Nonlinearity 17 (2004) 1965-1983

PII: S0951-7715(04)77367-3

Robust normal forms for saddles of analytic vector fields

Warwick Tucker

Department of Mathematics, Uppsala University, Box 480 Uppsala, Sweden

E-mail: warwick@math.uu.se

Received 9 March 2004, in final form 21 June 2004 Published 27 July 2004 Online at stacks.iop.org/Non/17/1965 doi:10.1088/0951-7715/17/5/020

Recommended by M Tsujii

Abstract

The aim of this paper is to introduce a technique for describing trajectories of systems of ordinary differential equations (ODEs) passing near saddle-fixed points. In contrast to classical linearization techniques, the methods of this paper allow for perturbations of the underlying vector fields. This robustness is vital when modelling systems containing small uncertainties, and in the development of numerical ODE solvers producing rigorous error bounds.

Mathematics Subject Classification: 34C20, 37M99, 65G30

1. Introduction

Consider a system of ordinary differential equations (ODEs) $\dot{x} = f(x)$, where $x \in \mathbb{R}^n$ and each component f_i of the vector field is analytic in x_1, \ldots, x_n . Suppose that f(0) = 0 and that $Df(0) = \Lambda$ is a diagonal matrix with non-vanishing, real entries λ_i , $i = 1, \ldots, n$, not all of the same sign. We then call the origin a *saddle-fixed point*, or simply a *saddle*. To emphasize the behaviour near the saddle, the differential equations can be expressed as

$$\dot{x} = \Lambda x + F(x),\tag{1}$$

where F contains only quadratic or higher-order terms. Thus, for small x, it is natural to expect the solutions of (1) to behave roughly like those of the linear system $\dot{y} = \Lambda y$.

In the past, significant effort has been made to provide explicit interpretations of the somewhat vague terms '*for small x*' and '*roughly like*' in the sentence above. Apart from the Hartman–Grobman theorem (see [Ha60, Ha64]), which is purely topological, the most comprehensive result is due to Siegel (see [Si52]), which essentially states that, if the

eigenvalues $\lambda_1, \ldots, \lambda_n$ satisfy a Diophantine condition¹, then there is an analytic change of coordinates *h* that takes trajectories of $\dot{x} = \Lambda x + F(x)$ to those of the linearized system $\dot{y} = \Lambda y$, whilst preserving their orientation with respect to time.

Note that, although the set of eigenvalues satisfying a Diophantine condition with $\tau > 1$ has full Lebesgue measure in the set of eigenvalues corresponding to a saddle, the set of resonant² eigenvalues is everywhere dense. This means that not even the *existence* of a *formal* linearizing change of coordinates is guaranteed if we allow for small perturbations of the eigenvalues of Λ . The *C*^r-linearization theorem by Sternberg (see [St57, St58]), as well as its variants (see [Ne64, Be78, Se85]), also share this sensitivity. Of course, the Diophantine condition required by Siegel's theorem is even more fragile.

If the differential equations (1) are obtained from experimental data, it is necessary to allow for small uncertainties in the eigenvalues of Λ , as well as in the coefficients of the Taylor series for *F*. From the discussion above, it is clear that we cannot hope to find a linearizing change of coordinates valid for such an open set of differential equations. Instead, in this situation, we must develop robust methods that allow for small perturbations of the underlying vector field.

For numerical applications, the problem at hand is to produce a *transfer map* $\Gamma: \Sigma_i \to \Sigma_j$, where Σ_i and Σ_j are different faces of the cube \mathfrak{B}_r with radius r, centred at the saddle point. Once r is fixed, we face the problem of describing the behaviour of the trajectories inside the cube \mathfrak{B}_r . This is not trivial for a nonlinear system such as (1), since on regions of Σ_i where some components of x are small (compared to r), the nonlinear part F(x) may very well dominate the linear part Λx . This problem remains even if we, by a polynomial change of variables, remove all nonlinear terms of F up to a high (but finite) degree. The choice of r is delicate: an upper bound is determined by the vector field itself, whereas a lower bound is governed by the integration method as well as the precision of the underlying floating point system. It is desirable to take r as large as possible, since any numerical solver breaks down in the vicinity of a fixed point due to the unbounded flow-times.

We close this introductory section by remarking that, in the past few years, several important results in dynamical systems have been proved using the so-called *validated numerics* for ordinary differential equations (see, e.g. [Be01,KZ03,ZM01]). The core technique utilized in these proofs is a means of numerically producing mathematically rigorous enclosures of solutions to systems of differential equations. In order to extend such methods to include the case of enclosing solutions passing near saddle points, estimates on the radius *r* and the transfer map Γ mentioned above are essential. In [Tu02], less general versions of the main results (see section 3) were successfully employed to prove the existence of a strange attractor for the Lorenz equations. For a different approach (due to Sil'nikov) to finding a description of Γ , see [De89] and references therein.

2. Normal forms

In what follows, we propose to locally find a close to identity change of coordinates $x = h(y) = y + \phi(y)$ which does *not* bring (1) into a completely linear system, but rather into

¹ We say that the eigenvalues $\lambda = (\lambda_1, ..., \lambda_n)$ satisfy a *Diophantine* condition of type (κ, τ) if there exists positive κ and τ such that for i = 1, ..., n we have $|m\lambda - \lambda_i| = \left|\sum_{k=1}^n m_k \lambda_k - \lambda_i\right| \ge \kappa |m|^{-\tau}$ for all natural numbers $m_1, ..., m_n$ with $|m| = \sum m_i \ge 2$.

² We say that the eigenvalues $\lambda = (\lambda_1, ..., \lambda_n)$ are *resonant* if there exist natural numbers $m_1, ..., m_n$ with $|m| \ge 2$ such that $m\lambda - \lambda_i = 0$ for some i = 1, ..., n. The number |m| is called the *order* of the resonance.

a system that, in some sense, is close to being linear:

$$\dot{x} = \Lambda x + F(x) \xrightarrow{x = y + \phi(y)} \dot{y} = \Lambda y + G(y).$$
 (2)

We call the resulting system $\dot{y} = \Lambda y + G(y)$ a *normal form*. There are of course many choices regarding the structure of *G*, and we will make a very careful selection. The first property we require from the particular normal form we have in mind is that its unstable and stable manifolds of the origin coincide with the appropriate coordinate axes. We then say the normal form is *rectified*.

In what follows, we will label the eigenvalues of Λ according to

$$\lambda_{s_q} < \cdots < \lambda_{s_1} < 0 < \lambda_{u_1} < \cdots < \lambda_{u_p}.$$

In order for the invariant manifolds to coincide with the coordinate axes, it is necessary that the axes are invariant under the flow. To ensure this, we need a change of variables which, in a fixed neighbourhood of the origin, transforms the original equations $\dot{x} = \Lambda x + F(x)$ into $\dot{y} = \Lambda y + G(y)$, where G satisfies the following conditions:

$$y_{u_1} = \dots = y_{u_p} = 0 \implies G_{u_i}(y) = 0 \qquad (i = 1, \dots, p)$$
 (3)

and

$$y_{s_1} = \dots = y_{s_q} = 0 \implies G_{s_i}(y) = 0 \qquad (i = 1, \dots, q).$$
 (4)

In these new coordinates, the unstable manifold coincides with the $(y_{u_1} \cdots y_{u_p})$ -plane, and the stable manifold coincides with the $(y_{s_1} \cdots y_{s_q})$ -plane, as desired. This will, however, *not* linearize the flow on the invariant manifolds. As an example, at a point y on the unstable manifold we have $y_{s_1} = \cdots = y_{s_q} = 0$, which brings the normal form into

$$\dot{y}_{u_i} = \lambda_{u_i} y_{u_i} + G_{u_i}(y)$$
 $(i = 1, ..., p)$
 $\dot{y}_{s_i} = 0$ $(i = 1, ..., q)$

which generally is nonlinear in the y_{u_i} -coordinates. An analogous statement can be made for points on the stable manifold. In order to guarantee linear behaviour on the invariant manifolds, we need to impose the additional condition that if a point y is close to the $(y_{u_1} \cdots y_{u_p})$ -plane (the unstable manifold) or the $(y_{s_1} \cdots y_{s_q})$ -plane (the stable manifold), then the perturbation G(y) is linearly small, i.e.

$$d(y) \stackrel{\text{def}}{=} \min\left\{\max_{i}\{|y_{u_i}|\}, \max_{i}\{|y_{s_i}|\}\right\} \Rightarrow |G_i(y)| = \mathcal{O}(d(y)) \qquad (i = 1, \dots, n).$$

Depending on the situation at hand, we may want to flatten the normal form even more. Flatness of order ℓ is given by requiring that $|G_i(y)| = O(d(y)^{\ell})$, for i = 1, ..., n.

In this case, it follows that the components of G can only contain terms of the form $y^m = y_1^{m_1} \cdots y_n^{m_n}$, where the multi-exponent $m \in \mathbb{N}^n$ satisfies both $\sum_{i=1}^p m_{u_i} \ge \ell$ and $\sum_{i=1}^q m_{s_i} \ge \ell$. For future reference, we define the sets

$$\mathbb{V}_{\ell}^{n} = \left\{ m \in \mathbb{N}^{n} \colon \sum_{i=1}^{p} m_{u_{i}} < \ell \quad \lor \quad \sum_{i=1}^{q} m_{s_{i}} < \ell \right\},$$
$$\mathbb{U}_{\ell}^{n} = \left\{ m \in \mathbb{N}^{n} \colon \sum_{i=1}^{p} m_{u_{i}} \ge \ell \quad \land \quad \sum_{i=1}^{q} m_{s_{i}} \ge \ell \right\};$$

or, equivalently

$$\mathbb{V}_{\ell}^{n} = \left\{ m \in \mathbb{N}^{n} : \min\left\{ \sum_{i=1}^{p} m_{u_{i}}, \sum_{i=1}^{q} m_{s_{i}} \right\} < \ell \right\},$$
$$\mathbb{U}_{\ell}^{n} = \left\{ m \in \mathbb{N}^{n} : \min\left\{ \sum_{i=1}^{p} m_{u_{i}}, \sum_{i=1}^{q} m_{s_{i}} \right\} \ge \ell \right\}.$$

In other words, writing G as a formal power series $G(y) = \sum g_m y^m$ (using multi-notation combined with vector notation), we require that

$$g_m \neq 0 \Rightarrow m \in \mathbb{U}_\ell^n$$

In what follows, we will sometimes omit the superscript *n* in \mathbb{U}_{ℓ}^{n} and \mathbb{V}_{ℓ}^{n} .

It is convenient to introduce the notion of *filters* for formal vector-valued power series: given any $f(y) = \sum_{|m| \ge 2} \alpha_m y^m$, we define

$$\langle f(y) \rangle_{\mathbb{U}_{\ell}} = \sum_{m \in \mathbb{U}_{\ell}} \alpha_m y^m, \qquad \langle f(y) \rangle_{\mathbb{V}_{\ell}} = \sum_{m \in \mathbb{V}_{\ell}} \alpha_m y^m.$$
 (5)

Note that we always have the decomposition $f(y) = \langle f(y) \rangle_{\mathbb{U}_{\ell}} + \langle f(y) \rangle_{\mathbb{V}_{\ell}}$, which splits f into its flat part and its non-flat part. It follows that the nonlinear part G of our normal form has flatness of order ℓ if $\langle G(y) \rangle_{\mathbb{U}_{\ell}} = G(y)$, or equivalently, $\langle G(y) \rangle_{\mathbb{V}_{\ell}} = 0$.

We stress the fact that flattening a function to order ℓ requires much more effort than simply linearizing it to the same order, i.e. removing all terms $\alpha_m y^m$ with $|m| < \ell$. As an example, in the three-dimensional case with e.g. $\lambda_3 < \lambda_2 < 0 < \lambda_1$, the term $y_1^3 y_2^{992} y_3^5$ is linear to order 1000, but only flat to order 3. In general, flattening a function to order ℓ requires the removal of infinitely many terms, as compared to a finite amount when linearizing to the same order.

3. Main results

Let S^n denote the space of all real-valued, diagonal $n \times n$ -matrices corresponding to the linearization at a saddle (i.e. strictly indefinite matrices), and let \mathcal{F}_{ℓ}^n denote the space of all such matrices whose diagonal elements $\lambda_1, \ldots, \lambda_n$ have no resonances for $m \in \mathbb{V}_{\ell}^n$:

$$\mathcal{F}_{\ell}^{n} = \{ \Lambda \in \mathcal{S}^{n} : m \in \mathbb{V}_{\ell}^{n} \Rightarrow m\lambda - \lambda_{i} \neq 0 \ (i = 1, \dots, n) \}.$$

We will use the following max norms:

$$|y| = \max\{|y_i|: i = 1, ..., n\}$$
 $||f||_r = \max\{|f(y)|: |y| \le r\}.$

Theorem 3.1. Given an integer $\ell \ge 2$ and a system $\dot{x} = \Lambda x + F(x)$, with $\Lambda \in \mathcal{F}_{\ell}^{n}$, and where $F(x) = \sum_{|m| \ge 2} a_m x^m$ is analytic, there exist positive constants r_0, r_1, K_0, K_1 and an analytic, close to identity change of variables $x = y + \phi(y)$ with

$$\|\phi\|_r \leqslant K_0 r^2 \qquad (r < r_0),$$

such that $\dot{x} = \Lambda x + F(x)$ is transformed into the normal form $\dot{y} = \Lambda y + G(y)$ satisfying $\langle G(y) \rangle_{\mathbb{U}_{\ell}} = G(y)$ and

$$\|G\|_r \leqslant K_1 r^{2\ell} \qquad (r < r_1).$$

This theorem tells us that the change of coordinates *and* the resulting normal form exist (as analytic functions) in a fixed neighbourhood of the origin.

Having established the change of coordinates, what can be said about the flow of the resulting normal form? In what follows, we will let \mathfrak{B}_r denote the closed ball (which in our norm looks like a box) centred at the origin, and having radius *r*. We will refer to the face $\{y \in \mathfrak{B}_r : y_{s_1} = r\}$ as the *lid* of the box \mathfrak{B}_r (recall that λ_{s_1} is the weakest contracting direction of the stable manifold). Within \mathfrak{B}_r , we let $\psi(y, t)$ denote the solution to the normal form $\dot{y} = \Lambda y + G(y)$.

We begin with the special case where Λ has only one positive eigenvalue λ_u . In this case, the saddle point has a unique unstable direction, and thus any trajectory starting from the lid of \mathfrak{B}_r (except points on the stable manifold of the origin) will exit through an unstable face $\{y \in \mathfrak{B}_r : |y_u| = r\}$. We would like to know how long a trajectory spends inside the box, and where it exits.

Theorem 3.2. If Λ has only one positive eigenvalue λ_u , then under the same conditions as in theorem 3.1, and given any $\kappa > 0$ sufficiently small, there exists r > 0 such that for any trajectory starting from the lid of \mathfrak{B}_r , we have the following enclosure of its point of exit:

$$\begin{split} \psi_u(y,\tau_e(y)) &= \operatorname{sign}(y_u)r;\\ r\left(\frac{|y_u|}{r}\right)^{(|\lambda_{s_1}|+\kappa)/(\lambda_u-\kappa)} \leqslant \psi_{s_1}(y,\tau_e(y)) \leqslant r\left(\frac{|y_u|}{r}\right)^{(|\lambda_{s_1}|-\kappa)/(\lambda_u+\kappa)} \end{split}$$

where $\tau_e(y)$ (the exit time) denotes the time spent inside \mathfrak{B}_r :

$$\frac{1}{\lambda_u + \kappa} \log \frac{r}{|y_u|} \leqslant \tau_e(y) \leqslant \frac{1}{\lambda_u - \kappa} \log \frac{r}{|y_u|}$$

If Λ has several negative eigenvalues $\lambda_{s_q} < \cdots < \lambda_{s_1} < 0$, and if we take $\ell > (|\lambda_{s_q}| + 1)/(|\lambda_{s_1}| - \kappa)$, then we also have the following enclosures:

$$(y_{s_i} - \kappa r) \left(\frac{|y_u|}{r}\right)^{|\lambda_{s_i}|/(\lambda_u - \sigma_1 \kappa)} \leq \psi_{s_i}(y, \tau_e(y)) \leq (y_{s_i} + \kappa r) \left(\frac{|y_u|}{r}\right)^{|\lambda_{s_i}|/(\lambda_u + \sigma_2 \kappa)}$$

where $\sigma_1 = \operatorname{sign}(y_{s_i} - \kappa r)$ and $\sigma_2 = \operatorname{sign}(y_{s_i} + \kappa r)$.

Remark 1. These additional enclosures can be made somewhat sharper, see lemma 8.5.

In the most general setting of this paper, we allow for Λ having several positive eigenvalues $0 < \lambda_{u_1} < \cdots < \lambda_{u_p}$. This situation adds the complication of determining through which unstable face of \mathfrak{B}_r a trajectory will exit. It is therefore more appropriate to provide enclosures of the trajectories within the box, and an enclosure of the required exit-time $\tau_e(y)$.

Theorem 3.3. Under the same conditions as in theorem 3.1, and given any $\kappa > 0$ sufficiently small, there exists r > 0 such that, for any trajectory starting from the lid of \mathfrak{B}_r , we have the following enclosures of the unstable components of its path throughout the box:

$$|\psi_{u_i}(y,t)-y_{u_i}e^{\lambda_{u_i}t}| \leqslant \frac{\kappa r}{\alpha_i}(1-e^{-\alpha_i t})e^{\lambda_{u_i}t} \qquad (i=1,\ldots,p),$$

for any α_i satisfying $0 < \lambda_{u_i} < \alpha_i \leq \lambda_{u_i} - \ell(\lambda_{s_1} + \kappa)$.

If we take $\ell > |\lambda_{s_q}|/(|\lambda_{s_1}| - \kappa)$, then for any α_i satisfying $0 < \alpha_i \leq \lambda_{s_i} - \ell(\lambda_{s_1} + \kappa)$, we also have similar enclosures of the stable components:

$$|\psi_{s_i}(y,t)-y_{s_i}e^{\lambda_{s_i}t}|\leqslant \frac{\kappa r}{\alpha_i}(1-e^{-\alpha_i t})e^{\lambda_{s_i}t} \qquad (i=1,\ldots,q).$$

As in theorem 3.2, there exist explicit bounds on the time spent inside \mathfrak{B}_r :

$$\tau_e^-(y) \leqslant \tau_e(y) \leqslant \tau_e^+(y),$$

where $\tau_e^{\pm}(y) \nearrow +\infty$ as $\max_{i=1,\dots,p}\{|y_{u_i}|\} \searrow 0$.

Remark 2. See corollary 8.11 for the explicit flow-time bounds $\tau_e^-(y)$ and $\tau_e^+(y)$.

These theorems have several strengths. First, the constants $r_0, r_1, K_0, K_1, \alpha, \kappa$ can be explicitly found, and are easy to obtain in terms of ℓ , Λ and F (naturally κ also depends on r). Second, the change of variables $x = y + \phi(y)$ is analytic for $|y| < r_0$, which means that explicit bounds on its inverse and derivatives can be obtained by Cauchy estimates. The same holds for Gwhen $|y| < r_1$. Furthermore, theorems 3.2 and 3.3 tell us that, inside \mathfrak{B}_r , solutions to the normal form act very much like those of the completely linearized system. This is *not* true for a system linearized up a certain high, but finite, order. Finally, the set \mathcal{F}_{ℓ}^{n} , viewed as a subset of \mathcal{S}^{n} , is open and has full Lebesgue measure. We call such a set *robust*: almost all members of \mathcal{S}^{n} belong to \mathcal{F}_{ℓ}^{n} , and any sufficiently small perturbation of an element in \mathcal{F}_{ℓ}^{n} remains in \mathcal{F}_{ℓ}^{n} . This allows us to perform the change of coordinates even when we only know the eigenvalues up to some finite degree of accuracy (see e.g. [Tu02]). In contrast to this, we point out that the theorems by Sternberg and Siegel fail on an everywhere dense subset of \mathcal{S}^{n} , and thus cannot be used in the situation at hand.

4. The change of variables

Returning to the normal form, we need to know how the vector field $\dot{x} = \Lambda x + F(x)$ is affected by the close to identity change of variables $x = y + \phi(y)$. We have the following identity:

$$\dot{x} = \Lambda(y + \phi(y)) + F(y + \phi(y)) = \Lambda y + \Lambda \phi(y) + F(y + \phi(y)).$$
(6)

On the other hand, we also have

$$\dot{x} = \frac{\mathrm{d}}{\mathrm{d}t}(y + \phi(y)) = (I + D\phi(y))\dot{y} = (I + D\phi(y))(\Lambda y + G(y))$$
$$= \Lambda y + D\phi(y)\Lambda y + G(y) + D\phi(y)G(y). \tag{7}$$

Comparing the two right-hand sides of (6) and (7) gives

$$D\phi(y)\Lambda y - \Lambda\phi(y) = F(y + \phi(y)) - D\phi(y)G(y) - G(y).$$
(8)
For shorthand, we will use the following notation

 $L_{\Lambda}\phi(y) = D\phi(y)\Lambda y - \Lambda\phi(y).$

The operator L_{Λ} is linear, and it acts on the space of formal vector fields. It leaves the spaces of homogeneous vector-valued polynomials of any degree invariant. Looking at (8) on the component level, we have

$$L_{\Lambda,i}\phi_{i}(y) = F_{i}(y + \phi(y)) - \sum_{j=1}^{n} \frac{\partial \phi_{i}}{\partial y_{j}}(y)G_{j}(y) - G_{i}(y) \qquad (i = 1, ..., n),$$
(9)

where

$$L_{\Lambda,i}\phi_i(y) = \sum_{j=1}^n \frac{\partial \phi_i}{\partial y_j}(y)\lambda_j y_j - \lambda_i \phi_i(y) \qquad (i = 1, \dots, n)$$

Note that

$$L_{\Lambda,i}(a_{i,m}y^m) = (m_1\lambda_1 + \dots + m_n\lambda_n - \lambda_i)a_{i,m_1,\dots,m_n}y_1^{m_1} \cdots y_n^{m_n} = (m\lambda - \lambda_i)a_{i,m}y^m$$

The crux is now to choose ϕ so that we produce only flat component functions in the normal form: $\langle G_i(y) \rangle_{U_\ell} = G_i(y)$. This means that $G_i(y)$ must *not* contain elements on the form $y^m = y_1^{m_1} \cdots y_n^{m_n}$ where the exponent *m* belongs to \mathbb{V}_ℓ . By (9), non-flat elements can only come from $F_i(y + \phi(y))$, and any such term can be absorbed by an appropriate choice of ϕ_i provided that the corresponding *divisor* $m\lambda - \lambda_i$ does not vanish. Thus the component functions ϕ_i need only consist of the non-flat terms appearing in the right-hand side of (9), which implies that we should choose ϕ_i such that $\langle \phi_i \rangle_{\mathbb{V}_\ell} = \phi_i$.

By filtering (9), we thus get

$$L_{\Lambda,i}\phi_i(y) = \langle F_i(y + \phi(y)) \rangle_{\mathbb{V}_\ell} \qquad (i = 1, \dots, n)$$
(10)

and

$$G_i(y) = \langle F_i(y + \phi(y)) \rangle_{\mathbb{U}_\ell} - \sum_{j=1}^n \frac{\partial \phi_i}{\partial y_j}(y) G_j(y) \qquad (i = 1, \dots, n).$$
(11)



Figure 1. The sets (*a*) $\{m \in \mathbb{N}^2 : |m| \ge 2\}$; (*b*) $\{m \in \mathbb{V}^2_4 : |m| \ge 2\}$.

We will begin by considering the existence and convergence of ϕ . The recursive scheme (10) can be formally solved by a power series

$$\phi_i(y) = \sum_{|m|=2}^{\infty} a_{i,m} y^m$$
 $(i = 1, ..., n)$

where the coefficients are determined by inserting this expression into (10). The existence of a solution ϕ is given by comparing both sides of (10): if $a_{i,m}y^m$ is a term of $\phi_i(y)$ with $|m| = m_1 + \cdots + m_n$, the comparison gives

$$(m\lambda - \lambda_i)a_{i,m} = \gamma,$$

where γ is a polynomial in the coefficients of the terms in ϕ_i (i = 1, ..., n) of degree less than |m|. Thus, the *existence* of ϕ is proved if we show that the divisors $m\lambda - \lambda_i$ do not vanish. As ϕ does not contain constant or linear terms, and since $\langle \phi \rangle_{\mathbb{V}_\ell} = \phi$, the only divisors we need to consider are of the form $m\lambda - \lambda_i$, where $m \in \mathbb{V}_\ell$ and $|m| \ge 2$ (see figure 1(*b*) for a two-dimensional example). In fact, the situation is generally more favourable than this: given an explicit system $\dot{x} = \Lambda x + F(x)$, we only have to consider elements of \mathbb{V}_ℓ that actually occur in the absorption process of the change of variables. These depend on the exact form of *F*, and may be very few compared to the total number of elements of \mathbb{V}_ℓ .

5. Small divisors and existence

In what follows, we let $\lceil x \rceil$ denote the *ceiling* of a real number x, i.e. $\lceil x \rceil = \min\{k \in \mathbb{Z} : x \le k\}$. We also introduce the numbers λ , λ and λ which denote the eigenvalue having the smallest modulus, the eigenvalue having the opposite sign of λ with largest modulus, and the eigenvalue of the same sign as λ with largest modulus, respectively:

$$\check{\lambda} = \begin{cases} \lambda_{s_1} : |\lambda_{s_1}| < |\lambda_{u_1}| \\ \lambda_{u_1} : \text{o.w.} \end{cases} \quad \hat{\lambda} = \begin{cases} \lambda_{u_p} : \check{\lambda} < 0 \\ \lambda_{s_q} : \text{o.w.} \end{cases} \quad \tilde{\lambda} = \begin{cases} \lambda_{s_q} : \check{\lambda} < 0 \\ \lambda_{u_p} : \text{o.w.} \end{cases}$$

Let us begin by stating a lemma that, together with its corollary, proves the existence of a formal series for ϕ for virtually every saddle fixed point.

Lemma 5.1. If the eigenvalues $\lambda = (\lambda_1, ..., \lambda_n)$ are non-resonant for $m \in \mathbb{V}_{\ell}$, then the divisors $m\lambda - \lambda_i$ are bounded away from zero. Furthermore, for all orders $|m| \ge \Gamma_{\Lambda,\ell} \equiv \ell - 1 + \lceil (\ell - 1) | \hat{\lambda} / \check{\lambda} \rceil + \check{\lambda} / \check{\lambda} \rceil$, we have the following sharp lower bound:

$$|m\lambda - \lambda_i| \ge |(|m| - (\ell - 1))\dot{\lambda} + (\ell - 1)\dot{\lambda} - \ddot{\lambda}| \qquad (i = 1, \dots, n).$$

Proof. Take |m| large. Since we are only considering $m \in \mathbb{V}_{\ell}$, this means that either $\sum_{i=1}^{q} m_{s_i}$ or $\sum_{i=1}^{p} m_{u_i}$ is large, but not both (since one of them must be less than ℓ). Although the corresponding eigenvalues have opposite signs, the modulus of the divisor $|m\lambda - \lambda_i|$ must then also be large. There are two cases to consider:

Case 1: $\sum_{i=1}^{p} m_{u_i} < \ell$. This means that $\sum_{i=1}^{q} m_{s_i}$ is large, i.e. the divisor $m\lambda - \lambda_i$ is large and negative. We clearly minimize the modulus of the divisor when $m_{u_p} = \ell - 1$ and $m_{s_1} = |m| - (\ell - 1)$ are the only non-zero components of m, and $\lambda_i = \lambda_{s_a}$, which gives

$$|m\lambda - \lambda_i| \ge |(\ell - 1)\lambda_{u_p} + (|m| - (\ell - 1))\lambda_{s_1} - \lambda_{s_q}| \qquad (i = 1, \dots, n).$$

Case 2: $\sum_{i=1}^{q} m_{s_i} < \ell$. This means that $\sum_{i=1}^{p} m_{u_i}$ is large, i.e. the divisor $m\lambda - \lambda_i$ is large and positive. We clearly minimize the modulus of the divisor when $m_{s_q} = \ell - 1$ and $m_{u_1} = |m| - (\ell - 1)$ are the only non-zero components of m, and $\lambda_i = \lambda_{u_p}$, which gives

$$|m\lambda - \lambda_i| \ge |(\ell - 1)\lambda_{s_q} + (|m| - (\ell - 1))\lambda_{u_1} - \lambda_{u_p}| \qquad (i = 1, \dots, n)$$

Combining both cases, we see that the lowest bound is given by

$$|m\lambda - \lambda_i| \ge |(|m| - (\ell - 1))\hat{\lambda} + (\ell - 1)\hat{\lambda} - \tilde{\lambda}| \qquad (i = 1, \dots, n),$$

which provides the sharp lower bound.

How large must |m| be for this bound to hold? Clearly, the bound is valid starting from the last sign change of $(|m| - (\ell - 1))\dot{\lambda} + (\ell - 1)\hat{\lambda} - \tilde{\lambda}$, which happens near the largest (in |m|) approximate zero:

$$(|m| - (\ell - 1))\dot{\lambda} + (\ell - 1)\hat{\lambda} - \tilde{\lambda} \approx 0.$$

Solving for |m| gives

$$|m| \approx \frac{1}{\lambda} ((\ell-1)\check{\lambda} - (\ell-1)\hat{\lambda} + \check{\lambda}) = \ell - 1 + (\ell-1) \left| \frac{\hat{\lambda}}{\check{\lambda}} \right| + \frac{\tilde{\lambda}}{\check{\lambda}}.$$

Rounding up to the nearest integer produces the desired bound:

$$|m| = \ell - 1 + \left[(\ell - 1) \left| \frac{\hat{\lambda}}{\check{\lambda}} \right| + \frac{\tilde{\lambda}}{\check{\lambda}} \right].$$

Beyond this order, the divisors will increase in modulus with |m|, and have the same sign as the eigenvalue of smallest modulus λ .

Remark 3. In the planar case (n = 2), we always have $\check{\lambda} = \tilde{\lambda}$, which gives the bound

$$|m\lambda - \lambda_i| \ge |(|m| - \ell)\check{\lambda} + (\ell - 1)\hat{\lambda}| \qquad (i = 1, 2),$$

which is valid for all $|m| \ge \Gamma_{\Lambda,\ell} \equiv \ell + \lceil (\ell-1)|\lambda/\lambda|\rceil$.

Remark 4. Note that the asymptotic growth of the divisors is given by

$$|m\lambda - \lambda_i| \sim |m||\dot{\lambda}|.$$



Figure 2. The resonant set in the planar case (n = 2) with $\lambda_2 < 0 < \lambda_1$ for $(a) \ell = 5$; $(b) \ell = 6$. (This figure is in colour only in the electronic version)

It might appear that requiring the eigenvalues to be non-resonant in \mathbb{V}_{ℓ} is a serious restriction. The following corollary, however, shows that this is in fact almost a completely void demand.

Corollary 5.2. For any integer $\ell \ge 2$, the set of eigenvalues

 $(\lambda_{s_1},\ldots,\lambda_{s_q},\lambda_{u_1},\ldots,\lambda_{u_p})\in\mathbb{R}^q_-\times\mathbb{R}^p_+$

that are resonant for $m \in \mathbb{V}_{\ell}^n$ form a closed set of n-dimensional Lebesgue measure zero.

The key word here is *closed*. This means that the non-resonant eigenvalues form an open set. Furthermore, this set has full measure. Recalling our wish to be allowed some uncertainty in the eigenvalues, this situation is ideal for our needs. The special ordering of the eigenvalues in the statement can be achieved by a simple permutation of the coordinates, and thus causes no loss of generality.

Proof. By lemma 5.1, there are only a finite number of orders |m| we need to consider. Since each order can give rise to at most a finite number of different resonances, it clearly suffices to show that each such resonance forms a closed set of measure zero in $\mathbb{R}^{q}_{-} \times \mathbb{R}^{p}_{+}$. But this is obvious: any resonance $m\lambda - \lambda_{i} = 0$ corresponds to a codimension-1 plane in \mathbb{R}^{p+q} passing through the origin (see figure 2 for the special case n = 2). A finite union of (n-1)-dimensional planes certainly forms a closed set of *n*-dimensional measure zero, as claimed.

Remark 5. As the order of flatness ℓ increases, so does the number of resonant planes. In the limit $\ell \to \infty$, the resonant set becomes everywhere dense in $\mathbb{R}^q_- \times \mathbb{R}^p_+$. This limiting case corresponds to completely linearizing the system, i.e. choosing $G \equiv 0$.

6. Majorants and convergence

Assuming, in what follows, that the formal power series for ϕ defined by (10) exists, we want to show that it also actually converges. To be able to talk about convergence, we need to specify

a norm. It is convenient to work in a complex neighbourhood of the origin, and we will use the appropriate max norms:

$$|y| = \max\{|y_i|: i = 1, ..., n\}$$
 $||f||_r = \max\{|f(y)|: |y| \le r\}$

In order to prove convergence, we follow [SM71, Hi76], and use the methods of majorants. If

$$f(\zeta) = \sum_{m} \alpha_{m_1,\dots,m_n} \zeta_1^{m_1} \cdots \zeta_n^{m_n}; \qquad g(\zeta) = \sum_{m} \beta_{m_1,\dots,m_n} \zeta_1^{m_1} \cdots \zeta_n^{m_n}$$

are two formal power series, g is said to be a *majorant* of f, which we denote $f \prec g$, if

 $|\alpha_{m_1,\ldots,m_n}| \leqslant \beta_{m_1,\ldots,m_n}$

holds for all the coefficients. Note that the coefficients of g must be real and non-negative, which implies that f must have *at least* as large a radius of convergence as g.

Suppose that we can find a function $\tilde{F}:\mathbb{C}^n \to \mathbb{C}$ such that $F_i \prec \tilde{F}$ (i = 1, ..., n) and, together with (10), consider the majorant system

$$\tilde{L}_{\Lambda}\tilde{\phi}_{i}(\zeta) = \langle \tilde{F}(\zeta + \tilde{\phi}(\zeta)) \rangle_{\mathbb{V}_{\ell}} \qquad (i = 1, \dots, n),$$
(12)

where $\tilde{L}_{\Lambda}(\zeta^m) = \tilde{\Omega}(m)\zeta^m$ and $\tilde{\Omega}: \mathbb{N}^n \to \mathbb{R}$ is defined by

$$\hat{\Omega}(m) = \min\{|m\lambda - \lambda_i|: i = 1, \dots, n\}.$$

This can be solved formally by a power series

$$\tilde{\phi}_i(\zeta) = \sum_{|m|=2}^{\infty} \tilde{a}_{i,m} \zeta^m \qquad (i = 1, \dots, n)$$
(13)

and it follows that $\tilde{\phi}_i$ is a majorant of ϕ_i . To see this, compare the two functional equations (10) and (12). In the latter, the divisors appearing on the left-hand side are positive and smaller than or equal to the modulus of those present in (10). Also, the coefficients of \tilde{F} , appearing on the right-hand side of (12), are positive and larger than or equal to the modulus of those of F. This implies that the coefficients satisfy $|a_{m_1,\ldots,m_n}| \leq \tilde{a}_{m_1,\ldots,m_n}$ for all m, as claimed.

Furthermore, since both \tilde{L}_{Λ} and the right-hand side of (12) are independent of *i*, we have $\tilde{\phi}_1 = \cdots = \tilde{\phi}_n$. If we set $\zeta_1 = \cdots = \zeta_n = z$, and find a new function $\check{F}: \mathbb{C} \to \mathbb{C}$ such that $\tilde{F}(z, \ldots, z) \prec \check{F}(z)$, we may, together with (12), consider the majorant system

$$\check{L}_{\Lambda}\check{\phi}(z) = \check{F}(z + \check{\phi}(z)), \tag{14}$$

where $\check{L}_{\Lambda}(z^k) = \check{\Omega}(k) z^k$ and $\check{\Omega}: \mathbb{N} \to \mathbb{R}$ is defined by

$$\hat{\Omega}(k) = \min \left\{ \hat{\Omega}(m) \colon |m| = k \land m \in \mathbb{V}_{\ell} \right\}.$$

Again, this can be solved formally by a power series

$$\check{\phi}(z) = \sum_{k=2}^{\infty} \check{a}_k z^k \tag{15}$$

and, from the same reasoning as above, it is clear that $\tilde{\varphi}(z, \ldots, z) \prec \check{\phi}(z)$. Note that this implies that $\|\phi\|_r \leq \check{\phi}(r)$ in the region of convergence. Thus it suffices to prove the convergence of $\check{\phi}$. We will now present explicit candidates for the above-mentioned majorants \tilde{F} and \check{F} .

Since we are assuming that F is analytic in a neighbourhood of the origin, we can therefore identify it with its power series

$$F(\zeta) = \begin{pmatrix} \sum_{|m|=2}^{\infty} c_{1,m} \zeta^m \\ \vdots \\ \sum_{|m|=2}^{\infty} c_{n,m} \zeta^m \end{pmatrix}$$

Thus, if we set

$$\tilde{F}(\zeta) = \sum_{|m|=2}^{\infty} \tilde{c}_m \zeta^m,$$

where $\tilde{c}_m = \max\{|c_{i,m}|: i = 1, ..., n\}$, we clearly have $F_i \prec \tilde{F}$ (i = 1, ..., n), and \tilde{F} has the same radius of convergence as F. Summing the coefficients of all terms having the same degree produces $\check{c}_k = \sum_{|m|=k} \tilde{c}_m$, and if we define

$$\check{F}(z) = \sum_{k=2}^{\infty} \check{c}_k z^k.$$

it follows that $\tilde{F}(z, ..., z) \prec \check{F}(z)$. Once again, \check{F} is analytic and has the same radius of convergence as F. Hence, if the solution to (14) converges, then we have $\|\phi\|_r \leq \check{\phi}(r)$ in the region of convergence.

Now, by lemma 5.1, there exists a positive constant A (depending only on Λ and ℓ) that satisfies $Ak \leq \check{\Omega}(k)$ for all k = 2, 3, ... Thus, we can replace the operator $\check{L}_{\Lambda}(z^k) = \check{\Omega}(k)z^k$, by the smaller operator $\hat{L}_{\Lambda}(z^k) = Akz^k$, which transforms (14) into the following functional equation:

$$L_{\Lambda}\hat{\phi}(z) = F(z + \hat{\phi}(z)), \tag{16}$$

where $\hat{L}_{\Lambda} z^k = Akz^k$, and $\hat{F} = \check{F}$. Substituting $\hat{F}(z) = \sum_{k=2}^{\infty} \hat{c}_k z^k$ and $\hat{\phi}(z) = \sum_{k=2}^{\infty} \hat{a}_k z^k$ gives the formal relation

$$\sum_{k=2}^{\infty} Ak \hat{a}_k z^k = \sum_{i=2}^{\infty} \hat{c}_i \left(z + \sum_{k=2}^{\infty} \hat{a}_k z^k \right)^i.$$
(17)

Note that the left-hand side of (17) is simply $Az\hat{\phi}'(z)$ (still only in a formal sense). Also note that all appearing coefficients are non-negative due to the majorization process. Therefore, the question regarding the convergence of $\hat{\phi}$ is reduced to that of the convergence of the solution to the real ordinary differential equation:

$$\hat{\phi}'(x) = (Ax)^{-1}\hat{F}(x+\hat{\phi}(x)), \qquad \hat{\phi}(0) = 0.$$
 (18)

Consider the partial sums $\hat{\phi}_d(x) = \sum_{k=2}^d \hat{a}_k x^k$. By (18), we have

$$0 \leqslant \hat{\phi}'_{d+1}(x) \leqslant (Ax)^{-1} \hat{F}(x + \hat{\phi}_d(x)) \qquad (0 \leqslant x),$$

which yields the following estimate

$$0 \leqslant \hat{\phi}_{d+1}(x) \leqslant x \hat{\phi}'_{d+1}(x) \leqslant A^{-1} \hat{F}(x + \hat{\phi}_d(x)).$$
(19)

Since $\hat{F}(x) = \hat{c}_2 x^2 + \cdots$, there are positive constants r_0 and B_0 such that $0 \leq \hat{F}(x) \leq B_0 x^2$ for all $0 \leq x \leq 2r_0$. Also, since $\hat{\phi}_d(x) = \hat{a}_2 x^2 + \cdots + \hat{a}_d x^d$, we can choose r_0 small enough to ensure that $0 \leq \hat{\phi}_d(x) \leq r_0$ for all $0 \leq x \leq r_0$. These estimates, combined with (19), give

$$0 \leqslant \hat{\phi}_{d+1}(x) \leqslant A^{-1} B_0 (r_0 + r_0)^2 = A^{-1} B_0 (2r_0)^2, \tag{20}$$

for all $0 \le x \le r_0$. By selecting $r_0 \le A/(4B_0)$, we have shown the induction step $\hat{\phi}_d(x) \le r_0 \Rightarrow \hat{\phi}_{d+1}(x) \le r_0$ for all $0 \le x \le r_0$. It follows that $\hat{\phi}(x) \le r_0$ for all $0 \le x \le r_0$, which settles the question of convergence of the change of variables $\zeta + \phi(\zeta)$.

7. Convergence of the normal form

All that remains is to prove the convergence of the nonlinear component G of the normal form. The aim of the proof is to give a lower bound of the radius of convergence r_1 appearing in theorem 3.1. Recall that G is recursively defined by

$$G_i(\zeta) = \langle F_i(\zeta + \phi(\zeta)) \rangle_{\mathbb{U}_\ell} - \sum_{j=1}^n \frac{\partial \phi_i}{\partial \zeta_j}(\zeta) G_j(\zeta) \qquad (i = 1, \dots, n).$$
(21)

As there are no small divisors to consider, the existence of a formal solution to (21) is immediate. The question of convergence, however, is complicated by the fact that the recursive formula is made up of two separate contributing terms. Following the spirit of the previous section, we will use majorization techniques to establish the convergence of *G*.

We begin by reducing the dimension of the range of the problem by considering the majorant system

$$\tilde{G}_i(\zeta) = \langle \tilde{F}(\zeta + \tilde{\phi}(\zeta)) \rangle_{\mathbb{U}_\ell} + \sum_{j=1}^n \frac{\partial \tilde{\phi}_i}{\partial \zeta_j}(\zeta) \tilde{G}_j(\zeta) \qquad (i = 1, \dots, n),$$
(22)

where $\tilde{\phi}$ solves (12), and $F_i \prec \tilde{F}$. This can be solved by a formal power series

$$\tilde{G}_i(\zeta) = \sum_{|m|=2\ell}^{\infty} \tilde{g}_{i,m} \zeta^m,$$
(23)

Note that, by our construction of the normal form, we know that the leading coefficients of \tilde{G} with $|m| < 2\ell$ are zero. Furthermore, since $\tilde{\phi}_1 = \cdots = \tilde{\phi}_n$, the right-hand side of (22) is independent of *i*, and we have $\tilde{G}_1 = \cdots = \tilde{G}_n$. Reducing the dimension of the domain of the problem is achieved by considering the one-dimensional functional equation

$$\hat{G}(z) = \hat{F}(z + \hat{\phi}(z)) + n\hat{\phi}'(z)\hat{G}(z),$$
(24)

where $\hat{\phi}$ solves (16). Again, this can be solved by a formal power series

$$\hat{G}(z) = \sum_{k=2\ell}^{\infty} \hat{g}_k z^k,$$
(25)

where the coefficients g_k can be explicitly solved for by rearranging the terms of (24) into

$$\hat{G}(z) = (1 - n\hat{\phi}'(z))^{-1}\hat{F}(z + \hat{\phi}(z)).$$
(26)

This expression is valid provided that $1 - n\hat{\phi}'(z)$ is invertible. But, since $\hat{\phi}'(z) = 2\hat{a}_2z + \cdots$, we can always arrange this by restricting ourselves to sufficiently small z. In other words, the radius of convergence of $\hat{G}(z)$ is at least as large as the smallest radius of convergence of $\hat{F}(z + \hat{\phi}(z))$ and $(1 - n\hat{\phi}'(z))^{-1}$, and is thus positive.

8. The solutions of the normal form

In this section, we will begin by proving a result on the structure of G using information obtained in section 7. We will use this result to show that the solutions of the normal form act very much like the solutions to the linearized system, as claimed in theorems 3.2 and 3.3.

Proposition 8.1. Under the same conditions as in theorem 3.1, and given $r_2 < r_1$, there exists a positive K_2 such that, in the open ball $B(0, r_2) = \{y: |y| < r_2\}$, we have

$$|G_i(y)| \leq K_2 \max_{i=1,\dots,p} \{|y_{u_i}|^\ell\} \max_{i=1,\dots,q} \{|y_{s_i}|^\ell\} \qquad (i=1,\dots,n)$$

Proof. Let $G_i(\zeta) = \sum_{m \in \mathbb{U}_\ell} g_{i,m} \zeta^m$, and consider the majorants

$$\tilde{G}(\zeta) = \sum_{m \in \mathbb{U}_{\ell}} \tilde{g}_m \zeta^m; \qquad \tilde{g}_m = \max_{i=1,\dots,n} \{|g_{i,m}|\}$$
$$\hat{G}(z) = \sum_{k \ge 2\ell} \hat{g}_k z^k; \qquad \hat{g}_k = \sum_{|m|=k} \tilde{g}_m.$$

We clearly have $G_i \prec \tilde{G} \prec \hat{G}$. Let $\beta(\zeta) = \max_{i=1,\dots,p} \{|\zeta_{u_i}|\} \max_{i=1,\dots,q} \{|\zeta_{s_i}|\}, \sigma_u(m) = \sum_{i=1,\dots,p} m_{u_i}, \sigma_s(m) = \sum_{i=1,\dots,q} m_{s_i}$, and suppose that $|\zeta| < r_2 < r_1$. Then we have

$$\begin{aligned} |G_{i}(\zeta)| &= \left| \sum_{m \in \mathbb{U}_{\ell}} g_{i,m} \zeta^{m} \right| \leq \sum_{m \in \mathbb{U}_{\ell}} |g_{i,m}||\zeta_{1}|^{m_{1}} \cdots |\zeta_{n}|^{m_{n}} \leq \sum_{m \in \mathbb{U}_{\ell}} \tilde{g}_{m} |\zeta_{1}|^{m_{1}} \cdots |\zeta_{n}|^{m_{n}} \\ &\leq \sum_{m \in \mathbb{U}_{\ell}} \tilde{g}_{m} \left(\max_{i=1,\dots,p} \{|\zeta_{u_{i}}|\} \right)^{\sum_{i=1}^{q} m_{u_{i}}} \left(\max_{i=1,\dots,q} \{|\zeta_{s_{i}}|\} \right)^{\sum_{i=1}^{q} m_{s_{i}}} \\ &= \max_{i=1,\dots,p} \{|\zeta_{u_{i}}|^{\ell}\} \max_{i=1,\dots,q} \{|\zeta_{s_{i}}|^{\ell}\} \sum_{m \in \mathbb{U}_{\ell}} \tilde{g}_{m} \max_{i=1,\dots,p} \{|\zeta_{u_{i}}|^{\sigma_{u}(m)-\ell}\} \max_{i=1,\dots,q} \{|\zeta_{s_{i}}|^{\sigma_{s}(m)-\ell}\} \\ &\leq \beta(\zeta)^{\ell} \sum_{m \in \mathbb{U}_{\ell}} \tilde{g}_{m} |\zeta|^{|m|-2\ell} = \beta(\zeta)^{\ell} \sum_{k \geqslant 2\ell} \hat{g}_{k} |\zeta|^{k-2\ell} = \beta(\zeta)^{\ell} |\zeta|^{-2\ell} \sum_{k \geqslant 2\ell} \hat{g}_{k} |\zeta|^{k}. \end{aligned}$$

Now we will use the fact that \hat{G} is analytic. Thus the coefficients \hat{g}_k satisfy $\hat{g}_k \leq DL^k$ for some positive constants D and L. Continuing the estimates, we have

$$\begin{split} |G_{i}(\zeta)| &\leq \beta(\zeta)^{\ell} |\zeta|^{-2\ell} \sum_{k \geqslant 2\ell} \hat{g}_{k} |\zeta|^{k} \leq \beta(\zeta)^{\ell} |\zeta|^{-2\ell} \sum_{k \geqslant 2\ell} DL^{k} |\zeta|^{k} \\ &= \beta(\zeta)^{\ell} |\zeta|^{-2\ell} D \sum_{k \geqslant 2\ell} (L|\zeta|)^{k} = \beta(\zeta)^{\ell} |\zeta|^{-2\ell} D \frac{(L|\zeta|)^{2\ell}}{1 - L|\zeta|} \\ &\leq \beta(\zeta)^{\ell} \frac{DL^{2\ell}}{1 - Lr_{2}} = K_{2} \beta(\zeta)^{\ell} = K_{2} \max_{i=1,...,p} \{|\zeta_{u_{i}}|^{\ell}\} \max_{i=1,...,q} \{|\zeta_{s_{i}}|^{\ell}\}, \end{split}$$

which completes the proof.

In what follows, we will let \mathfrak{B}_r denote the closed *n*-box centred at the origin, and having radius *r*. We will refer to the face { $\zeta \in \mathfrak{B}_r$; $\zeta_{s_1} = r$ } as the *lid* of the box \mathfrak{B}_r . Recall that s_1 is the index of the negative eigenvalue of the smallest modulus. We will also introduce the constant $\kappa = K_2 r^{2\ell-1}$, which should be thought of as being small compared to the minimal distance between the eigenvalues: $\kappa \ll \min\{||\lambda_i| - |\lambda_j||: i \neq j\}$. We also demand that κ be small compared to the minimal distance between the eigenvalues and the origin: $\kappa \ll \min\{|\lambda_{s_1}|, |\lambda_{u_1}|\}$. This can clearly be arranged by taking *r* sufficiently small or, if r < 1, by taking ℓ large. We begin by stating a lemma which establishes an important dominance property:

Lemma 8.2. For all trajectories $\psi(\zeta, t)$ of $\dot{y} = \Lambda y + G(y)$ starting from the lid of \mathfrak{B}_r , we have

$$\psi_{s_1}(\zeta, t) \ge |\psi_{s_i}(\zeta, t)| \qquad (i = 2, \dots, q)$$

throughout the entire box.

Proof. Using proposition 8.1, the differential equations for $\psi_{s_i}(\zeta, t)$ can be enclosed by the differential inequalities

$$\begin{aligned} |\psi_{s_{i}}(\zeta,t) - \lambda_{s_{i}}\psi_{s_{i}}(\zeta,t)| &= |G_{s_{i}}(\psi(\zeta,t))| \\ &\leqslant K_{2} \max_{i=1,\dots,p} \{|\psi_{u_{i}}(\zeta,t)|^{\ell}\} \max_{i=1,\dots,q} \{|\psi_{s_{i}}(\zeta,t)|^{\ell}\} \\ &\leqslant K_{2}r^{2\ell-1} \max_{i=1,\dots,q} \{|\psi_{s_{i}}(\zeta,t)|\} = \kappa \max_{i=1,\dots,q} \{|\psi_{s_{i}}(\zeta,t)|\}. \end{aligned}$$
(27)

Initially, we have $r = \psi_{s_1}(\zeta, 0) \ge |\psi_{s_i}(\zeta, 0)|$, and by the differential inequalities (27) it follows that, if $\psi_{s_1}(\zeta, 0) = |\psi_{s_i}(\zeta, 0)|$, then $|\psi_{s_i}|$ decreases faster than ψ_{s_1} . Now suppose that after some positive time t^* , we have the first occurrence of the situation $\psi_{s_1}(\zeta, t^*) = |\psi_{s_i}(\zeta, t^*)|$ for some i = 2, ..., q. Then, from (27), we have

$$\begin{aligned} |\psi_{s_1}(\zeta,t^\star) - \lambda_{s_1}\psi_{s_1}(\zeta,t^\star)| &\leqslant \kappa\psi_{s_1}(\zeta,t^\star), \\ |\dot{\psi}_{s_i}(\zeta,t^\star) - \lambda_{s_i}\psi_{s_i}(\zeta,t^\star)| &\leqslant \kappa|\psi_{s_i}(\zeta,t^\star)|. \end{aligned}$$

By the same reasoning as above, $|\psi_{s_i}|$ decreases faster than ψ_{s_1} . Hence $\psi_{s_1}(\zeta, t) \ge |\psi_{s_i}(\zeta, t)|$ for all i = 2, ..., q throughout the entire box.

It now follows that the ζ_{s_1} -component of the flow is monotonically decreasing within \mathfrak{B}_r .

Corollary 8.3. For all trajectories $\psi(\zeta, t)$ of $\dot{y} = \Lambda y + G(y)$ starting from the lid of \mathfrak{B}_r , we have

$$(\lambda_{s_1} - \kappa)\psi_{s_1}(\zeta, t) \leqslant \dot{\psi}_{s_1}(\zeta, t) \leqslant (\lambda_{s_1} + \kappa)\psi_{s_1}(\zeta, t)$$

and

$$re^{(\lambda_{s_1}-\kappa)t} \leqslant \psi_{s_1}(\zeta,t) \leqslant re^{(\lambda_{s_1}+\kappa)t}$$

throughout the entire box.

Proof. We simply note that $\max_{i=1,\dots,q}\{|\psi_{s_i}(\zeta, t)|\} = \psi_{s_1}(\zeta, t)$ in (27).

In what follows, we will repeatedly utilize the following lemma, which is easily proved by, e.g., the method of variation of parameters.

Lemma 8.4. The linear ODE $\dot{z} = \lambda z + \varepsilon e^{\mu t}$ has the following solution:

$$z(t) = z(0)e^{\lambda t} + \varepsilon \frac{e^{\lambda t} - e^{\mu t}}{\lambda - \mu}.$$

Regarding the remaining stable components of $\psi(\zeta, t)$, we have the following lemma.

Lemma 8.5. Given $\ell > |\lambda_{s_q}|/(|\lambda_{s_1}| - \kappa)$, there are positive α_i such that, for all trajectories $\psi(\zeta, t)$ of $\dot{y} = \Lambda y + G(y)$ starting from the lid of \mathfrak{B}_r , we have for all i = 2, ..., q

$$|\psi_{s_i}(\zeta,t)-\zeta_{s_i}e^{\lambda_{s_i}t}|\leqslant \frac{\kappa r}{\alpha_i}(1-e^{-\alpha_i t})e^{\lambda_{s_i}}$$

throughout the entire box.

Proof. Using lemma 8.2 and corollary 8.3 together with (27), we can enclose the differential equation for ψ_{s_i} by

$$\begin{aligned} |\dot{\psi}_{s_{i}}(\zeta,t) - \lambda_{s_{i}}\psi_{s_{i}}(\zeta,t)| &= |G_{s_{i}}(\psi(\zeta,t))| \\ &\leqslant K_{2} \max_{i=1,\dots,p} \{|\psi_{u_{i}}(\zeta,t)|^{\ell}\} \max_{i=1,\dots,q} \{|\psi_{s_{i}}(\zeta,t)|^{\ell}\} \\ &\leqslant K_{2}r^{\ell} \max_{i=1,\dots,q} \{|\psi_{s_{i}}(\zeta,t)|^{\ell}\} \leqslant K_{2}r^{\ell}|\psi_{s_{1}}(\zeta,t)|^{\ell} \\ &\leqslant K_{2}r^{\ell}(re^{(\lambda_{s_{1}}+\kappa)t})^{\ell} = \kappa re^{\ell(\lambda_{s_{1}}+\kappa)t}. \end{aligned}$$

Using lemma 8.4, we can explicitly solve for a bound on the perturbation from the linear flow:

$$|\psi_{s_i}(\zeta,t)-\zeta_{s_i}e^{\lambda_{s_i}t}|\leqslant \kappa r\left|\frac{e^{\lambda_{s_i}t}-e^{\ell(\lambda_{s_1}+\kappa)t}}{\lambda_{s_i}-\ell(\lambda_{s_1}+\kappa)}\right|$$

By our choice of ℓ , it is clear that there exist positive α_i satisfying $0 < \alpha_i \leq \lambda_{s_i} - \ell(\lambda_{s_1} + \kappa)$ (i = 2, ..., q). This implies that

$$|\psi_{s_i}(\zeta,t)-\zeta_{s_i}e^{\lambda_{s_i}t}| \leqslant \frac{\kappa r}{\alpha_i}|e^{\lambda_{s_i}t}-e^{\ell(\lambda_{s_1}+\kappa)t}| \leqslant \frac{\kappa r}{\alpha_i}(1-e^{-\alpha_i t})e^{\lambda_{s_i}t},$$

which completes the proof.

Remark 6. If we choose $\ell > (|\lambda_{s_q}| + 1)/(|\lambda_{s_1}| - \kappa)$, then we can take $\alpha_i > 1$, which is used in the estimates of theorem 3.2.

8.1. The proof of theorem 3.2

Assuming for now that Λ has only one positive eigenvalue λ_u , we bound the unstable component ψ_u in the following lemma.

Lemma 8.6. For all trajectories $\psi(\zeta, t)$ of $\dot{y} = \Lambda y + G(y)$ starting from the lid of \mathfrak{B}_r , we have

$$(\lambda_u - \sigma \kappa)\psi_u(\zeta, t) \leqslant \psi_u(\zeta, t) \leqslant (\lambda_u + \sigma \kappa)\psi_u(\zeta, t)$$

and

$$\zeta_u \mathrm{e}^{(\lambda_u - \sigma \kappa)t} \leqslant \psi_u(\zeta, t) \leqslant \zeta_u \mathrm{e}^{(\lambda_u + \sigma \kappa)t}.$$

throughout the entire box. Here $\sigma = \operatorname{sign}(\zeta_u)$.

Proof. As before, using proposition 8.1, the differential equation for $\psi_u(\zeta, t)$ can be enclosed by the differential inequality

$$\begin{aligned} |\dot{\psi}_{u}(\zeta,t) - \lambda_{u}\psi_{u}(\zeta,t)| &= |G_{u}(\psi(\zeta,t))| \leqslant K_{2}|\psi_{u}(\zeta,t)|^{\ell} \max_{i=1,\dots,q} \{|\psi_{s_{i}}(\zeta,t)|^{\ell}\} \\ &\leqslant K_{2}r^{2\ell-1}|\psi_{u}(\zeta,t)| = \kappa |\psi_{u}(\zeta,t)|. \end{aligned}$$
(28)

As an immediate consequence of this lemma, we can obtain bounds on the exit-time $\tau(\zeta)$, which is the time it takes the trajectory starting at ζ to leave the box \mathfrak{B}_r .

Corollary 8.7. For all trajectories $\psi(\zeta, t)$ of $\dot{y} = \Lambda y + G(y)$ starting from the lid of \mathfrak{B}_r , the flow-time required to exit the box \mathfrak{B}_r is enclosed by the bounds

$$\frac{1}{\lambda_u + \kappa} \log \frac{r}{|\zeta_u|} \leqslant \tau_e(\zeta) \leqslant \frac{1}{\lambda_u - \kappa} \log \frac{r}{|\zeta_u|}.$$

These bounds are attained by solving the equation $|\psi_u(\zeta, \tau_e(\zeta))| = r$, using lemma 8.6.

It is now a simple matter of substituting the enclosure of $\tau_e(\zeta)$ for t in the bounds of corollary 8.3 and lemma 8.5 to conclude the proof of theorem 3.2.



Figure 3. The uncertain regions of the lid of \mathfrak{B}_r with $0 < \lambda_{u_1} < \lambda_{u_2}$.

8.2. The proof of theorem 3.3

Turning to the case of having several unstable coordinates, the situation is slightly more delicate. As an example, it is not true that a trajectory will always exit the box through the face corresponding to the strongest expanding coordinate ζ_{u_p} . To illustrate this fact, let us consider the completely linear case $\dot{\zeta} = \Lambda \zeta$ with two unstable directions, ζ_{u_1} and ζ_{u_2} , and assume that a trajectory enters the lid of the box with

$$|\zeta_{u_1}| > r \left(\frac{|\zeta_{u_2}|}{r}\right)^{\lambda_{u_1}/\lambda_{u_2}}.$$
(29)

In this situation, even though $\lambda_{u_1} < \lambda_{u_2}$, the trajectory will exit through the face { $\zeta \in \mathfrak{B}_r: |\zeta_{u_1}| = r$ }. When both quantities of (29) are equal, the trajectory will exit through the intersection of both faces, i.e. through an edge of the box.

Returning to the nonlinear situation at hand $\dot{\zeta} = \Lambda \zeta + G(\zeta)$, the dividing lines become inflated as illustrated in figure 3. Trajectories starting from these *uncertain regions* may exit the box through any one of several faces of \mathfrak{B}_r , and with our limited knowledge of *G*, it is impossible to tell which. Any trajectory starting outside these regions, however, will have a well-defined face of exit.

Another complication is that, when $|\psi_{u_i}(\zeta, t)| \ll |\psi_{u_j}(\zeta, t)|$, we might very well have a situation where G_{u_i} is completely dominated by, e.g., a term of the form $a\zeta_{u_j}^k$, where $i \neq j$ and $k \ge \ell$. This means that the following situation could arise:

$$\begin{split} \psi_{u_i}(\zeta,t) &= \lambda_{u_i} \psi_{u_i}(\zeta,t) + G_{u_i}(\psi(\zeta,t)) \\ &\approx \lambda_{u_i} \psi_{u_i}(\zeta,t) + a \psi_{u_j}(\zeta,t)^k \approx a \psi_{u_j}(\zeta,t)^k \approx G_{u_i}(\psi(\zeta,t)), \end{split}$$

which shows that the u_i -coordinate of the normal form has no resemblance to its linear part. This makes a detailed analysis of the corresponding flow somewhat subtle.

A convenient concept in the forthcoming analysis is that of the *dominating unstable component*. This is simply the currently largest unstable component, which we label with the symbol \hat{i} :

$$|\psi_{u_i}(\zeta, t)| = \max\{|\psi_{u_i}(\zeta, t)|: i = 1, \dots, p\}.$$

Note that the dominating unstable component may change along an orbit flowing through the box \mathfrak{B}_r . We therefore have $\hat{i} = \hat{i}(\zeta, t)$.

We begin our treatment of the unstable components by noting that the dominating unstable component acts very much like its linear counterpart.

Lemma 8.8. While $\psi_{u_i}(\zeta, t) \in \mathfrak{B}_r$ is the dominating unstable component of the trajectory $\psi(\zeta, t)$ of $\dot{y} = \Lambda y + G(y)$, we have

$$(\lambda_{u_{\hat{i}}} - \sigma\kappa)\psi_{u_{\hat{i}}}(\zeta, t) \leqslant \dot{\psi}_{u_{\hat{i}}}(\zeta, t) \leqslant (\lambda_{u_{\hat{i}}} + \sigma\kappa)\psi_{u_{\hat{i}}}(\zeta, t)$$

and

$$\psi_{u_i}(\zeta, t_0) \mathrm{e}^{(\lambda_{u_i} - \sigma\kappa)(t - t_0)} \leqslant \psi_{u_i}(\zeta, t) \leqslant \psi_{u_i}(\zeta, t_0) \mathrm{e}^{(\lambda_{u_i} + \sigma\kappa)(t - t_0)}.$$

Here $\sigma = \text{sign}(\psi_{u_i}(\zeta, t))$, and t_0 is the first time ψ_{u_i} becomes dominating.

Proof. Using proposition 8.1, the differential equation for $\psi_{u_i}(\zeta, t)$ can be enclosed by the differential inequality

$$\begin{aligned} |\psi_{u_i}(\zeta, t) - \lambda_{u_i} \psi_{u_i}(\zeta, t)| &= |G_{u_i}(\psi(\zeta, t))| \\ &\leqslant K_2 \max_{i=1,\dots,p} \{ |\psi_{u_i}(\zeta, t)|^\ell \} \max_{i=1,\dots,q} \{ |\psi_{s_i}(\zeta, t)|^\ell \} \\ &\leqslant K_2 r^{2\ell-1} \max_{i=1,\dots,p} \{ |\psi_{u_i}(\zeta, t)| \} = \kappa |\psi_{u_i}(\zeta, t)|, \end{aligned}$$

which translates into

$$(\lambda_{u_{\hat{i}}} - \sigma\kappa)\psi_{u_{\hat{i}}}(\zeta, t) \leqslant \dot{\psi}_{u_{\hat{i}}}(\zeta, t) \leqslant (\lambda_{u_{\hat{i}}} + \sigma\kappa)\psi_{u_{\hat{i}}}(\zeta, t)$$

where $\sigma = \text{sign}(\psi_{u_i}(\zeta, t))$. The second statement of the lemma follows by integration.

By the same reasoning as in the proof of lemma 8.2, it follows that no weaker unstable component can ever overtake the dominating component:

$$i < \hat{i}(\zeta, t_0) \implies |\psi_{u_i}(\zeta, t)| < |\psi_{u_{\hat{i}(\zeta, t_0)}}(\zeta, t)| \qquad (t_0 \leqslant t).$$

Lemma 8.8 immediately gives a crude upper bound on the time required to exit the box \mathfrak{B}_r .

Corollary 8.9. Let \hat{i} be the dominating unstable component at time t_0 , i.e. let $|\psi_{u_i}(\zeta, t_0)| = \max\{|\psi_{u_i}(\zeta, t_0)|: i = 1, ..., p\}$. Then the flow-time required to exit the box \mathfrak{B}_r is bounded from above by

$$\tau_e \leqslant t_0 + \frac{1}{\lambda_{u_i} - \kappa} \log \frac{r}{|\psi_{u_i}(\zeta, t_0)|}$$

This bound is attained exactly when the dominating component remains dominating throughout the box.

In order to bound the unstable components of the normal form flow, we will make use of corollary 8.3, which provided an upper bound on the dominating stable component.

Lemma 8.10. There exist positive α_i such that, for all trajectories $\psi(\zeta, t)$ of $\dot{y} = \Lambda y + G(y)$ starting from the lid of \mathfrak{B}_r , we have for all i = 1, ..., p

$$|\psi_{u_i}(\zeta,t)-\zeta_{u_i}e^{\lambda_{u_i}t}|\leqslant \frac{\kappa r}{\alpha_i}(1-e^{-\alpha_i t})e^{\lambda_{u_i}t}$$

throughout the entire box.

Proof. Using proposition 8.1 and corollary 8.3, we can enclose the differential equation for ψ_{u_i} by

$$\begin{aligned} |\dot{\psi}_{u_i}(\zeta,t) - \lambda_{u_i}\psi_{u_i}(\zeta,t)| &= |G_{u_i}(\psi(\zeta,t))| \\ &\leqslant K_2 \max_{i=1,\dots,p} \{|\psi_{u_i}(\zeta,t)|^\ell\} \max_{i=1,\dots,q} \{|\psi_{s_i}(\zeta,t)|^\ell\} \\ &\leqslant K_2 r^\ell \max_{i=1,\dots,q} \{|\psi_{s_i}(\zeta,t)|^\ell\} \leqslant K_2 r^\ell |\psi_{s_1}(\zeta,t)|^\ell \\ &\leqslant K_2 r^{2\ell} e^{\ell(\lambda_{s_1}+\kappa)t} \leqslant \kappa r e^{\ell(\lambda_{s_1}+\kappa)t}. \end{aligned}$$

Using lemma 8.4, we can explicitly solve for a bound on the perturbation from the linear flow:

$$|\psi_{u_i}(\zeta,t)-\zeta_{u_i}e^{\lambda_{u_i}t}|\leqslant \kappa r\left|\frac{e^{\lambda_{u_i}t}-e^{\ell(\lambda_{s_1}+\kappa)t}}{\lambda_{u_i}-\ell(\lambda_{s_1}+\kappa)}\right|.$$

We have already chosen κ small enough to guarantee that $\lambda_{s_1} + \kappa < 0$. Therefore, we can find positive α_i such that $0 < \lambda_{u_i} < \alpha_i \leq \lambda_{u_i} - \ell(\lambda_{s_1} + \kappa)$ (i = 1, ..., p). This implies that

$$|\psi_{u_i}(\zeta,t)-\zeta_{u_i}e^{\lambda_{u_i}t}|\leqslant \frac{\kappa r}{\alpha_i}|e^{\lambda_{u_i}t}-e^{\ell(\lambda_{s_1}+\kappa)t}|\leqslant \frac{\kappa r}{\alpha_i}(1-e^{-\alpha_i t})e^{\lambda_{u_i}t}.$$

which completes the proof.

Using these results, we can enclose the time a trajectory starting from the lid of \mathfrak{B}_r spends inside the box.

Corollary 8.11. For all trajectories $\psi(\zeta, t)$ of $\dot{y} = \Lambda y + G(y)$ starting from the lid of \mathfrak{B}_r , the flow-time required to exit the box \mathfrak{B}_r is enclosed by the following inequalities:

$$\tau_e^-(\zeta) \leqslant \tau_e(\zeta) \leqslant \tau_e^+(\zeta),$$

where $\tau_e^-(\zeta)$ and $\tau_e^+(\zeta)$ are defined as follows:

1. Let $\hat{i} = \hat{i}(\zeta, 0)$ (in a tie, take the largest index), and define

$$\tau_{\hat{i}}^{\pm}(\zeta) = \frac{1}{\lambda_{u_{\hat{i}}} \mp \kappa} \log \frac{r}{|\zeta_{u_{\hat{i}}}|};$$

2. For all $i > \hat{i}$ compute

$$\tau_i^-(\zeta) = \frac{1}{\lambda_{u_i}} \log \frac{r}{|\zeta_{u_i}| + (\kappa r/\alpha_{u_i})} \quad \text{and} \quad \tau_i^+(\zeta) = \frac{1}{\lambda_{u_i}} \log \frac{r}{\max\{0, |\zeta_{u_i}| - (\kappa r/\alpha_{u_i})\}}.$$

Here the constants α_{u_i} are defined by $\alpha_{u_i} = \lambda_{u_i} - \ell(\lambda_{s_1} + \kappa)$; 3. Now define $\tau_e^-(\zeta) = \max\{\tau_i^-(\zeta): j \ge \hat{i}\}$ and $\tau_e^+(\zeta) = \min\{\tau_i^+(\zeta): j \ge \hat{i}\}$.

Note that the exit-time $\tau_e(\zeta)$ is infinite exactly when $\zeta_{u_1} = \cdots = \zeta_{u_p} = 0$, i.e. when we enter the box along the stable manifold.

It is now straightforward to obtain bounds on the trajectory when leaving the box \mathfrak{B}_r . Using corollary 8.11, we simply substitute the bounds on the exit-time $\tau_e(\zeta)$ into the enclosure bounds on the components of the flow. For the stable components, we use corollary 8.3 and lemma 8.5. For the unstable components, we use lemmas 8.8 and 8.10. This results in an interval enclosure I_i for each component, which we can possibly tighten by forming the intersection with the interval [-r, r], i.e. $\psi_i(\zeta, \tau_e(\zeta)) \in I_i \cap [-r, r]$. This concludes the proof of theorem 3.3.

Acknowledgment

This research was partially supported by NSF Award No DMS-0107242.

References

- [Be78] Belitskii G R 1978 Equivalence and normal forms of germs of smooth mappings Russ. Math. Surv. 33 107-77
- [Be01] Berz M, Makino K and Hoefkens J 2001 Verified integration of dynamics in the solar system Nonlinear Anal. 47 179–90
- [De89] Deng B 1989 Exponential expansion with Sil'nikov's saddle-focus J. Diff. Eqns 82 156-73
- [Ha60] Hartman P 1960 On local homeomorphisms of Euclidian spaces Bol. Soc. Math. Mexicana 5 220-41
- [Ha64] Hartman P 1964 Ordinary Differential Equations (New York: Wiley)
- [Hi76] Hille E 1976 Ordinary Differential Equations in the Complex Domain (New York: Wiley)
- [KZ03] Kapela T and Zgliczynski P 2003 The existence of simple choreographies for the N-body problem—a computer-assisted proof Nonlinearity 16 1899–918
- [Ne64] Nelson E 1964 Topics in Dynamics I: Flows (Princeton, NJ: Princeton University Press)
- [Po29] Poincaré H 1929 Sur les proprietes des fonctions definies par les equations aux differences partielles Oeuvres vol 1 (Paris: Gautiers-Villars)
- [Si52] Siegel C L 1952 Über die analytische Normalform analytischer Differentialgleichungen in der N\u00e4he einer Gleichgewichtsl\u00f6sung Nachr. Akad. Wiss. G\u00f6ttingen, Math. Phys. Kl. 21–30
- [SM71] Siegel C L and Moser J K 1971 Lectures on Celestial Mechanics (Berlin: Springer)
- [Se85] Sell G R 1985 Smooth linearization near a fixed point Am. J. Math. 107 1035–91
- [St57] Sternberg S 1957 Local contractions and a theorem of Poincaré Am. J. Math. 79 809–24
- [St58] Sternberg S 1958 On the structure of local homeomorphisms of Euclidian *n*-space II Am. J. Math. 80 623–31
- [Tu02] Tucker W 2002 A rigorous ODE solver and Smale's 14th problem Found. Comput. Math. 2 53-117
- [ZM01] Zgliczynski P and Mischaikow K 2001 Rigorous numerics for partial differential equations: the Kuramoto– Sivashinsky equation Found. Comput. Math. 1 255–88