



HAL
open science

A Posteriori Validation of Generalized Polynomial Chaos Expansions

Maxime Breden

► **To cite this version:**

Maxime Breden. A Posteriori Validation of Generalized Polynomial Chaos Expansions. *SIAM Journal on Applied Dynamical Systems*, 2023, 22 (2), pp.765-801. 10.1137/22M1493197 . hal-04296607

HAL Id: hal-04296607

<https://hal.science/hal-04296607>

Submitted on 2 Feb 2024

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

A posteriori validation of generalized polynomial chaos expansions

Maxime Breden *

March 7, 2022

Abstract

Generalized polynomial chaos expansions are a powerful tool to study differential equations with random coefficients, allowing in particular to efficiently approximate random invariant sets associated to such equations. In this work, we use ideas from validated numerics in order to obtain rigorous a posteriori error estimates together with existence results about gPC expansions of random invariant sets. This approach also provides a new framework for conducting validated continuation, i.e. for rigorously computing isolated branches of solutions in parameter-dependent systems, which generalizes in a straightforward way to multi-parameter continuation. We illustrate the proposed methodology by rigorously computing random invariant periodic orbits in the Lorenz system, as well as branches and 2-dimensional manifolds of steady states of the Swift-Hohenberg equation.

Keywords: generalized polynomial chaos; validated numerics; validated continuation; uncertainty quantification

1 Introduction

Most of the mathematical models that are used nowadays to try and describe the world we live in, or at least some very specific region or aspect of it, include some stochastic component. This randomness can have various sources: sometimes we do not fully know or understand the mechanisms underlying the phenomenon we are trying to describe, sometimes we need to account for the influence of events occurring at much smaller scales than that of the full system, for which we cannot afford to solve too accurately, and sometimes our model contains crucial parameters whose value can only be known up to some uncertainty level.

A common mathematical framework to study this last situation is the one of random differential equations, say a random ODE described by a nonlinear vector field f

$$X' = f(X, p), \tag{1}$$

or more generally a random PDE, where p denotes a parameter whose value is not known precisely, and is therefore represented by a random variable. In this work we assume that the probability distribution of p is known, for instance through some preliminary statistical inference. In this situation, one would like to understand and quantify as precisely as possible how the uncertainty in p affects the output of the system [40].

If we want to study the global behavior of (1), one option is to use a Monte-Carlo type approach: sample p according to its known distribution, and study for each sampled value p_i the deterministic system $X' = f(X, p_i)$. Of course, even in the deterministic case, understanding the global dynamics of a system of nonlinear ODEs can already be a daunting task. Numerical simulations can then be of great help to get some insights, in particular in order to study invariant sets (equilibria, periodic orbits, invariant manifolds, connecting orbits, etc), which typically act as building blocks of the global dynamics.

*CMAP, École Polytechnique, route de Saclay, 91120 Palaiseau, France. maxime.breden@polytechnique.edu

Another option to study problems with random parameters like (1), which has risen in popularity in the last decades, is the usage of generalized polynomial chaos (gPC) expansions [20, 52]. The main idea is to expand the random quantity of interest as a series with a well chosen basis, namely polynomials in p which are orthogonal with respect to probability distribution of p . One is then left with computing the (deterministic!) coefficients of this series expansion, and as in the Monte-Carlo approach we recover a deterministic problem, and the ability to use existing algorithms for it. This strategy has proven very effective in various contexts [27, 53], and in particular the recent work [8] showcases that gPC can be used to approximate some random invariant sets generated by random ODEs of the form (1).

Once the gPC expansion has been computed, it readily provides quantitative information about the way the randomness in p influences the solutions of the system. In order to be more concrete, let us focus for instance on periodic orbits. With a gPC representation, we directly have access to the mean and the variance of the period (which typically depends in a non-explicit, nonlinear way on p), and we can also do cheaper Monte-Carlo simulations to estimate the full probability distribution of the period, or to quantify the shape of the orbit in phase space, etc.

The above discussion exemplifies why gPC is a very powerful tool to quantify uncertainties, at least if we had access to the *exact* gPC representation of the object of interest. However, in practice, the fact that we heavily rely on numerical computations introduces an extra level of uncertainty. The two main sources of approximations in the above procedure are: the fact that the (theoretically infinite) gPC series expansion is truncated, because we can only compute finitely many coefficients, and the fact that the deterministic algorithms used to compute these gPC coefficients also contain truncation errors (indeed even for a fully deterministic nonlinear ODE, one cannot hope to compute exactly periodic orbits, or more complicated invariant sets).

Regarding the truncation of the gPC expansion, it is known a priori that the truncation error decays quickly (spectral convergence) when X depends smoothly on p [11, 16], and some tight convergence results were even obtained recently in a non-smooth case [6]. However, when X is not known a priori, these estimates cannot give any quantitative information about the truncation error. A posteriori error estimators for gPC expansions have also been developed, especially in the context of random linear elliptic PDEs [13, 15, 5], but also for more general random PDEs [10, 31, 32]. Yet again, for nonlinear problems these estimators typically still contain some approximations and cannot provide fully rigorous error bounds between the approximate solution and the exact one, if only because the existence of an exact solution is not always readily available.

The purpose of this work is to quantify in a very explicit way all the errors involved in the computation of some random invariant sets using gPC. For instance, if \bar{X} is a gPC representation of an approximate random periodic solution that we obtained numerically, we are going to provide guaranteed a posteriori estimates stating that there exists an exact random periodic solution X^* , with $\|\bar{X} - X^*\| \leq r$ in some well chosen norm, where the error bound r will be explicit. In the context of deterministic dynamical systems, such guaranteed a posteriori error estimates which also provide existence results go back at least to the proof of the Feigenbaum conjecture [14, 26] (see also [45] for an even earlier work) and have become more and more popular since then, mostly under the name of *validated/rigorous numerics* or *computer-assisted proofs*. We will recall some of the main ideas behind these techniques in this work, and refer to the survey papers [21, 22, 23, 39, 47] and books [34, 44] for a more in-depth overview of the field. The main contribution of this work is to show that these ideas can be extended, in a computationally efficient way, to dynamical systems with random coefficients.

Before proceeding further, let us present an alternate viewpoint for the techniques we develop in this paper. The important starting observation is that, in any kind of gPC expansion, the probability distribution of p only influences the choice of the expansion basis. Once a basis has been selected, one can forget the random character of p , and simply view (1) as a deterministic parameter-dependent problem. In that context, numerical continuation methods can be used, for instance to approximate a curve of periodic orbits. This is exactly what we do with a gPC expansion, the only difference being that traditional continuation methods would typically proceed by computing points close to one another

along the curve, and glue them together in a low-order (say piece-wise linear) fashion, whereas here we directly compute a larger chunk of curve at once, by looking for a higher order parameterization.

Numerical continuation can be used together with validated numerics to prove the existence of curves of solutions and to get tight and explicit error bounds (see e.g. [3, 9, 48, 50, 51]), but up to now this has mostly been done with the piece-wise linear approximations provided by usual predictor-corrector techniques. The only exceptions seems to be the recent works [1, 2], where Taylor expansions in the parameter are used to compute and validate larger pieces of curve at once. The approach proposed in this paper is very similar, but we generalize it to other kind of expansions bases, which proves to be sometimes more efficient than using Taylor expansions. This framework also generalizes in a completely straightforward way to rigorous multi-parameter continuation, which again provides a higher-order and more global alternative to the existing techniques [17], which also rely on local piece-wise linear approximations.

The remainder of the paper is organized as follows. In Section 2, we introduce some of the tools that will be required in this work, in particular well chosen sequence spaces provided with a discrete convolution and a type of Newton-Kantorovich Theorem, and start with a basic example (Section 2.6) describing how these tools can be combined to rigorously validate gPC expansions. We then explain in Section 3 how this framework can be applied to random invariant sets, via the example of random periodic orbits in the Lorenz system. This section ends with some comparisons regarding the performance of several choices of polynomial bases. We continue with a different example in Section 4, namely the Swift-Hohenberg equation, for which we rigorously compute parameter-dependent families of steady states. With this example we focus more on the validation continuation viewpoint, and on how the proposed technique interacts with bifurcations, and also discuss how to handle multiple parameters at once. We wrap up in Section 5, where we summarize our work, and discuss the current limitations and possible extensions of the proposed approach. All the codes associated with this work are available at [7].

2 Background material, notations and a basic example

In this section, we introduce some of the objects and tools that we make use of in this work. Most of the material presented here is not original, and mainly included for the convenience of the reader, and for the sake of fixing some notations. We discuss weighted ℓ^1 spaces of Fourier coefficients in Section 2.1, and an extension where each Fourier coefficient is itself written as a gPC expansions together with associated generalized convolutions structures in Section 2.2. We introduce notations for finite dimensional projections in Section 2.3, and state a useful lemma for studying the norm of linear operators on Schauder spaces in Section 2.4. We then recall a specific variation of the Newton-Kantorovich theorem, which is a cornerstone of many computer-assisted technique, in Section 2.5, and then present a very easy example where we use this theorem to validate a gPC expansion in Section 2.6.

2.1 Fourier coefficients and ℓ^1 spaces

It will be convenient to represent several of the solutions we look for in this work as Fourier series, such as

$$u(t) = \sum_{k \in \mathbb{Z}} u_k e^{ikt}.$$

A natural function space to work with is then to consider the set of Fourier coefficients having some prescribed decay rate.

Definition 2.1. *Let $(X, \|\cdot\|)$ be a normed vector space and $\nu \geq 1$. We define*

$$\ell_\nu^1(\mathbb{Z}, X) = \left\{ u = (u_n)_{n \in \mathbb{Z}} \in X^{\mathbb{Z}}, \|u\|_{\ell_\nu^1(\mathbb{Z}, X)} := \sum_{n \in \mathbb{Z}} \|u_n\| \nu^{|n|} < \infty \right\}.$$

In the sequel, we sometimes shorten $\ell_\nu^1(\mathbb{Z}, X)$ into ℓ_ν^1 when knowledge about the set of indices and X is not relevant or clear from context.

Remark 2.2. As soon as $\nu > 1$, the coefficients of an element of ℓ_ν^1 decay at least geometrically, which means the associated function has analytic regularity. This might be seen as a strong requirement, but it should rather be thought of as a precise information: if it happens that the solutions we are dealing with have analytic regularity, by choosing such weighted spaces for the a posteriori analysis we will be able to prove that they do have said regularity. If we had to deal with less smooth solutions, we could use different spaces [28].

We recall that the discrete convolution makes weighted ℓ^1 spaces into Banach algebras, as soon as the weights are submultiplicative. This Banach algebra property is going to be very useful for obtaining the validation estimates.

Lemma 2.3. Let $(X, \|\cdot\|)$ be a Banach algebra, with multiplication denoted by $*$, and $\nu \geq 1$. Given u and v in $\ell_\nu^1(\mathbb{Z}, X)$, we can define their convolution product $u \otimes v$ by

$$(u \otimes v)_k = \sum_{l \in \mathbb{Z}} u_l * v_{k-l} \quad \forall k \in \mathbb{Z},$$

and we have

$$\|u \otimes v\|_{\ell_\nu^1} \leq \|u\|_{\ell_\nu^1} \|v\|_{\ell_\nu^1},$$

i.e., $\ell_\nu^1(\mathbb{Z}, X)$ is a Banach algebra for the multiplication \otimes .

For a deterministic periodic solution, each coefficient u_k is simply a complex number and we will therefore use the above definition with $X = \mathbb{C}$. However, in the presence of a random parameter p in the system, each u_k will also be random, and therefore expressed using a gPC expansion. In that case, u_k will itself be a sequence of coefficients and X an associated sequence space.

In the next subsection, we recall the necessary ingredients for equipping spaces of gPC coefficients with a Banach algebra structure, so as to be able to use the above Lemma.

2.2 Linearization formulas and generalized convolution products

Most of the material presented in this subsection about the relationships between orthogonal polynomials and Banach algebras can be found (in a different context) in the lecture notes [42]. We also refer to the appendix of [8] for discussions related to implementation issues.

In this work, any quantity x (a Fourier coefficient, the period of a periodic orbit, etc) which depends on a parameter p will be written using gPC expansions, i.e.

$$x(p) = \sum_{n \in \mathbb{N}} x_n \phi_n(p),$$

where $(\phi_n)_{n \in \mathbb{N}}$ is basis of polynomials.

The basis of gPC is that, when p is a random variable having a density function ϱ_p with finite moments, one should use for the basis $(\phi_n)_{n \in \mathbb{N}}$ orthogonal polynomials with respect to ϱ_p , i.e. such that $\int \phi_m \phi_n \varrho_p = 0$ as soon as $m \neq n$.

Example 2.4. Here are a couple of examples, which are particular cases of Jacobi polynomials, that which we make use of in this work

- The Legendre polynomials P_n , which correspond to $\varrho_p(t) = \frac{1}{2} \mathbf{1}_{(-1,1)}(t)$;
- The Chebyshev polynomials of the first kind T_n , which correspond to $\varrho_p(t) = \frac{1}{\pi \sqrt{1-t^2}} \mathbf{1}_{(-1,1)}(t)$;

- The Chebyshev polynomials of the second kind U_n , which correspond to $\varrho_p(t) = \frac{2}{\pi}\sqrt{1-t^2}\mathbf{1}_{(-1,1)}(t)$;
- The Gegenbauer or ultraspherical polynomials C_n^μ , $\mu > -\frac{1}{2}$, $\mu \neq 0$, which correspond to $\varrho_p(t) = \frac{2^{2\mu-1}\mu B(\mu,\mu)}{\pi}(1-t^2)^{\mu-\frac{1}{2}}\mathbf{1}_{(-1,1)}(t)$.

We point out that, for a given ϱ_p , each orthogonal polynomial ϕ_n is only defined up to a multiplicative constant, and a normalization condition is required in order to uniquely characterize them. In this work, we choose the condition $\phi_n(1) = 1$ for all n . With this normalization, we recover the traditional definition of the Legendre and Chebyshev polynomials of the first kind, but the usual Chebyshev polynomials of the second kind and Gegenbauer polynomials have to be renormalized. The reason behind this normalization choice is explained in Lemma 2.10.

Finally, we will also use the monomial basis $\phi_n(p) = p^n$, in order to compare the performances of gPC expansions with the one of Taylor expansions.

For a more complete description of gPC choices and their relations to the Askey scheme, see [54].

Definition 2.5. Let $\eta \geq 1$. We define

$$\ell_\eta^1(\mathbb{N}, \mathbb{C}) = \left\{ x = (x_n)_{n \in \mathbb{N}} \in \mathbb{C}^{\mathbb{N}}, \|x\|_{\ell_\eta^1(\mathbb{N}, \mathbb{C})} := \sum_{n \in \mathbb{N}} |x_n| \eta^n < \infty \right\}.$$

In the sequel, we always use η as a weight when we consider a space of gPC coefficients, and ν for the Fourier coefficients, in order to better know at a glance which type of object we are currently dealing with.

Remark 2.6. As in the previous subsection, these spaces encode regularity properties. Indeed, having $\eta \geq 1$ will be sufficient to ensure that, for any x in ℓ_η^1 , the corresponding function

$$p \mapsto \sum_{n \in \mathbb{N}} x_n \Phi_n(p),$$

is at least continuous (see Lemma 2.10), and even analytic when $\eta > 1$ [43, Theorem 8.2]. Thereby, we will often not distinguish between a sequence $x = (x_n)_{n \in \mathbb{N}}$ in ℓ_η^1 and the corresponding function $x : p \mapsto \sum_{n \in \mathbb{N}} x_n \Phi_n(p)$, and use the same symbol to denote both.

Given two functions $x(p)$ and $y(p)$ written as gPC expansions, we now want to define a product on the sequence space $\ell_\eta^1(\mathbb{N}, \mathbb{C})$ corresponding to the multiplication $x(p)y(p)$.

Definition 2.7. Let $(\phi_n)_{n \in \mathbb{N}}$ be a family of (univariate) real polynomials such that ϕ_n is of degree n for all n . The linearization coefficients $(\alpha_k^{m,n})_{k,m,n \in \mathbb{N}}$ for this family are the real numbers such that

$$\phi_m \phi_n = \sum_{k=0}^{n+m} \alpha_k^{m,n} \phi_k, \quad \forall k, n, m \in \mathbb{N}, \quad (2)$$

with $\alpha_k^{m,n} = 0$ for all $k > m + n$.

Definition 2.8. Given linearization coefficients, we define the generalized convolution product $*$ (associated to the linearization coefficients, or equivalently to the polynomial basis) of two sequences of complex numbers $x = (x_n)_{n \in \mathbb{N}}$ and $y = (y_n)_{n \in \mathbb{N}}$ by

$$(x * y)_k = \sum_{m=0}^{\infty} \sum_{n=0}^{\infty} x_m y_n \alpha_k^{m,n}, \quad \forall k \in \mathbb{N}.$$

The generalized convolution product of coefficients corresponds to the pointwise product of functions, at least formally. The Lemma below gives sufficient conditions for the generalized convolution product to be well defined, and for this identification to be justified.

Lemma 2.9. Let $\eta \geq 1$ and $(\alpha_k^{m,n})_{k,m,n \in \mathbb{N}}$ be linearization coefficients such that

$$\sum_{k=0}^{m+n} |\alpha_k^{m,n}| = 1, \quad \forall m, n \in \mathbb{N}. \quad (3)$$

Then, for any x and y in $\ell_\eta^1(\mathbb{N}, \mathbb{C})$, the generalized convolution product $x * y$ (associated to the linearization coefficients) is well defined, belongs to $\ell_\eta^1(\mathbb{N}, \mathbb{C})$, and

$$\|x * y\|_{\ell_\eta^1} \leq \|x\|_{\ell_\eta^1} \|y\|_{\ell_\eta^1},$$

i.e., $\ell_\eta^1(\mathbb{N}, \mathbb{C})$ is a Banach algebra for $*$.

Proof. We simply use the triangle inequality and exchanges sums:

$$\begin{aligned} \|x * y\|_{\ell_\eta^1} &= \sum_{k \in \mathbb{N}} |(x * y)_k| \eta^k \\ &\leq \sum_{m \in \mathbb{N}} \sum_{n \in \mathbb{N}} |x_m| \eta^m |y_n| \eta^n \sum_{k \in \mathbb{N}} |\alpha_k^{m,n}| \eta^{k-n-m}, \end{aligned}$$

which allows us to conclude since $\alpha_k^{m,n} = 0$ for $k > m + n$ and $\eta \geq 1$. □

Lemma 2.10. Let $(\phi_n)_{n \in \mathbb{N}}$ be either:

- the monomial basis,
- the Legendre polynomials P_n ,
- the Chebyshev polynomials of the first kind T_n ,
- the Chebyshev polynomials of the second kind U_n (normalized so that $U_n(1) = 1$),
- the Gegenbauer polynomials C_n^μ (normalized so that $C_n^\mu(1) = 1$).

Then, the associated linearization coefficients satisfy (3). In particular, the corresponding generalized convolution product provides $\ell_\eta^1(\mathbb{N}, \mathbb{C})$ with a Banach algebra structure.

Moreover, for any $\eta \geq 1$ and any $x = (x_n)_{n \in \mathbb{N}}$ in $\ell_\eta^1(\mathbb{N}, \mathbb{C})$, the associated function $x(p) = \sum_{n \in \mathbb{N}} x_n \phi_n(p)$ satisfies

$$\|x\|_{C^0} := \sup_{p \in [-1, 1]} |x(p)| \leq \|x\|_{\ell_\eta^1}.$$

Proof. The first part of the Lemma is known more generally, for a large class of Jacobi polynomials (see [18, 19]), and is based on the nonnegativity of the linearization coefficients. In our context, we have explicit formula for those coefficients in each case [35], which are indeed nonnegative. It then suffices to evaluate (2) at 1 to get (thanks to the normalization condition)

$$1 = \sum_{k=0}^{m+n} \alpha_k^{m,n},$$

and therefore (3).

The second part of the Lemma is a direct consequence of the fact that, with our choice of normalization, $|\phi_n(p)| \leq 1$ for all p in $[-1, 1]$, see [35]. □

2.3 Finite dimensional projections

In practice, we approximate elements in $\ell_\eta^1(\mathbb{N}, \mathbb{C})$ by finite dimensional vectors (or equivalently truncated series).

Definition 2.11. Given N in \mathbb{N} , we define the projector $\Pi_N : \mathbb{C}^{\mathbb{N}} \rightarrow \mathbb{C}^{\mathbb{N}}$ as follows:

$$(\Pi_N x)_n = \begin{cases} x_n & n < N, \\ 0 & n \geq N, \end{cases}$$

and

$$\Pi_N \ell_\eta^1(\mathbb{N}, \mathbb{C}) := \{x \in \ell_\eta^1(\mathbb{N}, \mathbb{C}), \Pi_N x = x\}.$$

We use similar projectors for Fourier series, i.e. ℓ^1 spaces of sequences indexed by \mathbb{Z} . In order not to use too many different notations, we keep the same letter Π to also denote these projectors, but make sure to always use the letter N to refer to gPC indices, and K to refer to Fourier indices.

Definition 2.12. Given K in \mathbb{N} and a vector space X , we define the projector $\Pi_K : X^{\mathbb{Z}} \rightarrow X^{\mathbb{Z}}$ as follows:

$$(\Pi_K u)_k = \begin{cases} u_k & |k| < K, \\ 0 & |k| \geq K, \end{cases}$$

and

$$\Pi_K \ell_\eta^1(\mathbb{Z}, X) := \{u \in \ell_\eta^1(\mathbb{Z}, X), \Pi_K u = u\}.$$

In the case where $X = \ell_\eta^1(\mathbb{N}, \mathbb{C})$, given N in \mathbb{N} we denote by $\Pi_{K,N}$ the composition of Π_K with Π_N (applied component-wise). That is, for $u = (u_k)_{k \in \mathbb{Z}}$ in $(\ell_\eta^1(\mathbb{N}, \mathbb{C}))^{\mathbb{Z}}$,

$$\Pi_{K,N} u = (\dots, 0, 0, \Pi_N u_{-K+1}, \Pi_N u_{-K+2}, \dots, \Pi_N u_{K-2}, \Pi_N u_{K-1}, 0, 0, \dots),$$

and

$$\Pi_{K,N} \ell_\eta^1(\mathbb{Z}, \ell_\eta^1(\mathbb{N}, \mathbb{C})) := \{u \in \ell_\eta^1(\mathbb{Z}, \ell_\eta^1(\mathbb{N}, \mathbb{C})), \Pi_{K,N} u = u\}.$$

2.4 Operator norms

Controlling operator norms will be crucial in our work, and we are often going to rely on the following statement.

Lemma 2.13. Let X be a Banach space with a Schauder basis $(e_j)_{j \in \mathbb{N}}$ and a norm $\|\cdot\|$ of the form

$$\left\| \sum_{j \in \mathbb{N}} x_j e_j \right\| = \sum_{j \in \mathbb{N}} |x_j| w_j, \quad (4)$$

for some prescribed weights $w_j > 0$. Then, for any bounded linear operator B on X ,

$$\|B\| = \sup_{j \in \mathbb{N}} \frac{1}{w_j} \|B e_j\|.$$

Moreover, for any disjoint subsets I and J of \mathbb{N} such that $I \cup J = \mathbb{N}$, if we denote by X_I and X_J the subspaces of X having $(e_j)_{j \in I}$ and $(e_j)_{j \in J}$ as Schauder bases,

$$\|B\| = \max \left(\sup_{\substack{x \in X_I \\ x \neq 0}} \frac{\|Bx\|}{\|x\|}, \sup_{\substack{x \in X_J \\ x \neq 0}} \frac{\|Bx\|}{\|x\|} \right). \quad (5)$$

Proof. For any $x = \sum_{j \in \mathbb{N}} x_j e_j$ in X , we simply use the triangle inequality to get

$$\|Bx\| \leq \sum_{j \in \mathbb{N}} |x_j| \|Be_j\| \leq \left(\sup_{j \in \mathbb{N}} \frac{1}{w_j} \|Be_j\| \right) \sum_{j \in \mathbb{N}} |x_j| w_j,$$

hence $\|B\| \leq \sup_{j \in \mathbb{N}} \frac{1}{w_j} \|Be_j\|$. However, for any j in \mathbb{N} ,

$$\frac{1}{w_j} \|Be_j\| = \frac{\|Be_j\|}{\|e_j\|} \leq \|B\|,$$

which proves (4). The identity (5) then simply amounts to

$$\sup_{j \in \mathbb{N}} \frac{1}{w_j} \|Be_j\| = \max \left(\sup_{j \in I} \frac{1}{w_j} \|Be_j\|, \sup_{j \in J} \frac{1}{w_j} \|Be_j\| \right). \quad \square$$

2.5 A kind of Newton-Kantorovich theorem

As is the case in many works on validated numerics, a crucial tool in our argument is a kind of Newton-Kantorovich theorem [36], which allows us to validate a posteriori a numerically obtained solution.

Given a map \mathcal{F} defined on a Banach space \mathcal{X} , and an *approximate* zero \bar{x} of \mathcal{F} , this theorem provides us with sufficient conditions guaranteeing the existence of a *genuine* zero x^* of \mathcal{F} near \bar{x} , together with explicit error bounds between \bar{x} and x^* .

Theorem 2.14. *Let \mathcal{X} and \mathcal{Y} be Banach spaces, \mathcal{F} be a C^1 map from \mathcal{X} to \mathcal{Y} , \bar{x} an element of \mathcal{X} , A a linear injective map from \mathcal{Y} to \mathcal{X} , and r^* in $(0, +\infty]$. Assume there exist nonnegative constants Y , Z_1 and Z_2 such that*

$$\|A\mathcal{F}(\bar{x})\|_{\mathcal{X}} \leq Y \quad (6a)$$

$$\|I - A D\mathcal{F}(\bar{x})\|_{\mathcal{X}} \leq Z_1 \quad (6b)$$

$$\|A(D\mathcal{F}(x) - D\mathcal{F}(\bar{x}))\|_{\mathcal{X}} \leq Z_2 \|x - \bar{x}\|_{\mathcal{X}} \quad \forall x \in \mathcal{B}_{\mathcal{X}}(\bar{x}, r^*), \quad (6c)$$

where $D\mathcal{F}$ denotes the Fréchet derivative of \mathcal{F} , $\|\cdot\|_{\mathcal{X}}$ simultaneously denotes the norm on \mathcal{X} and the associated operator norm, and $\mathcal{B}_{\mathcal{X}}(\bar{x}, r^*)$ is the closed ball of center \bar{x} and radius r^* in \mathcal{X} . If these constants satisfy

$$Z_1 < 1 \quad (7a)$$

$$2YZ_2 < (1 - Z_1)^2, \quad (7b)$$

then, for any r satisfying

$$\frac{1 - Z_1 - \sqrt{(1 - Z_1)^2 - 2YZ_2}}{Z_2} \leq r < \min \left(\frac{1 - Z_1}{Z_2}, r^* \right) \quad (8)$$

there exists a unique zero x^* of \mathcal{F} in $\mathcal{B}_{\mathcal{X}}(\bar{x}, r)$.

As previously mentioned, similar results already appeared many times, especially in the computer-assisted proof literature (see, e.g., [4, 12, 37, 55]), and we refer to [46] for a detailed proof, which merely consists in applying the contraction mapping theorem to $x \mapsto x - D\mathcal{F}(\bar{x})^{-1}\mathcal{F}(x)$.

Remark 2.15. *In practice, applying this theorem requires two main ingredients: a good enough approximate solution \bar{x} so that Y is small enough, and a good enough approximate inverse A of $D\mathcal{F}(\bar{x})$ so that (7a) holds. For a given Z_1 satisfying (7a) and Z_2 , the condition (7b) tells us in a quantitative way how small Y has to be (and therefore, in some sense, how good of an approximate solution \bar{x} has to be),*

for the existence of a nearby true solution to be guaranteed. Let us also mention that the injectivity assumption for A is usually automatically satisfied as soon as (6b) and (7a) hold, thanks to some structural properties of A , as we will see whenever we apply Theorem 2.14 in this work.

Conceptually, defining a suitable A , which is not only a good approximate inverse of $DF(\bar{x})$ but also simple enough that all the estimates (6) can be obtained, is often the crucial part. For deterministic problems, say a periodic orbit in (1) for a given value of p , this A is often defined as a finite rank perturbation of a somewhat simple operator (for instance a diagonal operator). In this work, we will see how to generalize this construction to parameter dependent problems.

Finally, let us point out that the codomain \mathcal{Y} of \mathcal{F} is inconsequential, as it does not appear anywhere in the estimates (6). The only thing that really matters is that the composition $A\mathcal{F}$ does map \mathcal{X} into itself.

2.6 A basic example

In this subsection, we showcase on a very simple example how the Banach algebra structure of gPC expansions presented in Section 2.2 and Theorem 2.14 can be combined to provide a fully rigorous uncertainty quantification for an algebraic problem.

Let p be a uniform random variable on $[-1, 1]$. Assume we are given a function g as a truncated Legendre series:

$$g(p) = \sum_{n=0}^{N_g-1} g_n P_n(p), \quad (9)$$

with some given N_g and coefficients $(g_n)_{0 \leq n \leq N_g}$ such that g is positive on $[-1, 1]$, and that we want to compute $x = x(p) = \sqrt{g(p)}$. For a given g , and a given truncated Legendre series $\bar{x} = \bar{x}(p)$ approximating $\sqrt{g(p)}$, we are going to derive a fully computable error bound between this approximate gPC representation and the true object $\sqrt{g(p)}$. While this example in itself is of limited interest, the techniques we use to solve it generalize very well to more complicated and interesting problems, in particular when we do not have a closed-form expression for x in terms of p , as we will see in the remaining sections of the paper.

Proposition 2.16. *Consider the function g of the form (9) with $N_g = 6$, $g_0 = 2$, $g_1 = -1$, $g_5 = -1/2$ and $g_n = 0$ otherwise. Let $\bar{x} = \bar{x}(p) = \sum_{n=0}^5 \bar{x}_n \phi_n(p)$, where the coefficients \bar{x}_n are given in Table 1. Both g and \bar{x} are represented in Figure 1.*

There exists a unique x^ in $\ell_1^1(\mathbb{N}, \mathbb{R})$ satisfying $x^*(p) = \sqrt{g(p)}$ for all p in $[-1, 1]$ and $\|\bar{x} - x^*\|_{\ell_1^1} \leq 0.074$.*

Proof. We consider $\eta = 1$, and the map $F : \ell_\eta^1(\mathbb{N}, \mathbb{R}) \rightarrow \ell_\eta^1(\mathbb{N}, \mathbb{R})$ defined by

$$F(x) = x * x - g \quad \forall x \in \ell_\eta^1,$$

where we identify the function g with its sequence of Legendre coefficients, and $*$ is the generalized convolution product associated to the Legendre basis.

We are going to apply Theorem 2.14 to the map F , with $\mathcal{X} = \mathcal{Y} = \ell_\eta^1(\mathbb{N}, \mathbb{R})$, and \bar{x} as an approximate solution (in practice \bar{x} was obtained by using Newton's method on the finite dimensional projection $\Pi_{N_g} F \Pi_{N_g}$ of F). In order to do so, we first need to define a linear map A on $\ell_\eta^1(\mathbb{N}, \mathbb{R})$, which should be a *good enough* approximate inverse of $DF(\bar{x})$, in the sense that (6b) should be satisfied with $Z_1 < 1$. Since $DF(\bar{x})$ is nothing but the multiplication operator $x \mapsto (2\bar{x}) * x$, we compute numerically an element a in $\Pi_{N_g} \ell_\eta^1$ such that $2\bar{x} * a \approx 1$ (see Table 1), and define A as the multiplication operator by a , i.e.

$$Ax = a * x \quad \forall x \in \ell_\eta^1.$$

We are now ready to derive bounds Y , Z_1 and Z_2 satisfying assumption (6).

- For the Y bound, we simply have

$$AF(\bar{x}) = a * (\bar{x} * \bar{x} - g).$$

Since all the elements involved, namely a , \bar{x} and g , belong to $\Pi_{N_g} \ell_\eta^1$, or equivalently are polynomials of degree at most $N_g - 1$, $AF(\bar{x})$ is a polynomial of degree at most $3(N_g - 1)$. Its coefficients, and therefore its norm, can thus be computed explicitly, and we can take $Y = \|a * (\bar{x} * \bar{x} - g)\|_{\ell_\eta^1}$.

- For the Z_1 bound, notice that $I - ADF(\bar{x})$ is simply the multiplication operator by $1 - a * (2\bar{x})$. Therefore its operator norm is equal to the norm of $1 - a * (2\bar{x})$, which can also be computed explicitly, and we can take $Z_1 = \|1 - 2a * \bar{x}\|_{\ell_\eta^1}$.
- Finally, for the Z_2 bound, $A(D\mathcal{F}(x) - D\mathcal{F}(\bar{x}))$ is the multiplication operator by $2a * (x - \bar{x})$, therefore we can take $Z_2 = \|2a\|_{\ell_\eta^1}$ and $r^* = +\infty$.

In principle, one could evaluate the obtained bounds Y , Z_1 and Z_2 for a , \bar{x} and g given in Table 1 by hand, and then check whether the conditions (7) hold. To do it by hand would of course be a waste of time, and become very impractical for higher dimensional problems. Therefore, we evaluate these bounds with a computer, but with interval arithmetic [33, 44] rather than the usual floating point arithmetic, so as to control rounding errors and ensure that the numbers we obtain for Y , Z_1 and Z_2 do satisfy (6). We obtain

$$Y = 0.06509919865113, \quad Z_1 = 0.07544522585888, \quad Z_2 = 1.18378588092726.$$

Hence (7a) holds, and we check that (7b) holds as well. Moreover, the definition of Z_1 together with (7a) imply that $a * \bar{x}$ is invertible, therefore so is a (because ℓ_η^1 is a commutative ring), and we do have that A is injective. We can thus apply Theorem 2.14, and according to (8), there exists a unique zero x^* of F in $\mathcal{B}_{\ell_\eta^1}(\bar{x}, r)$ for all $r \in [r_{min}, r_{max})$, with

$$r_{min} = 0.073908425226395, \quad r_{max} = 0.781015206412929.$$

The computational parts of the proof can be reproduced by downloading the Matlab code at [7], and running `script_BasicExample.m` (in order to get a rigorous proof including rounding error control, you also need Intlab [38]). \square

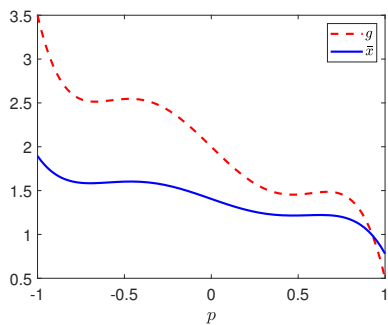


Figure 1: The input function g and the approximate solution \bar{x} which is validated a posteriori in Proposition 2.16.

n	\bar{x}	a
0	1.397466142483791	0.367312683971841
1	-0.361966926543100	0.100133468883877
2	-0.034437671606020	0.030295726801934
3	-0.010225992059514	0.014887942068709
4	-0.025424456095587	0.023944696446413
5	-0.184395889766637	0.055318422290860

Table 1: The coefficients of \bar{x} and a used in Proposition 2.16 and in its proof.

Remark 2.17. Let us assume again that p is a uniform random variable on $[-1, 1]$. If we had only computed \bar{x} numerically as an approximation of $\sqrt{g(p)}$, we would “know” for instance that

$$\mathbb{E}\left(\sqrt{g(p)}\right) \approx \mathbb{E}(\bar{x}(p)) = \bar{x}_0 = 1.397466\dots,$$

but without any guarantee regarding the precision of the approximation. With Proposition 2.16, we get a fully guaranteed (even rounding errors are accounted for) a posteriori estimate controlling the distance between $\sqrt{g(p)}$ and \bar{x} , which proves that

$$\left|\mathbb{E}\left(\sqrt{g(p)}\right) - \bar{x}_0\right| \leq r_{\min} \approx 0.074.$$

It turns out that this estimate is rather conservative, but more precise results can be obtained by looking for \bar{x} (and then for a), in a larger dimensional subspace. For instance, if we look for \bar{x} in $\Pi_N \ell_\eta^1$ with $N = 20$ instead of $N = 6$, we already get an error bound r_{\min} of the order of 10^{-4} , and with $N = 50$ we get again a much more precise approximation, together with a validation radius r_{\min} of less than 10^{-8} , which tells us in particular that

$$\mathbb{E}\left(\sqrt{g(p)}\right) = 1.39729844 \pm 10^{-8}.$$

For more details, simply run `script_BasicExample.m` with different values of N .

3 A posteriori validation of parameter-dependent periodic orbits for the Lorenz system

In this section, we study periodic orbits in the Lorenz system

$$\begin{cases} x' = \sigma(y - x) \\ y' = \rho x - y - xz \\ z' = -\beta z + xy, \end{cases} \quad (10)$$

with the *usual* parameter values $\sigma = 10$, $\beta = 8/3$, but assuming ρ is of the form

$$\rho = \bar{\rho} + \delta p, \quad (11)$$

where $\bar{\rho} = 28$, $\delta \geq 0$ is a given constant, and p varies in $[-1, 1]$. We can either think of p as being a random variable taking values in $[-1, 1]$, in which case the p -depending periodic orbits are random periodic orbits, or consider a deterministic continuation problem in p .

3.1 Setting

We follow the framework developed in [8] for the approximation of random periodic orbits based on Fourier×gPC expansions, and extend it to include the a posteriori validation of the obtained truncated expansions. We look for x , y and z of the form

$$x(t, p) = \sum_{k \in \mathbb{Z}} x_k(p) e^{i\Omega(p)t}, \quad (12)$$

where each Fourier coefficient $x_k(p)$ as well as the unknown frequency $\Omega(p)$ are also expanded as series, this time with the appropriate gPC basis

$$x_k(p) = \sum_{n \in \mathbb{N}} x_{k,n} \phi_n(p), \quad \Omega(p) = \sum_{n \in \mathbb{N}} \Omega_n \phi_n(p). \quad (13)$$

A natural space in which to apply Theorem 2.14 so as to validate an approximate solution $\bar{X} = (\bar{\Omega}, \bar{x}, \bar{y}, \bar{z})$ is then given by

$$\mathcal{X} = \ell_\eta^1(\mathbb{N}, \mathbb{C}) \times (\ell_\nu^1(\mathbb{Z}, \ell_\eta^1(\mathbb{N}, \mathbb{C})))^3, \quad (14)$$

for some $\nu, \eta \geq 1$ to be specified later. For $X = (\Omega, x, y, z)$ in \mathcal{X} , we consider the norm

$$\|X\|_{\mathcal{X}} := |\Omega| + \|x\|_{\ell_\nu^1(\mathbb{Z}, \ell_\eta^1(\mathbb{N}, \mathbb{C}))} + \|y\|_{\ell_\nu^1(\mathbb{Z}, \ell_\eta^1(\mathbb{N}, \mathbb{C}))} + \|z\|_{\ell_\nu^1(\mathbb{Z}, \ell_\eta^1(\mathbb{N}, \mathbb{C}))}.$$

We assume that the linearization coefficients of the chosen gPC basis satisfy (3), so that Lemma 2.10 and then Lemma 2.3 apply, i.e. $\ell_\eta^1(\mathbb{N}, \mathbb{C})$ is a Banach algebra for the generalized convolution product $*$, and then $\ell_\nu^1(\mathbb{Z}, \ell_\eta^1(\mathbb{N}, \mathbb{C}))$ is itself a Banach algebra for the convolution product \otimes .

Introducing the linear operator \mathfrak{K} , which to any sequence $x \in \ell_\nu^1(\mathbb{Z}, \ell_\eta^1(\mathbb{N}, \mathbb{C}))$ associates the sequence $\mathfrak{K}x$ defined as

$$(\mathfrak{K}x)_k = kx_k \quad \forall k \in \mathbb{Z}, \quad (15)$$

and plugging the Ansatz (12)-(13) into the Lorenz system (10), we can rewrite the resulting set of equations on the Fourier \times gPC coefficients in the form

$$F(X) = 0,$$

where $X = (\Omega, x, y, z)$ belongs to \mathcal{X} , and $F(X) = (F^{(x)}(X), F^{(y)}(X), F^{(z)}(X))$ with

$$\begin{cases} F^{(x)}(X) = -i\mathfrak{K}(\Omega \otimes x) - \sigma x + \sigma y \\ F^{(y)}(X) = -i\mathfrak{K}(\Omega \otimes y) + \rho \otimes x - y - (x \otimes z) \\ F^{(z)}(X) = -i\mathfrak{K}(\Omega \otimes z) - \beta z + (x \otimes y). \end{cases} \quad (16)$$

In the above equations, we identify an element of $\ell_\eta^1(\mathbb{N}, \mathbb{C})$ like Ω or ρ with its natural injection in $\ell_\nu^1(\mathbb{Z}, \ell_\eta^1(\mathbb{N}, \mathbb{C}))$, which allows us to write for instance $\Omega \otimes x$, which is nothing but the sequence $(\Omega * x_k)_{k \in \mathbb{Z}}$, where each $\Omega * x_k$ is an element of $\ell_\eta^1(\mathbb{N}, \mathbb{C})$.

This F is almost the one to which we will apply Theorem 2.14, but we first need to add a phase condition to get rid of time-translation invariance and allow F to have isolated zeros. Here we depart slightly from the framework introduced in [8], in which we used a Poincaré phase (or transversality) condition, and instead impose

$$G(X) := \sum_{k \in \mathbb{Z}} ik(x_k * \text{conj}(\tilde{x}_k) + y_k * \text{conj}(\tilde{y}_k) + z_k * \text{conj}(\tilde{z}_k)) = 0, \quad (17)$$

where $(\tilde{x}, \tilde{y}, \tilde{z})$ is an approximate solution previously computed. This is inspired from the integral phase condition

$$\int \langle u, \tilde{u}' \rangle = 0,$$

which has proven to be more robust numerically [24], and turns out to also be more efficient regarding the a posteriori validation.

Given truncation levels K and N in \mathbb{N} , and an approximate periodic solution $\bar{X} = (\bar{\Omega}, \bar{x}, \bar{y}, \bar{z})$ in $\Pi_{K,N}\mathcal{X}$, we are going to try and validate this approximation using Theorem 2.14 for the map

$$\mathcal{F} = (G, F) \quad (18)$$

and the space \mathcal{X} . In order to do so, we first need to derive a suitable approximate inverse A of $D\mathcal{F}(\bar{X})$, and then to obtain bounds satisfying (6). We accomplish these tasks in the next two subsections.

3.2 The approximate inverse A

If we were considering a deterministic periodic orbit (for a given value of p), and therefore working with the space $\mathcal{X}_{deter} = \mathbb{C} \times (\ell_\nu^1(\mathbb{Z}, \mathbb{C}))^3$ rather than $\mathcal{X} = \ell_\eta^1(\mathbb{N}, \mathbb{C}) \times (\ell_\nu^1(\mathbb{Z}, \ell_\eta^1(\mathbb{N}, \mathbb{C})))^3$, a typical way to construct A would be as follows. One would split the space into a *finite part* and a *tail part*

$$\mathcal{X}_{deter} = \Pi_K \mathcal{X}_{deter} \oplus (I - \Pi_K) \mathcal{X}_{deter},$$

where $\Pi_K \mathcal{X}_{deter} = \mathbb{C} \times (\Pi_K \ell_\nu^1(\mathbb{Z}, \mathbb{C}))^3$, and define A separately on both subspaces. For the finite part, we would simply compute numerically an inverse A_K of $\Pi_K D\mathcal{F}(\bar{X})\Pi_K$ which can be represented as a $(6K-2) \times (6K-2)$ matrix with complex entries. For the tail part, i.e. the higher order modes, the parts of \mathcal{F} corresponding to the differential operator would be the most important one, and we would neglect the rest to define A :

$$AX_k := -\frac{1}{ik\Omega} X_k, \quad \forall |k| \geq K,$$

where $X_k = (x_k, y_k, z_k)$. This is but an example of a general strategy for defining approximate inverses in the context of computer-assisted proofs, which consists in choosing A as a finite rank perturbation of some leading order operator that can be inverted by hand, the finite rank part being directly related to a finite dimensional projection of the map \mathcal{F} .

In this work, for a random periodic orbit, we adopt this strategy with a slight twist, by trying to mimic as much as possible the situation in the deterministic case, but replacing \mathbb{C} by $\ell_\eta^1(\mathbb{N}, \mathbb{C})$ and the multiplication on \mathbb{C} by the generalized convolution product $*$.

We consider the same splitting as above, based only on the Fourier modes

$$\mathcal{X} = \Pi_K \mathcal{X} \oplus (I - \Pi_K) \mathcal{X},$$

where $\Pi_K \mathcal{X} = \ell_\eta^1(\mathbb{N}, \mathbb{C}) \times (\Pi_K \ell_\nu^1(\mathbb{Z}, \ell_\eta^1(\mathbb{N}, \mathbb{C})))^3$. We will still refer informally to both subspaces as the *finite part* and the *tail part* respectively, but we emphasize that $\Pi_K \mathcal{X}$ is only “finite” in terms of Fourier modes, but remains an infinite dimensional subspace because we did not truncate anything in the gPC components.

Remark 3.1. *It is really crucial to take the “finite” part, i.e. the part on which the inverse will be computed accurately, as $\Pi_K \mathcal{X}$ and not as $\Pi_{K,N} \mathcal{X}$, otherwise the resulting A will not be a good enough approximate inverse. Indeed, we are allowed to truncate in Fourier because of the regularizing properties of the equation in t (which corresponds to the Fourier expansion), but there is no such regularization in p (which corresponds to the gPC expansion).*

Regarding the tail part, we first compute numerically an approximate inverse Υ of $\bar{\Omega}$ in $\Pi_N \ell_\eta^1(\mathbb{N}, \mathbb{C})$, i.e. such that $\Upsilon * \bar{\Omega} \approx 1$. Then, we define A in the tail in a similar way as above, except $\frac{1}{\Omega} X_k$ now becomes $\Upsilon * X_k$.

For the finite part, we will also compute numerically an approximate inverse A_K of $\Pi_K D\mathcal{F}(\bar{X})\Pi_K$. However, $\Pi_K D\mathcal{F}(\bar{X})\Pi_K$ is no longer finite dimensional, but can be identified with a linear operator on $(\ell_\eta^1(\mathbb{N}, \mathbb{C}))^{6K-2}$. To make things slightly more concrete, this means we can still represent $\Pi_K D\mathcal{F}(\bar{X})\Pi_K$ as a $(6K-2) \times (6K-2)$ matrix, except each entry is now a linear operator on $\ell_\eta^1(\mathbb{N}, \mathbb{C})$ rather than a complex number. The key point here is that each of these linear operators is not any linear operator, but a multiplication operator, and can therefore be represented compactly by an element of $\ell_\eta^1(\mathbb{N}, \mathbb{C})$ (analogously to the way complex numbers in the deterministic case actually represent multiplication operators on \mathbb{C}). Therefore, we compute an approximate inverse A_K of $\Pi_K D\mathcal{F}(\bar{X})\Pi_K$ under the form of a $(6K-2) \times (6K-2)$ matrix of multiplication operators on $\ell_\eta^1(\mathbb{N}, \mathbb{C})$, represented by $(6K-2)^2$ elements of $\Pi_N \ell_\eta^1(\mathbb{N}, \mathbb{C})$. Each of these multiplication operator is still of infinite rank, but the fact they are multiplications (generalized convolutions) with elements of $\Pi_N \ell_\eta^1(\mathbb{N}, \mathbb{C})$ means that A_K can be represented and stored on a computer.

To summarize, we define the linear operator A by

$$\begin{cases} A\Pi_K(X) = A_K\Pi_K X \\ AX_k = \frac{1}{-ik}\Upsilon * X_k, & |k| \geq K, \end{cases} \quad (19)$$

where $\Upsilon * X_k$ must be understood as $(\Upsilon * x_k, \Upsilon * y_k, \Upsilon * z_k)$, and A_K is a linear operator on $(\ell_\eta^1(\mathbb{N}, \mathbb{C}))^{6K-2}$ which takes the form of $(6K-2)^2$ multiplication operators with elements of $\Pi_N \ell_\eta^1(\mathbb{N}, \mathbb{C})$ (which are computed numerically so that $A_K \approx (\Pi_K D\mathcal{F}(\bar{X})\Pi_K)^{-1}$).

Remark 3.2. *In practice, one of the main limiting factors for computer-assisted proofs like the ones we are using here is the dimension of the finite part of A , and the computing power and memory requirement associated to it. Given the type of expansion we are using, namely bi-infinite series, it is remarkable that the number of complex numbers needed to represent this finite part scales like $K^2 N$ (it is actually equal to $(6K-2)^2 N$), rather than like $K^2 N^2$. This is possible because we take advantage of the multiplication operator structure.*

3.3 Bounds for the a posteriori validation

Now that A has been defined in (19), we are left with deriving estimates Y , Z_1 and Z_2 satisfying (6). Once a proper framework has been obtained, including an appropriate definition of A and a suitable choice of sequence space, the derivation of these estimates is by now standard in the computer-assisted proof literature. Therefore, we only go into the details when they are specific to the new structure of A that is used in this work.

We recall that the map \mathcal{F} we are considering is defined in (16)-(18), the space \mathcal{X} in (14), and that the approximate solution \bar{X} belongs to $\Pi_{K,N}\mathcal{X}$, i.e. is a trigonometric polynomial of degree at most $N-1$, whose coefficients are polynomials of degree at most $K-1$. Similarly, we assume that \tilde{x} , \tilde{y} and \tilde{z} involved in the phase condition (17) all belong to $\Pi_{K,N}\ell_\nu^1(\mathbb{Z}, \ell_\eta^1(\mathbb{N}, \mathbb{C}))$.

3.3.1 The bound Y

As was the case in the example of Section 2.6, obtaining a Y bound satisfying (6a) is rather straightforward, as we can simply take

$$Y := \|A\mathcal{F}(\bar{X})\|_{\mathcal{X}}.$$

The only thing to notice is that $A\mathcal{F}(\bar{X})$ has only finitely non-zero coefficients, hence it can be computed exactly on a computer, up to rounding errors which are taken care of by the use of interval arithmetic. Indeed, having \bar{X} in $\Pi_{K,N}\mathcal{X}$ means $\mathcal{F}(\bar{X})$ belongs to $\Pi_{2K-1,2N-1}\mathcal{X}$, and that $A\mathcal{F}(\bar{X})$ belongs to $\Pi_{2K-1,3N-2}\mathcal{X}$, hence the above defined Y is computable in finitely many operations.

3.3.2 The bound Z_1

In order to obtain a Z_1 estimate, we need to bound the operator norm of $B := I - A D\mathcal{F}(\bar{X})$. To that end, we use the following splitting (see Lemma 2.13)

$$\|B\|_{\mathcal{X}} = \max \left(\sup_{\substack{X \in \Pi_{2K-1}\mathcal{X} \\ X \neq 0}} \frac{\|BX\|_{\mathcal{X}}}{\|X\|_{\mathcal{X}}}, \sup_{\substack{X \in (I - \Pi_{2K-1})\mathcal{X} \\ X \neq 0}} \frac{\|BX\|_{\mathcal{X}}}{\|X\|_{\mathcal{X}}} \right). \quad (20)$$

In order to handle the first part, we will use the following lemma, which is a direct consequence of the Banach algebra property of Lemma 2.10 and of the usual computation of ℓ^1 operator norms.

Lemma 3.3. *Let B be a linear operator on $\ell_\nu^1(\mathbb{Z}, \ell_\eta^1(\mathbb{N}, \mathbb{C}))$, represented as an infinite matrix $(B_{k,l})_{k,l \in \mathbb{Z}}$ of linear operators on $\ell_\eta^1(\mathbb{N}, \mathbb{C})$, and assume that each of those is in fact a multiplication operator by an element $b_{k,l}$ in $\ell_\eta^1(\mathbb{N}, \mathbb{C})$. Then, denoting by $\|b\|_{\ell_\eta^1}$ the infinite matrix of real numbers $(\|b_{k,l}\|_{\ell_\eta^1})_{k,l \in \mathbb{Z}}$, we have that*

$$\|B\|_{\ell_\nu^1(\mathbb{Z}, \ell_\eta^1)} = \left\| \|b\|_{\ell_\eta^1} \right\|_{\ell_\nu^1} = \sup_{l \in \mathbb{Z}} \frac{1}{\nu^{|l|}} \sum_{k \in \mathbb{Z}} \|b_{k,l}\|_{\ell_\eta^1} \nu^{|k|}.$$

This lemma easily generalizes to a linear operator defined on \mathcal{X} (or on a subspace of \mathcal{X}), and allows us to get a computable upper-bound Z_1^{finite} of the first supremum in (20). Indeed, for X in $\Pi_{2K-1}\mathcal{X}$ and $B = I - ADF(\bar{X})$, BX belongs to $\Pi_{3K-2}\mathcal{X}$, which means we only have finitely many multiplication operators on ℓ_η^1 whose norm we need to compute in order to control this supremum over $\Pi_{2K-1}\mathcal{X}$. Finally, since we assumed that all the multiplication operators in A were with elements of $\Pi_N \ell_\eta^1(\mathbb{N}, \mathbb{C})$, each of those multiplication operators in B are with elements of $\Pi_{2N-1} \ell_\eta^1(\mathbb{N}, \mathbb{C})$, which have only finitely many non-zero coefficients, which makes their norm fully computable.

Regarding the second part of the splitting, if X belongs to $(I - \Pi_{2K-1})\mathcal{X}$, then $DF(\bar{X})X$ is in $(I - \Pi_K)\mathcal{X}$, and therefore so is $ADF(\bar{X})X$. Hence, still assuming X belongs to $(I - \Pi_{2K-1})\mathcal{X}$, BX is in $(I - \Pi_K)\mathcal{X}$, we can write BX rather explicitly since it does not involve A_K : for $|k| \geq K$ we get

$$\begin{aligned} (BX)_k &= X_k - ADF_k(\bar{X})X \\ &= \begin{pmatrix} x_k \\ y_k \\ z_k \end{pmatrix} + \frac{1}{ik} \Upsilon * \begin{pmatrix} -ik\bar{\Omega} * x_k - \sigma x_k + \sigma y_k \\ -ik\bar{\Omega} * y_k + \rho * x_k - y_k - (\bar{x} \otimes z + x \otimes \bar{z})_k \\ -ik\bar{\Omega} * z_k - \beta z_k + (\bar{x} \otimes y + x \otimes \bar{y})_k \end{pmatrix} \\ &= (1 - \Upsilon * \bar{\Omega}) * \begin{pmatrix} x_k \\ y_k \\ z_k \end{pmatrix} + \frac{1}{ik} \Upsilon * \begin{pmatrix} -\sigma x_k + \sigma y_k \\ \rho * x_k - y_k - (\bar{x} \otimes z + x \otimes \bar{z})_k \\ -\beta z_k + (\bar{x} \otimes y + x \otimes \bar{y})_k \end{pmatrix}, \end{aligned}$$

which yields

$$\begin{aligned} \|BX\|_{\mathcal{X}} &\leq \|1 - \Upsilon * \bar{\Omega}\|_{\ell_\eta^1} \|X\|_{\mathcal{X}} \\ &\quad + \frac{1}{K} \left((\sigma \|\Upsilon\|_{\ell_\eta^1} + \|(\bar{z} - \rho) \otimes \Upsilon\|_{\ell_\nu^1(\ell_\eta^1)} + \|\bar{y} \otimes \Upsilon\|_{\ell_\nu^1(\ell_\eta^1)}) \|x\|_{\ell_\nu^1(\ell_\eta^1)} \right. \\ &\quad \left. + ((\sigma + 1) \|\Upsilon\|_{\ell_\eta^1} + \|\bar{x} \otimes \Upsilon\|_{\ell_\nu^1(\ell_\eta^1)}) \|y\|_{\ell_\nu^1(\ell_\eta^1)} + (\|\bar{x} \otimes \Upsilon\|_{\ell_\nu^1(\ell_\eta^1)} + \beta \|\Upsilon\|_{\ell_\eta^1}) \|z\|_{\ell_\nu^1(\ell_\eta^1)} \right) \\ &\leq \left(\|1 - \Upsilon * \bar{\Omega}\|_{\ell_\eta^1} + \frac{1}{K} \max \begin{pmatrix} \sigma \|\Upsilon\|_{\ell_\eta^1} + \|(\bar{z} - \rho) \otimes \Upsilon\|_{\ell_\nu^1(\ell_\eta^1)} + \|\bar{y} \otimes \Upsilon\|_{\ell_\nu^1(\ell_\eta^1)} \\ (\sigma + 1) \|\Upsilon\|_{\ell_\eta^1} + \|\bar{x} \otimes \Upsilon\|_{\ell_\nu^1(\ell_\eta^1)} \\ \|\bar{x} \otimes \Upsilon\|_{\ell_\nu^1(\ell_\eta^1)} + \beta \|\Upsilon\|_{\ell_\eta^1} \end{pmatrix} \right) \|X\|_{\mathcal{X}}. \end{aligned}$$

Therefore, we can define

$$Z_1^{tail} = \|1 - \Upsilon * \bar{\Omega}\|_{\ell_\eta^1} + \frac{1}{K} \max \begin{pmatrix} \sigma \|\Upsilon\|_{\ell_\eta^1} + \|(\bar{z} - \rho) \otimes \Upsilon\|_{\ell_\nu^1(\ell_\eta^1)} + \|\bar{y} \otimes \Upsilon\|_{\ell_\nu^1(\ell_\eta^1)} \\ (\sigma + 1) \|\Upsilon\|_{\ell_\eta^1} + \|\bar{x} \otimes \Upsilon\|_{\ell_\nu^1(\ell_\eta^1)} \\ \|\bar{x} \otimes \Upsilon\|_{\ell_\nu^1(\ell_\eta^1)} + \beta \|\Upsilon\|_{\ell_\eta^1} \end{pmatrix},$$

and finally $Z_1 := \max(Z_1^{finite}, Z_1^{tail})$ which satisfies (6b).

3.3.3 Z_2

As in the example of Section 2.6, our map \mathcal{F} is quadratic, which allows us to take $r^* = +\infty$. In particular, for any X and X' in \mathcal{X} we have

$$D^2G(\bar{X})(X, X') = 0$$

$$D^2F(\bar{X})(X, X') = i\mathfrak{K} \left(\Omega \circledast \begin{pmatrix} x' \\ y' \\ z' \end{pmatrix} + \Omega' \circledast \begin{pmatrix} x \\ y \\ z \end{pmatrix} \right) + \begin{pmatrix} 0 \\ -x \circledast z' - x' \circledast z \\ x \circledast y' + x' \circledast y \end{pmatrix},$$

where \mathfrak{K} was introduced in Section 3.1. Therefore

$$\begin{aligned} \|AD^2\mathcal{F}(\bar{X})(X, X')\|_{\mathcal{X}} &\leq \|A\mathfrak{K}\|_{\mathcal{X}} (|\Omega| \|X'\|_{\mathcal{X}} + |\Omega'| \|X\|_{\mathcal{X}}) \\ &\quad + \|A\|_{\mathcal{X}} \left(\|x\|_{\ell_v^1(\ell_\eta^1)} \left(\|y'\|_{\ell_v^1(\ell_\eta^1)} + \|z'\|_{\ell_v^1(\ell_\eta^1)} \right) + \|x'\|_{\ell_v^1(\ell_\eta^1)} \left(\|y\|_{\ell_v^1(\ell_\eta^1)} + \|z\|_{\ell_v^1(\ell_\eta^1)} \right) \right) \\ &\leq \max(\|A\mathfrak{K}\|_{\mathcal{X}}, \|A\|_{\mathcal{X}}) \|X\|_{\mathcal{X}} \|X'\|_{\mathcal{X}}, \end{aligned}$$

and $Z_2 := \max(\|A\mathfrak{K}\|_{\mathcal{X}}, \|A\|_{\mathcal{X}})$ satisfies (6c) (with $r^* = +\infty$).

Remark 3.4. *To be precise, in the above definition of Z_2 we replace the norm of A and $A\mathfrak{K}$ by easily computable upper bounds of their norms, obtained using Lemma 3.3.*

This estimate could also be made slightly sharper, by noticing that some ‘‘columns’’ of A and $\mathfrak{K}A$ are always multiplied by zero in the above computation of $AD^2\mathcal{F}(\bar{X})(X, X')$, and can therefore be excluded from the norm computation.

3.4 Results

We are now ready to rigorously validate approximate periodic solutions of (10)-(11) represented as truncated Fourier \times gPC series, by proving the existence of a true solution within a distance at most r of the approximate one, for an explicit value of r .

Approximate solutions using truncated Fourier \times gPC series were already obtained in [8, Section 6], but without guarantee regarding their accuracy, which is what we add in this paper, thanks to Theorem 2.14 and the estimates derived up to now in this section. Here is an example of the kind of results we can obtain with this approach.

Theorem 3.5. *Take the parameter values $\bar{\rho} = 28$, $\delta = 10$ in (11), the generalized convolution product (Definition 2.8) associated to the Legendre polynomials P_n , the weights $\nu = \eta = 1$ in the definition of \mathcal{X} (14), and the truncation parameters $K = 100$ and $N = 15$. Consider the approximate Fourier \times Legendre solution \bar{X} in $\Pi_{K,N}\mathcal{X}$ of (10)-(11), which can be downloaded at [7], and for which a couple of orbits are represented on Figure 2.*

There exists a periodic solution X^ in \mathcal{X} of (10)-(11), such that $\|X^* - \bar{X}\|_{\mathcal{X}} \leq r_{min} = 1.3724 \times 10^{-4}$. This is the unique solution in the open ball of center \bar{X} and radius $r_{max} = 1.382 \times 10^{-3}$ in \mathcal{X} .*

Proof. We consider \mathcal{F} as in (16)-(18), A as in (19), and evaluate the bounds Y , Z_1 and Z_2 (with $r^* = +\infty$) obtained in Section 3.3. We get

$$Y = 3.62368... \times 10^{-5} \quad Z_1 = 0.72214... \quad Z_2 = 201.067...,$$

hence assumptions (7) are satisfied. From (6b) and (7a) we know that A must be surjective, and that the tail part of A is bijective ($Z_1^{tail} < 1$ yields $\|1 - \Upsilon * \bar{\Omega}\|_{\ell_\eta^1} < 1$, therefore Υ is invertible). The finite part A_K of A is only surjective a priori, but since A_K can be represented as a (finite) matrix over the commutative ring ℓ_η^1 , surjectivity implies injectivity and we do have that A is injective. Theorem 2.14 then yields the announced results, with

$$r_{min} = \frac{1 - Z_1 - \sqrt{(1 - Z_1)^2 - 2YZ_2}}{Z_2} \quad \text{and} \quad r_{max} = \frac{1 - Z_1}{Z_2}.$$

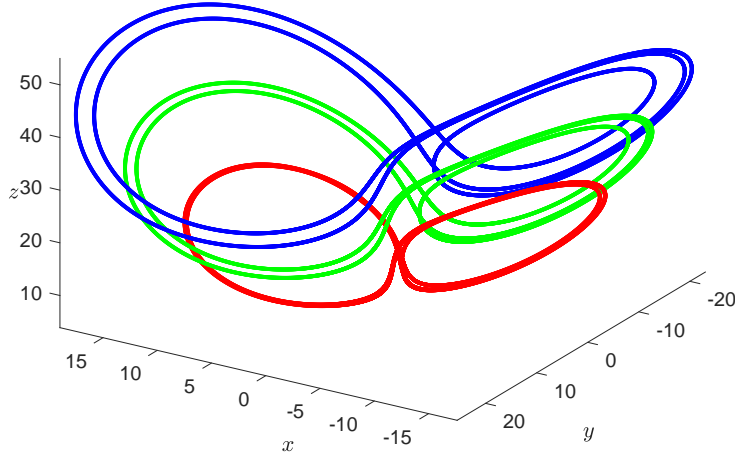


Figure 2: Several approximate periodic orbits of the Lorenz system (10)-(11), for $\rho = 18$ ($p = -1$) in red, $\rho = 28$ ($p = 0$) in green, and $\rho = 38$ ($p = 1$) in blue, all encoded in the Fourier \times gPC expansion \bar{X} , and validated in Theorem 3.5.

The computational parts of the proof, namely the computation of the finite part A_K of A and the evaluation of the bounds, can be reproduced using `script_Lorenz.m` available at [7] (with Intlab [38] for the required interval arithmetic computations). \square

Remark 3.6. *Since we used a Legendre expansion in p , the solution X^* described in Theorem 3.5 would be a natural gPC representation of a random periodic orbit of (10)-(11) where p is a uniform random variable in $[-1, 1]$. We would then directly get statistics about the random periodic orbit, for instance an approximation of the expectation of its frequency*

$$\mathbb{E}(\Omega^*) \approx \bar{\Omega}_0 = 1.5993\dots,$$

together with a guaranteed error bound

$$|\mathbb{E}(\Omega^*) - \bar{\Omega}_0| \leq r_{min}.$$

Moreover, the obtained solution contains a precise description of a periodic orbit for each p in $[-1, 1]$, therefore it could also be used to compute statistics of a random periodic orbit assuming a different distribution for p . For instance, if p has a density ϱ_p , then we get an approximation of the expectation of its frequency

$$\mathbb{E}(\Omega^*) \approx \mathbb{E}(\bar{\Omega}) = \sum_{n=0}^N \bar{\Omega}_n \int_{-1}^1 P_n(s) \varrho_p(s) ds,$$

and an error bound

$$\begin{aligned} |\mathbb{E}(\Omega^*) - \mathbb{E}(\bar{\Omega})| &= \left| \sum_{n=0}^{\infty} \Omega_n^* \int_{-1}^1 P_n(s) \varrho_p(s) ds - \sum_{n=0}^N \bar{\Omega}_n \int_{-1}^1 P_n(s) \varrho_p(s) ds \right| \\ &\leq \sum_{n=0}^{\infty} |\Omega_n^* - \bar{\Omega}_n| \int_{-1}^1 \varrho_p(s) ds \\ &= \|\Omega^* - \bar{\Omega}\|_{\ell_n^1} \\ &\leq r_{min}, \end{aligned}$$

since $\bar{\Omega}_n = 0$ for $n \geq N$ and each $|P_n|$ is bounded by 1 on $[-1, 1]$.

If p does not have a uniform distribution, the approximate solution obtained using Legendre polynomials will be less accurate (at least in L^2 norm) than the one obtained with the gPC basis associated to ϱ_p , and one then has to do extra computations a posteriori, like the integrals $\int_{-1}^1 P_n(s)\varrho_p(s)ds$. Nonetheless, since the cost of the validation can change significantly from one choice of basis to the other (see the discussion below), the natural choice of gPC basis (i.e. the one associated to ϱ_p) might not always be the cheapest option.

In the above discussion, we mostly adopted the viewpoint of random periodic orbits, but Theorem 3.5 also provides us with a deterministic continuation result, namely the existence (and precise description) of a branch of periodic orbits for ρ going from $\bar{\rho} - \delta = 18$ to $\bar{\rho} + \delta = 38$. If there is no underlying random distribution for p , we are completely free from the gPC paradigm, and should try to chose the *best* expansion basis, where of course one has to specify in which sense we mean best. In the following we investigate two criteria:

- For each basis, what is the smallest value of N for which the validation is successful?
- For a fixed N , what is the minimal validation radius r_{min} obtained with each basis?

The first criterion is related to the cost of the validation, both in terms of computational time and memory requirement (see Remark 3.2). The second one assesses the accuracy of the obtained approximation, or at least the accuracy that can be guaranteed.

The output of these comparisons is described in Table 2 regarding the cost of the validation, and in Table 3 regarding the accuracy. In both cases we considered (10)-(11) with $\bar{\rho} = 28$ and $\delta = 10$, a fixed truncation level $K = 100$ for the Fourier modes, and weights $\nu = \eta = 1$ in the norm on \mathcal{X} . These experiments can be reproduced using `script_Lorenz.m` available at [7] (with Intlab [38] for the required interval arithmetic computations).

Polynomial basis	Legendre	Chebyshev	Chebyshev 2nd kind	Gegenbauer $\mu = 20$	Taylor
Minimal value of N	14	13	15	22	28

Table 2: We validate an approximate solution of (10)-(11) with $\bar{\rho} = 28$ and $\delta = 10$, for several choices of polynomial bases for the expansion in p . As soon as N is taken strictly smaller than the value indicated here, the validation fails, typically because (7b) is no longer satisfied.

Polynomial basis	Legendre	Chebyshev	Chebyshev 2nd kind	Gegenbauer $\mu = 20$	Taylor
Error bound	8.7×10^{-7}	7.6×10^{-7}	9.7×10^{-7}	1.5×10^{-5}	7.1×10^{-4}

Table 3: We validate an approximate solution of (10)-(11) with $\bar{\rho} = 28$ and $\delta = 10$, for several choices of polynomial bases for the expansion in p , but this time with a fixed truncation level $N = 28$. In each case the validation is successful (meaning that the assumptions of Theorem 2.14 can be checked using the estimates of Section 3.3), and the displayed error bound corresponds to $r_{min} = \frac{1-Z_1-\sqrt{(1-Z_1)^2-2YZ_2}}{Z_2}$.

The Chebyshev polynomials (of the first kind) prove to be the best choice in both metrics (cost and accuracy), which is not surprising given their remarkable approximation properties [43], although the difference with the Legendre polynomials or the Chebyshev polynomials of the second kind is barely noticeable. On the other hand, there seems to be a significant difference between using a Chebyshev expansion, and a Gegenbauer expansion (with $\mu = 20$) or a Taylor expansion, in particular since the latter require significantly more modes for the validation to be successful (Table 2), which is related to the fact that they yield less accurate approximations (Table 3), at least in the \mathcal{X} norm which is used for the proof. While Taylor expansions were already used successfully to obtain impressive results about

validated branches of stationary and periodic solutions of PDEs [1, 2], the present comparison suggests that replacing the Taylor expansion by a Chebyshev expansion in the continuation variable would prove even more efficient.

4 Validated continuation of steady states of the Swift-Hohenberg equation

We now concentrate fully on the parameter continuation viewpoint, and consider as an example the Swift-Hohenberg equation [41]

$$\partial_t u = -(1 + \Delta)^2 u + \rho u - \beta u^3, \quad (21)$$

where $u = u(t, x)$ is a scalar function, for which we compute and validate branches of equilibria. We focus on the 1 dimensional case, where the spatial variable x belongs in $(0, L)$ together with homogeneous Neumann boundary conditions, which is already very rich, as is illustrated by the bifurcation diagram of steady states represented in Figure 3. For the moment we keep β fixed, and take ρ as the continuation parameter, which we again normalize by writing

$$\rho = \bar{\rho} + \delta p, \quad (22)$$

where $\bar{\rho}$ and $\delta \geq 0$ are given constants, and p varies in $[-1, 1]$.

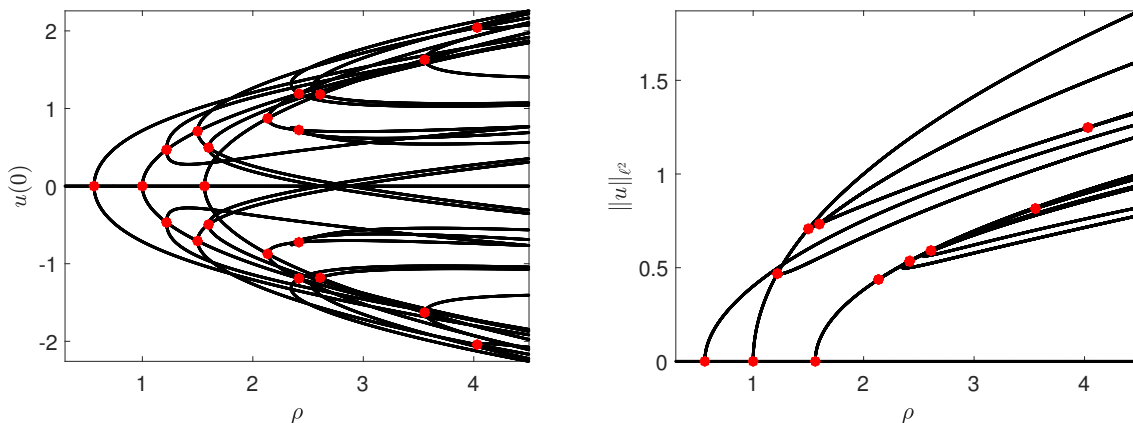


Figure 3: A numerical bifurcation diagram of steady states of (21), for $\beta = 1$ and $L = 2\pi$, represented using two different projections: $u(0)$ on the left, and $\|u\|_{\ell^2}$ on the right. The red dots indicate numerically detected bifurcation.

This problem has already been used as a test case for rigorous continuation methods. We re-emphasize that the main novelty of our work in that regard is the fact that we also expand the solution in the continuation parameter, with several choices of bases, including Chebyshev polynomials, which allows us to represent and validate “in one go” large portions of the curves of solutions.

4.1 Setup for the validation

The validation setup is very similar to the one used in Section 3, so we go over it more briefly.

The homogeneous Neumann boundary conditions make it natural to expand the solution in Fourier series in the x variable:

$$u(x, p) = u_0(p) + 2 \sum_{k=1}^{\infty} u_k(p) \cos\left(\frac{k\pi}{L} x\right) = \sum_{k \in \mathbb{Z}} u_k(p) e^{i \frac{k\pi}{L} x}, \quad (23)$$

with $u_{-k} = u_k$, and we use a gPC expansion for the p variable

$$u_k(p) = \sum_{n \in \mathbb{N}} u_{k,n} \phi_n(p). \quad (24)$$

We look for solutions $u = (u_{k,n})_{\substack{k \in \mathbb{Z} \\ n \in \mathbb{N}}}$ in the space $\mathcal{X} = \ell_\nu^1(\mathbb{Z}, \ell_\eta^1(\mathbb{N}, \mathbb{R}))$ where we impose that $u_{-k,n} = u_{k,n}$ for all k and n , and consider the norm

$$\|u\|_{\mathcal{X}} = \|u\|_{\ell_\nu^1(\mathbb{Z}, \ell_\eta^1(\mathbb{N}, \mathbb{R}))}.$$

The zero-finding map \mathcal{F} is defined as

$$\mathcal{F}(u) = - \left(I - \left(\frac{\pi}{L} \mathfrak{K} \right)^2 \right)^2 u + \rho \otimes u - \beta u \otimes u \otimes u,$$

where the operator \mathfrak{K} is defined in equation (15). Given an approximate zero $\bar{u} \in \Pi_{K,N} \mathcal{X}$, we construct an approximate inverse A of $\mathcal{DF}(\bar{u})$ as in Section 3.2, the main difference being that the $\frac{-1}{ik} \Upsilon$ factor in the tail part of A is now taken as $\frac{-1}{\lambda_k}$, where

$$\lambda_k := \left(1 - \frac{\pi k}{L} \right)^2,$$

and we make sure to take the truncation level K large enough to ensure that λ_k cannot vanish for $k \geq K$.

Regarding the bounds Y , Z_1 and Z_2 needed to apply Theorem 2.14, Y can again be computed by simply evaluating $A\mathcal{F}(\bar{u})$ (with interval arithmetic). For Z_1 , we again introduce $B = I - A\mathcal{DF}(\bar{u})$ and separate its norm using Lemma 2.13

$$\|B\|_{\mathcal{X}} = \max \left(\sup_{\substack{X \in \Pi_{3K-2} \mathcal{X} \\ X \neq 0}} \frac{\|BX\|_{\mathcal{X}}}{\|X\|_{\mathcal{X}}}, \sup_{\substack{X \in (I - \Pi_{3K-2}) \mathcal{X} \\ X \neq 0}} \frac{\|BX\|_{\mathcal{X}}}{\|X\|_{\mathcal{X}}} \right).$$

As in Section 3.3.2, we can again compute explicitly an upper-bound Z_1^{finite} for the first supremum, and control the second one by

$$Z_1^{tail} = \frac{\|\rho - 3\beta \bar{u} \otimes \bar{u}\|_{\mathcal{X}}}{\min_{k \geq K} \lambda_k}.$$

Finally, we can take $Z_2 = 6|\beta| \|A\|_{\mathcal{X}} (\|\bar{u}\|_{\mathcal{X}} + r^*)$ for the last bound.

4.2 Results in the 1-parameter case

Here is an example of the type of results that can be obtained with this approach.

Theorem 4.1. *Consider the 1D Swift–Hohenberg equation (21) with $\beta = 1$, $L = 2\pi$, $\bar{\rho} = 2.9$ and $\delta = 1.6$ in (22), and the branch of approximate steady states \bar{u} represented in blue in Figure 4, whose precise description in terms of Fourier×Chebyshev coefficients can be downloaded at [7]. Take also $\nu = \eta = 1$.*

With the notations introduced in Section 4.1, there exists a zero u^ of \mathcal{F} in \mathcal{X} such that $\|\bar{u} - u^*\|_{\mathcal{X}} \leq r_{min} = 3.8 \times 10^{-4}$, and which is unique among all u in \mathcal{X} such that $\|\bar{u} - u^*\|_{\mathcal{X}} \leq r_{max} = 6.5 \times 10^{-3}$. This u^* corresponds to an isolated branch of steady states of (21) with $\beta = 1$ and $L = 2\pi$, for ρ in $[1.3, 4.5]$.*

Proof. We again evaluate the bounds Y , Z_1 and Z_2 , obtained in Section 4.1, check that assumptions (7) are satisfied, and apply Theorem 2.14, which yields the existence and uniqueness statement for a zero u^* of \mathcal{F} near \bar{u} . The computational parts of the proof, namely the computation of the finite part A_K of A

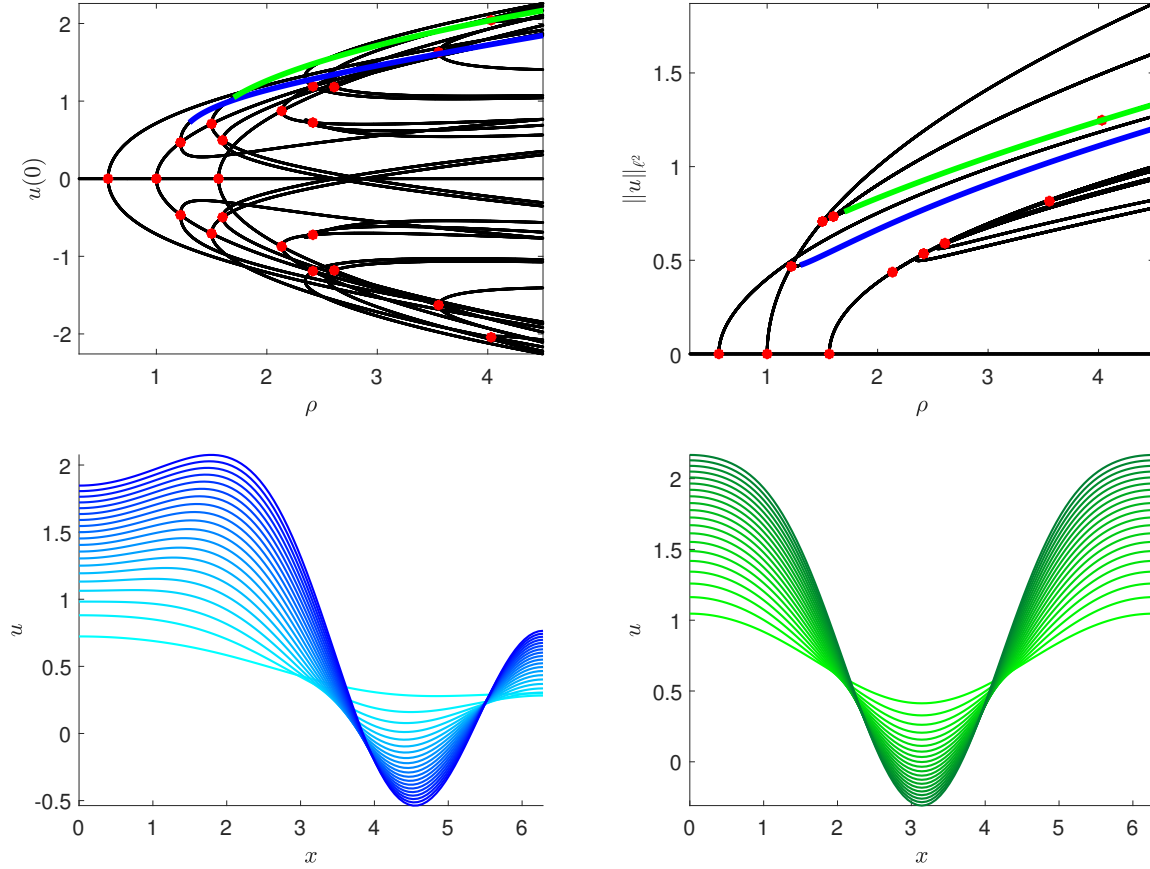


Figure 4: Two specific portions of branches, in blue and green on the bifurcation diagrams at the top, for which several solutions along the branch are represented at the bottom (left for the blue branch and right for the green one). Regarding the solutions at the bottom, the lighter the color the smaller the corresponding value of ρ . The whole blue branch has been validated in Theorem 4.1. The validation of green branch failed, which is coherent with the fact that numerics suggest the presence of a bifurcation crossing that branch.

and the evaluation of the bounds, can be reproduced using `script_SwiftHohenberg.m` available at [7] (with Intlab [38] for the required interval arithmetic computations).

It remains to be proven that the branch of steady states corresponding to u^* is isolated. This is essentially due to the fact that the ℓ_η^1 norm controls the \mathbb{C}^0 norm (Lemma 2.10), hence we can apply Theorem 2.14 uniformly in p . To be more precise, for any p in $[-1, 1]$ and $u = (u_{k,n})$ in $\mathcal{X} = \ell_\nu^1(\mathbb{Z}, \ell_\eta^1(\mathbb{N}, \mathbb{R}))$, we can consider $u(p) = (u_k(p))$ which now belongs to $\ell_\nu^1(\mathbb{Z}, \mathbb{R})$, with $u_k(p)$ as in (24). Thanks to Lemma 2.10, we have that, for any p in $[-1, 1]$:

$$\|u(p)\|_{\ell_\nu^1(\mathbb{Z}, \mathbb{R})} \leq \|u\|_{\ell_\nu^1(\mathbb{Z}, \ell_\eta^1(\mathbb{N}, \mathbb{R}))}. \quad (25)$$

We also consider the map \mathcal{F}_p , which is defined as \mathcal{F} but with p fixed, and only acts on elements of $\ell_\nu^1(\mathbb{Z}, \mathbb{R})$. Similarly, we recall that A can be represented as an *infinite matrix* $(A_{k,l})_{k,l \in \mathbb{Z}}$ of operators on $\ell_\eta^1(\mathbb{N}, \mathbb{C})$, and that each $A_{k,l}$ is in fact a multiplication operator on $\ell_\eta^1(\mathbb{N}, \mathbb{R})$, represented by an element

$a_{k,l}$ in $\ell_\eta^1(\mathbb{N}, \mathbb{R})$. Hence we can consider $A(p) = (a_{k,l}(p))_{k,l \in \mathbb{Z}}$, which now acts on $\ell_\nu^1(\mathbb{Z}, \mathbb{R})$, where

$$a_{k,l}(p) = \sum_{n \in \mathbb{N}} (a_{k,l})_n \phi_n(p).$$

Since the generalized convolution product of coefficients corresponds to the pointwise product of functions, we have that $A(p)\mathcal{F}_p(\bar{u}(p)) = (A\mathcal{F}(\bar{u}))(p)$, and by (25)

$$\|A(p)\mathcal{F}_p(\bar{u}(p))\|_{\ell_\nu^1(\mathbb{Z}, \mathbb{R})} \leq \|A\mathcal{F}(\bar{u})\|_{\mathcal{X}} \leq Y,$$

for all p in $[-1, 1]$. Similarly,

$$\|I_{\ell_\nu^1(\mathbb{Z}, \mathbb{R})} - A(p)D\mathcal{F}_p(\bar{u}(p))\|_{\ell_\nu^1(\mathbb{Z}, \mathbb{R})} \leq \|I_{\mathcal{X}} - AD\mathcal{F}(\bar{u})\|_{\mathcal{X}} \leq Z_1,$$

for all p in $[-1, 1]$. Indeed, for any linear operator $B = (B_{k,l})_{k,l \in \mathbb{Z}}$ acting on \mathcal{X} , where each $B_{k,l}$ is a multiplication operator on $\ell_\eta^1(\mathbb{N}, \mathbb{R})$ represented by an element $b_{k,l}$ in $\ell_\eta^1(\mathbb{N}, \mathbb{R})$, we have (see Lemma 3.3)

$$\|B(p)\|_{\ell_\nu^1(\mathbb{Z}, \mathbb{R})} = \left\| (b_{k,l}(p))_{k,l \in \mathbb{Z}} \right\|_{\ell_\nu^1} \leq \left\| \left(\|b_{k,l}\|_{\ell_\eta^1(\mathbb{N}, \mathbb{R})} \right)_{k,l \in \mathbb{Z}} \right\|_{\ell_\nu^1} = \|B\|_{\mathcal{X}}.$$

Finally, for any p in $[-1, 1]$ and any v in $\mathcal{B}_{\ell_\nu^1(\mathbb{Z}, \mathbb{R})}(\bar{u}(p), r^*)$, we get

$$\begin{aligned} \|A(p)(D\mathcal{F}_p(v) - D\mathcal{F}(\bar{u}(p)))\|_{\ell_\nu^1(\mathbb{Z}, \mathbb{R})} &\leq \|A(p)\|_{\ell_\nu^1(\mathbb{Z}, \mathbb{R})} \|D\mathcal{F}_p(v) - D\mathcal{F}(\bar{u}(p))\|_{\ell_\nu^1(\mathbb{Z}, \mathbb{R})} \\ &\leq \|A(p)\|_{\ell_\nu^1(\mathbb{Z}, \mathbb{R})} 6|\beta| \left(\|\bar{u}(p)\|_{\ell_\nu^1(\mathbb{Z}, \mathbb{R})} + r^* \right) \\ &\leq Z_2. \end{aligned}$$

Hence, for each p in $[-1, 1]$ we can apply Theorem 2.14 to the map \mathcal{F}_p and the approximate solution $\bar{u}(p)$, which proves that the steady state $u^*(p)$ of (21) is locally unique, and in particular there cannot be a another branch of steady states of (21) bifurcating from u^* . \square

Remark 4.2. *We used a Chebyshev expansion in p in Theorem 4.1 because we expect it to be the most efficient choice to represent the branch of solutions. Indeed, while we could for instance have gotten a similar result with a Taylor expansion, we would have needed to take at least $N = 47$ for assumption (7b) to be satisfied, and $N = 72$ if we wanted to get an error estimate r_{min} which is as small as in Theorem 4.1, whereas the current proof with a Chebyshev expansion uses only $N = 15$.*

A remarkable part of Theorem 4.1 is that it guarantees that the portion of the branch that is validated is isolated, i.e. we have a proof that there is no other branch of steady states connected to this part. On the other hand, this means that we cannot hope to validate a part of a branch that goes through a bifurcation. Indeed, if we try to validate the branch of steady states represented in green in Figure 4, the proof fails because Z_1 (in fact Z_1^{finite}) remains larger than 1, no matter how large we take K and N . This does not prove, but strongly suggests, that there is indeed a bifurcation on this part of the branch. Computer-assisted proofs of the existence of bifurcations are possibles, but require more work, see for instance [2, 3, 30, 49] and the references therein. If the parameter ρ is modeled by a random variable, one may want to try and quantify how these possible bifurcations impact the behavior of the system, which is for instance discussed in the recent work [25].

4.3 Extension to multi-parameter validated continuation

Let us now consider both ρ and β as varying parameters in (21), normalized as

$$\rho = \bar{\rho} + \delta_\rho p_1, \quad \beta = \bar{\beta} + \delta_\beta p_2, \tag{26}$$

where $\bar{\rho}, \bar{\beta}$ and $\delta_\rho, \delta_\beta \geq 0$ are given constants, and p_1, p_2 vary in $[-1, 1]$. Away from bifurcation points, we expect to get a 2-dimensional manifold of steady states parametrized by $p = (p_1, p_2)$. Using a bi-variate gPC expansion, we can approximate and then rigorously validate such manifold of steady states. It is remarkable that this generalization from the 1-parameter case requires only very minor modifications, both in terms of the estimates and in terms of the code. The only other work we are aware of in which validated multi-parameter continuation is studied is [17], in which the transition from the 1-parameter case requires a significant effort.

Starting back from (23), we now consider a bi-variate expansion for each Fourier coefficient

$$u_k(p) = u_k(p_1, p_2) = \sum_{n \in \mathbb{N}^2} u_{k,n} \phi_n(p),$$

where the new basis is simply obtained by taking the tensor product of two univariate bases:

$$\phi_n(p) := \phi_{n_1}^{(1)}(p_1) \phi_{n_2}^{(2)}(p_2) \quad \forall n = (n_1, n_2) \in \mathbb{N}^2.$$

In the sequel we take the Chebyshev polynomials of the first kind for both $\phi_{n_1}^{(1)}$ and $\phi_{n_2}^{(2)}$, but all the bases mentioned up to now could be combined here. A generalized convolution product associated to such bi-variate expansion can be defined in a straightforward way from the generalized convolution products associated to each univariate basis, see e.g. [8, Appendix].

Up to changing the space to $\mathcal{X} = \ell_\nu^1(\mathbb{Z}, \ell_{\eta_1}^1(\mathbb{N}, \ell_{\eta_2}^1(\mathbb{N}, \mathbb{R}))) \simeq \ell_\nu^1(\mathbb{Z}, \ell_\eta^1(\mathbb{N}^2, \mathbb{R}))$, to taking $\beta = \bar{\beta} + \delta_\beta p_2$ instead of β constant in \mathcal{F} , and to replacing $|\beta|$ by $\|\beta\|_{\ell_{\eta_2}^1(\mathbb{N}, \mathbb{R})}$ in the Z_2 estimate, we can use exactly the same setup as in Section 4.1 to validate an approximate 2-dimensional manifold of steady states.

Theorem 4.3. *Consider the 1D Swift–Hohenberg equation (21) with $L = 2\pi$ and the manifold of approximate steady states \bar{u} represented in Figure 5, whose precise description in terms of Fourier×gPC coefficients can be downloaded at [7].*

There exists a zero u^ of \mathcal{F} in \mathcal{X} such that $\|\bar{u} - u^*\|_{\mathcal{X}} \leq r_{min} = 4.4 \times 10^{-3}$, and which is unique among all u in \mathcal{X} such that $\|\bar{u} - u^*\|_{\mathcal{X}} \leq r_{max} = 8 \times 10^{-3}$. This u^* corresponds to an isolated manifold of steady states of (21) with $L = 2\pi$, for (ρ, β) in $[2, 4] \times [0.25, 1.75]$.*

Proof. The proof again amounts to checking the assumptions of Theorem 2.14. The computational parts of the proof, namely the computation of the finite part A_K of A and the evaluation of the bounds, can be reproduced using `script_SwiftHohenberg_2para.m` available at [7] (with Intlab [38] for the required interval arithmetic computations). \square

5 Conclusion

In this work, we introduced a new methodology to obtain fully rigorous a posteriori error bounds for several types of gPC expansions (Legendre, Chebyshev of the first and the second kind, and Gegenbauer expansions). We showcased via several examples that this strategy can be used in the context of random invariant sets generated by random ODEs or PDEs, allowing to get a very precise and certified description of random periodic orbits of ODEs and of random steady states of parabolic PDEs.

These techniques can also be seen through the lens of rigorous/validated numerics, and in this context they provide a new way of rigorously computing curves or higher-dimensional manifolds of solutions in parameter-dependent systems, generalizing an approach introduced recently in [2]. It is remarkable that the memory requirements associated to this approach can be made to scale linearly with the dimension of the gPC projection (see Remark 3.2).

We finish by mentioning possible generalizations but also current limitations and open questions related to this work that we believe to be of interest.

- We only considered ODEs or PDEs with polynomial nonlinearities, which is a particularly convenient framework to work in with spectral techniques. Yet, some non-polynomial nonlinearities can be handled in a similar way, making use of ideas from automatic differentiation, see e.g. [29].

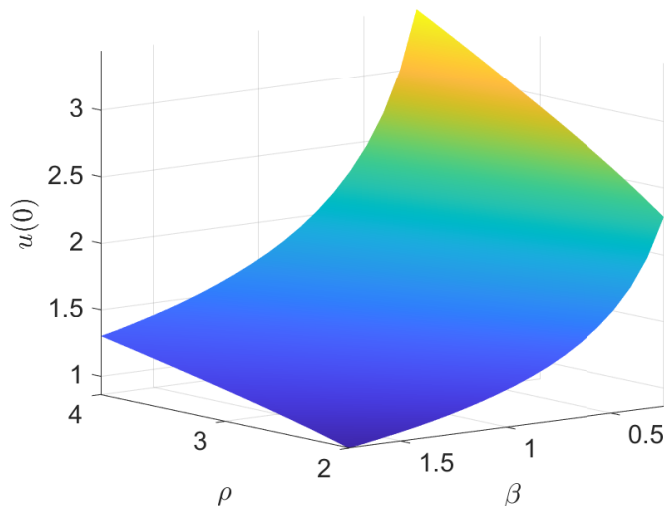


Figure 5: A validated manifold of steady states of (21), for $L = 2\pi$, and (ρ, β) in $[2, 4] \times [0.25, 1.75]$, represented using $u(0)$ as a projection. The approximate manifold of steady states \bar{u} is composed of $K = 20$ Fourier modes, $N_1 = 5$ Chebyshev modes in ρ and $N_2 = 10$ Chebyshev modes in β .

- While we only studied random steady states and random periodic orbits in this work, the proposed approach generalizes in a straightforward way to rigorously compute other types of random invariant sets, as soon as we already have the tools to rigorously compute them in the deterministic case, which is for instance the case for invariant manifolds or connecting orbits.
- We restricted our attention to random parameters having somewhat classical distributions (namely uniform distributions or at least symmetric beta distributions). For more exotic distributions, in particular distributions that are obtained from data and have no analytic expression, one of the main difficulty with our approach is that we require an explicit knowledge of the linearization coefficients. We believe that generalizing the techniques of this paper to a wider class of random parameters (maybe making use of a probability transform to recover a uniform distribution) would be of interest.
- Even if we stick with classical distributions, for which the linearization coefficients are known analytically, our approach can currently only handle bounded random parameters, and in particular excludes Gaussian or exponential distributions. The main reason is that the corresponding orthogonal polynomials, namely Hermite and Laguerre polynomials, do not readily give rise to a discrete convolution structure like the one we could make use of in this work (Lemma 2.10). Finding a way to rigorously compute gPC expansions with those bases, which occur very naturally in many problems, would also be of great interest.
- We conclude with a comment about the implementation. Because we wanted to handle several different expansions in a uniform way, we did not take advantage of the fact that for some expansions (namely Chebyshev and Taylor expansions), the corresponding convolutions can be very efficiently computed using FFT (or DCT) algorithms. If one wanted to focus solely on Chebyshev expansions, which we would for instance recommend if one is only interested in the deterministic parameter-continuation viewpoint, making use of the FFT could improve the performances of the code significantly, especially for higher dimensional problems.

References

- [1] G. Arioli. Computer assisted proof of branches of stationary and periodic solutions, and Hopf bifurcations, for dissipative PDEs. *Communications in Nonlinear Science and Numerical Simulation*, 105:106079, 2022.
- [2] G. Arioli, F. Gazzola, and H. Koch. Uniqueness and bifurcation branches for planar steady Navier–Stokes equations under Navier boundary conditions. *Journal of Mathematical Fluid Mechanics*, 23(3):1–20, 2021.
- [3] G. Arioli and H. Koch. Computer-assisted methods for the study of stationary solutions in dissipative systems, applied to the Kuramoto-Sivashinski equation. *Arch. Ration. Mech. Anal.*, 197(3):1033–1051, 2010.
- [4] G. Arioli, H. Koch, and S. Terracini. Two novel methods and multi-mode periodic solutions for the Fermi-Pasta-Ulam model. *Communications in mathematical physics*, 255(1):1–19, 2005.
- [5] A. Bespalov, C. E. Powell, and D. Silvester. Energy norm a posteriori error estimation for parametric operator equations. *SIAM Journal on Scientific Computing*, 36(2):339–363, 2014.
- [6] F. Bourgey, E. Gobet, and C. Rey. A comparative study of polynomial-type chaos expansions for indicator functions. *HAL preprint, hal-03199734*, 2021.
- [7] M. Breden. Matlab code for “A posteriori validation of generalized polynomial chaos expansions”. https://github.com/MaximeBreden/gPC_expansions, 2022.
- [8] M. Breden and C. Kuehn. Computing invariant sets of random differential equations using polynomial chaos. *SIAM Journal on Applied Dynamical Systems*, 19(1):577–618, 2020.
- [9] M. Breden, J.-P. Lessard, and M. Vanicat. Global bifurcation diagrams of steady states of systems of PDEs via rigorous numerics: a 3-component reaction-diffusion system. *Acta applicandae mathematicae*, 128(1):113–152, 2013.
- [10] T. Butler, C. Dawson, and T. Wildey. A posteriori error analysis of stochastic differential equations using polynomial chaos expansions. *SIAM Journal on Scientific Computing*, 33(3):1267–1291, 2011.
- [11] C. Canuto, M. Y. Hussaini, A. Quarteroni, and T. A. Zang. *Spectral methods in fluid dynamics*. Springer Science & Business Media, 2012.
- [12] S. Day, J.-P. Lessard, and K. Mischaikow. Validated continuation for equilibria of PDEs. *SIAM J. Numer. Anal.*, 45(4):1398–1424, 2007.
- [13] M. K. Deb, I. M. Babuška, and J. T. Oden. Solution of stochastic partial differential equations using galerkin finite element techniques. *Computer Methods in Applied Mechanics and Engineering*, 190(48):6359–6372, 2001.
- [14] J.-P. Eckmann and P. Wittwer. A complete proof of the Feigenbaum conjectures. *Journal of statistical physics*, 46(3):455–475, 1987.
- [15] M. Eigel, C. J. Gittelsohn, C. Schwab, and E. Zander. Adaptive stochastic galerkin fem. *Computer Methods in Applied Mechanics and Engineering*, 270:247–269, 2014.
- [16] D. Funaro. *Polynomial approximation of differential equations*, volume 8. Springer Science & Business Media, 2008.
- [17] M. Gameiro, J.-P. Lessard, and A. Pugliese. Computation of smooth manifolds via rigorous multi-parameter continuation in infinite dimensions. *Foundations of Computational Mathematics*, 16(2):531–575, 2016.

- [18] G. Gasper. Linearization of the product of Jacobi polynomials. I. *Canadian Journal of Mathematics*, 22(1):171–175, 1970.
- [19] G. Gasper. Linearization of the product of Jacobi polynomials. II. *Canadian Journal of Mathematics*, 22(3):582–593, 1970.
- [20] R. G. Ghanem and P. D. Spanos. Stochastic Finite Element Method: Response Statistics. In *Stochastic Finite Elements: A Spectral Approach*, pages 101–119. Springer, 1991.
- [21] J. Gómez-Serrano. Computer-assisted proofs in PDE: a survey. *SeMA Journal*, 76(3):459–484, 2019.
- [22] T. Kapela, M. Mrozek, D. Wilczak, and P. Zgliczyński. CAPD:: DynSys: a flexible C++ toolbox for rigorous numerical analysis of dynamical systems. *Communications in nonlinear science and numerical simulation*, 101:105578, 2021.
- [23] H. Koch, A. Schenkel, and P. Wittwer. Computer-assisted proofs in analysis and programming in logic: a case study. *SIAM review*, 38(4):565–604, 1996.
- [24] B. Krauskopf, H. M. Osinga, and J. Galán-Vioque. *Numerical continuation methods for dynamical systems*, volume 2. Springer, 2007.
- [25] C. Kuehn and K. Lux. Uncertainty quantification of bifurcations in random ordinary differential equations. *SIAM Journal on Applied Dynamical Systems*, 20(4):2295–2334, 2021.
- [26] O. E. Lanford III. A computer-assisted proof of the Feigenbaum conjectures. *Bulletin of the American Mathematical Society*, 6(3):427–434, 1982.
- [27] O. Le Maître and O. M. Knio. *Spectral methods for uncertainty quantification: with applications to computational fluid dynamics*. Springer Science & Business Media, 2010.
- [28] J.-P. Lessard and J. D. Mireles James. Computer assisted fourier analysis in sequence spaces of varying regularity. *SIAM Journal on Mathematical Analysis*, 49(1):530–561, 2017.
- [29] J.-P. Lessard, J. D. Mireles James, and J. Ransford. Automatic differentiation for Fourier series and the radii polynomial approach. *Physica D: Nonlinear Phenomena*, 334:174–186, 2016.
- [30] J.-P. Lessard, E. Sander, and T. Wanner. Rigorous continuation of bifurcation points in the diblock copolymer equation. *Journal of Computational Dynamics*, 4(1&2):71, 2017.
- [31] L. Mathelin and O. Le Maître. Dual-based a posteriori error estimate for stochastic finite element methods. *Communications in Applied Mathematics and Computational Science*, 2(1):83–115, 2007.
- [32] F. Meyer, C. Rohde, and J. Giesselmann. A posteriori error analysis for random scalar conservation laws using the stochastic galerkin method. *IMA Journal of Numerical Analysis*, 40(2):1094–1121, 2020.
- [33] R. E. Moore. *Methods and applications of interval analysis*. SIAM, 1979.
- [34] M. T. Nakao, M. Plum, and Y. Watanabe. *Numerical Verification Methods and Computer-Assisted Proofs for Partial Differential Equations*, volume 53 of *Springer Series in Computational Mathematics*. Springer Singapore, 2019.
- [35] F. W. Olver, D. W. Lozier, R. F. Boisvert, and C. W. Clark. *NIST handbook of mathematical functions*. Cambridge University Press, 2010.
- [36] J. M. Ortega. The Newton-Kantorovich theorem. *The American Mathematical Monthly*, 75(6):658–660, 1968.

- [37] M. Plum. Explicit H^2 -estimates and pointwise bounds for solutions of second-order elliptic boundary value problems. *Journal of Mathematical Analysis and Applications*, 165(1):36–61, 1992.
- [38] S. M. Rump. INTLAB - INTerval LABoratory. *Developments in Reliable Computing*, Kluwer Academic Publishers, Dordrecht, pp, pages 77–104, 1999.
- [39] S. M. Rump. Verification methods: Rigorous results using floating-point arithmetic. *Acta Numer*, 19:287–449, 2010.
- [40] T. J. Sullivan. *Introduction to uncertainty quantification*, volume 63. Springer, 2015.
- [41] J. Swift and P. C. Hohenberg. Hydrodynamic fluctuations at the convective instability. *Physical Review A*, 15(1):319, 1977.
- [42] R. Szwarc. Orthogonal polynomials and Banach algebras. *Inzell Lectures on Orthogonal Polynomials. Advances in the Theory of Special Functions and Orthogonal Polynomials*, Nova Science Publishers, 2:103–139, 2005.
- [43] L. N. Trefethen. *Approximation theory and approximation practice*, volume 128. Siam, 2013.
- [44] W. Tucker. *Validated numerics: a short introduction to rigorous computations*. Princeton University Press, 2011.
- [45] M. Urabe. Galerkin’s procedure for nonlinear periodic systems. *Archive for Rational Mechanics and Analysis*, 20(2):120–152, 1965.
- [46] J. B. van den Berg, M. Breden, J.-P. Lessard, and L. van Veen. Spontaneous periodic orbits in the Navier–Stokes flow. *Journal of Nonlinear Science*, 31(2):1–64, 2021.
- [47] J. B. van den Berg and J.-P. Lessard. Rigorous numerics in dynamics. *Notices Amer. Math. Soc.*, 62(9), 2015.
- [48] J. B. van den Berg, J.-P. Lessard, and K. Mischaikow. Global smooth solution curves using rigorous branch following. *Mathematics of computation*, 79(271):1565–1584, 2010.
- [49] J. B. Van den Berg, J.-P. Lessard, and E. Queirolo. Rigorous verification of hopf bifurcations via desingularization and continuation. *SIAM Journal on Applied Dynamical Systems*, 20(2):573–607, 2021.
- [50] J. B. van den Berg and E. Queirolo. A general framework for validated continuation of periodic orbits in systems of polynomial ODEs. *Journal of Computational Dynamics*, 8(1):59, 2021.
- [51] T. Wanner. Computer-assisted bifurcation diagram validation and applications in materials science. In *Proc. Sympos. Appl. Math. Rigorous Numerics in Dynamics.*, volume 74, pages 123–174. Amer. Math. Soc., 2018.
- [52] N. Wiener. The homogeneous chaos. *American Journal of Mathematics*, 60(4):897–936, 1938.
- [53] D. Xiu. Fast numerical methods for stochastic computations: a review. *Communications in computational physics*, 5(2-4):242–272, 2009.
- [54] D. Xiu and G. E. Karniadakis. The Wiener–Askey polynomial chaos for stochastic differential equations. *SIAM journal on scientific computing*, 24(2):619–644, 2002.
- [55] N. Yamamoto. A numerical verification method for solutions of boundary value problems with local uniqueness by Banach’s fixed-point theorem. *SIAM J. Numer. Anal.*, 35:2004–2013, 1998.