $\{c_j\}$ and $\{\widehat{h}_j\}$ respectively, and thus

$$C_0 \mathbf{x} = \widehat{\mathbf{h}}_0 \quad (\text{A.10})$$

The solution then follows as

$$\mathbf{x} = C_0^\dagger \widehat{\mathbf{h}}_0 \quad (\text{A.11})$$

where C_0^\dagger is the pseudo-inverse of C_0 , as defined in section 3.2.

Appendix B

Implementation

All of the experiments described in this thesis have been carried out using a piece of software designed specifically for that purpose. The software comprises four separate programs: a ‘parent’ program and three subprograms which are run under the parent’s control when required. A fair amount of time and effort has gone into the development of these programs, which are written in C and run on a network of multi-processor Unix workstations. Various subroutines have been adapted from Numerical Recipes in C [34]. In addition to this set of programs, a number of other programs have been written. These produce time series data by numerically iterating or integrating the maps described in section 2 and perform various transformations on those time series. In this appendix we will describe in some detail the structure and operation of the suite of programs which make up the main piece of software.

There are two available interfaces to the parent program: an X Window interface and a text-only batch interface, both of which allow the user to manipulate the same set of underlying data objects. The relationships between these objects will be discussed in the next section, and the two interfaces will be described in the following sections. To take advantage of this dual functionality, the parent may be compiled into two distinct forms: the first, called *ts*, is compiled with both interfaces; the second, called *tsb* for no satisfactorily explored reason, is compiled with only the batch interface, resulting in a much smaller binary.

B.1 Structure

The structure of the program takes the form of a hierarchy of objects, implemented as linked lists of C structures. This hierarchy descends from a single *base* object. The *base* serves primarily as an anchor for its descendants, an arbitrarily long list of *root* objects. Each *root* contains a distinct time series (which

may be multi-valued). It also stores variously calculated estimates of that time series, along with their associated per-point error magnitudes (see section 3.2.4) and power spectra (if required).

Descended from each *root* is an arbitrarily long list of *trajectory* objects. Each of these holds the trajectory resulting from the application of an n, τ -delay window to the (possibly FIR filtered) time series stored by its parent, with the added wrinkle that if the time series in question is multi-valued then the resulting delay vectors are formed by concatenating the vectors arising from each individual time series—this latter feature also enables us to work with general multi-valued data sets, such as the circle and torus data of chapter 4, merely by setting $n = 1$. A *trajectory* may also store a principal component basis, calculated from a segment of its own trajectory, as a *PCA* object. This trajectory may then be further modified by projection onto a *PCA* basis subset stored either locally or, given the appropriate dimensions, in any other *trajectory* object. (As an implementational nicety, to minimise storage requirements, if no projection has been made, the time series is single-valued and unfiltered, and $n = \tau = 1$, then the *trajectory* merely references, with the appropriate offsets, the time series stored in its parent *root*.) A *trajectory* may also store various trajectory estimates, along with their per-point error magnitudes.

Descended, in their turn, from each *trajectory*, are two lists of *phi* and *RBF* objects, respectively. Each *phi* implements the nonlinear part of an RBF map, storing the appropriate centers and basis function parameters (if any). Each *RBF* implements a general form of the RBF map, storing the linear part locally and referencing zero, one or two *phi* objects as required. That is, an *RBF* may be either a linear map or a composition of one or two nonlinear maps with a linear one, following chapter 3. If no *phi*'s are referenced then the *RBF* is linear. If a single *phi* is referenced, either from the *RBF*'s parent *trajectory* or from any other, then it is applied to the domain, or range, respectively, resulting in the 'classical' RBF of section 3.2. If a *phi* is referenced from both the parent and another (possibly the same) *trajectory* then we have the 'symmetrical' RBF of section 3.3. These *RBF*'s can be trained to approximate the relationship between trajectory segments in any two *trajectory* objects, between any two segments of the same trajectory, or between a *trajectory* and any *root* time series segment, with a variable prediction offset. Trajectory segments may then be mapped through a trained *RBF* to any other *trajectory* or *root* data set of the appropriate dimension, thus encompassing all of the fitting tasks required for this thesis.

These objects form the primary structure of the program. With the exception of the *base*, objects may be created or deleted at will during the program's execution: if an object is deleted then so are all of its descendants. Deletion of the *base* terminates the program. Certain classes of object—specifically, the *PCA*, *phi* and *RBF*—may also be saved and loaded between runs, to eliminate unnecessary repetition of time-consuming calculations.

B.2 Process control

Among the various functions built into the parent program, there are three which deserve a special mention: these are the three which are sufficiently complex and time consuming to be implemented as subprograms, rather than subroutines. All three are controlled from a *trajectory*. The first is called *pca* and returns the principal components, in the form of a *PCA*, of the distribution of delay vectors in a specified interval of the trajectory. The second is called *fit* and returns an *RBF*, fitted between the specified interval of the delay reconstruction and a (possibly different) interval in another *trajectory* or the time series in a *root*. The last is called *iterate* and, not surprisingly, iterates an *RBF* with a one-dimensional range as a time series predictor.

Each of these subprograms is run by ‘forking’ a copy of the parent program onto another (or the same) processor and then replacing that copy with the appropriate subprogram. The necessary data is then sent to the subprogram via a Unix ‘pipe’ and operation in the parent program is free to continue while the subprogram finishes its calculations. On return, the results are piped back to the parent, which interrupts whatever it’s currently doing to take the appropriate action and allow the subprogram to exit.

For simplicity, the *pca* and *iterate* processes are just forked onto the same processor as the parent, and only one is allowed to run at a time from a given *trajectory*. However, the *fit* process is implemented in a more flexible manner, so that several may be run in parallel by the same *trajectory*. To this end, a list of ‘virtual’ processors is specified by the user and maintained by the parent. Each element of this list can represent a physically distinct processor, or several elements can refer to the same physical processor. While a *fit* process is running on one of these virtual processors, that element of the list is flagged as busy. When the user requests a new *fit* process the list is checked: if all of the elements are busy then the parent goes into a ‘wait’ state until one becomes free; otherwise the process is forked onto the first free element.

B.3 The batch interface

In batch mode, the program is driven by a list of commands contained in an ASCII file. Every program function is accessible as a unique command, usually followed by a list of optional parameters and a matching *End* keyword. Default values for these parameters can be set from within the batch file and from the command line. The flow of execution within the batch file can be controlled by grouping adjacent commands together using the keywords *Loop*, *Seq* and *Par*, each of which also has a matching *End*.

The *Loop* . . . *End* construction takes parameters which define an integer loop variable with a specific range: the group of commands within this construction is repeated for each value of the loop variable. The variable itself may be substituted within the loop for any parameter, either directly or mapped onto a real-valued range specified by a further pair of parameters after the *Loop* command.

The *Seq . . . End* and *Par . . . End* constructions, conceptually borrowed from the OCCAM language, determine the mode of execution for certain commands within those groups: commands immediately enclosed by *Seq . . . End* are executed sequentially; those immediately enclosed by *Par . . . End* are executed in parallel. At present, this distinction is only valid for the forking of a *fit* process. Within a sequential group, once a *fit* process has been forked, execution continues from the next pending command outside that group, returning only after the relevant process has exited. Within a parallel group, multiple *fit* processes are forked in parallel as long as there are free virtual processors available for them to run on.

All three of these group types can be nested indefinitely, so that it is the innermost group which affects the mode of execution at a given level. However, a sequential group cannot be nested inside another sequential group unless there is an intervening parallel group, and vice versa. Inside a sequential group, a parallel group and its subgroups are just treated as a single sequential command; inside a parallel group, a sequential group is similarly treated as a parallel command. By default, the batch file itself is treated as an outer sequential group. Sensible use of these constructions allows efficient use to be made of a list of virtual processors in batch mode: commands which depend on the results of *fit* processes can be grouped with them inside sequential groups, while the use of outer loop and parallel groups enables multiple instances of these processes to run in parallel, making good use of the available resources.

An example batch file is illustrated in figure B.1. This particular file was used to generate one of the error curves in figure 5.1, used in the detection of periodic orbits in the Ikeda system. The first command, "Root", results in the creation of a new *root* object, containing the first 5000 elements of the one-dimensional time series specified by the "name" variable; the "Traj" command causes a single *trajectory* object to be derived from the *root* object. There follows a parallel group, enclosing a loop whose (scaled) variable is varied from -1 to 1 , in steps of 0.02 . A sequential group is defined inside this loop because each subsequent, enclosed command depends on the outcome of its predecessor. The first of this sequence is the "Window" command, which operates on the specified *trajectory* object to construct a trajectory from a 5-delay, unit-lag filtered reconstruction, through the FIR filter $(1, a_1, 1)$, where a_1 is substituted by the (scaled) loop variable. The "Phi" command then constructs a *phi* object, comprising 200 centers selected from the specified portion of trajectory and the "Fit" command forks a *fit* process, returning an *RBF* object trained to relate first 2000 elements of trajectory, transformed through the *phi* indicated by the (unscaled) loop variable, to the corresponding portion of time series in the *root* object; the errors and other indicators of the quality of the fit are labelled by the scaled loop variable and stored in a log file. Finally, the "Test" command applies that *RBF* to the specified test set, logging the resulting errors in a similar fashion. Although only a subset of the applicable parameters have actually been specified in this file, the remainder are taken from an optional "Default" command, otherwise defaulting to hard-wired values.

Once it has finished executing a batch file the *tsb* version of the parent program will exit automatically, but the *ts* version can be instructed to switch then into window mode, described in the next section, whilst

```

Root
  name "/masp/data/ikeda/20000/x"
  length 5000
  dimension 1
End
Traj 1
Par
  Loop 1 101 range -1 1
    Seq
      Window 1 1
        embed 5
        lag 1
        coeffs 3
        filter 1 #1 1
      End
      Phi 1 1
        number 200
        start 1
        stop 2000
      End
      Fit 1 1
        label #1
        offset 0
        start 1
        stop 2000
        target 1 0
        phi_X $1
        id 1
      End
      Test 1 1
        label #1
        start 2001
        stop 4000
        target 1 0
        rbf 1 1 1
      End
    End
  End
End

```

Figure B.1 Example command file for batch-mode time series analysis. Generates error curves for the detection of periodic orbits in the Ikeda system by executing 101 parallel copies of a sequence of commands comprising a filtered embedding followed by the construction and testing of an RBF map.

retaining its current state. In this way a batch file can be used to bring the program into a particular state automatically, at which time the user is able to continue to operate the program through the window interface.

B.4 The graphical user interface

The graphical user interface is for real-time operation of the program, with immediate visual feedback. In this mode, the *base* and each of the *root* and *trajectory* objects have a display window associated with them. In the case of the *base* this just contains a single button, which creates a new *root* containing the time series specified in its associated input field. However, the windows associated with *root* and *trajectory* objects also display the data associated with those objects in graphical form. Operations on this data are initiated with the appropriate menu buttons and other controls, and any necessary parameters can be specified on their associated pop-up property windows.

A *root* window displays the entire time series contained in that *root*. A specific portion of the time series can also be displayed in close-up, in a graphical subwindow of the *root* window. The contents of this window can be altered by selecting an interval of the time series in the *root* window, and defines the default portion of time series on which any subsequent operations, carried out by descendants of that *root*, are performed. In addition, a *root* window has a second graphical subwindow in which the FFT of that portion of the time series can be displayed.

A *trajectory* window displays a perspective view of the path described by the delay vectors corresponding to its particular delay reconstruction of the portion of time series currently displayed in close-up; the viewpoint can be varied at will. The view is represented either as a single two-dimensional projection or as a stereo pair. It also has graphical subwindows displaying the singular spectra obtained from SVD of the distribution of delay vectors before and after transformation through *phi* transformations, the latter being useful for examining the rank of the linear part of an *RBF*.