# THE 19TH INTERNATIONAL SYMPOSIUM ON SCIENTIFIC COMPUTING, COMPUTER ARITHMETIC, AND VERIFIED NUMERICAL COMPUTATIONS

Volume of abstract papers

# SCAN-2020

Organized by the Institute of Informatics, University of Szeged

September 13 – September 15, 2021
Szeged, Hungary

**Scientific Committee:**

G. Alefeld (Karlsruha, Germany)
A. Bauer (Ljubljana, Slovenia)
J. B. van den Berg (Amsterdam, the Netherlands)
G. F. Corliss (Milwaukee, USA)
T. Csendes (Szeged, Hungary)
R. B. Kearfott (Lafayette, USA)
V. Kreinovich (El Paso, USA)
J.-P. Lessard (Montreal, Canada)
W. Luther (Duisburg, Germany)
S. Markov (Sofia, Bulgaria)
G. Mayer (Rostock, Germany)
J.-M. Muller (Lyon, France)
M. Nakao (Tokyo, Japan)
T. Ogita (Tokyo, Japan)
S. Oishi (Tokyo, Japan)
K. Ozaki (Tokyo, Japan)
M. Plum (Karlsruhe, Germany)
A. Rauh (Brest, France)
N. Revol (Lyon, France)
J. Rohn (Prague, Czech Republic)
S. Rump (Hamburg, Germany/Tokyo, Japan)
S. Shary (Novosibirsk, Russia)
W. Tucker (Uppsala, Sweden)
W. Walter (Dresden, Germany)
J. Wolff von Gudenberg (Wuerzburg, Germany)
N. Yamamoto (Tokyo, Japan)

**Organizing Committee:**
Balázs Bánhelyi, University of Szeged, Hungary (chair)
Tibor Csendes, University of Szeged, Hungary
Boglárka G.-Tóth, University of Szeged, Hungary
Viktor Homolya, University of Szeged, Hungary
Tamás Vinkó, University of Szeged, Hungary
Dániel Zombori, University of Szeged, Hungary

**Address of the Organizing Committee**
c/o. Balázs Bánhelyi
University of Szeged, Institute of Informatics
H-6701 Szeged, P.O. Box 652, Hungary
Phone: +36 62 544 810
E-mail: `scan2020@inf.szte.hu`
URL: `http://www.inf.u-szeged.hu/scan2020/`

**Sponsors**

HU-MATHS-IN, the Hungarian Service Network for Mathematics in Industry and Innovations
http://hu-maths-in.hu/en/
http://hu-maths-in.hu/en/

University of Szeged, Institute of Informatics
https://www.inf.u-szeged.hu/en/about-us

City of Szeged Mayor's Office
https://www.szegedvaros.hu/
http://szegedtourism.hu/en/

# Preface

The 19th International Symposium on Scientific Computing, Computer Arithmetic and Verified Numerical Computation was planned to be organized in Szeged, Hungary in the year 2020. Due to the pandemic situation the Scientific Committee of SCAN decided to have the meeting in fully online version. More than 50 submissions arrived for the call, out of which 45 regular talks will be held and 4 plenary presentation. The plenary speakers will be Jay Mireles James, Fabienne Jézéquel and Kazuaki Tanaka, together with the Moore Prize winners Marko Lange and Siegfried Rump. The papers emerging from the talks can be submitted to the special issues of the journals Acta Cybernetica and Reliable Computing.

The Organizing Committee: Balázs Bánhelyi (chair), Tibor Csendes, Boglárka G.-Tóth, and Tamás Vinkó wishes a fruitful meeting and a memorable event for all participants.

# Contents

# Program

## Monday, September 13

**09:00 – 09:20**    Opening session

**09:20 – 10:00**    Plenary session (Kazuaki Tanaka)

**10:00 – 10:20**    Coffee break

**10:20 – 12:00**    Parallel sessions (2 * 4 talks)

**12:00 – 13:00**    Lunch break

**13:00 – 13:50**    Parallel sessions (2 * 2 talks)

**13:50 – 14:10**    Coffee break

**14:10 – 15:25**    Parallel sessions (2 * 3 talks)

**16:00 – 16:30**    Cultural Program

## Tuesday, September 14

**09:00 – 09:40**    Plenary session (Fabienne Jézéquel)

**09:40 – 10:00**    Coffee break

**10:00 – 11:40**    Parallel sessions (2 * 4 talks)

**11:40 – 13:00**    Lunch break

**13:00 – 14:15**    Parallel sessions (2 * 3 talks)

**14:15 – 14:35**    Coffee break

**14:35 – 15:25**    Parallel sessions (2 * 2 talks)

**18:00 – 18:30**    Cultural Program

## Wednesday, September 15

**09:00 – 09:50**    Parallel sessions (2 * 2 talks)

**09:50 – 10:10**    Coffee break

**10:10 – 11:25**    Session (1 * 3 talks)

**11:40 – 13:00**    Lunch break

**13:00 – 13:50**    Session (1 * 2 talks)

**13:50 – 14:10**    Coffee break

**14:10 – 14:50**    Plenary session (Jason Mirelles James)

**14:50 – 15:50**    Moore Prize laudation and talk

**16:00 – 16:30**    Cultural Program and closing session

# Detailed program

## Monday, September 13

09:20 - 10:00 – **Plenary session**, chair: Shin'ichi Oishi

- Kazuaki Tanaka: Verification of the sign of solutions to elliptic partial differential equations

10:00 - 10:20 – **Coffee break**

10:20 - 12:00 – **Parallel sessions (2 * 4 talks)**

---

**Session A: PDE**, chair: Shin'ichi Oishi

- Shin'ichi Oishi and Kouta Sekine: *Inverse Bifurcation Diagram Problem of Forced El Nino Equation*

- Yuuki Saito, Naoki Takamatsu, Shin'ichi Oishi and Kouta Sekine: *Inverse bifurcation diagram problem for delayed van der Pol-Duffing equation*

- Jonathan Wunderlich: *Computer-assisted Existence Proofs for Navier-Stokes Equations on an Unbounded Strip with Obstacle*

- Akitoshi Takayasu and Jean-Philippe Lessard: *A rigorous forward integration method for time-dependent PDEs*

**Session B: Arithmetic and Implementation**, chair: Shinya Miyajima

- Shinya Miyajima: *Verified bounds for matrix gamma function*

- Tomoaki Okayama and Shota Ogawa: *Improvement of selection formulas of mesh size and truncation number for the DE-Sinc approximation and its theoretical error bound*

- Naoya Yamanaka and Takeo Uramoto: *Verified algorithm for high-order partial derivatives using nilpotent matrix*

- Shinya Miyajima: *Computing enclosure for matrix real powers*

---

12:00 - 13:00 – **Lunch break**

13:00 - 13:50 – **Parallel sessions (2 * 2 talks)**

---

**Session A: Dynamic Systems**, chair: Tibor Csendes

- Alexander Morozov and Dmitry Reviznikov: *Kd-tree based adaptive interpolation algorithm for modeling dynamic systems with interval parameters*

- Anna Gierzkiewicz and Piotr Zgliczynski: *The Sharkovskii Theorem for multidimensional maps with attracting periodic orbits*

**Session B: Application and Software**, chair: Ekaterina Auer

- Ekaterina Auer and Wolfram Luther: *Assessing Uncertainty in Hereditary Risk Models for BRCA1/2 Related Cancer*

- Ekaterina Auer, Lorenz Gillner, Wolfram Luther and Andreas Rauh: *VERICOMP 2.0: Comparing and Recommending Verified IVP Solvers in a Flexible Way*

---

13:50 - 14:10 – **Coffee break**

14:10 - 15:25 – **Parallel sessions (2 * 3 talks)**

**Session A: Optimization**, chair: Ralph Baker Kearfott

- Mihály Csaba Markót: *Interval methods for packing problems on the sphere*

- Dun Liu and Ralph Baker Kearfott: *On Convexity Density and Difficulty of Global Optimization Problems*

- David Sanders and Valentin Churavy: *Interval constraint propagation and branch-and-bound-type methods on the GPU using Julia*

**Session B: Arithmetic and Implementation**, chair: Vladik Kreinovich

- Tamás Dózsa: Inverses of Rational Functions

- Mantas Mikaitis: A Trick for an Accurate $e^{(-|x|)}$ Function in Fixed-Point Arithmetics

- Vladik Kreinovich, Olga Kosheleva and Victor Selivanov: *Kinematic Metric Spaces Under Interval Uncertainty: Towards an Adequate Definition*

09:00 - 09:40 – **Plenary session**, chair: Nathalie Revol

- Fabienne Jézéquel: *Benefits of stochastic arithmetic in high performance simulations and arbitrary precision codes*

09:40 - 10:00 – **Coffee break**

10:00 - 11:40 – **Parallel sessions (2 * 4 talks)**

---

**Session A: Dynamic Systems**, chair: Nobito Yamamoto

- Taisei Asai, Kazuaki Tanaka, Kouta Sekine and Shin'ichi Oishi: *Computer-assisted analysis for bifurcation diagrams of the one-dimensional Henon equation*

- Naoki Takamatsu, Yuuki Saito, Shin'ichi Oishi and Kouta Sekine: *Numerical verification of existence for subharmonic solutions to delayed van der Pol-Duffing equation*

- Nobito Yamamoto and Koki Nitta: *A numerical verification method on time-global solutions of autonomous systems of complex functions*

- Kaname Matsue: *Rigorous numerics of blow-up separatrix in autonomous ODEs*

**Session B: Algorithms**, chair: Andreas Rauh

- Takehiko Kinoshita, Yoshitaka Watanabe and Mitsuhiro T. Nakao: *On some convergence properties for finite element approximations to the inverse of linear elliptic operators*

- Xuefeng Liu: *Rigorous maximum norm estimation for polynomial systems*

- Ekaterina Auer and Andreas Rauh: *Parameter Identification for Cooperative SOFC Models on the GPU*

- Marco De Angelis: *Linear-time algorithm for interval uncertainty propagation through the discrete Fourier transform*

---

11:40 - 13:00 – **Lunch break**

13:00 - 14:15 – **Parallel sessions (2 * 3 talks)**

---

**Session A: Dynamic Systems**, chair: Tibor Krisztin

- Andreas Rauh and Rachid Malti: *Quantification of Time-Domain Truncation Errors for the Reinitialization of Fractional Integrators*

- Tibor Krisztin and János Dudás: *Global stability for the three-dimensional logistic map*

- Ferenc Ágoston Bartha, Tibor Krisztin and Alexandra Vígh: *Stable periodic orbits for the Mackey–Glass equation*

**Session B: Algorithms**, chair: Elena Chausova

- Imre Fekete: *Local error estimation and step size control in adaptive linear multistep methods*

- Elena Chausova: *The inventory control problem for a supply chain with a mixed type of demand uncertainty*

- Auguste Bourgois, Amine Chaabouni, Andreas Rauh and Luc Jaulin: *Proving the stability of navigation cycles*

---

14:15 - 14:35 – **Coffee break**

14:35 - 15:25 – **Parallel sessions (2 * 2 talks)**

**Session A: Artificial Intelligence**, chair: Tibor Csendes

- Tibor Csendes, Nándor Balogh, Balázs Bánhelyi, Dániel Zombori, Richárd Tóth and István Megyeri: *Adversarial Example Free Zones for Specific Inputs and Neural Networks*

- Jonatan Contreras, Martine Ceberio and Vladik Kreinovich: *Why rectified linear neurons: a possible interval-based explanation*

**Session B: Optimization**, chair: Boglárka G.-Tóth

- Leocadio G. Casado, Boglárka G.-Tóth, Frédéric Messine and E.M.T. Hendrix: *Directional derivative bounds and border facets in simplicial B&B monotonicity tests*

- Bartlomiej Kubica: *How many constraints are satisfied? An approach to solving classification and regression problems*

**Wednesday, September 15**

09:00 - 09:50 – **Parallel sessions (2 * 2 talks)**

**Session A: Numerical Linear Algebra**, chair: Katsuhisa Ozaki

- Katsuhisa Ozaki: *Error-free transformation of matrix multiplication for multi-precision computations*

- Matyáš Lorenc: *B-matrices and their generalizations in the interval setting*

**Session B: Artificial Intelligence**, chair: Balázs Bánhelyi

- Dániel Zombori, Tamás Szabó, János Horváth, Attila Szász, Tibor Csendes and Balázs Bánhelyi: *Verification of artificial neural networks via Taylor models of INTLAB*

- Dániel Zombori, Tamás Szabó, János Horváth, Attila Szász, Tibor Csendes and Balázs Bánhelyi: *Verification of artificial neural networks via MIPVerify and SCIP*

09:50 - 10:10 – **Coffee break**

10:10 - 11:25 – **Session (1 * 3 talks)**

**Session A: Arithmetic and Implementation**, chair: Nathalie Revol

- Nathalie Revol: *Convergent Real Matrix Powers with Divergent Results in Interval Arithmetic*

- Massimiliano Fasi and Mantas Mikaitis: CPFloat: *A C library for emulating low-precision arithmetic*

- Sergey Kumkov: *Information Sets of Rebuilding Dependence Parameters for Criteria of Strong and Weak Compatibility under Heavy Two-Dimensional Measuring Errors*

11:40 - 13:00 – **Lunch break**

13:00 - 13:50 – **Session (1 * 2 talks)**

**Session A: Algorithms**, chair: Sergey Shary

- Sergey Shary: *Variability measures for estimates in interval data fitting*

- Vladik Kreinovich and Sergey Shary: *How probabilistic methods for data fitting deal with interval uncertainty: a more realistic analysis*

13:50 - 14:10 – **Coffee break**

14:10 - 14:50 – **Plenary session**, chair: Jean-Philippe Lessard

- Jason Mirelles James: *Computer assisted proofs for connecting orbits in infinite dimensions*

14:50 - 15:50 – **Moore Prize laudation and talk**, chair: Vladik Kreinovich

- Marko Lange and Siegfried M. Rump: *Verified inclusions of a nearest matrix of specified rank via a generalization of Wedin's $\sin(\theta)$ theorem*

# Benefits of stochastic arithmetic in high performance simulations and arbitrary precision codes

**Fabienne Jézéquel**
**Sorbonne University, CNRS, Paris, France**
**Université Panthéon-Assas, Paris, France**

**Keywords:** approximate GCD, BLAS, Discrete Stochastic Arithmetic, floating-point arithmetic, Newton method, numerical validation, polynomial roots, rounding errors

DSA (Discrete Stochastic Arithmetic) [Vig04] enables one to estimate the rounding error propagation which occurs with floating-point arithmetic. This probabilistic method uses a random rounding mode: at each elementary operation, the result is rounded up or down with the same probability. Therefore, the computer's deterministic arithmetic is replaced by a stochastic arithmetic, where each arithmetic operation is performed several times before the next one is executed. With DSA, temporary results that are actually numerical noise can be discarded and iterative algorithms can be stopped in an optimal way that does not rely on any parameter. DSA can be used to control the accuracy of programs in half, single, double and/or quadruple precision via the CADNA library [CAD, EBFJ15, GJP+18, JSHH21], and also in arbitrary precision via the SAM library [SAM, GJWZ11].

We present an algorithm that takes benefits of DSA to efficiently and accurately compute polynomial roots, in particular multiple roots. Thanks to a stochastic version of the polynomial GCD algorithm and the polynomial Euclidean division, the proposed algorithm provides a low-degree polynomial with single roots. Then Newton method can be applied to get fast and accurate approximations of the roots in arbitrary precision [GJQMS21].

Thanks to DSA, the accuracy estimation and the detection of numerical instabilities can be performed in parallel codes on CPU and on GPU [EBFJ15, EBFJ16, ELB+18]. However its performance overhead may be large compared with the standard floating-point operations. We show that with perturbed data it is possible to use standard floating-point arithmetic instead of DSA for the purpose of numerical validation. For instance, for codes including matrix multiplications, we can directly utilize the matrix multiplication routine (GEMM) of level-3 BLAS that is performed with standard floating-point arithmetic. Consequently, we can achieve a significant performance improvement by avoiding the performance overhead of DSA operations as well as by exploiting the speed of highly-optimized BLAS implementations [JGM+20].

## References

[CAD]  The CADNA library. http://cadna.lip6.fr.

[EBFJ15]  P. Eberhart, J. Brajard, P. Fortin, and F. Jézéquel. High performance numerical validation using stochastic arithmetic. *Reliable Computing*, 21:35–52, 2015.

[EBFJ16] P. Eberhart, J. Brajard, P. Fortin, and F. Jézéquel. Estimation of Round-off Errors in OpenMP Codes. In *12th International Workshop on OpenMP (IWOMP)*, volume 9903 of *LNCS*, pages 3–16. Springer, 2016.

[ELB+18] P. Eberhart, B. Landreau, J. Brajard, P. Fortin, and F. Jézéquel. Improving CADNA Performance on GPUs. In *19th IEEE International Workshop on Parallel and Distributed Scientific and Engineering Computing (PDSEC) in conjunction with the 32nd International Parallel and Distributed Processing Symposium (IPDPS)*, pages 1016–1025, Vancouver, Canada, 2018.

[GJP+18] S. Graillat, F. Jézéquel, R. Picot, F. Févotte, and B. Lathuilière. Numerical validation in quadruple precision using stochastic arithmetic. In *TNC'18. Trusted Numerical Computations*, volume 8 of *Kalpa Publications in Computing*, pages 38–53, 2018.

[GJQMS21] S. Graillat, F. Jézéquel, E. Queiros Martins, and M. Spyropoulos. Computing multiple roots of polynomials in stochastic arithmetic with Newton method and approximate GCD. working paper: https://hal.archives-ouvertes.fr/hal-03274453, June 2021.

[GJWZ11] S. Graillat, F. Jézéquel, S. Wang, and Y. Zhu. Stochastic arithmetic in multiprecision. *Mathematics in Computer Science*, 5(4):359–375, 2011.

[JGM+20] F. Jézéquel, S. Graillat, D. Mukunoki, T. Imamura, and R. Iakymchuk. Can we avoid rounding-error estimation in HPC codes and still get trustworthy results? In *NSV'20, 13th International Workshop on Numerical Software Verification*, Los Angeles, CA, United States, July 2020.

[JSHH21] F. Jézéquel, S. Sadat Hoseininasab, and T. Hilaire. Numerical validation of half precision simulations. In *1st Workshop on Code Quality and Security (CQS 2021) in conjunction with WorldCIST'21 (9th World Conference on Information Systems and Technologies)*, Terceira Island, Azores, Portugal, March 2021.

[SAM] The SAM library. http://www-pequan.lip6.fr/~jezequel/SAM.

[Vig04] J. Vignes. Discrete Stochastic Arithmetic for validating results of numerical software. *Numerical Algorithms*, 37(1–4):377–390, December 2004.

# Verified inclusions of a nearest matrix of specified rank via a generalization of Wedin's $\sin(\theta)$ theorem

**Marko Lange and Siegfried M. Rump**
**Hamburg University of Technology, Hamburg, Germany**

**Keywords:** verified error bounds, rank deficiency, $\sin(\theta)$ theorem, separation of singular vector subspaces, unitarily invariant norms, ill-posedness, INTLAB

In the scientific computing community it is well know that proving singularity of a matrix A is an ill-posed problem. Naturally, this implies that validating a certain rank deficiency of a matrix is out of the scope of verification methods. But where is the boundary between these ill-posed problems and closely related but well-posed problems? In this talk we immerse ourselves in this question. We present verification methods to solve two closely related problems. In this context, we further prove a generalization of Wedin's $\sin(\theta)$ theorem [1]. The corresponding result is the singular vector space counterpart to Davis and Kahan's generalized $\sin(\theta)$ theorem [2] for eigenspaces.

**References**

[1] Per-Åke Wedin. Perturbation bounds in connection with singular value decomposition. *BIT Numer. Math.*, 12(1):99–111, 1972.

[2] Chandler Davis and William M. Kahan. The rotation of eigenvectors by a perturbation. III. *SIAM J. Numer. Anal.*, 7(1):1–46, 1970.

# Computer assisted proofs for connecting orbits in infinite dimensions

**J.D. Mireles James**
**Florida Atlantic University, United States**

**Keywords:** Connecting orbits, Invariant Manifolds, Banach Spaces, Computer Assisted Proofs

Connecting orbits occupy a special place in dynamical systems theory, appearing as hypotheses in theorems forcing the existence of complex behavior. Important examples date back to the work of Poincare, who used a transverse connecting orbit to show the non-existence of conserved quantities in the gravitational three body problem. The homoclinic tangency of Shilnikov, and the Smale tangle theorem provide further illustrations of this principle. In infinite dimensions, recall that the proof of the Feigenbaum conjectures studies the intersection of invariant manifolds attached to fixed points of certain renormalization operators [1, 2].

In a given example problem it may be difficult to establish the existence of a connecting orbit using classical pen and paper techniques, and numerical simulations provide much needed insights into the dynamics of strongly nonlinear systems. Moreover, a number of researchers have devoted substantial energy into developing constructive, mathematically rigorous numerical techniques for proving the existence of connecting orbits. This is a fascinating area of research which involves invariant manifold theory, boundary value problems, implicit function theory, and computational mathematics.

In a typical problem, proving the existence of a connecting orbit involves a "tower" of arguments. One recipe is as follows.

- **Step 1:** Prove the existence of the underlying fixed point.

- **Step 2:** Prove theorems about the stability the fixed point: hyperbolicity, validated morse index, etcetera.

- **Step 3:** Obtain validated bounds on the stable/unstable manifolds attached to the fixed point.

- **Step 4:** Prove the existence of an orbit which starts on the unstable, and ends on the stable manifold.

- **Step 5:** (If possible) establish transversality of the connection.

Completing this outline for a finite dimensional problem is already a challenge, and a large number of works are devoted to this topic. A few examples include [3, 4, 5]. One could give a similar outline for connecting orbits in finite dimensional ordinary differential equations, and again many authors have worked on this problem. See for example the works of [6, 7, 8].

Applying an outline like the one above to an infinite dimensional problem, where even steps one and two are nontrivial, presents substantial challenges. In this talk I will focus on some infinite dimensional problems where the outline has been successfully implemented in full. I'll highlight infinite dimensional obstacles, and indicate how they are overcome using validated numerical methods. I'll focus on ideas for compact infinite dimensional discrete time dynamics developed in [9, 10], and discuss some recent applications to delay differential equations as in [11, 12, 13]. Time permitting, I will make some comments about how these ideas extend to partial differential equations.

# References

[1] Oscar E. Lanford, III. A computer-assisted proof of the Feigenbaum conjectures. *Bull. Amer. Math. Soc. (N.S.)*, 6(3):427–434, 1982.

[2] Jean-Pierre Eckmann and Peter Wittwer. A complete proof of the Feigenbaum conjectures. *J. Statist. Phys.*, 46(3-4):455–475, 1987.

[3] Arnold Neumaier and Thomas Rage. Rigorous chaos verification in discrete dynamical systems. *Phys. D*, 67(4):327–346, 1993.

[4] Marian Gidea and Piotr Zgliczyński. Covering relations for multidimensional dynamical systems. *J. Differential Equations*, 202(1):59–80, 2004.

[5] Piotr Zgliczyński. Covering relations, cone conditions and the stable manifold theorem. *J. Differential Equations*, 246(5):1774–1819, 2009.

[6] Gianni Arioli and Piotr Zgliczyński. The Hénon-Heiles Hamiltonian near the critical energy level—some rigorous results. *Nonlinearity*, 16(5):1833–1852, 2003.

[7] Daniel Wilczak and Piotr Zgliczy´nski. Heteroclinic connections between periodic orbits in planar restricted circular three-body problem—a computer assisted proof. *Comm. Math. Phys.*, 234(1):37–75, 2003.

[8] Gianni Arioli and Hans Koch. Existence and stability of traveling pulse solutions of the FitzHugh-Nagumo equation. *Nonlinear Anal.*, 113:51–70, 2015.

[9] J. D. Mireles James. Fourier-Taylor approximation of unstable manifolds for compact maps: numerical implementation and computer-assisted error bounds. *Found. Comput. Math.*, 17(6):1467–1523, 2017.

[10] R. de la Llave and J. D. Mireles James. Connecting orbits for compact infinite dimensional maps: computer assisted proofs of existence. *SIAM J. Appl. Dyn. Syst.*, 15(2):1268–1323, 2016.

[11] J. P. Lessard and J. D. Mireles James. A functional analytic approach to validated numerics for eigenvalues of delay equations. *J. Comput. Dyn.*, 7(1):123–158, 2020.

[12] Olivier Hénot, J.P. Lessard, and J. D. Mireles James. Parameterization of unstable manifolds for ddes: Formal series solutions and validated error bounds. *Journal of Dynamics and Differential Equations*, 2021.

[13] J.P. Lessard and J.D. Mireles James. A rigorous implicit $C^1$ chebyshev integrator for delay equations. *Journal of Dynamics and Differential Equations*, pages 1–30, 2020.

# Verification of the sign of solutions to elliptic partial differential equations

**Kazuaki Tanaka**
**Institute for Mathematical Science**
**Waseda University, Tokyo, Japan**

**Keywords:** Computer-assisted proof, Numerical verification, Elliptic differential equations, Sign-change structure, Positive solutions

The objective of this talk is the elliptic differential equation

$$-\Delta u(x) = f(u(x)), \;\; x \in \Omega, \tag{1}$$

where $\Omega \subset \mathbb{R}^N$ ($N = 1, 2, 3, \cdots$) is a bounded domain, and $f : \mathbb{R} \to \mathbb{R}$ is a given nonlinear map. We consider solutions of (1) with one of the three types of homogeneous boundary conditions: Dirichlet, Neumann, and mixed.

This talk surveys the methods proposed in [1, 2] that verify the information on the sign of solutions to (1). In [1], we proposed a method for verifying the positivity of a weak solution $u$ of (1) with the Dirichlet boundary condition assuming $H_0^1$-error estimation $\|u - \hat{u}\|_{H_0^1} \leq \rho$ given some numerical approximation $\hat{u}$ and an explicit error bound $\rho$. Subsequently, in [2], we rigorously analyzed the sign-change structure of solutions to (1) with one of the three boundary conditions and applying this to the Dirichlet problem of the Allen–Cahn equation.

Recent developments on this topic will be introduced.

## References

[1] K. TANAKA: Numerical verification method for positive solutions of elliptic problems *Journal of Computational and Applied Mathematics*, 370 (2020), 1112647.

[2] K. TANAKA: A posteriori verification for the sign-change structure of solutions of elliptic partial differential equations *Japan Journal of Industrial and Applied Mathematics*, (2021). https://doi.org/10.1007/s13160-021-00456-0

# Computer-assisted analysis for bifurcation diagrams of the one-dimensional Hénon equation

Taisei Asai[1], Kazuaki Tanaka[2], Kouta Sekine[3] and
Shin'ichi Oishi[4]

[1] Graduate School of Fundamental Science and Engineering,
Waseda University, Tokyo, Japan
[2] Institute for Mathematical Science, Waseda University, Tokyo, Japan
[3] Faculty of Information Networking for Innovation and Design,
Toyo University, Tokyo, Japan
[4] Faculty of Science and Engineering, Waseda University, Tokyo, Japan
captino@fuji.waseda.jp

We consider the one-dimensional Hénon equation which is the two-point boundary value problem

$$\begin{cases} -u'' = |x|^l u^p, \quad x \in (-1, 1), \\ u(-1) = u(1) = 0, \end{cases} \tag{1}$$

where $l \geq 0$, $2 \leq p < \infty$. The parameter $l$ is the potential index, and the parameter $p$ is the polytropic index. It is known that if $l = 0$, then there is no asymmetric positive solution, and if $l > 0$ is sufficiently large, then there are some asymmetric solutions. The importance of the Hénon equation has led to an active mathematical study on it over the last decade. In particular, its symmetry-breaking phenomena are attracting a lot of attention. S. Tanaka [1] proved that if $l(p - 1) \geq 4$, the Morse index of the positive least energy solution equals 1 and the Morse index of the positive symmetric solution equals 2, and hence the positive least energy solution is asymmetric and symmetry-breaking phenomena occur. It is also shown that if $l$ and $p$ are sufficiently small, then there is no positive asymmetric solution and the Morse index of the symmetric positive solution equals 1. However, still only sufficient conditions for symmetry–breaking bifurcation have been clarified, and the existence of multiple solutions near the bifurcation point and the structure of the bifurcation are not known completely.

The purpose of our study is to verify the existence of multiple solutions of (1) near the bifurcation point, and tracking the bifurcation diagrams by computer assistance.

Due to the variable coefficient $|x|^l$ in the problem (1), the solution $u$ has a singularity at $x = 0$. We will show the numerical verification method that follows such a singularity. In addition, the bifurcation structure with respect to $l$ will be discussed. Numerical examples of some symmrtric or asymmetric solutions and solution curves will be presented as follows.

## References

[1] S. TANAKA: Morse index and symmetry-breaking for positive solutions of one-dimensional Hénon type equations, *Journal of Differential Equations*, 255:7 (2013), 1709–1733.

Figure 1: Solution curves ($p = 3$).

# VERICOMP 2.0: Comparing and Recommending Verified IVP Solvers in a Flexible Way

Ekaterina Auer[1], Lorenz Gillner[1], Wolfram Luther[2], and Andreas Rauh[3]

[1] University of Applied Sciences Wismar, Wismar, Germany
[2] University of Duisburg-Essen, Duisburg, Germany
[3] Lab-STICC, ENSTA Bretagne, Brest, France
`ekaterina.auer@hs-wismar.de`

**Keywords:** verified IVP solvers, VERICOMP, recommender systems

Methods with result verification, for example, interval analysis [1], have been applied in engineering since the 1970s at the latest. They help not only to prove automatically that computer-obtained results are correct, but also to represent bounded uncertainty and propagate it through system models in an easy-to-understand, deterministic manner. After over half a century of research, many up-to-date libraries are available (and still emerging), implementing the concepts for a variety of programming languages and computer algebra systems such as C++, Python, Matlab, Julia. Regardless of the merits of such methods and their general accessibility, they are rarely used outside of the university context or cooperations.

Aside from the necessity to learn new material, an impediment on the way of a larger application of verified techniques in industrial engineering practice is the lack of information which of the available tools to choose for a given task. Sometimes, making an inappropriate choice for a program or approach can lead to too conservative results discouraging the use of the whole branch of methods. To improve the situation at least in one area, we have been working on a web-based platform VERICOMP [2] for promoting, comparing and recommending verified initial value problem software for ordinary differential equations (IVPS) for over a decade. Almost as a by-product, VERICOMP offers developers of new IVPS a possibility to compare their solvers with the established ones. Here, VERICOMP can be of use for facilitating such projects as ARCH-COMP (`cps-vo.org/group/ARCH/FriendlyCompetition`), a competition on verifying continuous and hybrid systems.

Differential equations are fundamental in many applied areas of science as a mathematical model for dynamic systems or processes, with techniques for comparing traditional, non-verified software available since the seventies [3] (`archimede.dm.uniba.it/~testset/testsetivpsolvers/`). Some of the challenging tasks on the way to highlight advantages of various tools are to develop

- a standard set of problems,

- a set of fair criteria and testing conditions,

- a means to easily incorporate new software into the comparison, and,

- last but not least, a means to present and visualize the gathered information.

The goal is to allow a user to obtain knowledge easily and to grasp it immediately. Moreover, the maintenance of these sets and means should be as flexible as possible facilitating, for example, entry of new (specialized or standard) problems or replacement of criteria if the need arises.

In this talk, we focus on the current version of VERICOMP available at `vericomp.fiw.hs-wismar.de` that enhances the previous one especially from the point of view of flexibility. Additionally, we demonstrate how ideas from the area of recommender systems can be employed to produce an automatic suggestion about the right tool to use for a given application

(cf. [4]). Moreover, we discuss in general what kind of metadata, data, quality criteria, metrics, and visualization are required to be able to compare and recommend IVPS. Finally, we give an outlook on the possibility of easier management and comparability of different solvers inside VERICOMP based on containerization. This task requires standardisation of data flow and IVPS interfaces (which are quite non-uniform at the moment).

**References**

[1] R. E. MOORE, R. B. KEARFOTT, M. J. CLOUD: *Introduction to Interval Analysis*, SIAM, Philadelphia, 2009.

[2] E. AUER AND A. RAUH, VERICOMP: A System to Compare and Assess Verified IVP Solvers, *Computing* 94, 163–172, 2012.

[3] T. E. HULL ET AL., Comparing Numerical Methods for Ordinary Differential Equations, *SINUM*, 9, 603–637, 1972.

[4] E. AUER AND W. LUTHER, Recommender techniques for software with result verification, UNCECOMP 2019, M. Papadrakakis, V. Papadopoulos, G. Stefanou (eds.), Crete, Greece, 24-26 June, 2019, 25p.

# Assessing Uncertainty in Hereditary Risk Models for *BRCA1/2* Related Cancer

Ekaterina Auer[1] and Wolfram Luther[2]

[1] University of Applied Sciences Wismar, Germany
[2] University of Duisburg-Essen, Germany
`ekaterina.auer@hs-wismar.de,wolfram.luther@uni-due.de`

**Keywords:** interval analysis, Dempster-Shafer theory, *BRCA1/2* related risk assessment and genetic counseling

Public websites making general recommendations about preventive services for major diseases are becoming increasingly important in the healthcare area. Government agencies, leading universities, and independent foundations [1, 2] provide them nationally or internationally for the purpose of informing the population about the available possibilities. Additionally, such websites allow individuals and their families to assess their risk of contracting a particular disease based on various factors such as their age, origin, or genetic predisposition. With the help of this risk assessment, concrete recommendations can be made for individual prevention and risk mitigation.

In this contribution, we focus on *BRCA1/2* related cancer. Complex prediction software relying on various kinds of mathematical stochastic models plays an essential role in the process of genetic counseling with the goal of determining either the gene mutation probability of a patient or their lifetime risk of cancer. Using personal information about individuals and their relatives as well as standardized case data from medical databases, a recommendation can be calculated and communicated via a suitable output interface (e.g., as a graph or a report) [3, 4, 5, 6]. Often, such recommendations can be augmented based on opinions of a multidisciplinary team of experts who collaborate in a final meta-study on issues concerning, for example, benefits/harms of counseling or clinical treatment for specific disease patterns [7]. Among the questions most meta-studies are raising is the issue of reliability of the generated recommendations, since they are considerably influenced by uncertainty. It is noted that existing verification and validation approaches usually account only for aleatory uncertainty and tend to disregard other kinds.

We consider genetic risk assessment and genetic counseling for *BRCA1/2* related breast cancer from the point of view of reliable uncertainty handling. First, we provide a short overview of existing risk models, software tools as well as family history interfaces and repositories. We show how missing or conflicting information on mutation probabilities can be improved using Dempster-Shafer theory. Based on multi-criteria binary decision trees and interval analysis, we combine the referral screening tool RST [8] designed to determine patients at risk of breast cancer with three further widely spread risk assessment tools for this purpose. The combined method has the advantage of assigning individuals to appropriate risk classes depending on their family history, taking into account epistemic uncertainty in the information about such factors as the age of onset in a relative, the degree of kinship or the relative's origin.

## References

[1] U.S. PREVENTIVE SERVICES TASK FORCE: https://www.uspreventiveservicestaskforce.org/uspstf/topic_search_results?topic_status=P

[2] HEALTH CARE DATA EXCHANGE: https://www.hl7.org/fhir/riskassessment.html

[3] J. A. CINTOLO-GONZALEZ ET AL.: Breast cancer risk models: a comprehensive overview of existing models, validation, and clinical applications. *Breast Cancer Res. Treat.*, 164 (2017), 263—284

[4] BRCAPRO WEB SERVICE CLIENT: http://bayesmendel.dfci.harvard.edu/risk/

[5] CENTRE FOR CANCER GENETIC EPIDEMIOLOGY: https://ccge.medschl.cam.ac.uk/boadicea/

[6] BRCA CHALLENGE PROJECT: https://brcaexchange.org/

[7] H. D. NELSON ET AL.: https://www.ncbi.nlm.nih.gov/books/NBK545867/pdf/Bookshelf_NBK545867.pdf, 2019

[8] D. K. OWENS ET AL.: Risk assessment, genetic counseling, and genetic testing for BRCA-related cancer, *JAMA*, 322 (2019), 652–665.

# Parameter Identification for Cooperative SOFC Models on the GPU

Ekaterina Auer[1] and Andreas Rauh[2]

[1] University of Applied Sciences Wismar, Wismar, Germany
[2] Lab-STICC, ENSTA Bretagne, Brest, France
ekaterina.auer@hs-wismar.de

**Keywords:** cooperative ODE systems, GPU, parameter identification, SOFC

Over the last decade, using graphic processing units (GPUs) for scientific computations has given a significant boost to such varied fields as neural networks, bioinformatics, medical imaging, or cryptography. Especially in control engineering, this paradigm can open up possibilities which remained unexplored because of the lack of cheap computing power.

One of such fields is modeling, parameter identification, simulation, and control of solid oxide fuel cells (SOFCs). Models for SOFC temperature are based on partial-differential equations, which are usually discretized wrt. space and time into algebraic equations. A disadvantage of this technique is the lack of flexibility: Only stationary states of SOFC systems can be simulated this way, which is unsuitable for control. Using the same finite volume method without discretization in time, it is possible to arrive at *dynamic* system models consisting of a set of ordinary differential equations which can be shown to be cooperative [2]. A challenge here is to identify a large number of parameters based on uncertain measured values from the SOFC test rig.

Our special focus is on the property of cooperativity [3]. For a cooperative system with uncertain but bounded parameters, two bracketing systems with crisp parameters can be defined to catch the bulk of uncertainty. (These bounding systems might be coupled with each other if lower and upper interval bounds for the system parameters to be identified appear simultaneously in an equation.) A brute force approach using the GPU would be to partition the parameter search space and evaluate the system over the subintervals in parallel, eliminating the regions inconsistent with available measurements. To avoid a prohibitively large number of system evaluations due to naive interval multi-sectioning schemes, we propose to employ a set of additional simple consistency tests. They represent knowledge from physics such as non-negativity and strict monotonicity of heat capacities and reaction enthalpies that occur multiple times in the dynamic SOFC model (common subexpressions). Such constraints in combination with inequalities reflecting physically meaningful temporal variation rates of the measured SOFC stack temperature help to reduce the number of parameter subintervals. This preprocessing stage is carried out prior to the parallelized evaluation of the system model.

In this contribution, we extend the GPU-based technique described in [1] to deal with the reaction phase of a dynamic SOFC model. Considering this phase separately from the electrochemically idle heating phase represents a further possibility to cope with the high dimensionality of the parameter space. Finally, we show how using Bernstein polynomials helps to reduce the amount of data for controlled and measured system inputs and outputs.

## References

[1] Auer, E., Rauh, A. and J. Kersten. Experiments-Based Parameter Identification on the GPU for Cooperative Systems. In *J. of Comp. and App. Math.*, 371, 2020.

[2] Ifqir, S., Rauh, A., Kersten, J., Ichalal, D., Ait-Oufroukh, N. and S. Mammar. Interval Observer-Based Controller Design for Systems with State Constraints: Application to SOFC

Stacks. In *24th International Conference on Methods and Models in Automation and Robotics (MMAR)*, pages 372–377, 2019.

[3] H.L. Smith. Monotone Dynamical Systems: An Introduction to the Theory of Competitive and Cooperative Systems. *Mathematical Surveys and Monographs, American Mathematical Soc.*, 41, 1995.

# Stable periodic orbits for the Mackey?Glass equation

Ferenc Ágoston Bartha Tibor Krisztin, and Alexandra Vígh

Bolyai Institute, University of Szeged, Szeged, Hungary
`barfer@math.u-szeged.hu`

We study the classical Mackey–Glass delay differential equation

$$x'(t) = -ax(t) + bf_n(x(t-1))$$

where $a, b, n$ are positive reals, and $f_n(\xi) = \xi/[1 + \xi^n]$ for $\xi \geq 0$. As a limiting ($n \to \infty$) case we also consider the discontinuous equation

$$x'(t) = -ax(t) + bf(x(t-1))$$

where $f(\xi) = \xi$ for $\xi \in [0, 1)$, $f(1) = 1/2$, and $f(\xi) = 0$ for $\xi > 1$. First, for certain parameter values $b > a > 0$, an orbitally asymptotically stable periodic orbit is constructed for the discontinuous equation. Then it is shown that for large values of $n$, and with the same parameters $a, b$, the Mackey–Glass equation also has an orbitally asymptotically stable periodic orbit near to the periodic orbit of the discontinuous equation.

Although the obtained periodic orbits are stable, their projections $\mathbb{R} \ni t \mapsto (x(t), (x(t-1))) \in \mathbb{R}^2$ can be complicated.

# Proving the stability of navigation cycles

Auguste Bourgois[1,2], Amine Chaabouni[3], Andreas Rauh[1] and Luc Jaulin[1]

[1] Ensta-Bretagne, Lab-STICC, Brest, France
[2] FORSSEA, Paris, France
[3] Ecole polytechnique, Palaiseau, France

**Keywords:** Interval analysis, Poincaré sections, Hybrid systems, Stable cycles

When navigating in an unstructured environment a robot may not be able to geolocalize due to the absence of marks or to the fact that it has no access to the GPS. Now, using ancestral navigation methods, it is possible to move without getting lost. The principle is to find a discrete sequence of control such that the robot converges to a limit cycle [2]. The transition from one discrete state to another is triggered by an event such as a timer has reached a given time value or the robot has crossed an isobath line. Consider for instance the robot described by the state equations

$$\begin{cases} \dot{x}_1 &= \cos x_3 \\ \dot{x}_2 &= \sin x_3 \\ \dot{x}_3 &= u, \end{cases}$$

where $(x_1, x_2)$ is the position of the robot and $x_3$ is its heading. The heading control has the form $u = \sin(\bar{\psi} - x_3)$, where the desired heading $\bar{\psi}$ obeys to the automaton (or Petri net) of Figure 1. The variable $q \in \{0, 1, 2\}$ is discrete and $c$ is a continuous clock initialized to 0 each time $q$ changes.



Figure 1: Automaton deciding the desired heading

Using a simulation, we observe (see Figure 2) that the state $\mathbf{x}(t)$ converges to a stable limit cycle. The left figure is a simulation for $t \in [0, 18]$, and $\mathbf{x}(0) = (-3, 1, -1)$. The figure in the center shows the same simulation with $t \in [0, \infty]$. In the right, the robot is now represented by points in order to see the cycle. The colors blue, red, greed are associated to $q = 0, 1, 2$, respectively.

The goal of this paper is to provide a rigorous method to show that such a cycle is stable. It will combine interval analysis [1] and Poincaré section concepts [3].

**References**

[1]  R. E. MOORE (1966). *Interval Analysis*. Prentice-Hall.

[2]  A. BOURGOIS, L. JAULIN (2020). *Interval centred form for proving stability of non-linear discrete-time system*, SNR'2020 , Vienna, 2020.

Figure 2: The robot converges to a stable limit cycle in the $(x_1, x_2, x_3)$ space. The frame box is $[-4, 4] \times [-1, 7]$

[3] T. KAPELA, D. WILCZAK, P. ZGLICZYŃSKI (2021), *Recent advances in rigorous computation of Poincaré maps*, arXiv

# Directional derivative bounds and border facets in simplicial B&B monotonicity tests [1]

L.G. Casado[1], B.G.- Tóth[2], F. Messine[3], and E.M.T. Hendrix[4]

[1]University of Almería (CeiA3), Almería, Spain, [2]University of Szeged, Szeged, Hungary, [3]University of Toulouse, Toulouse, France, and [4]Universidad de Málaga, Málaga, Spain
leo@ual.es

**Keywords:** simplex, branch and bound, derivative, interval arithmetic

The concept of monotonicity tests and dimension reduction to a facet of an interval partition set are well elaborated concepts in IA branch and bound algorithms. The extension of those concepts to simplicial branch and bound implies several challenges.

- For a simplicial feasible set, it is much harder to verify a face of the partition set is on the boundary of the feasible set.
- For a simplicial feasible set, when dealing with reduced dimension partition sets, the monotonicity test relies on directional derivative bounds.
- The border faces to reduce to, depend on directions that are more or less perpendicular to facets.

We consider the minimization of a continuously differential function $f : \mathbb{R}^n \to \mathbb{R}$. Initially, the partition sets may consist of $n-$simplices, so they have $n + 1$ affinely independent vertices $\mathcal{V} := \{v_0, \ldots, v_n\} \subset \mathbb{R}^n$. However, after dimension reduction due to a monotonicity test, we are dealing with $m-$simplices on the boundary of the feasible set. Given a simplicial partition set $S$, a facet $F$ of $S$ is called border if there exists a face $\mathcal{F}$ of the feasible set of the same dimension, such that $F \subseteq \mathcal{F}$.

Given a simplicial subset $S$, its centroid $c = \frac{1}{n} \sum_{i=0}^{n} v_i$, $\boldsymbol{f}$ the interval extension of $f$, $\boldsymbol{f}'$ the interval extension of its gradient, $\Box S$ the interval hull of a simplex $S$, and a facet $F_j$ obtained by removing vertex $v_j$, we define the directional vector $d_j = (v_j - c)$ from the centroid of $S$ to $v_j$ and $\boldsymbol{g}_j = d_j^T \boldsymbol{f}'(\Box S)$ the directional interval gradient from $c$ to $v_j$. Basically, if $\underline{\boldsymbol{g}}_j > 0$, the minimum over $S$ is on $F_j$ as depicted in blue arrow in Figure 1.

As example, consider $f = x_1^2 + x_2^2$ on simplex $\text{conv}(\{v_0, v_1, v_2\})$ with $v_0 = (0.1, 0.9), v_1 = (0.8, 0.2)$ and $v_2 = (1, 1)$ where the minimum point can be found in $(0.5, 0.5)$. Blue, red, and black arrows correspond to $\underline{\boldsymbol{g}}_j > 0$, $\overline{\boldsymbol{g}}_j < 0$ and $0 \in \boldsymbol{g}_j$, respectively. In a dimension reduction, we can drop vertex $v_1$ and go further with facet $F_1$ as partition set. From there we proceed with bisection combined with rejection due to monotonicity when the reached line segment partition sets have not border facets. Notice that a red arrow in one direction has a blue counterpart in the opposite direction and vice versa.

In the presentation and paper, we will discuss the complete algorithm and all variants related to implement the monotonicity tests and border checking for simplicial partition sets.

## References

[1] B. G.-Tóth, L. G. Casado, E. M. T. Hendrix, and F. Messine. On new methods to construct lower bounds in simplicial branch and bound based on Interval Arithmetic. *Journal of Global Optimization*, Forthcoming article DOI: https://doi.org/10.1007/s10898-021-01053-8.

---

18

Figure 1: Example minimizing $f = x_1^2 + x_2^2$ over a simplicial feasible set

[2] E. M. T. Hendrix, B. G.-Tóth, F. Messine and L. G. Casado. On derivative based bounding for simplicial branch and bound. *RAIRO-Operations Research*, vol 55, n.3, pp. 2023-2034, DOI: https://doi.org/10.1051/ro/2021081.

---

**Algorithm 1** Monotonicity-Reduction test $(S, 0 \in \boldsymbol{f}'(\square S))$

---

1: **if** $S$ is non-reduced **then**
2:     **if** $S$ has not border facets **then**
3:         **return** Reject $S$
4:     **end if**
5:     Evaluate $\boldsymbol{g}_j$ for border facets $F_j$
6:     **if** $\exists \, \underline{\boldsymbol{g}}_j > 0$ **then**
7:         **return** Reduce $S$ to just one face by removing all $v_j$ at once
8:     **end if**
9:     **if** $\exists \, 0 \in \boldsymbol{g}_j$ **then**
10:        **return** Reduce $S$ to each facet $F_j$
11:     **end if**
12:     **return** Error                      $\triangleright$ All $\underline{\boldsymbol{g}}_j > 0$ is not possible
13: **end if**
14: Evaluate $\boldsymbol{g}_j$ for all facets $F_j$             $\triangleright$ $S$ is a reduced simplex
15: **if** $S$ has not border facets **then**
16:     **if** $\exists \, \overline{\boldsymbol{g}_j} < 0$ **Or** $\exists \, \boldsymbol{g}_j > 0$ **then**
17:         **return** Reject $S$
18:     **end if**
19:     Return Divide $S$                     $\triangleright$ Only $0 \in \boldsymbol{g}_j$ occurs
20: **end if**
21: **if** $\exists \, \underline{\boldsymbol{g}}_j < 0$ **and** $F_j$ is (are) border **then**
22:     **return** Reduce $S$ to just one face by removing all $v_j$ at once
23: **end if**
24: **if** $\exists \, \overline{\boldsymbol{g}_j} < 0$ **then**
25:     **if** $\exists \, 0 \in \boldsymbol{g}_k$ **then**
26:         **return** Reduce $S$ to each border facet $F_k$
27:     **end if**
28:     **return** Error                   $\triangleright$ All $\underline{\boldsymbol{g}}_j > 0$ is not possible
29: **end if**
30: **return** Divide $S$                        $\triangleright$ Only $0 \in \boldsymbol{g}_j$ occurs

---

# The inventory control problem for a supply chain with a mixed type of demand uncertainty

Elena Chausova[1],

[1] Tomsk State University, Tomsk, Russia
chauev@mail.ru

**Keywords:** inventory control, supply chain, network model, interval-stochastic uncertainty, model predictive control

The paper develops the results of [1] for the case of mixed uncertainty with not only interval assigned variables but also random inputs. We use the method of model predictive control to derive the control strategy. This method is widely applied in the practice of control and allows solving complex control problems (for example, [2]). To handle with interval uncertainty we use the interval analysis tools and operate according to the interval analysis theory.

We consider a dynamic inventory control system with a network structure (supply chain). The evolution of the network is described by the equation:

$$x(k+1) = x(k) + Bu(k) + C(d(k) + w(k)), \quad k = 0, 1, 2, \ldots. \tag{1}$$

Here $x(k) \in \mathbb{R}^n$ is the system state whose components represent storage levels in the network nodes at the moment $k$; $u(k) \in \mathbb{R}^m$ means control representing controllable resource flows between the network nodes at the moment $k$; $d(k), w(k) \in \mathbb{R}^l$ are uncertain noncontrollable flows at the moment $k$ describing demand in the network nodes; the matrices $B \in \mathbb{R}^{n \times m}$ and $C \in \mathbb{R}^{n \times l}$ describe the network structure.

Interval uncertainty in the system is represented by the vector $d(k)$. We only know that $d(k)$ takes its values within an assigned interval but otherwise unknown:

$$d(k) \in \mathbf{D}, \quad k = 0, 1, 2, \ldots,$$

where $\mathbf{D} \in \mathbb{IR}^l$, $\mathbf{D} \geq 0$, $\mathbb{IR}$ is the set of the real intervals $\mathbf{x} = [\underline{x}, \overline{x}]$, $\underline{x} \leq \overline{x}$, $\underline{x}, \overline{x} \in \mathbb{R}$. The uncertain vector $w(k)$ is the vector of white noises with zero mean and covariance matrix $E\{w(k)w^T(k)\} = W$. This is a stochastic uncertainty.

We assume as well, that both expected storage levels and controls must be nonnegative and bounded

$$\mathsf{E}\{x(k+1)|x(k)\} \in \mathbf{X}, \quad k = 0, 1, 2, \ldots, \tag{2}$$
$$u(k) \in \mathbf{U}, \quad k = 0, 1, 2, \ldots, \tag{3}$$

where $\mathsf{E}\{\cdot|\cdot\}$ denotes the conditional expectation; $\mathbf{X} \in \mathbb{IR}^n$, $\mathbf{X} = \left[0, \overline{X}\right]$; $\mathbf{U} \in \mathbb{IR}^q$, $\mathbf{U} = \left[0, \overline{U}\right]$ ($\underline{X} = 0$ means that backlogged demand is not desirable in the system).

For the system (1) we synthesize the strategies with a predictive model according to the following rule. At each time $k$ we minimize the criterion with a sliding control horizon:

$$J(k+p|k) = \mathsf{E}\Bigg\{ \sum_{i=1}^{p} \Big( (x(k+i|k) - x_0)^T Q(x(k+i|k) - x_0) -$$

$$-Q_1(x(k+i|k) - x_0) + u(k+i-1|k)^T R u(k+i-1|k) \Big) \Big| x(k) \Bigg\}.$$

on the trajectories of the system (1) over a sequence of predictive controls $u(k|k), u(k+1|k)$, $\ldots, u(k+p-1|k)$, which depend on system state at moment $k$, under the constraints (2), (3), where $Q \in \mathbb{R}^{n \times n}, Q = Q^T, Q_1 \in \mathbb{R}^{1 \times n}, R \in \mathbb{R}^{m \times m}, R = R^T$, are the given weight matrices; $x(k+i|k)$ is the state of the system at time $k+i$ derived at time $k$ by applying the sequence of predictive controls $u(k|k), u(k+1|k), \ldots, u(k+i-1|k)$ on the system (1), and $x(k|k) = x(k)$ is the state of the system measured at time $k$; $x_0$ is the target storage level; $p$ is the prediction horizon. We reduce the problem to a quadratic programming problem with constraints for whose solution the efficient techniques exist.

Although more than one control move is calculated, only the first one is used. So, we get the feedback control $u(k) = u(k|k)$ as a function of state $x(k)$. Then the state $x(k+1)$ is measured and the optimization is repeated at the next sampling time $k+1$, etc. As a result, we have the feedback inventory control strategy with high service level.

The developed results are illustrated with a numerical example.

## References

[1] E. V. CHAUSOVA: Dynamic Network Inventory Control Model with Interval Nonstationary Demand Uncertainty // *Numerical Algorithms*, 37 (2004), 71–84.

[2] V. V. DOMBROVSKII, E. V. CHAUSOVA: Model Predictive Control for Linear Systems with Interval and Stochastic Uncertainties // *Reliable computing*, 19 (2014), 351–360.

# Why rectified linear neurons: a possible interval-based explanation

Jonatan Contreras, Martine Ceberio, and Vladik Kreinovich

University of Texas at El Paso, USA
`jmcontreras2@utep.edu, mceberio@utep.edu, vladik@utep.edu`

**Keywords:** neural networks, rectified linear neurons, interval uncertainty

**What are rectified linear neurons.** At present, the most efficient machine learning techniques are deep neural networks; see, e.g., [1]. In general, in a neural network, a signal repeatedly undergoes two types of transformations: linear combination, and a non-linear transformation of each value $v \to s(v)$. The corresponding nonlinear function $s(v)$ is called an *activation function*. In deep neural networks, most nonlinear layers use the function $s(v) = \max(0, v)$ which is called the *rectified linear (ReLU) activation function*.

*Comment.* Taking into account that we also have linear layers, what can be represented by the ReLU function can also be represented if we use any piece-wise linear activation function.

**Why rectified linear neurons?** Empirically, rectified linear activation functions work the best. There are some partial explanations for this empirical success (see, e.g., [2]), but none of them is fully convincing, so yet another explanation is always welcome.

**What we do.** In this paper, we analyze this why-question from the viewpoint of uncertainty propagation, and we show that some reasonable uncertainty-related arguments indeed lead to a possible (partial) explanation.

**Need to take interval uncertainty into account.** The activation function transforms the input $v$ into the output $y = s(v)$. The input $v$ comes either directly from measurements, or from processing measurement results. Measurements are never absolutely accurate: the measurement result $\widetilde{v}$ is, in general, different from the actual (unknown) value of the quantity $v$. In many practical situations, all we know about the measurement error $\Delta v \stackrel{\text{def}}{=} \widetilde{v} - v$ is the upper bound $\Delta$ on its absolute value: $|\widetilde{v} - v| \leq \Delta$. In this case, possible values of $v$ form an interval $[\widetilde{v} - \Delta, \widetilde{v} + \Delta]$.

**First natural requirement.** A first natural requirement is that the output $y$ should not be too much affected by inaccuracy with which we know the input. Ideally, this inaccuracy should not increases after data processing, i.e., we should have $|s(\widetilde{v}) - s(v)| \leq |\widetilde{v} - v|$. In mathematical terms, this means that the function $s(v)$ should be 1-Lipschitz – so its derivative (or generalized derivative) should be limited by 1: $|s'(v)| \leq 1$.

**Second natural requirement: first try.** On the other hand, we do not want to lose information about the signal, so we must be able to reconstruct the input signal from the output as accurately as possible. This idea can be naturally described as $|\widetilde{v} - v| \leq |s(\widetilde{v}) - s(v)|$. Together with the first requirement, this means that $|\widetilde{v} - v| = |s(\widetilde{v}) - s(v)|$. Taking into account that we want to uniquely reconstruct $v$ from $s(v)$, this implies that either $s(v) = v$ or $s(v) = -v$. However, we wanted the function $s(v)$ to be nonlinear, since otherwise we will only be able to represent linear dependencies.

**Second natural requirement made realistic.** Since we cannot accurately reconstruct the input $v$ from $s(v)$, a natural idea is to use *two* activation functions $s_1(v)$ and $s_2(v)$ so that for each $v$, we can accurately reconstruct the signal from at least one of the two outputs $s_i(v)$.

**What we can conclude.** A natural conclusion is that for (almost) all values $v$, we must have either $|s_1'(v)| = 1$ or $|s_2'(v)| = 1$. In other words, the real line – the set of all possible values $v$ – is divided into two subsets: on one of them $s_1(v) = \pm v$, on another one $s_2(v) = \pm v$.

22

**Third natural requitement.** Since many real-life dependencies are linear, it is desirable to require that a linear function – e.g., the function $f(v) = v$ – can be represented as a linear combination of the two activation functions, i.e., that $v = c_0 + c_1 \cdot s_1(v) + c_2 \cdot s_2(v)$.

**What we can now conclude.** For values $v$ for which $s_1(v) = \pm v$, we conclude that $s_2(v) = c_2^{-1} \cdot (v - c_0 - c_1 \cdot s_1(v))$ is linear. Similarly, for remaining values $v$ – for which $s_2(v) = \pm v$ – we can conclude that the function $s_1(v)$ is linear. Thus, both activation functions $s_1(v)$ and $s_2(v)$ are piecewise linear – which is exactly what we wanted to explain.

## References

[1] I. GOODFELLOW, Y. BENGIO, A. COURVILLE: *Deep Learning*, MIT Press, Cambridge, Massachusetts, 2016.

[2] V. KREINOVICH, O. KOSHELEVA: Optimization under uncertainty explains empirical success of deep learning heuristics, In: P. PARDALOS, V. RASSKAZOVA, M. N. VRAHATIS (EDS.): *Black Box Optimization, Machine Learning and No-Free Lunch Theorems*, Springer, Cham, Switzerland, 2021, pp. 195–220.

# Adversarial Example Free Zones for Specific Inputs and Neural Networks

Tibor Csendes[1], Nándor Balogh[2], Balázs Bánhelyi[1], Dániel Zombori[1], Richárd Tóth[1], and István Megyeri[1]

[1] University of Szeged, Szeged, Hungary
[2] Redink Ltd., Szeged, Hungary
csendes@inf.szte.hu

Recent machine learning models are highly sensitive to adversarial input perturbation. That is, an attacker may easily mislead a well-performing image classification system by altering some pixels. However, proving that a network will have correct output when changing some regions of the images, is quite challenging – mostly due to the high dimensionality and/or the nonlinearity of the problem. Because of this, only a few works targeted this problem, and some of these verification tools are not reliable [4]. Although there are an increasing number of studies on this field, reliable robustness evaluation is still an open issue. We will present interval arithmetic based algorithms to provide adversarial example free image patches for trained artificial neural networks [2]. The method is based on an earlier interval technique to bound level sets of parameter estimation problems [1].

The obtained results are illustrated on Figure 1 for some of the studied images from the MNIST dataset. The calculated number of pixels to be changed arbitrarily were between 88 and 190 (compare it with the $28 \times 28 = 784$ pixels in the images). The combined running time for the second round of 10 test images was 1971.87 second, i.e. closely half an hour.



Figure 1: Original pictures and proven rectangles where we can change *everything* without having an adversarial example.

We are still in the phase when we explore the capabilities of interval arithmetic based algorithms for describing the sensitivity of trained natural neural networks to changes in object to be classified, but we find our present results encouraging enough to continue our research project.

**References**

[1] T. CSENDES: An interval method for bounding level sets of parameter estimation problems, *Computing*, 41 (1989), 75–86.

[2] T. CSENDES, N. BALOGH, B. BÁNHELYI, D. ZOMBORI, R. TÓTH, I. MEGYERI: Adversarial Example Free Zones for Specific Inputs and Neural Networks. *Proc. ICAI*, Eger, Hungary, 2020, 76–84, http://ceur-ws.org/Vol-2650/paper9.pdf

[3] D. ZOMBORI, B. BÁNHELYI, T. CSENDES, I. MEGYERI, M. JELASITY: Fooling a Complete Neural Network Verifier. *Proc. ICLR*, 2021, https://openreview.net/pdf?id=4IwieFS44l

# Linear-time algorithm for interval uncertainty propagation through the Fourier transform

Marco De Angelis

University of Liverpool, United Kingdom
`marco.de-angelis@liverpool.ac.uk`

**Keywords:** Interval extension, United extension, Convex hull, Zonotopes

The use of the Fourier transform in practical applications is nearly ubiquitous. Noisy signals can be decomposed into simple harmonics to filter them, compress them or further study their features. However, real data is often subject to uncertainties due to various reasons, such as sensor inaccuracy or measurement errors. When only error bounds can be known, no assumptions on dependency nor on the error distribution should be made.

In [1] the authors have developed an algorithm to propagate interval uncertainty through the discrete Fourier transform to obtain the exact bounds on the amplitude spectrum in polynomial time $O(n^2 \log n)$. In this paper we show that it is possible to obtain the same results but in linear time $O(n)$. In [1] , it has been shown that the interval propagation has to map the boundary of the united set onto the Fourier domain in order for the dependency to be fully tracked. It has also been shown that the united set that encodes the functional dependence between interval real and imaginary components of the Fourier signal is always a disk. Furthermore, the boundary of this disk is always inscribed in the box arising from the transform's natural extension. This is because such natural extension is also the optimal one, due to the absence of repeating variables. With this premise, the linear time algorithm follows immediately; from the box that results from evaluating the natural extension, the united set can be derived by computing the disk inscribed in the obtained box.

This result can have a significant impact on applications as measurement error can be propagated through the Fourier transform with nearly no additional computational cost and without artificial distributional assumptions.

## References

[1] M. De Angelis, M. Behrendt, L. Comerford, Y. Zhang, M. Beer: *Forward interval propagation through the discrete Fourier transform*, In: 9th International Workshop on Reliable Engineering Computing, (2020), 39–52.

# Inverses of Rational Functions

Tamás Dózsa[1]

[1] Eötvös Loránd University, Faculty of Informatics, Department of Numerical Analysis, Budapest, Hungary

`dotuaai@inf.elte.hu`

**Keywords:** Rational functions, Blaschke-products, Newton-iteration

In this presentation we consider numerical methods to produce the global inverses of analytic functions on the closed disk. The inverses coincide with the solutions $z$ of the implicit equation

$$f(z) = w = Re^{i\alpha}, \ (z \in \overline{\mathbb{D}}, R \in \mathbb{R}, \alpha \in [-\pi, \pi]). \tag{1}$$

The investigated problem relates closely to the identification of transfer functions [1, 4], the construction of discrete orthogonal and biorthogonal rational systems [3] and the description of electrostatic equilibrium [5].

We begin by restricting the rational function $f$ to the Torus, (or in other words considering its Nyquist-diagram) and determining the $z = e^{it_j}, \ (j = 1, \ldots, N \in \mathbb{N})$ solutions of (1). The number of solutions $N$ depends on the number of zeros and poles of $f$.

We refer to the zeros of $f'$ and their image eastablished by $f$ as critical points and critical values, respectively. One can prove that if the line segment $S_\alpha := \{re^{i\alpha} : 0 \leq r \leq R\}$ does not contain any critical values, then a continuous solution curve $z_j : [0, R] \to \overline{\mathbb{D}} \ (j = 1, \ldots, N)$, for which $f(z_j(r)) = r \cdot e^{i\alpha}$ and $z_j(R) = z_j = e^{it_j} \ (t_j \in [-\pi, \pi])$ exists uniquely. Furthermore once $z_j = e^{it_j}$ is known, any inverse point $z_j(r)$ on this curve can be found by applying a finite number of Newton-iterations.

In the presentation we pay special attention to the cases when $f$ is a polynomial or a Blaschke-product. We give several variations of the proposed algorithm including ones which do not depend explicitly on the derivatives of $f$. In addition, we aim to make the proposed algorithms more robust by the introduction of adaptive step sizes. We also discuss previous results such as the interpretation of the inverse branches of Blaschke-products, which can be regarded as a generalization of the $n$-th root function, as well as considering possible physical interpretations of the inverses [2]. As a potential new application, we show how the results can be used to find and interpret the roots of classical orthogonal polynomials.

## References

[1] BOKOR J., SCHIPP F., SOUMELIDIS A. *Pole structure estimation from Laguerre representation using hyperbolic metric on the unite disc.* 50th IEEE Conf. on Decision and Control an European Control Conf., Orlando, Florida, December 12-15, 2136-2141 (2011)

[2] DÓZSA T., SCHIPP F., *A generalization of the root function* (2020), Manuscript submitted for publication.

[3] FRIDLI S., SCHIPP F. *Discrete rational biorthogonal systems on the disc.* Annales Univ. Sci. Budapest, Sect. Comp. 50 (2020) 127-134.

[4] HEUBERGER P., VAN DEN HOF P., WAHLBERG B. *Modelling and Identification with Rational Orrthogonal Basis Functions.* Springer, 2005.

[5] PAP M., SCHIPP F. *Equilibrium conditions for the Malmquist-Takanaka systems.* Acta Sci. Math.(Szeged) 81 (2015), 169-182.

# Global stability for the three-dimensional logistic map

János Dudás, and Tibor Krisztin

Bolyai Institute, University of Szeged, Szeged, Hungary

`krisztin@math.u-szeged.hu`

For the delayed logistic equation $x_{n+1} = ax_n(1 - x_{n-2})$ it is well known that the nontrivial fixed point is locally stable for $1 < a \leq \left(\sqrt{5} + 1\right)/2$, and unstable for $a > \left(\sqrt{5} + 1\right)/2$. We prove that for $1 < a \leq \left(\sqrt{5} + 1\right)/2$ the fixed point is globally stable, in the sense that it is locally stable and attracts all points of $S$, where $S$ contains those $(x_0, x_1, x_2) \in \mathbb{R}_+^3$ for which the sequence $(x_n)_{n=0}^{\infty}$ remains in $\mathbb{R}_+$. The proof is a combination of analytical and reliable numerical methods. The novelty is an explicit construction of a relatively large attracting neighborhood of the nontrivial fixed point of the 3-dimensional logistic map by using center manifold techniques and the Neimark–Sacker bifurcational normal form. The results appeared in [1].

## References

[1] J. DUDÁS, T. KRISZTIN: Global stability for the three-dimensional logistic map, *Nonlinearity* 34 (2021), no. 2, 894–938.

# CPFloat: A C library for emulating low-precision arithmetic

Massimiliano Fasi[1] and Mantas Mikaitis[2]

[1] Örebro University, Örebro, Sweden
S-701 82 Örebro, Sweden
`massimiliano.fasi@oru.se`

[2] The University of Manchester, Manchester, UK
Oxford Road, M13 9PL Manchester, UK `mantas.mikaitis@manchester.ac.uk`

**Keywords:** low-precision floating-point arithmetic, binary16, bfloat16

From a software perspective, it is straightforward to simulate a given floating-point format on a machine whose hardware supports a format with a larger exponent range and more bits of precision: one can perform each arithmetic operation in hardware, using the available precision, and then round the result to the desired number of significant binary digits.

If the underlying floating-point arithmetic is IEEE compliant, rounding can be performed by using only the definition of floating-point numbers and a few standard mathematical library functions. This approach, however, is not necessarily robust: handling subnormals, underflow, and overflow requires special attention, and numerical errors can creep in and cause a mathematically correct formula to perform incorrectly in finite arithmetic. Moreover, the ensuing implementations are not necessarily efficient, as the library functions these techniques build upon are typically designed to handle a broad range of cases and may not be optimised for the specific needs of floating-point rounding algorithms.

CPFloat [1] is a header-only C library that offers efficient routines for rounding arrays of binary32 or binary64 numbers to lower precision. The library is distributed as free software, is fully documented, is accompanied by a comprehensive test suite, and is hosted on GitHub.[1] The repository also contains a MEX interface for MATLAB and Octave, and an object-oriented interface for C++. The library implements a variety of rounding modes: round-to-nearest with several tie-breaking rules, directed rounding, two variants of stochastic rounding, and round-to-odd. Any format that can fit into binary64 is supported, but if round-to-nearest is used then only formats with up to 26 digits of precision are recommended, as these are immune from double rounding [4]. The underlying techniques exploit the bit level representation of these formats, and perform only low-level bit manipulations and integer arithmetic without relying on costly library calls.

In numerical experiments the new techniques bring a considerable speedup (typically one order of magnitude or more) over existing C/C++ alternatives, such as GNU MPFR [3] and the FloatX library [2]. To the best of our knowledge, CPFloat is currently the most efficient and complete library for rounding floating-point numbers to a custom low-precision format.

## References

[1] Massimiliano Fasi and Mantas Mikaitis. CPFloat: A C library for emulating low-precision arithmetic. MIMS EPrint 2020.22, Manchester Institute for Mathematical Sciences, The University of Manchester, UK, October 2020.

[2] Goran Flegar, Florian Scheidegger, Vedran Novaković, Giovani Mariani, Andrés E. Tomás, A. Cristiano I. Malossi, and Enrique S. Quintana-Ortí. FloatX: A C++ library for customized floating-point arithmetic. *ACM Trans. Math. Software*, 45(4):1–23, December 2019.

---

[1] `https://github.com/mfasi/cpfloat`

[3] Laurent Fousse, Guillaume Hanrot, Vincent Lefèvre, Patrick Pélissier, and Paul Zimmermann. MPFR: A multiple-precision binary floating-point library with correct rounding. *ACM Trans. Math. Software*, 33(2):13:1–13:15, June 2007.

[4] Siegfried M. Rump. IEEE754 precision-$k$ base-$\beta$ arithmetic inherited by precision-$m$ base-$\beta$ arithmetic for $k < m$. *ACM Trans. Math. Software*, 43(3):1–15, January 2017.

# Local error estimation and step size control in adaptive linear multistep methods

Carmen Arévalo[1], Gustaf Söderlind[1], Yiannis Hadjimichael[2], and <u>Imre Fekete</u>[3,4]

[1] Centre for Mathematical Sciences, Lund University, Sweden
[2] Weierstrass Institute for Applied Analysis and Stochastics, Germany
[3] Institute of Mathematics, Eötvös Loránd University, Hungary
[4] MTA-ELTE NUMNET Research Group, Hungary

`imre.fekete@ttk.elte.hu`

**Keywords:** adaptivity, variable step size, linear multistep methods

The talk is based on the recently published open access paper [1].

In a $k$-step adaptive linear multistep methods the coefficients depend on the $k-1$ most recent step size ratios. In a similar way, both the actual and the estimated local error will depend on these step ratios.

The classical error model has been the asymptotic model, $ch^{p+1}y^{(p+1)}(t)$, based on the constant step size analysis, where all past step sizes simultaneously go to zero. This does not reflect actual computations with multistep methods, where the step size control selects the next step, based on error information from previously accepted steps and the recent step size history. In variable step size implementations the error model must therefore be dynamic and include past step ratios, even in the asymptotic regime.

In this talk we derive dynamic asymptotic models of the local error and its estimator, and show how to use dynamically compensated step size controllers that keep the asymptotic local error near a prescribed tolerance TOL. The new error models enable the use of controllers with enhanced stability, producing more regular step size sequences. Numerical examples illustrate the impact of dynamically compensated control, and that the proper choice of error estimator affects efficiency.

## References

[1] C. ARÉVALO, G. SÖDERLIND, Y. HADJIMICHAEL, I. FEKETE: Local error estimation and step size control in adaptive linear multistep methods, *Numerical Algorithms*, 86 (2021), 537–563.

# The Sharkovskii Theorem for multidimensional maps with attracting periodic orbits

Anna Gierzkiewicz[1], and Piotr Zgliczyński[2]

[1] University of Agriculture in Krakow, Poland
[2] Jagiellonian University, Krakow, Poland
anna.gierzkiewicz@urk.edu.pl

**Keywords:** computer-assisted proof, Sharkovskii theorem, periodic orbit, covering relation, Roessler system

The Sharkovskii Theorem [1] is a powerful tool for proving the existence of periodic orbits and chaotic phenomena in discrete one-dimensional dynamical systems:

**Theorem 1** (Sharkovskii)**.** *Define an ordering '$\triangleleft$' of natural numbers:*

$$3 \triangleleft 5 \triangleleft 7 \triangleleft \cdots \triangleleft 2 \cdot 3 \triangleleft 2 \cdot 5 \triangleleft \cdots \triangleleft 2^2 \cdot 3 \triangleleft 2^2 \cdot 5 \triangleleft \cdots \triangleleft 2^k \triangleleft 2^{k-1} \triangleleft \cdots \triangleleft 2^2 \triangleleft 2 \triangleleft 1. \tag{1}$$

*Let $f : I \to \mathbb{R}$ be a continuous map of an interval. If $f$ has an $n$-periodic point and $n \triangleleft m$, then $f$ also has an $m$-periodic point.*

In general, the above Theorem is not valid for multidimensional maps. However, we show that the methods used by Burns and Hasselblatt [2] to prove Sharkovskii's Theorem can be generalized to the case of higher-dimensional maps with an attracting periodic orbit. One can find a connection between the one-dimensional covering of intervals and the multidimensional covering of h-sets, which is a well-known tool in computer-assisted proofs in dynamics. The result is [4]:

**Theorem 2.** *Consider a continuous map $F : I \times \overline{B}(0, R) \to \mathrm{int}\left(I \times \overline{B}(0, R)\right)$, where $I \subset \mathbb{R}$ is a closed interval and $\overline{B}(0, R) \subset \mathbb{R}^{n-1}$ a closed ball of radius $R$. Let us denote by $(x, y)$ points in $I \times \overline{B}(0, R)$.*

*Suppose that $F$ has an $n$-periodic point $(x_0, y_0) \in \mathbb{R} \times \mathbb{R}^{n-1}$ with least period $n$ and denote its orbit by $\{(x_0, y_0), (x_1, y_1) = F(x_0, y_0), \ldots, (x_{n-1}, y_{n-1}) = F^{n-1}(x_0, y_0), (x_n, y_n) = (x_0, y_0)\} \subset \mathrm{int}\, I \times \overline{B}(0, R).$*

*Suppose that there exist $\delta_0, \delta_1, \ldots, \delta_{n-1} > 0$ such that*

$$\forall\, i \in \{0, \ldots, n-1\} \qquad F\left([x_i \pm \delta_i] \times \overline{B}(0, R)\right) \subset (x_{i+1} \pm \delta_{i+1}) \times B(0, R).$$

*Then for every natural number $m$ succeeding $n$ in the Sharkovskii order (1), $F$ has a point with the least period $m$.*

As an application, we prove the existence of $n$-periodic orbits for almost all $n \in \mathbb{N}$ in the Rössler system with an attracting periodic orbit, for four sets of parameters. The proof is computer-assisted with the use of CAPD library for C++ [3].

## References

[1] A.N. SHARKOVSKII: Co-existence of cycles of a continuous mapping of the line into itself, *Ukrainian Math. J.*, 16 (1964) 61–71 (in Russian, English translation in *J. Bifur. Chaos Appl. Sci. Engrg.*, 5 (1995) 1263–1273).

[2] K. BURNS AND B. HASSELBLATT: The Sharkovsky Theorem: A Natural Direct Proof, *The American Mathematical Monthly*, vol. 118, No. 3 (2011), 229–244.

[3] 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 (2021), 105578.

[4] A. GIERZKIEWICZ AND P. ZGLICZYŃSKI: From the Sharkovskii theorem to periodic orbits for the Rössler system, submitted to *J. Diff. Eq*, (2021).

# Rigorous numerics of blow-up separatrix in autonomous ODEs

Kaname Matsue[1,2]

[1]Institute of Mathematics for Industry / [2]International Institute for Carbon-Neutral Energy Research
(WPI-I$^2$CNER),
Kyushu University,
744 Motooka, Fukuoka, 819-0395, Japan

kmatsue@imi.kyushu-u.ac.jp

**Keywords:** Blow-up Solutions, Rigorous numerics, Separatrix.

My talk is concerned with special blow-up solutions for ordinary differential equations (ODEs) separating sets of initial points into two regions, one of which admits solutions global in time, and another admits blow-up solutions. There are several examples of differential equations such that the amplitude of initial points is not essential to determine asymptotic behavior of solutions including blow-ups. The present topic addresses such a dynamical nature with computer-assisted proof. Our approach used here is based on the combination of machineries in dynamical systems (e.g. phase space compactifications, time-scale desingularizations of vector fields) and computer-assisted proofs (e.g. rigorous integrators, local Lyapunov functions and parameterization method for invariant manifolds), which enables us to characterize blow-up solutions by means of (un)stable manifolds of "invariant sets at infinity". In particular, (un)stable manifolds of saddle-type invariant sets at infinity can characterize blow-up solutions which are unstable under small perturbations of initial points and separations of phase spaces mentioned above. The above invariant manifolds, which we shall call *blow-up separatrices*, provide singular behavior of trajectories from the viewpoints of not only asymptotic behavior but also maximal existence times (blow-up times) of trajectories. One example of such a singular nature in an autonomous ODE is exhibited with various concretely enclosed results. The present talk is based on the joint work with Jean-Philippe Lessard and Akitoshi Takayasu [1].

## References

[1] LESSARD, J.-P. AND MATSUE, K. AND TAKAYASU, A.: A geometric characterization of unstable blow-up solutions with computer-assisted proof, *arXiv:2103.12390*, 2021.

# Error-free transformation of matrix multiplication for multi-precision computations

## Katsuhisa Ozaki[1]

[1] Shibaura Institute of Technology, Saitama, Japan
ozaki@sic.shibaura-it.ac.jp

Multi-precision computations are used if a numerical result of standard floating-point (FP) arithmetic is inaccurate. GMP [1], MPFR [2], and exflib [3] are examples of multi-precision arithmetic libraries. Hardware-supported FP arithmetic, such as binary32 and binary64 in IEEE 754, can be performed rapidly in modern computers. However, the performance of software-emulated multi-precision arithmetic is slow compared with binary32 and binary64. This study emulates and accelerates multi-precision computations for numerical linear algebra problems using only FP arithmetic.

Ozaki et al. proposed an error-free transformation (EFT) of matrix multiplication [4]. For FP matrices $A$ and $B$, EFT algorithm splits $A$ and $B$ into

$$
\begin{aligned}
A &= A^{(1)} + A^{(2)} + \cdots + A^{(k)}, \\
B &= B^{(1)} + B^{(2)} + \cdots + B^{(\ell)},
\end{aligned}
$$

to avoid a rounding error in the FP evaluation of $A^{(i)}B^{(j)}$ for all $(i, j)$ pairs. Then, $AB$ can be transformed into an unevaluated sum of $k\ell$ FP matrices using only FP arithmetic. If we employ accurate summation algorithms, an accurate numerical result can be obtained. Mukunoki et al. [5] applied the EFT into the emulation of binary128.

We employ this technique for multi-precision computations in numerical linear algebra, especially, focusing on matrix multiplication and Cholesky decomposition. Using diagonal scaling, we exploit the EFT of matrix multiplication with binary64. In addition, we employ the EFT for block Cholesky decomposition. The block Cholesky decomposition consists of (i) Cholesky decomposition for diagonal blocks, (ii) solving triangular systems, and (iii) block matrix multiplications. For (i) and (ii), we straightforwardly use multi-precision computations. For (iii), we use the EFT. Because (iii) is the most computational-intensive, using EFT for matrix multiplication accelerates the performance of the block Cholesky decomposition.

Finally, numerical examples were used to illustrate the efficiency of the proposed method. We compared the computation times and accuracy of numerical results of the proposed method and Advanpix Multiprecision Computing Toolbox for MATLAB (MCT) [6], a user-friendly and well-turned toolbox. We used Core i7-86665U and MATLAB2020a as the computational environment. We set 34 as digits (comparable to binary128) in MCT. As a result, the proposed method for matrix multiplication is about 8–20 times faster than MCT's matrix multiplication in the best case. In addition, the EFT-based block Cholesky decomposition was approximately 1.5–2.5 times faster than the MCT's Cholesky decomposition. More detailed numerical examples, including results from other CPUs and precision, will be shown in the presentation.

## References

[1] The GNU Multiple Precision Arithmetic Library, https://gmplib.org/.

[2] The GNU MPFR Library, https://www.mpfr.org/.

[3] exflib - extend precision floating-point arithmetic library, http://www-an.acs.i.kyoto-u.ac.jp/~fujiwara/exflib/.

[4] K. Ozaki, T. Ogita, S. Oishi, S. M. Rump: Error-Free Transformation of Matrix Multiplication by Using Fast Routines of Matrix Multiplication and its Applications, *Numerical Algorithms*, 59 (2012), 95-118.

[5] D. MUKUNOKI, K. OZAKI, T. OGITA, T. IMAMURA: Accurate Matrix Multiplication on Binary128 Format Accelerated by Ozaki Scheme, *in Proc. of 50th International Conference on Parallel Processing*, to appear.

[6] Multiprecision Computing Toolbox for MATLAB, https://www.advanpix.com/.

# On some convergence properties for finite element approximations to the inverse of linear elliptic operators

Takehiko Kinoshita[1], Yoshitaka Watanabe[2], and Mitsuhiro T. Nakao[3]

[1] Department of Information Science, Saga University, Saga 840-8502, Japan
[2] Research Institute for Information Technology, Kyushu University, 744 Motooka, Nishi-ku, Fukuoka 819-0395, Japan
[3] Faculty of Science and Engineering, Waseda University, 3-4-1 Okubo, Shinjuku-ku, Tokyo 169-8555, Japan
kinosita@cc.saga-u.ac.jp

This talks deals with convergence theorems of the Galerkin finite element approximation for the second order elliptic boundary value problems. Under some quite general setting, we show not only the pointwise convergence but also prove that the norm of approximate operator converges to the corresponding norm for the inverse of linear elliptic operator. Since the approximate norm estimates of linearized inverse operator plays an essential role in the numerical verification method of solutions for nonlinear elliptic problems, our result is also important in terms of guaranteeing its validity.

## References

[1] T. KINOSHITA, Y. WATANABE, MITSUHIRO T. NAKAO: Some lower bound estimates for resolvents of a compact operator on an infinite-dimensional Hilbert space, *Journal of Computational and Applied Mathematics*, 369 (2020).

[2] M. T. NAKAO, K. HASHIMOTO, Y. WATANABE: A numerical method to verify the invertibility of linear elliptic operators with applications to nonlinear problems, *Computing*, 75 (2005) 1–14.

[3] M. T. NAKAO, M. PLUM, Y. WATANABE: Numerical Verification Methods and Computer-Assisted Proofs for Partial Differential Equations, Springer Series in Computational Mathematics 53 *Springer Nature Singapore*, (2019).

# Kinematic Metric Spaces Under Interval Uncertainty: Towards an Adequate Definition

Vladik Kreinovich[1], Olga Kosheleva[1], and Victor Selivanov[2]

[1] University of Texas at El Paso, USA
[2] A. P. Ershov Institute of Informatics Systems, Novosibirsk, Russia
vladik@utep.edu, olgak@utep.edu, vseliv@iis.nsk.su

**What is a kinematic metric: physical introduction.** In the physical space, we can define the distance $d(a, b)$ between two points as the length of the shortest possible path between them. Thus defined distance is symmetric ($d(a, b) = d(b, a)$) and satisfies the usual triangle inequality $d(a, c) \leq d(a, b) + d(b, c)$. The mathematical notion of a metric is a natural generalization of this physical notion.

From the viewpoint of space-time, physical space corresponds to the situation when we take space-time points ("events") $(a, t_0)$, $(b, t_0)$, etc. corresponding to the same moment of time $t_0$. In relativity theory, such events cannot causally influence each other.

When an event $a$ *can* causally influence an event $b$ (we will denote this strict order – i.e., irreflexive transitive – relation by $a < b$), this influence is implemented by a particle or particles whose trajectories start at $a$ and end up at $b$. For each such trajectory, we can measure the proper time of the corresponding particle. In principle, particles can travel as close to the speed of light as possible, in which case the proper time can be as close to 0 as possible – so the *smallest* proper time over all trajectories is always 0. Interestingly, there is the *largest* proper time $\tau(a, b)$ – which corresponds to inertial motion. The corresponding function $\tau(a, b)$ – defined only when $a < b$ – satisfies the "anti-triangle" inequality $\tau(a, c) \geq \tau(a, b) + \tau(b, c)$.

This inequality describes the known *twins paradox* of relativity theory: when a twin brother who traveled to the stars comes back to Earth, he will be younger than his twin who stayed on Earth: the biological age of the stay-home brother is $\tau(a, c)$, while the biological age of the astronaut brother is $\tau(a, c) + \tau(c, b)$, where $c$ is the moment when the brother reached a faraway star.

A natural generalization of this function is a notion of *kinematic metric*.

**Kinematic metric: definition.** Let $(X, <)$ be an ordered set. A function $\tau(a, b)$ – defined for all pairs for which $a < b$ – is called a *kinematic metric* if all its values are non–negative and it satisfies the anti-triangle inequality.

**Need for interval uncertainty.** All information about the values of a physical quantity $v$ – including the values of the kinematic metric – comes from measurements. Measurements are never absolutely accurate, so the measurement result $\widetilde{v}$ is, in general, different from the actual (unknown) value $v$: there is a measurement error $\Delta v \stackrel{\text{def}}{=} \widetilde{v} - v$. Often, the only information that we have about the measurement error is an upper bound $\Delta$ on its absolute value. In this case, the only information that we have about the actual value $v$ is that this value belongs to the interval $[\underline{v}, \overline{v}] \stackrel{\text{def}}{=} [\widetilde{v} - \Delta, \widetilde{v} + \Delta]$.

**Natural question.** Suppose that we have, for all pairs $a < b$, intervals $[\underline{\tau}(a, b), \overline{\tau}(a, b)]$, wth $\underline{\tau}(a, b) \geq 0$, obtained from measurement. If all the upper bounds $\Delta(a, b)$ are correct, then there is a kinematic metric $\tau(a, b)$ for which $\tau(a, b) \in [\underline{\tau}(a, b), \overline{\tau}(a, b)]$ for all $a < b$. However, if we – as happens – underestimated the measurement errors, we may not have such a function.

So, a natural question is: what is the condition on the intervals $[\underline{\tau}(a, b), \overline{\tau}(a, b)]$ under which such a function $\tau(a, b)$ exists?

**A seemingly natural idea does not work.** Anti-triangle inequality implies that $\overline{\tau}(a, c) \geq \underline{\tau}(a, b) + \underline{\tau}(b, c)$ for all $a < b < c$. So, it may seem that this inequality is the right condition for the existence of the desired kinematic metric $\tau(a, b)$. However, this inequality does not guarantee the existence of $\tau(a, b)$: e.g., for $X = \{a_1 < a_2 < a_3 < a_4\}$ and $[\underline{\tau}(a_i, a_j), \overline{\tau}(a_i, a_j)] = [1, 2]$ for all $i < j$, this inequality is satisfied, but the desired function $\tau(a, b)$ is not possible: indeed, if it existed, we would have $2 \geq \tau(a_1, a_4) \geq \tau(a_1, a_2) + \tau(a_2, a_3) + \tau(a_3, a_4) \geq 3$, i.e., $2 \geq 3$.

**Main result.** For an interval-valued function $[\underline{\tau}(a,b), \overline{\tau}(a,b)]$ defined for all $a < b$, the existence of the kinematic metric $\tau(a,b)$ for which always $\tau(a,b) \in [\underline{\tau}(a,b), \overline{\tau}(a,b)]$ is equivalent to the condition that $\overline{\tau}(a_1, a_n) \geq \sum\limits_{i=1}^{n-1} \underline{\tau}(a_i, a_{i+1})$ for all sequences $a_1 < \ldots < a_n$.

**Proof: main idea.** If $\tau(a,b)$ exists, then this inequality is clearly satisfied. Vice versa, if the above condition is satisfied, then we can take $\tau(a,b) = \sup \left\{ \sum\limits_{i=1}^{n-1} \underline{\tau}(a_i, a_{i+1}) \right\}$, where the supremum is taken overall all the chains $a = a_1 < a_2 < \ldots < a_n = b$ that connect $a$ and $b$.

*Comment.* We need the above condition for *all* natural numbers $n$: if we only require it only for $n \leq n_0$, this does not guarantee the existence of $\tau(a,b)$.

this may not guarantee the existence of a kinematic metric: example is $X = \{a_1 < \ldots < n_0 + 2\}$ and $[\underline{\tau}(a_i, a_j), \overline{\tau}(a_i, a_j)] = [1, n_0]$ for all $i < j$.

### References

[1] R. I. PIMENOV: *Kinematic Spaces: Mathematical Theory of Space-Time*, Consultants Bureau, New York, 1970.


Mahdokht Afravi, Vladik Kreinovich, and Thongchai Dumrongpokaphan, "Metric Spaces Under Interval Uncertainty: Towards an Adequate Definition", In: Grigori Sidorov and Oscar Herrera-Alcántara (Eds.), *Advances in Computational Intelligence: Proceedings of the 15th Mexican International Conference on Artificial Intelligence MICAI'2016, Cancun, Mexico, October 25–29, 2016, Part I*, Springer Lecture Notes in Artificial Intelligence, Cham, Switzerland, 2017, Part I, Vol. 10061.

# How probabilistic methods for data fitting deal with interval uncertainty: a more realistic analysis

Vladik Kreinovich[1] and Sergey P. Shary[2]

[1] University of Texas at El Paso, USA
[2] Novosibirsk University, Novosibirsk, Russia
vladik@utep.edu, shary@ict.nsc.ru

**General motivation.** When processing data, most practitioners use probabilistic methods. It is therefore desirable to study how, for the case of interval uncertainty, these methods compare with interval techniques.

**Data fiting problem.** In many situations, we know the general form $y = F(x, c)$ of the dependence of a quantity $y$ on quantities $x = (x_1, \ldots, x_n)$, but we do not know the exact values of the parameters $c = (c_1, \ldots, c_m)$. These values must be determined from the measurement results. For this purpose, several ($K$) times, we measure $x_i$ and $y$. Based on the measurement results $\widetilde{x}_k = (\widetilde{x}_{k1}, \ldots, \widetilde{x}_{kn})$ and $\widetilde{y}_k$, we need to estimate the values of the parameters. This problem is also called *problem of parameter estimation*.

Measurements are never absolutely accurate. Because of this, we need to take into account that the measurement results $\widetilde{v}$ are, in general, different from the actual (unknown) values of the corresponding quantity $v$, i.e., that there is a non-zero measurement error $\Delta v := \widetilde{v} - v$.

**Known probability distributions.** In many cases, we know the probability distributions $f_i(\Delta x_i)$ and $f(\Delta y)$ of the measurement errors. In this case, we can use the Maximum Likelihood (ML) approach — i.e., select the *most probable* values $c$ (and $x_{ki}$) for which the product $\prod\limits_{k=1}^{K} \left( f(\widetilde{y}_k - F(x_k, c)) \cdot \prod\limits_{i=1}^{n} f_i(\widetilde{x}_{ki} - x_{ki}) \right)$ is the largest. Usually, the logarithm of this product, known as *log-likelihood*, is maximized for computational convenience.

**Interval uncertainty.** In many practical situations, we do not know the probability distributions, all we know is that the measurement errors $\Delta v$ are located on the given interval $[-\Delta_v, \Delta_v]$. In such situations, a usual probabilistic approach is to select, on this interval, the distribution with maximal entropy – which turns out to be the uniform distribution.

**Simplest case.** The simplest – and rather frequent – case is when the values $x_i$ are measured very accurately, so we can safely ignore the corresponding measurement errors and conclude that $\widetilde{x}_{ik} = x_{ik}$ for all $i$ and $k$. In this case, the ML approach selects all possible values $c$ for which, for all $k$, we have $F(x_k, c) \in [\widetilde{y}_k - \Delta_y, \widetilde{y}_k + \Delta_y]$; see, e.g., [1]. Interestingly, in this case, the probabilistic approach leads to the same answer as the interval techniques.

**General case.** If we also know the values $x_{ki}$ with interval uncertainty, then the ML approach selects the set of all the values $c$ for which $F(x_k, c) \in \boldsymbol{y}_k = [\widetilde{y}_k - \Delta_y, \widetilde{y}_k + \Delta_y]$ for some values $x_{ki} \in \boldsymbol{x}_{ki} = [\widetilde{x}_{ki} - \Delta_{x_i}, \widetilde{x}_{ki} + \Delta_{x_i}]$. This is exactly the *united solution set* to the interval equation system constructed from interval data [1, 2]. Thus, the united solution set has a natural probabilistic meaning.

**A more realistic description of the practical problem.** Often, when we get a measurement result, this does not mean that there was only one measurement: it means that there were several different measurements leading to the same result – e.g., same intervals.

**How probabilistic techniques deal with this situation.** For each $k$, instead of a single combination $x_k$, we have several $x_{k\ell}$ for different $\ell$. For each combination of values $x_{k\ell i} \in \boldsymbol{x}_{ki}$, we can form the log-likelihood $\sum\limits_{k=1}^{K} \sum\limits_{\ell} \sum\limits_{i=1}^{n} \ln(f_i(\widetilde{y}_k - F(x_{k\ell}, c)))$. We do not know the actual values $x_{k\ell i}$; following the maximum entropy idea, we assume that they are uniformly distributed on the corresponding intervals $\boldsymbol{x}_{ki}$. For a large number of constituent measurement $\ell$, the sum over $\ell$ is proportional to the expected value. Thus, a reasonable idea is to maximize the expected value of the log-likelihood over this uniform distribution.

**What is the resulting estimate.** We show that, as a result, we return all values of $c$ for which $f(x_k, c) \in \boldsymbol{y}_k$ for *all* $x_{ki} \in \boldsymbol{x}_{ki}$ — which is exactly the *tolerable solution set* to the interval equation system constructed from data, a solution set that has many useful properties; see, e.g., [2]. So, the tolerable solution set also makes sense in the probabilistic setting.

## References

[1] V. KREINOVICH, S. P. SHARY: Interval methods for data fitting under uncertainty: a probabilistic treatment, *Reliable Computing*, 23 (2016), 105–141.

[2] S. P. SHARY: Weak and strong compatibility in data fitting problems under interval uncertainty, *Advances in Data Science and Adaptive Analysis*, 12 (2020), No. 1, Paper 2050002.

# How many constraints are satisfied? An approach to solving classification and regression problems

Bartłomiej Jacek Kubica[1],

[1] Department, of Information Systems, Institute of Information Technology, Warsaw University of Life Sciences SGGW, Poland
bartlomiej_kubica@sggw.edu.pl

**Keywords:** interval methods, constraint satisfaction, classification, regression, machine learning, fuzzy numbers

Problems of classification and regression are ubiquitous in several branches of science, technology, and economy. They can both be formulated in a similar manner: We have a training set of pairs $(x_i, y_i)$, $i = 1, \ldots, N$, where $x_i$'s are some attribute vectors and $y_i$'s are their labels – either from a discrete (classification) or continuous (regressions) space.

Provided this training set, we want to parameterize some model, i.e., a function $f(x, p)$, by providing the value of $p$ such that: $f(x_i, p) \approx y_i$ for all (or sometimes most) $i = 1, \ldots, N$.

By obtaining the proper value of $p$, we shall be able to predict/assign the labels $y$ also to arguments $x$ from outside of the training set.

The above formulation refers, in particular, to several kinds of regression (linear, logistic, and more advanced ones), but also to many machine learning tools: support vector machines, and various kinds of neural networks, including deep neural networks.

It is worth noting that, in many practical situations, the training set is not faultless. It can have two kinds of mistakes:

- inaccuracies: values of $x_i$ or $y_i$ may be somewhat noisy,

- outliers: examples that are completely wrong and should be ignored, when identified.

Interval methods are usually well at bounding noisy values, using intervals: $\mathbf{x}_i = [\underline{x}_i, \overline{x}_i]$, and $\mathbf{y}_i = [\underline{y}_i, \overline{y}_i]$.

How to formulate the above problem in precise terms, i.e., how to express the relation $\approx$? As indicated in [1], we can do it in two manners: either as an optimization problem or a constraint satisfaction problem (CSP).

The optimization problem can look as follows:

$$\min_p \|y - f(x, p)\|, \tag{1}$$

where the minimized norm is usually the least-squares: $\sum_{i=1}^{N} \left( y_i - f(x_i, p) \right)^2$, but other norms can be used as well; in deep learning applications, the Kullback-Leibler divergence [2] has been recognized as particularly useful.

The CSP can be formulated like:

$$\text{Find the set } \{p \mid f(\mathbf{x}_i, p) \subseteq \mathbf{y}_i \; i = 1, \ldots, n\}. \tag{2}$$

Both approaches have their advantages and drawbacks. Formulation (1) may be more familiar to researchers from outside of the interval community. Also, this formulation allows us to find a solution, even if the model is not strictly correct, or if there are some outliers. Unfortunately, its main advantage is also the main drawback: it does not allow us to check whether the model is correct, hence the results might be considered non-reliable (they are verified solutions to a non-precisely formulated problem!).

On the other hand, using (2), we can verify if the model is correct, i.e., if it is consistent with all the measurements. But in the presence of outliers, the solution to (2) will often be the empty set.

Fortunately, this formulation (2) can be enhanced to take the possible outliers into account. Firstly, we can assume *a priori* a number $\tau$ of faulty examples in the training set and seek the set:

$$\left\{ p \mid f(\mathbf{x}_i, p) \subseteq \mathbf{y}_i \text{ for at least } (N - \tau) \text{ of } i \in \{1, \ldots, n\} \right\}, \tag{3}$$

instead of (2).

A more advanced, but also more computationally intensive approach, would be to find the fuzzy set of solutions. Its membership function at $p$ would be the percentage of equations that are satisfied for $p$: from zero to (possibly) one, although all equations are likely not to be satisfied anywhere (this would correspond to a situation with no outliers).

These approaches are going to be presented in the paper.

**References**

[1] B. J. KUBICA: *Interval Methods for Solving Nonlinear Constraint Satisfaction, Optimization and Similar Problems: From Inequalities Systems to Game Solutions*, Series: Studies in Computational Intelligence, 805, Springer, 2019.

[2] O. KOSHELEVA, V. KREINOVICH: Why deep learning methods use KL divergence instead of least squares: A possible pedagogical explanation *Mathematicheskiye struktury y modelirovaniye*, 2 (2018).

# Information Sets of Rebuilding Dependence Parameters for Criteria of Strong and Weak Compatibility under Heavy Two-Dimensional Measuring Errors

Kumkov S. I.

N. N. Krasovskii Institute of Mathematics and Mechanics,
Ural Branch of Russian Academy of Sciences, Ekaterinburg, Russia
sikumkov@gmail.com

An experimental process sample is investigated that contains 10 "heavy" noised measurements with two-dimensional measuring errors (Fig. 1). The function approximating the process is $F(t) = At + B$ , where $t > 0$ is the argument, $A > 0$ and $B > 0$ are the parameters (as in the middle part of the complete sample from [1]). Uncertainty box of each measurement is constrained by errors of large values $|e_t| \leq e_{t,max} \leq 0.42$ and $|e_F| \leq e_{F,max} \leq 0.3$, correspondingly. Note interesting aspect: step in $t$ is 0.04, but the constraint on it exceeds the step and is $e_{t,max} = 0.042$. As a result, uncertainty boxes (Fig. 1) overlap each other. Probability characteristics of both errors are unknown.



Figure 1: Noised measurements (crosses) and their uncertainty boxes (rectangles); for illustration, the least squares mean line is shown in dots.

**Problem:** it is necessary to build and compare the information sets of compatible values of parameters $A, B$ for two compatibility criteria — weak and strong [2].

Due to simple form of the approximating function and using the Partial Information Sets approach [3], it was succeeded to build exact information set for the strong compatibility criterion.

But for the weak compatibility criterion, corresponding information set was peeled (as described in [4]) by small boxes, satisfying the experimentalist (Fig. 2).

Obviously, the second information set is significantly larger than the first one and completely contains the latter inside.

**Conclusion**

The shown fact is very important for choosing the criterion of processing the noised experimental information with interval uncertainty.

## References

[1] S. V. GLADCOVSKY, S. I. KUMKOV: *Application of approximation methods to analysis of peculiarities of breaking-up and forecasting the break-resistibility of high-strength steel*, Matematicheskoe modelirovanie sistem i protsessov. Sbornik nauchnykh trudov, Permskii Gos. Tekhnicheskii Universitet, Perm. (1997), no. 5, pp. 26–34. (in Russian)

[2] S. P. SHARY: *Weak and Strong Compatibility in Data Fitting Problems Under Interval Uncertainty*, Advances in Data Science and Adaptive Analysis Vol. 12, No. 1 (2020).

Figure 2: a) Exact information set (for strong compatibility) and b) peeled set (for weak compatibility); minimal outer box-estimates of both sets are marked in dots.

[3] S. I. KUMKOV, YU. V. MIKUSHINA: *Interval approach to identification of catalytical process parameters*, Reliable Computing, 2013, Vol. 19, pp. 197–214.

[4] L. JAULIN, M. KIEFER, O. DIDRIT, E. WALTER: *Applied Interval Analysis*, Springer-Verlag, London, 2001.

# On Convexity and Difficulty of Global Optimization Problems

Dun Liu[1] and R. Baker Kearfott[2]

[1] University of Louisiana at Lafayette, Louisiana, United States
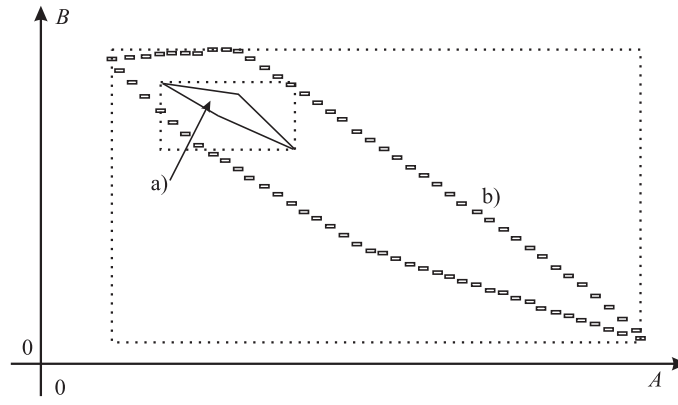[2] University of Louisiana at Lafayette, Louisiana, United States
dun.liu1@louisiana.edu
rbk@lusfiber.net

**Keywords:** global optimization, convexity, interval analysis, branch and bound

Several papers from the early-2000's to the present have considered the standard quadratic optimization problem, which can be defined as

$$\begin{aligned} \min \quad & \phi(x) = x^T C x \\ \text{subject to} \quad & x \in \mathcal{S}, \end{aligned} \tag{1}$$

where $C \in \mathbb{R}^{n+1 \times n+1}$, and $\mathcal{S}$ is the standard simplex defined by $\mathcal{S} = \{x = (x_1, \ldots, x_{n+1}) | \sum_{i=1}^{n+1} x_i = 1, x_i \geq 0\}$, see [3, 2, 4, 1]. In these works, the concept of *convexity density* is proposed, which is defined as the ratio of the number of convex edges to the total number of edges of the simplex thereof. In particular, some of these works can be construed as generalizations of theory of linear programming based on a measure of the "amount" of convexity on the boundary of the simplex, while other results seem to indicate that, the more circumstantial evidence of convexity there is on the boundary, the more difficult the problem can be; this second type of result possibly leads to an a-priori way of detecting the difficulty of global optimization problems. We will briefly review the results.

Based on the preliminary works, we generalized the concept of *convexity density* to more general objective functions using interval arithmetic. Finally, we test the power of the generalized concept to predict difficulty by running test sets on various significantly different branch-and-bound solvers.

## References

[1] A. SCOZZARI AND F. TARDELLA: A clique algorithm for standard quadratic programming *Discrete Applied Mathematics*, 156 (2008), 2439–2448.

[2] G. LIUZZI, M. LOCATELLI AND V. PICCIALLI: A new branch-and-bound algorithm for standard quadratic programming problems *Optimization Methods and Software*, 34 (2019), 79–97.

[3] I.M. BOMZE: On Standard Quadratic Optimization Problems *Journal of Global Optimization*, 13 (1998), 369–387.

[4] I. NOWAK: A new semidefinite programming bound for indefinite quadratic forms over a simplex *Journal of Global Optimization*, 14 (1999), 357–364.

# Rigorous maximum norm estimation for polynomial systems

Xuefeng Liu[1],

[1] Niigata University, Japan
xfliu@math.sc.niigata-u.ac.jp

**Keywords:** Maximum norm, finite element method, rigorous norm estimation

The maximum norm estimation for general functions is not an easy task, even in the sense of an approximate estimation. The developing team of the popular finite element method library FEniCS discussed whether to provide the maximum norm estimation for finite element solutions, but gave up the idea to provide such a feature due to the problem of complexity and trustability of computed results [1].

In this research, a new idea is proposed to utilize the convex hull property for Bernstein polynomials (see, e.g., [2]) to evaluate the maximum norm of polynomial systems, including the function of finite element solutions, which is usually represented by piecewise polynomials over triangulation of domains. The convex hull property says that for a Bernstein polynomial function of degree $n$ over interval $[0, 1]$,

$$f(x) = \sum_{i=0}^{n} c_i \binom{n}{i} x^i (1-x)^{(n-i)} \ ,$$

the following inequality holds

$$\|f\|_\infty \le \max |c_i|.$$

Since a direct application of the convex hull property will lead to a raw estimation of the range of function $f$, we apply De Casteljau's algorithm (see, e.g., [2]) to recursively subdivide the objective domain to have piecewise and exact representation of the objective polynomial function. The application of the convex hull property for subdivided polynomial over small sub-interval will have a sharp estimation of the maximum norm of the function.

In the Verified Finite Element Method (VFEM) library for function over 2D and 3D domains, the above idea has been successfully implemented to provide efficient and sharp maximum norm estimation for functions from FEM function spaces [3].

The new efficient algorithm for maximum norm estimation has also been applied to the rigorous root computation for nonlinear polynomial systems. Compared with the existing approaches which estimate the polynomial function range by naive interval arithmetic over sub-domains, the newly proposed algorithm demonstrates dramatically improved efficiency in rigorous solution estimation.

## References

[1] Discussion of FEniCS developing team on the $L^\infty$ norm feature: url: https://bitbucket.org/fenics-project/dolfin/issues/663/add-linf-norm-for-functions

[2] FARIN, GERALD & HANSFORD, DIANNE (2000). The Essentials of CAGD. Natic, MA: A K Peters, Ltd. ISBN 1-56881-123-3

[3] XUEFENG LIU. Verified Finite Element Method library as MATLAB package, URL: https://ganjin.online/xfliu/NS.SolutionVerification, May, 2021

# B-matrices and their generalizations in the interval setting

Matyáš Lorenc[1]

[1] Charles University, Prague, Czech Republic
Charles University, Faculty of Mathematics and Physics, Department of Applied Mathematics,
Malostranské nám. 25, 11800, Prague, Czech Republic
Mafi412@seznam.cz

In 1968 mathematicians Cottle and Dantzig proposed the linear complementarity problem, denoted $LCP(M, q)$, where $M$ is a matrix and $q$ a vector. It was shown that $LCP(M, q)$ has a unique solution for every vector $q$ if and only if $M$ is a P-matrix, i.e. all its principal minors are positive. However, this class of matrices is computationally complex to recognize, the task of verifying given matrix on being a P-matrix is co-NP-complete. This leads us to try and define several classes of P-matrices that are easily recognizable. Such classes are e.g. B-matrices (introduced by Peña in [1]), doubly B-matrices (introduced also by Peña in [2]) or $B_\pi^R$-matrices (introduced by Neumann, Peña and Pryporova in [3]). The B-matrices and doubly B-matrices found their use in localization of eigenvalues, as shown by Peña in [1] and [2].

Ever since the beginning of rigorous measurements, mathematicians had to deal with inaccuracy in data or any form of uncetainty they may encounter. And one of the answers to this problem is interval analysis, sometimes called interval computing or interval mathematics as well.

In our work we generalize our special subclasses of P-matrices, those being B-matrices, doubly B-matrices and $B_\pi^R$-matrices, into interval settings, thus interconnecting these two topics. We lay grounds to recognizing the interval variants through characterizations, both through some property they posses or via reduction to finite number of real instances ($n$ real matrices for B-matrices, $n^3$ instances for doubly B-matrices and $n$ matrices for $B_\pi^R$-matrices), and we derive necessary conditions and sufficient ones. Also we take a closer look at some of their properties, be it some fundamental ones (e.g. some conditions that the matrix entries must fulfill), or closure properties. What is interesting is that whereas the complexity of characterizations of interval B-matrices and interval $B_\pi^R$-matrices is the same as that of the real cases, in this case $O(n^2)$, for interval doubly B-matrices it is $O(n^4)$ compared to $O(n^2)$ for the real case.

## References

[1] Peña, J. M. A class of $P$-matrices with applications to the localization of the eigenvalues of a real matrix. *SIAM J. Matrix Anal. Appl.*, 22(4):1027–1037, 2001.

[2] Peña, J. M. On an alternative to Gerschgorin circles and ovals of Cassini. *Numer. Math.*, 95(2):337–345, 2003.

[3] Neumann, Michael and Peña, J.M. and Pryporova, Olga. Some classes of nonsingular matrices and applications. *Linear Algebra Appl.*, 438(4):1936–1945, 2013.

# Interval methods for packing problems on the sphere

Mihály Csaba Markót[1]

[1] University of Vienna, Vienna, Austria
Oskar-Morgenstern-Platz 1, A-1090 Vienna, Austria
`mihaly.markot@univie.ac.at`

**Keywords:** interval arithmetic, global optimization, computer-assisted proof, sphere packing

The subject of the talk is the development of computer-assisted proofs for solving optimal point packing problems on the sphere. During the last two decades we established numerous tools for solving point packing problems in two dimensions [1, 2]. Although some of these tools can be generalized to three dimensions, the development of the basic interval data structures and the key algorithmic components requires a whole new way of thinking.

In the talk a proposed interval representation of the subsets of the search space is introduced. These sets will be called spherical intervals. Then a new method for eliminating suboptimal regions from these spherical intervals will be given. This method is actually a constraint propagation technique applied to the distance constraints written for the pairs of spherical intervals.

The power of the new methodology is demonstrated by solving the famous 3-dimensional kissing number problem (also known as the 13 spheres problem) on a computer. During the solution process an important symmetry breaking idea, known from two-dimensional packing methods, is applied, which divides the original search space into regions such that each region can contain at most one point of an optimal configuration.

The developed methods provide a good starting point for solving instances of various related optimization problems. These include the Tammes problems (maximizing the minimal pairwise distance of $n$ points on the sphere), the Fekete point problems (minimizing a potential energy function computed from the pairs of points on the sphere), the 'double kissing' problem (a generalization of the original kissing number problem, involving two touching, identical spheres), and even various atom cluster and molecular modeling problems (involving multiple, not necessarily identical spheres).

## References

[1] M. C. MARKÓT AND T. CSENDES: A New Verified Optimization Technique for the "Packing Circles in a Unit Square" Problems, *SIAM J. Optimization*, 16 (2005), 193–219.

[2] M. C. MARKÓT: Improved Interval Methods for Solving Circle Packing Problems in the Unit Square, to appear in *J. Global Optimization*.

# A Trick for an Accurate $e^{-|x|}$ Function in Fixed-Point Arithmetics

## Mantas Mikaitis

Department of Mathematics, University of Manchester, Manchester, UK
mantas.mikaitis@manchester.ac.uk

We demonstrate a method of improving the accuracy of the exponential decay function in fixed-point arithmetic.

Consider two standard fixed-point number formats from the ISO/IEC 18037:2008 embedded C standard[2]: `s16.15` (sign, 16 integer bits, 15 fractional bits) and `s0.31` (sign, zero integer bits, 31 fractional bits—maximized accuracy of values in $(-1, 1)$). Assume we have a standard implementation of $y = e^x$ where $x$ and $y$ are `s16.15`. The input domain is $x \in (\log_e(2^{-15}) = -10.397..., \log_e(2^{16}-2^{-15}) = 11.09...)$. Outside of this range the output underflows or overflows. If we take $x$ to be `s16.15` but $y$ `s0.31`, the input domain changes to $x \in (\log_e(2^{-31}) = -21.487..., 0)$.

Consider implementing the exponential decay which requires the computation of $e^{-|x|}$. Since inputs are always negative, excluding zero as a special case, the output domain of the function will be $[0, 1)$ (including zero produced on underflow). We noticed that no integer bits are required for representing this range of values and thus in the `s16.15` outputs from the exponential function the top 16 bits after the binary point are not used. In this case it would be better to have inputs as `s16.15` and outputs as `s0.31`. This can be achieved by a new implementation of $e^x$. However, in some cases it might not be possible, for example if it is already a hardware routine that outputs `s16.15`. The following method allows to generate more accurate `s0.31` outputs with the `s16.15` exponential function without modifying it.

To obtain the exponential decay output $e^x$ as `s0.31`, we need to arrive at $2^{16}e^x$ (exponential in `s16.15` shifted 16 places left, which gives the same value when the bits are interpreted as `s0.31`). However, this cannot be done by simply running the arithmetic algorithm for calculating $e^x$ and then shifting the output as that will place 0's at the bottom part and no accuracy improvement will be achieved. The following exponential function property is of interest: $2^{16}e^x = e^{\log_e(2^{16})}e^x = e^{16 \times \log_e(2)+x}$. Now we have $2^{16}e^x$ as `s16.15` or $e^x$ as `s0.31` when the binary point location is interpreted to be located at 16 bits to the left. Note that this is achieved not by shifting but by manipulating the exponent $x$ (add the constant $\log_e(2) \times 16$) to get the same effect as first calculating $e^x$ and then shifting 16 places left, without propagating 0's at the bottom end but letting the exponential algorithm fill in the bottom bits as best as it can.[3] This method is illustrated in Figure 1.



Figure 1: Illustration of the standard one precision use of `s16.15` exponential (top) and the required input transformation for obtaining `s0.31` exp decay using the `s16.15` exponential function (bottom).

We have implemented this method in MATLAB and measured errors (compared with MATLAB's `exp`) of both the `s16.15` exponential decay and the method of obtaining a `s0.31` decay with the same `s16.15` function. (Figure 2). The results show that in the original input domain the exponential function

---

[2] https://www.iso.org/standard/51126.html
[3] Note that the format `u0.32` was not chosen because $17 \times \log_e(2)+x$ term would cause saturation in the `s16.15` exponential function when $x \geq \log_e(0.5)$.

Figure 2: Errors of `s16.15` $e^x$ (black) and the mixed-precision method introduced here (blue).

output is more accurate. Furthermore, the input domain in which the outputs do not suffer underflow is wider. As per Figure 1, this is achieved by adding a constant to the input, using the `s16.15` exponential function unchanged, and interpreting the `s16.15` output as `s0.31`.

# Verified bounds for matrix gamma function

## Shinya Miyajima

Iwate University, Iwate, Japan
4-3-5 Ueda, Morioka, Iwate 020-8551, Japan
`miyajima@iwate-u.ac.jp`

**Keywords:** matrix gamma function, verified block diagonalization, verified numerical computation

For $z \in \mathbb{C}$ with positive real part, the gamma function is defined by

$$\Gamma(z) := \int_0^\infty e^{-t} t^{z-1} dt,$$

and otherwise by analytic continuation. It is well known that $\Gamma(z)$ is analytic everywhere in $\mathbb{C}$, with the exception of non-positive integer numbers $\mathbb{Z}_-$. Therefore, the general theory of primary matrix functions ensures that the matrix gamma function $\Gamma(A)$ is well defined for $A \in \mathbb{C}^{n \times n}$ having no eigenvalues on $\mathbb{Z}_-$. If all eigenvalues of $A$ have positive real parts, then we have the representation
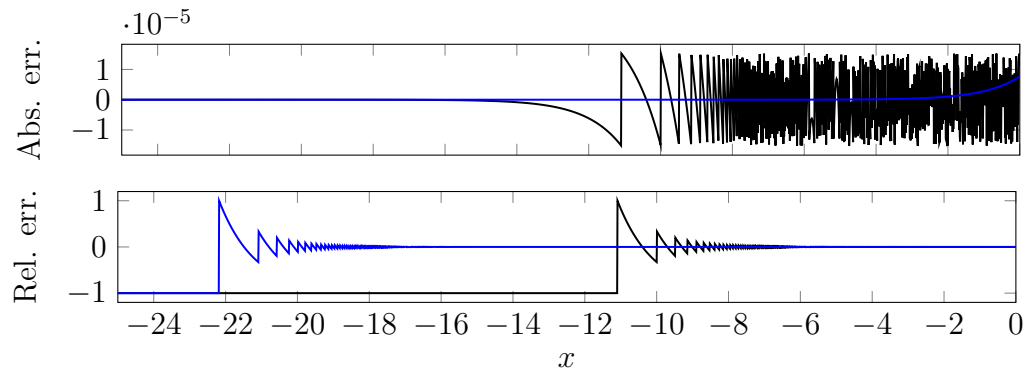
$$\Gamma(A) = \int_0^\infty e^{-t} t^{A-I_n} dt,$$

where $t^{A-I_n} := e^{(A-I_n)\log(t)}$ and $I_n$ denotes the $n \times n$ identity matrix.

The function $\Gamma(A)$ has connections with other special functions, which play an important role in solving certain matrix differential equations [1]. Two of these special functions are the matrix beta and Bessel functions. In [1], mathematical properties of $\Gamma(A)$ are elegantly clarified, and fast and accurate algorithms for computing $\Gamma(A)$ are proposed.

The work presented in this talk addresses the problem of verified computations for $\Gamma(A)$, specifically, numerically computing interval matrices which are guaranteed to contain $\Gamma(A)$. There are many sophisticated verification algorithms for matrix functions. To the author's best knowledge, on the other hand, a verification algorithm designed specifically for $\Gamma(A)$ has not yet appeared in the literature. A possible method is to use the VERSOFT routine `vermatfun`. This routine is applicable not only to the matrix gamma function but also to other matrix functions, and computes the interval matrices by enclosing all the eigenvalues and eigenvectors of $A$ via the INTLAB routine `verifyeig`. The routine `vermatfun` fails when $A$ is defective or close to defective, and requires $\mathcal{O}(n^4)$ operations.

The purpose of this talk is to propose two verification algorithms for $\Gamma(A)$. In [2], algorithms for enclosing all the eigenvalues and basis of invariant subspaces of $A$ are presented. As byproducts of these algorithms, we can obtain interval matrices containing blocks whose spectrums are included in that of $A$. Moreover, the spectrum of $A$ is also contained in the union of the spectrums of the blocks. In this talk, we interpret the interval matrices containing the basis and blocks as a result of verified block diagonalization (VBD), and develop the verification algorithms. To achieve enclosure for the gamma function of the blocks, we derive normwise computable perturbation bounds. Here, the word "computable" means that we can numerically obtain a rigorous upper bound which takes rounding and truncation errors into account. We can apply the derived perturbation bounds if disks containing the spectrums of input matrices lie in the open right half plane. We incorporate matrix argument reductions (ARs) to force the input matrices to have this property, and develop theories for accelerating the ARs. The first algorithm uses the VBD based on a numerical spectral decomposition, and involves only $\mathcal{O}(n^3)$ operations if the total computational cost of the accelerated ARs is $\mathcal{O}(n^3)$. The second algorithm adopts the VBD based on a numerical Jordan decomposition, is applicable even when $A$ is defective, and requires $\mathcal{O}(n^4)$ operations. We present a theory for verifying that $A$ has no eigenvalues on $\mathbb{Z}_-$. By the aid of this theory, these algorithms do not assume but prove that $A$ has no eigenvalues on $\mathbb{Z}_-$. The first and second algorithms require intervals containing $\Gamma^{(0)}(z)/0!, \ldots, \Gamma^{(\ell)}(z)/\ell!$, where $\ell$ is a non-negative integer and $z \in \mathbb{C}$. To the author's best knowledge, an algorithm for computing such intervals is not available in literature, whereas there are well-established algorithms for computing intervals containing *real* scalar gamma functions. We thus present a way for computing such intervals, which is based on the Spouge approximation. Numerical results show efficiency and robustness of the algorithms.

## References

[1] J. R. CARDOSO, A. SADEGHI: Computation of matrix gamma function, *BIT*, 59 (2019), 343–370.

[2] S. MIYAJIMA: Fast enclosure for all eigenvalues and invariant subspaces in generalized eigenvalue problems, *SIAM J. Matrix Anal. Appl.*, 35 (2014), 1205–1225.

# Computing enclosure for matrix real powers

## Shinya Miyajima

Iwate University, Iwate, Japan
4-3-5 Ueda, Morioka, Iwate 020-8551, Japan
`miyajima@iwate-u.ac.jp`

**Keywords:** matrix real power, verified block diagonalization, verified numerical computation

Let $\alpha \mathbb{R}$, and $A \in \mathbb{C}^{n \times n}$ have no non-positive real eigenvalues. In this talk, we are concerned with accuracy of numerically computed result for $A^\alpha$, which is uniquely defined as follows:

**Definition 1** (Higham and Lin [1, Definition 1.1]). *Let $\alpha \mathbb{R}$, and $A \in \mathbb{C}^{n \times n}$ have no non-positive real eigenvalues. Then, $A^\alpha = e^{\alpha \log(A)}$, where $\log(A)$ is the principal logarithm of A.*

If $\alpha = 1/d$, with $d$ a positive integer, then $A^\alpha$ reduces to the principal $d$th root of $A$. If $A$ has no non-positive real eigenvalues and $\alpha \in (0, 1)$, then we have the representation

$$A^\alpha = \frac{\sin(\alpha\pi)}{\alpha\pi} A \int_0^\infty (t^{1/\alpha} I_n + A)^{-1} dt,$$

where $I_n$ denotes the $n \times n$ identity matrix.

Computing $A^\alpha$ has many applications such as Markov chain models in health-care and finance, fractional partial differential equations, discrete representations of norms corresponding to finite element discretization of fractional Sobolev spaces, and the computation of geodesic-midpoints in neural networks. Efficient numerical algorithms for computing $A^\alpha$ are proposed (see [1], e.g.).

In this talk, we consider verified computations for $A^\alpha$, specifically, numerically computing interval matrices which are guaranteed to contain $A^\alpha$. Although there are many well-established verification algorithms for matrix functions, a verification algorithm designed specifically for $A^\alpha$ seems to be unavailable in the literature. If $\alpha$ is rational, then $\alpha$ can be written as $\alpha = c/d$, where $c, d \in \mathbb{Z}$ and $d > 0$, so that $A^\alpha = (A^{1/d})^c$. Therefore, the verification algorithm for the matrix principal $d$th root is applicable. However, we cannot adopt this approach when $\alpha$ is irrational. One of possible methods is to use the VERSOFT routine `vermatfun`. This routine is applicable not only to the matrix real power but also to other matrix functions, and computes the interval matrices by enclosing all the eigenvalues and eigenvectors of $A$ via the INTLAB routine `verifyeig`. The routine `vermatfun` fails when $A$ is defective or close to defective, and requires $\mathcal{O}(n^4)$ operations. Since $A^\alpha = e^{\alpha \log(A)}$, we can compute the interval matrices by combining the algorithm for enclosing the matrix exponentials and principal logarithms. On the other hand, this combination tends to give interval matrices whose radii are not small. This seems to be because an interval matrix containing $\alpha \log(A)$ is computed in the first stage, and its radius is enlarged in the process of enclosing the matrix exponential.

The purpose of this talk is to propose two verification algorithms for $A^\alpha$. In [2], algorithms for enclosing all the eigenvalues and basis of invariant subspaces of $A$ are presented. As byproducts of these algorithms, we can obtain interval matrices containing blocks whose spectrums are included in that of $A$. Moreover, the spectrum of $A$ is also contained in the union of the spectrums of the blocks. In this talk, we interpret the interval matrices containing the basis and blocks as a result of verified block diagonalization (VBD), and develop the verification algorithms. To achieve enclosure for the real power of the blocks, we derive componentwise computable perturbation bounds. Here, the word "computable" means that we can numerically obtain a rigorous upper bound which takes all possible errors into account. The first algorithm uses the VBD based on a numerical spectral decomposition, and involves only $\mathcal{O}(n^3)$ operations if $|\alpha|$ is not too large. The second algorithm adopts the VBD based on a numerical Jordan decomposition, is applicable even when $A$ is defective, and requires $\mathcal{O}(n^4)$ operations. We present a theory for verifying that $A$ has no non-positive real eigenvalues, which is based on checking nonsingularity of a matrix. By the aid of this theory, these algorithms do not assume but prove that $A$ has no non-positive real eigenvalues. Numerical results show effectiveness of the algorithms.

# References

[1] N. J. HIGHAM, L. LIN: A Schur-Padé algorithm for fractional powers of a matrix, *SIAM J. Matrix Anal. Appl.*, 32 (2011), 1056–1078.

[2] S. MIYAJIMA: Fast enclosure for all eigenvalues and invariant subspaces in generalized eigenvalue problems, *SIAM J. Matrix Anal. Appl.*, 35 (2014), 1205–1225.

# Kd-tree based adaptive interpolation algorithm for modeling dynamic systems with interval parameters

Alexander Y. Morozov[1,2]
Dmitry L. Reviznikov[1,2]

[1]Federal Research Center Computer Science and Control of the Russian Academy of Sciences (FRC CSC RAS), Russia, Moscow
[2]Moscow Aviation Institute (National Research University) (MAI), Russia, Moscow
morozov@infway.ru, reviznikov@mai.ru

When solving applied and research problems, there often arise situations when certain parameters are not exactly known, but there is information about their ranges [1]. For such problems, it is necessary to obtain an interval estimate of the solution based on interval values of parameters. For dynamic systems such problems are formulated in the form of a Cauchy problem for a system of ordinary differential equations (ODEs) with interval initial conditions or interval parameters [2]. There are many works devoted to the interval methods. Some methods are based on the representation of the set of solutions through geometric primitives, for example, parallelepipeds and ellipses; there are also methods based on interval arithmetic; methods based on symbolic calculations; stochastic methods such as the Monte Carlo method. The adaptive interpolation algorithm [3, 4] is an alternative to existing methods. Though it lacks the guarantee property inherent in interval methods it obtains the boundaries of solutions with a given accuracy, has a high degree of parallelization and works faster than most of other methods.

The main idea of the algorithm is to build a dynamic structured grid based on kd-tree over the set formed by the interval parameters. Each vertex of the tree is a regular interpolation grid corresponding to a given degree of interpolation polynomial. A solution that is found using the parameters determined by the position of the node in the parametric space, is associated with the node. In the computational process, at each step of the ODE system integration, a piecewise polynomial function that interpolates the dependence of the solution on the interval parameters is constructed. The calculation of the interval estimation of the solution comes down to computation of the range of interpolating function values.

One iteration of the algorithm consists of two stages. At the first stage, all solutions associated with nodes are recalculated to the next "temporary" layer using some numerical integration method. At the second stage, the kd-tree is rebuilt according to the principle of minimizing the interpolation error. The interpolation error of the solution is estimated for each vertex. If at some leaf vertex the interpolation error exceeds some given value, then the vertex is divided into two ones. In this case, the construction of a new partition is carried out at the previous step, and the integration step is repeated. This is due to the fact that when a new vertex is formed, the values associated with the nodes of the new grids are interpolated, and they need to be calculated at time when the error is still valid. If the error becomes acceptable for the vertex and all its descendants, then the descendants are deleted, and the vertex itself becomes a leaf.

Statements regarding the adaptive interpolation algorithm properties are formulated and proved. It is shown that the global error estimate is directly proportional to the height of the kd-tree. Testing the algorithm on a representative series of problems of different dimensions with a different number of interval parameters (including problems of chemical kinetics, gas dynamics, problems with bifurcations and dynamic chaos) demonstrates its effectiveness.

## References

[1] S. P. SHARY: *Maximum compatibility method for data fitting under interval uncertainty*, Journal of Computer and Systems Sciences International, 2017, Vol. 56, No. 6, pp. 897–913. DOI: 10.1134/S1064230717050100

[2] R. E. MOORE: *Interval Analysis*, Englewood Cliffs: Prentice Hall, 1966.

[3] A. Y. Morozov, D. L. Reviznikov: *Adaptive Interpolation Algorithm Based on a kd-Tree for Numerical Integration of Systems of Ordinary Differential Equations with Interval Initial Conditions*, Differential Equations, 2018, Vol. 54, No. 7, pp. 945-956., DOI: 10.1134/S0012266118070121

[4] A. Yu. Morozov, D. L. Reviznikov, V. Yu. Gidaspov: *Adaptive Interpolation Algorithm Based on a KD-Tree for the Problems of Chemical Kinetics with Interval Parameters*, Mathematical Models and Computer Simulations, 2019, Vol. 11, No. 4, pp. 622–633., DOI: 10.1134/S2070048219040100

# Inverse Bifurcation Diagram Problem of Forced El Niño Equation

Shin'ichi Oishi[1], and Kouta Sekine[2]

[1] Department of Applied Mathematics, Faculty of Science and Engineering, Waseda University, and JST CREST, Tokyo 169-8555 Japan
[2] Faculty of Information Networking for Innovation and Design, Toyo University, Tokyo 115-0053 Japan
oishi@waseda.jp

**Keywords:** Subharmonic solutions, Almost periodic solutions

In 1988, Suarez and Schopf [1] have introduced a *delayed action oscillator* equation

$$\frac{dx(t)}{dt} - x(t) + x^3(t) + \alpha x(t - \tau) = 0 \tag{1}$$

as a simple model of El Niño. This paper proposes to add a seasonal forcing term to SS equation and consider the following forced delay action oscillator equation

$$\frac{dx(t)}{dt} - x(t) + x^3(t) + \alpha x(t - \tau) - \beta \cos \omega t = 0, \tag{2}$$

which we will call forced Suarez and Schopf's equation, or fSS equation in short. The variable $x$ represent a deviation of a sea surface temperature near Peru from average. The term $f(x) = -x(t) + x^3(t)$ represents an effect of energy exchange between sea surface and atmosphere. The term $\alpha x(t - \tau)$ represents an effect of delay by wave propagation on the equator from this area to the east end of Asia (the eastward Kelvin wave), reflected at there, and reflecting back to near Peru (the westward Rossby wave). The parameter $\tau (> 0)$ expresses a turn around time of these Kelvin and Rossby waves propagation. Here, the parameters $\beta (\geqq 0)$ and $\omega (> 0)$ express the strength of the effect of the seasonal force, and the angular frequency of the seasonal force, respectively.

One of the main purpose of this paper is to show that fSS equation exhibits various complicated dynamics. Fig. 1 shows a stroboscopic bifurcation diagram obtained by numerical integration. We consider a problem of finding what kind of solutions generate such a bifurcation diagram. We will call this problem as an inverse bifurcation diagram problem. The 1st, the 4th and 5th figures of Fig. 1 display a solution of this problem, which we will call an inverse bifurcation diagram. This diagram mainly consists of Galerkin's approximate solution branches of periodic solutions. *Throughout this paper, in the inverse bifurcation diagram, the odd symmetric 1-periodic solution[4] branch is labeled by 's' or '1s'. On the other hand, the branch consisting of asymmetric 1-periodic solutions is labeled by 'a' or '1a'. The label '1/n', or simply 'n' indicates an $1/n$ subharmonic solution. If further the symbol 's' or 'a' is concatenated, such a periodic solution is odd-symmetric or not odd-symmetric, respectively.* We choose a representative from each solution branch and prove the existence of exact periodic solutions nearby such representatives via the method of computer assisted existence proof of periodic solutions proposed in Ref [2]. The red line of the 2nd and 4th figures represent time average $l_2$ norm of solutions. The final two figures in Fig. 1 show Galerkin's approximations of almost periodic solutions of fSS equation. In these cases, the red lines represent Poincaré plots.

## References

[1] M. J. SUAREZ AND P. S. SCHOPF: A delayed action oscillator for ENSO, *Journal of the Atmospheric Sciences*, **45(21)** (1988) pp. 3283-3287.

[2] S. OISHI AND K. SEKINE: Inclusion of periodic solutions for forced delay differential equation modeling El Niño *Nonlinear Theory and Its Applications, IEICE*, vol. **12**, no. 3, pp.1-36.

---

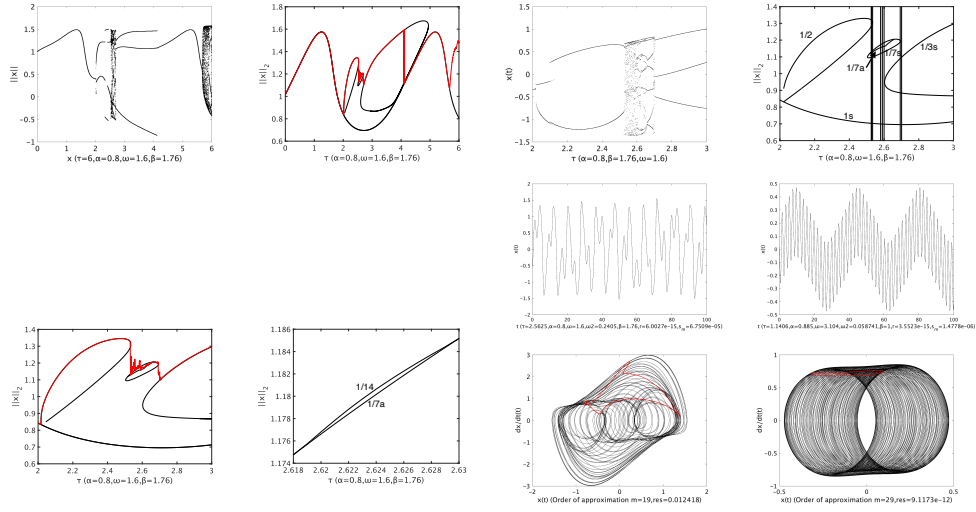[4] We call a $2\pi$ periodic solution as 1-periodic.

Figure 1: Stroboscopic bifurcation diagram and inverse bifurcation diagram

# Improvement of selection formulas of mesh size and truncation number for the DE-Sinc approximation and its theoretical error bound

Tomoaki Okayama[1] and Shota Ogawa[2]

[1] Hiroshima City University, Hiroshima, Japan
[2] Systemers Inc., Tokyo, Japan
okayama@hiroshima-cu.ac.jp

**Keywords:** Sinc approximation, double-exponential transformation, error bound

The Sinc approximation is based on Shannon's sampling formula

$$F(x) \approx \sum_{k=-\infty}^{\infty} F(kh) \operatorname{sinc}\left(\frac{x-kh}{h}\right), \tag{1}$$

where $\operatorname{sinc} x = \sin(\pi x)/(\pi x)$, and $h$ is a mesh size. If $F$ is analytic and absolutely integrable on $\mathcal{D}_d = \{z \in \mathbb{C} : |\operatorname{Im} z| < d\}$ $(d > 0)$, this approximation achieves exponential convergence: $\mathrm{O}(\mathrm{e}^{-\pi d/h})$. Furthermore, if $|F(x)|$ decays exponentially as $x \to \pm\infty$, i.e., $|F(x)| = \mathrm{O}(\mathrm{e}^{-\alpha|x|})$ $(\alpha > 0)$, the infinite sum in (1) may be truncated at some truncation number $N$ as

$$\sum_{k=-\infty}^{\infty} F(kh) \operatorname{sinc}\left(\frac{x-kh}{h}\right) \approx \sum_{k=-N}^{N} F(kh) \operatorname{sinc}\left(\frac{x-kh}{h}\right), \tag{2}$$

where its error rate is $\mathrm{O}(\mathrm{e}^{-\alpha N h})$. In view of the approximation error in (1) (called discretization error) and that in (2) (called truncation error), the optimal formula of the mesh size $h$ with respect to the truncation number $N$ is given by

$$\mathrm{e}^{-\pi d/h} = \mathrm{e}^{-\alpha N h} \quad \Longleftrightarrow \quad h = \sqrt{\frac{\pi d}{\alpha N}},$$

which gives $\mathrm{O}(\mathrm{e}^{-\sqrt{\pi d \alpha N}})$ as a final (overall) error rate. The Sinc approximation for exponentially decaying functions is called the SE-Sinc approximation, which has extensively been developed and analyzed by Stenger [1].

Recently, it has been reported [2] that if $|F(x)|$ decays double-exponentially as $x \to \pm\infty$, i.e., $|F(x)| = \mathrm{O}(\mathrm{e}^{-(\pi/2)\alpha \exp(|x|)})$, the approximation error in (2) is $\mathrm{O}(\mathrm{e}^{-(\pi/2)\alpha \exp(Nh)})$, which is far smaller than that in Stenger's results. The Sinc approximation for double-exponentially decaying functions is called the DE-Sinc approximation. In this case, the optimal formula of the mesh size $h$ with respect to the truncation number $N$ should be given by the equation

$$\mathrm{e}^{-\pi d/h} = \mathrm{e}^{-(\pi/2)\alpha \exp(Nh)}, \tag{3}$$

but the solution $h$ cannot be expressed in terms of elementary functions of $N$. Instead, several authors have employed a near-optimal formula

$$h = \frac{\log(2dN/\alpha)}{N}.$$

By the formula, the left-hand side in (3) becomes $\mathrm{e}^{-\pi dN/\log(2dN/\alpha)}$, and the right-hand side in (3) becomes $\mathrm{e}^{-\pi dN}$. Although those two rates are approximately equal, the former rate is a bit lower than the latter rate, which gives $\mathrm{O}(\mathrm{e}^{-\pi dN/\log(2dN/\alpha)})$ as a final error rate of the DE-Sinc approximation. Furthermore, aiming for a computation with guaranteed accuracy, computable error bounds for both SE- and DE-Sinc approximations has been already given [3].

The main contribution of this study is to propose the optimal formulas of $h$ and $N$ for the DE-Sinc approximation, and to give its computable error bound. We introduce an integer $n$ as a parameter, and select $h$ and $N$ as

$$h = \frac{\operatorname{arcsinh}(dn/\alpha)}{n}, \quad N = \left\lceil n \cdot \frac{\operatorname{arcsinh}(q(dn/\alpha))}{\operatorname{arcsinh}(dn/\alpha)} \right\rceil,$$

where $q(x) = x/\operatorname{arcsinh} x$. Then, both discretization and truncation errors become exactly the same rate: $\mathrm{O}(\mathrm{e}^{-\pi dn/\operatorname{arcsinh}(dn/\alpha)})$. Numerical comparison with the existing formulas will also be given in this talk.

## References

[1] F. STENGER: *Handbook of Sinc Numerical Methods*, CRC Press, Boca Raton, FL, 2011.

[2] M. SUGIHARA, T. MATSUO: Recent developments of the Sinc numerical methods, *Journal of Computational and Applied Mathematics*, 164/165 (2004), 673–689.

[3] T. OKAYAMA, T. MATSUO, M. SUGIHARA: Error estimates with explicit constants for Sinc approximation, Sinc quadrature and Sinc indefinite integration, *Numerische Mathematik*, 124 (2013), 361–394.

# Quantification of Time-Domain Truncation Errors for the Reinitialization of Fractional Integrators

Andreas Rauh[1] and Rachid Malti[2]

[1] ENSTA Bretagne, Lab-STICC, 29806 Brest, France
[2] IMS Laboratory, University of Bordeaux, 33405 Talence, France
Andreas.Rauh@interval-methods.de, rachid.malti@ims-bordeaux.fr

FDEs are powerful modeling tools in many engineering applications in which non-standard dynamics, characterized by infinite horizon states, can be observed. An example for such applications is modeling the charging and discharging dynamics of batteries [1]. Previous work for an interval-based state estimation of such systems has accounted for a cooperativity preserving or cooperativity enforcing design of observers [1, 3]. These interval observers exploit specific monotonicity properties of positive dynamic systems and provide lower and upper bounding trajectories for all pseudo-state variables as soon as suitable initialization functions are specified. Moreover, interval-valued iteration procedures have been developed [5] for a verified simulation of such systems. The latter, based on Mittag-Leffler function parameterizations of the pseudo-state enclosures, are not *a priori* restricted to cooperative models but are applicable also to nonlinear systems with interval parameters.

However, the evaluation of observer-based pseudo-state estimation procedures for continuous-time fractional models supposes that measurements are also available in a continuous-time form or at least at each sampling period [3]. For many practical applications, this is not the case, so that continuous-time pseudo-state predictions need to be performed between the discrete time instants at which measurements are available. Then, the measured pseudo-state information (described by intervals to represent bounded measurement errors) can be intersected with the predicted state information to enhance the knowledge of the actual system dynamics. However, this intersection demands reinitializing the integration of the fractional model. Similar requirements are discussed in [5], where temporal sub-slices were considered to reduce the overestimation of interval-based simulation approaches.

Due to the infinite horizon memory property of fractional systems, the reinitialization of time-domain simulations requires a rigorous consideration of the arising truncation errors. Guaranteed outer bounds for these errors were derived in [4]. These bounds are the basis for a novel error refinement strategy between discrete reinitialization points in an observer-based setting. In this contribution, we discuss the following aspects:

1. Expressing non-constant pseudo-state initializations from a bounded past time window in terms of uncertain initial conditions at a single point with a conservative interval-valued correction of the FDE model;

2. Implementation of an observer-based quantification of truncation errors for simulations with periodic reinitialization (e.g. based on the floating point MATLAB routines from [2]);

3. Performing an interval contractor-based state estimation of a continuous-time battery model [1] with discrete-time measurements;

4. Describing possible interfaces with the verified simulation routines from [5] as an outlook for future work.

## References

[1] E. HILDEBRANDT, J. KERSTEN, A. RAUH, H. ASCHEMANN: Robust interval observer design for fractional-order models with applications to state estimation of batteries, *IFAC-PapersOnLine*, 53,2 (2020), 3683–3688.

[2] R. GARRAPPA: Predictor-corrector PECE method for fractional differential equations, *MATLAB Central File Exchange*, https://www.mathworks.com/matlabcentral/fileexchange/32918-predictor-corrector-pece-method-forfractional-differential-equations, accessed: May 09, 2021.

[3] G. BEL HAJ FREJ, R. MALTI, M. AOUN, T. RAÏSSI: Fractional interval observers and initialization of fractional systems, *Communications in Nonlinear Science and Numerical Simulation*, 82 (2020), 105030.

[4] I. PODLUBNY: *Fractional Differential Equations: An Introduction to Fractional Derivatives, Fractional Differential Equations, to Methods of Their Solution and Some of Their Applications*, Mathematics in Science and Engineering, Academic Press, London, UK, 1999.

[5] A. RAUH, L. JAULIN: Novel techniques for a verified simulation of fractional-order differential equations, *Fractal Fract*, 5 (2021), 17.

# Convergent Real Matrix Powers with Divergent Results in Interval Arithmetic

## Nathalie Revol[1]

[1] University of Lyon, INRIA, LIP-ENS de Lyon, France
46 allée d'Italie, 69364 Lyon Cedex 07, France
Nathalie.Revol@inria.fr

**Keywords:** interval arithmetic, matrix powers, wrapping effect

In control theory, as well as in other fields, one regularly encounters iterated powers of a given matrix $A$. The precise definition of the problem we study is given by Mayer and Warnke in [2]. While these iterations converge, when exact arithmetic over the reals is considered, that is, $A^k \to 0$ as $k \to +\infty$, divergence may occur when floating-point or interval arithmetic is used. This phenomenon and some ways to circumvent it are described by Lohner in [1]. We will explore, experiment and as much as possible revisit some of the solutions proposed in [1]. More precisely, we will explore the use of a preconditioner based on an approximate SVD of $A$, and the use of a step $k_0$ such that the iterations using $A^{k_0}$ converge, even when floating-point or interval arithmetic is employed.

## References

[1] R. J. LOHNER: On the Ubiquity of the Wrapping Effect in the Computation of Error Bounds, *Perspectives on Enclosure Methods*, Springer, Vienna, (2001), 201–216.

[2] G. MAYER, I. WARNKE: On the fixed points of the interval function $[f]([x]) = [A][x] + [b]$, *Linear Algebra and Applications*, 363 (2003), 201–216.

# Inverse bifurcation diagram problem for delayed van der Pol-Duffing equation

Yuuki Saito[1], Naoki Takamatsu[1] Shin'ichi Oishi[2], and
Kouta Sekine[3]

[1] Graduate School of Fundamental Science and Engineering,
Waseda University
[2] Department of Applied Mathematics, Faculty of Science and Engineering, Waseda University, Japan
[3] Faculty of Information Networking for Innovation and Design,
Toyo University, Japan
yuki-swim7.wu@toki.waseda.jp

**Keywords:** van der Pol-Duffing equation, Computer-assisted proof

In this talk, we are concerned with the bifurcation problem of the delay van der Pol-Duffiing equation

$$\frac{d^2x(t)}{dt^2} + k(x^2(t)-1)\frac{dx(t)}{dt} + \mu x(t) + \gamma x^3(t) - \alpha x(t-\tau) - \beta\cos\omega t = 0, \tag{1}$$

where $k > 0, \tau \geq 0, \mu \geq 0, \gamma > 0, \alpha > 0$ and $\beta > 0$ are parameters. It is known that subharmonic solutions and chaotic aperiodic solutions have been observed numerically for the equation (1) (see, e.g., [1]).

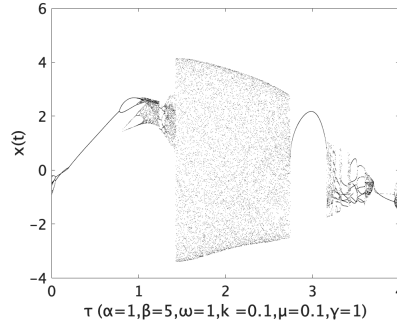Fig. 1 shows a stroboscopic bifurcation diagram obtained by numerical integration. We consider a



Figure 1: Bifurcation diagram

problem of finding what kind of solutions generate such a bifurcation diagram. We will call this problem as an inverse bifurcation diagram problem.

Figure 1 exhibits complicated dynamics, *i.e.*, we can observe $2\pi$-periodic solutions branches, subharmonic solution branches, and aperiodic solution branches. To solve the inverse bifurcation problem, in the first place, we have calculated approximate $2\pi$-periodic and subharmonic solutions using Galerkin's method. Then, using the continuation method, we have calculated approximate periodic solution branches. For each approximate periodic solution branch we have picked upped a subharmonic solution. Then, extending a verification method presented in Ref [2], the existence of an exact subharmonic solution is verified nearby the approximated subharmonic solution via verified numerical computations. Fig.2 is a solution of the inverse bifurcation diagram problem associated with the bifurcation diagram shown in Fig.1. In this figure, the label '$1/n$' indicates an $1/n$ subharmonic solution (branch). If further the symbol 's' or 'a' is concatenated, such a periodic solution is odd-symmetric or not odd-symmetric, respectively.

Figure 2: Inverse bifurcation diagram

## References

[1] YOSHIHIRO TSUDA, HIDEYUKI TAMURA, ATSUTO SUEOKA AND TSUYOSHI FUJII: Chaotic Behavior of a Nonlinear Vibrating System with a Retarded Argument (Characteristics in the Region of Subharmonic Resonance) *JSME international journal. Ser. 3, Vibration, control engineering, engineering for industry*, 35(2) (1992) pp. 259-267.

[2] SHIN'ICHI OISHI, KOUTA SEKINE: Inclusion of Periodic Solutions for Forced Delay Differential Equation Modeling El Niño, *Nonlinear Theory and Its Applications, IEICE*, vol. 12, no. 3, (2021), pp.1-36.

# Branch-and-bound interval methods and constraint propagation on the GPU using Julia

David P. Sanders[1,2,4] and Valentin Churavy[3,4]

[1] Facultad de Ciencias, Universidad Nacional Autónoma de México
[2] Department of Mathematics, Massachusetts Institute of Technology
[3] Department of Electrical Engineering and Computer Science, Massachusetts Institute of Technology
[4] Julia Lab, Computer Science and Artificial Intelligence Laboratory, Massachusetts Institute of Technology
dpsanders@gmail.com

**Keywords:** GPU, root finding, constraint propagation

`IntervalArithmetic.jl` [3] is a library for interval arithmetic written entirely in the Julia language [6]. As such it is both usable interactively, while still being competitive in performance with state-of-the-art C++ libraries. The `JuliaIntervals` suite of libraries [2] contains implementations of branch-and-bound-type algorithms built on `IntervalArithmetic.jl`, including `IntervalRoot-Finding.jl` and constraint propagation in `IntervalConstraintProgramming.jl` [4].

These libraries were originally designed to be run on a single core of the CPU. However, over the last few years the parallel programming capabilities available within Julia have significantly improved. They are designed so that the *same code* may be run on different platforms, including multi-threading with shared memory, on distributed processors, and on the GPU (Graphics Processing Unit).

In this talk we will focus on GPUs, which offer the tantalising possibility of large performance gains for those algorithms that can take advantage of their highly-parallel design. For example, the latest NVIDIA CPUs, such as the V100 that we use, have over 5,000 cores on a single code. The ideal paradigm is to run the same program on multiple data (SPMD) where each "thread" (piece of work) carries out the same instructions on their own piece of data. Branch-and-bound-type methods are clearly a fruitful playground for applying GPUs – although with some caveats.

The `CUDA.jl` library [5, 1] provides facilities for automatically compiling to the GPU functions that are written in Julia code. But it is also possible to run vectorised computations on vectors of objects living on the GPU without writing a single line of GPU-specific code, using broadcasting. We will demonstrate vectorised implementations of branch-and-bound-type methods, such as root finding and global optimisation, that are **generic**, i.e. they are able to run on either the CPU or the GPU, with no change in code.

As an application, we introduce a new benchmark of computing and verifying $10^6$ roots of the 2D Griewank function in one second, representing a 240x speed-up over a single CPU core, with room for further code optimisation by writing GPU kernels by hand. We also show significant speed-ups for the standard benchmark problem of computing the ground state of a 5-atom Lennard-Jones cluster.

Finally, we will demonstrate what we believe to be the first implementation of the forward-backward (HC4revise) interval constraint propagation algorithm running on the GPU, achieved using a code-generation approach.

## References

[1] CUDA.jl, https://github.com/JuliaGPU/CUDA.jl

[2] JuliaIntervals, https://github.com/JuliaIntervals/

[3] JuliaIntervals/IntervalArithmetic.jl, https://github.com/JuliaIntervals/IntervalArithetic.jl

[4] JuliaIntervals/IntervalConstraintProgramming.jl , https://github.com/JuliaIntervals/IntervalConstraintProgramming.jl

[5] Besard, T., Foket, C., De Sutter, B.: Effective extensible programming: Unleashing Julia on GPUs. IEEE Transactions on Parallel and Distributed Systems (2018).

[6] Bezanson, J., Edelman, A., Karpinski, S., Shah, V.B.: Julia: A Fresh Approach to Numerical Computing. SIAM Review **59**(1), 65–98 (Jan 2017).

# Variability measures for estimates in interval data fitting

Sergey P. Shary

Novosibirsk State University, Novosibirsk, Russia
shary@ict.nsc.ru

We consider the following data fitting problem: given results of measurements or observations, it is required to construct a functional dependence of a fixed type that "best fit" these data. Specifically, we need to determine the parameters $\beta_1, \beta_2, \ldots, \beta_n$ of a linear function

$$y = \beta_1 x_1 + \beta_2 x_2 + \ldots + \beta_n x_n \tag{1}$$

from a number of values of the independent variables $x_1, x_2, \ldots, x_n$ (also called *predictor* or *input* variables), and the corresponding values of the dependent variable $y$ (also called *criterion* or *output* variable). Both $x_1, x_2, \ldots, x_n$ and $y$ are not known precisely, and, in the $i$-th measurement, we only know intervals of their possible values, that is, $x_1 \in \boldsymbol{x}_{i1}, x_2 \in \boldsymbol{x}_{i2}, \ldots, x_n \in \boldsymbol{x}_{in}$ and $y \in \boldsymbol{y}_i$, $i = 1, 2, \ldots, m$. Overall, the data of our data fitting problem form an interval $m \times n$-matrix $\boldsymbol{X} = (\boldsymbol{x}_{ij})$ and an interval $m$-vector $\boldsymbol{y} = (\boldsymbol{y}_i)$.

To find estimates of the coefficients $\beta_1, \beta_2, \ldots, \beta_n$, several techniques have been developed. We focus on the so-called maximum compatibility method elaborated in [1] and other works. After the estimates for all $\beta_i$ are found, we need to somehow evaluate their possible variability and non-uniqueness. Our work presents the construction of two quantitative measures of variability for parameter estimates in the data fitting problem under interval uncertainty. They show the degree of variability and ambiguity of the estimate, and the need for their introduction is dictated by the fact that the results of processing interval data are typically non-unique. These measures can serve, in a certain sense, as an analog of the variance of the estimate in traditional probabilistic statistics.

In the maximum compatibility method, the estimate $\hat{\beta} = (\hat{\beta}_1, \hat{\beta}_2, \ldots, \hat{\beta}_n)$ of the parameters $\beta_1, \beta_2, \ldots, \beta_n$ is taken as the argument of the maximum of a special "recognizing functional", a function $\mathrm{Tol} : \mathbb{R}^n \to \mathbb{R}$, constructed from the interval data of the problem (see [1]. In other words,

$$\hat{\beta} = (\hat{\beta}_1, \hat{\beta}_2, \ldots, \hat{\beta}_n) \ = \ \arg\max_{x \in \mathbb{R}^n} \mathrm{Tol}\,(x).$$

To quantitatively characterize the variability of the estimate of the parameter vector $\beta = (\beta_1, \beta_2, \ldots, \beta_n)$ in the linear function (1), which is obtained by the maximum compatibility method from the data $\boldsymbol{X}, \boldsymbol{y}$, we propose the values

$$\mathrm{IVE}\,(\boldsymbol{X}, \boldsymbol{y}) \ = \ \sqrt{n} \ \max_{\mathbb{R}^n} \mathrm{Tol} \cdot \left( \min_{X \in \mathrm{vert}\,\boldsymbol{X}} \mathrm{cond}_2 X \right) \cdot \frac{\left\| \arg\max_{\mathbb{R}^n} \mathrm{Tol} \right\|_2}{\|\hat{\boldsymbol{y}}\|_2}$$

and

$$\mathrm{IDE}\,(\boldsymbol{X}, \boldsymbol{y}) \ = \ \sqrt{n} \ \max_{\mathbb{R}^n} \mathrm{Tol} \cdot \mathrm{cond}_2 \, \mathcal{X} \cdot \frac{\left\| \arg\max_{\mathbb{R}^n} \mathrm{Tol} \right\|_2}{\|\hat{\boldsymbol{y}}\|_2},$$

where

$n$ is the dimension of the parameter vector of function (1);

$\| \cdot \|_2$ is the Euclidean norm (2-norm) of vectors from $\mathbb{R}^n$;

$\mathrm{vert}\,\boldsymbol{X}$ is the set of corner matrices for the interval matrix $\boldsymbol{X}$, i.e., the set of such point matrices $X = (x_{ij})$ that $x_{ij} \in \{\underline{\boldsymbol{x}}_{ij}, \mathbb{1}^{\mathrm{v}}_{\boldsymbol{x}ij}\}$ for every $i$ and $j$;

$\mathcal{X}$ is a special matrix, called *endpoint combination matrix*, of the size $N \times n$ with $N \leq m \cdot 2^n$, made up of combinations of endpoints of the interval elements along each row of the data matrix $\boldsymbol{X}$;

$\mathrm{cond}_2$ means the spectral condition number of the matrix, defined as

$$\mathrm{cond}_2 A \;=\; \sigma_{\mathrm{max}}(A) \,/\, \sigma_{\mathrm{min}}(A),$$

the ratio of its maximal ($\sigma_{\mathrm{max}}$) and minimal ($\sigma_{\mathrm{min}}$) singular values;

$\hat{\boldsymbol{y}}$ is a certain "most representative" point from the interval vector $\boldsymbol{y}$, which is taken as

$$\hat{\boldsymbol{y}} \;=\; \tfrac{1}{2}\big(|\mathrm{mid}\,\boldsymbol{y} + \mathrm{rad}\boldsymbol{y}| + |\mathrm{mid}\,\boldsymbol{y} - \mathrm{rad}\boldsymbol{y}|\big)$$

and the operations "mid" and "rad" are applied componentwise.

The rationale for new variability measures is given, their motivation and applications are discussed.

## References

[1] S.P. SHARY: Weak and strong compatibility in data fitting under interval uncertainty, *Advances in Data Science and Adaptive Analysis* 12 (2020), 34 p.

# Numerical verification of existence for subharmonic solutions to delayed van der Pol-Duffing equation

Naoki Takamatsu[1], Yuuki Saito[1] Shin'ichi Oishi[2],
and Kouta Sekine[3]

[1] Graduate School of Fundamental Science and Engineering, Waseda University
[2] Department of Applied Mathematics, Faculty of Science and Engineering, Waseda University, Japan
[3] Faculty of Information Networking for Innovation and Design, Toyo University, Japan
naokiy@suou.waseda.jp

**Keywords:** Computer-assisted proof, Existence of subharmonic solutions

We are concerned with a problem of numerical verification of the existence for solutions to the delayed van der Pol-Duffing equation

$$\frac{d^2x(t)}{dt^2} + k(x^2(t) - 1)\frac{dx(t)}{dt} + \mu x(t) + \gamma x^3(t) - \alpha x(t - \tau) = \beta \cos t, \tag{1}$$

where $k > 0, \tau \geq 0, \mu \geq 0, \gamma > 0$, and $\beta > 0$. It is known that subharmonic solutions and chaotic aperiodic solutions have been observed numerically for the equation (1) (see, e.g., [1]). In this paper, we present an algorithm for proving the existence of $1/n$ subharmonic solutions with the periods $n$ times than that of the external forcing term. Via the time scaling, we reduce the problem to a problem of seeking $2\pi$ periodic solutions of the following equation:

$$\frac{d}{dt}\begin{pmatrix} x(t) \\ y(t) \end{pmatrix} + n\begin{pmatrix} -y(t) \\ k(x(t)^2 - 1)y(t) + \mu x(t) + \gamma x^3(t) - \alpha x\left(t - \frac{\tau}{n}\right) - \beta \cos nt \end{pmatrix} = 0. \tag{2}$$

We first compute a Galerkin approximate solution of (2) using the truncated Fourier series
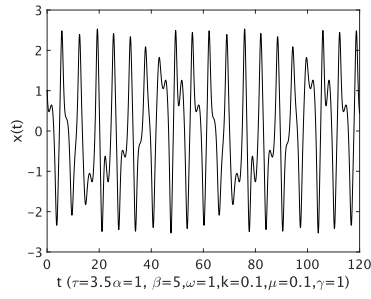
$$x_0(t) = \frac{a_0}{2} + \sum_{i=1}^{m} (a_i \cos it + b_i \sin it).$$

We then prove the existence of an exact subharmonic solution close to the approximate solution $x_0$ using the extended Newton-Kantorovich theorem. For the evaluation of the norm of the inverse operator of the Fréchet derivative, we have used the theory of asymptotic diagonally dominant matrix developed in Ref. [2].
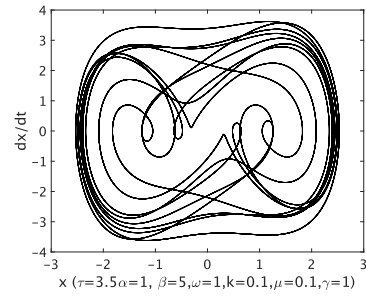
As an example, Figure 1 shows an approximate 1/9 subharmonic solution obtained for (2) via Galerkin's method. The numerical verification result of the existence of the exact solution around this approximate solution is shown in Table 1. The verification theory of Ref. [2] declares that $b(r_0) < 1$ implies the existence of a 1/9 subharmonic solution nearby the approximate solution $x_0$. Here, $M_n$, $M$, $\eta$, $r_0$ and $\|x_0\|_\infty$ are bounds for the norm of the inverse of Jacobian matrix of the Galerkin's equation at $x_0$, the norm of the inverse of the Frechét derivative at $x_0$, the residual, the radius of the ball in $H^1$-Sobolev space centered at $x_0$ including an exact solution, and the maximum norm of the approximate solution, respectively. For detail, we will explain in the presentation.

Table 1: Result for verification of the existence

| $m$ | $M_n$ | $M$ | $\eta$ | $r_0$ | $b(r_0)$ | $\|x_0\|_\infty$ |
|-----|-------|-----|--------|-------|----------|------------------|
| 500 | 372 | 9782 | $3.28e^{-7}$ | $6.57e^{-7}$ | $8.42e^{-2}$ | 2.57 |

a) $x_0(t)$          b) Phase diagram

Figure 1: Approximate solution (1/9 subharmonic)

**References**

[1] Y.TSUDA, H.TAMURA, A.SUEOKA, T.FUJII: Chaotic Behaviour of a Nonlinear Vibrating System with a Retarded Argument, *JSME International Journal, Series III*, **35** No.2 pp. 259-267 (1992).

[2] SHIN'ICHI OISHI, KOUTA SEKINE: Inclusion of Periodic Solutions for Forced Delay Differential Equation Modeling El Niño, *Nonlinear Theory and Its Applications, IEICE*, vol. 12, no. 3, (2021), pp.1-36.

# A rigorous forward integration method for time-dependent PDEs

Akitoshi Takayasu[1] and Jean-Philippe Lessard[2]

[1] University of Tsukuba, Ibaraki, Japan
[2] McGill University, Montreal, Canada
takitoshi@risk.tsukuba.ac.jp

**Keywords:** rigorous forward integration, time-dependent PDEs, evolution operator

This talk provides a numerical method of rigorous integration in forward time for a class of time-dependent partial differential equations. Our tool of rigorous integration is based on a fixed point form via semigroup theory (the evolution operator generated on a Banach space). A uniform bound of the evolution operator, which is a solution map of the linearized problem at an approximate solution, provides the existence of such evolution operator. Checking the hypothesis of local inclusion theorem, we obtain rigorous error bounds of the solution in a time interval. The proof of local inclusion theorem is based on Banach's fixed-point theorem. Furthermore, a time stepping scheme extends the local inclusion of the solution, which is designed to avoid the propagation of errors. As applications of this method, we introduce results of computer-assisted proofs for complex-valued nonlinear heat equation [1], nonlinear Schrödinger equation [2], and Swift-Hohenberg equation.

## References

[1] A. TAKAYASU, J.-P. LESSARD, J. JAQUETTE, H. OKAMOTO: *Rigorous numerics for nonlinear heat equations in the complex plane of time*, submitted 2019. (arXiv:1910.12472)

[2] J. JAQUETTE, J.-P. LESSARD, A. TAKAYASU: *Global dynamics in nonconservative nonlinear Schrödinger equations*, submitted 2020. (arXiv:2012.09734)

# Computer-assisted Existence Proofs for Navier-Stokes Equations on an Unbounded Strip with Obstacle

Jonathan Wunderlich

Karlsruhe Institute of Technology
Department of Mathematics
Englerstraße 2, 76131 Karlsruhe, Germany
Jonathan.Wunderlich@kit.edu

**Keywords:** Computer-assisted proof, Navier-Stokes, existence, enclosure

The incompressible stationary 2D Navier-Stokes equations

$$\left.\begin{aligned} -\Delta v + Re\left[(v \cdot \nabla)v + \nabla q\right] = f \\ \operatorname{div} v = 0 \end{aligned}\right\} \text{ in } \Omega$$

$$v = 0 \quad \text{on } \partial\Omega$$

are considered on an unbounded strip domain $\Omega \subseteq \mathbb{R}^2$ perturbed by a compact obstacle $D$, i.e., $\Omega = \mathbb{R} \times (0,1) \setminus D$. Here, $Re$ denotes the Reynolds number and $f$ models external forces acting on the fluid.

With $U$ denoting the Poiseuille flow and $P$ its associated pressure we are interested in solutions of the form $v = U + \bar{u}$ where $\bar{u}(x,y) \to 0$ as $|x| \to \infty$ and $q = P + p$.

Since such functions $\bar{u}$ do not satisfy the Dirichlet boundary conditions anymore we perform a second transformation using a solenoidal vector field $V$ with

$$V = 0 \text{ on } \partial\Omega \setminus \partial D, \quad V = U \text{ on } \partial D \quad \text{and} \quad V(x,y) \to 0 \text{ as } |x| \to \infty$$

which finally leads to the transformed Navier-Stokes equations

$$\left.\begin{aligned} -\Delta u + Re\left[(u \cdot \nabla)u + (u \cdot \nabla)(U - V) + ((U - V) \cdot \nabla)u + \nabla p\right] = g \\ \operatorname{div} u = 0 \end{aligned}\right\} \text{ in } \Omega$$

$$u = 0 \quad \text{on } \partial\Omega$$

with the right-hand side $g = f - \Delta V - Re\left[(V \cdot \nabla)V - (V \cdot \nabla)U - (U \cdot \nabla)V\right]$.

Modeling the divergence free condition in the solution space

$$H(\Omega) := \left\{u \in H_0^1(\Omega, \mathbb{R}^2) \colon \operatorname{div} u = 0\right\}$$

we can eliminate the pressure from the first equation which leads to the following weak formulation for the velocity:

Find $u \in H(\Omega)$ such that

$$\int_\Omega \left(\nabla u \cdot \nabla\varphi + Re\left[(u \cdot \nabla)u + (u \cdot \nabla)(U + V) + ((U + V) \cdot \nabla)u\right] \cdot \varphi\right) d(x,y)$$

$$= \int_\Omega g \cdot \varphi \, d(x,y) \quad (\varphi \in H(\Omega)).$$

Applying computer-assisted techniques to this problem, we are able to prove existence of a (non-trivial) solution $u^* \in H(\Omega)$ to the weak formulation (with $f \equiv 0$) for different Reynolds numbers and several domains $\Omega$. We point out that our methods do not use a stream function formulation which allows us to handle domains which are not simply connected as well.

Starting from an approximate solution (computed with divergence-free finite elements), we determine a bound for its defect, and a norm bound for the inverse of the linearization at the approximate solution. For the latter, bounds for the essential spectrum and for eigenvalues play a crucial role, especially for the eigenvalues "close to" zero. Therefor we use the Rayleigh-Ritz method, a corollary of

the Temple-Lehmann theorem and a homotopy method to get enclosures of the eigenvalues below the essential spectrum.

With these data at hand, we can use a fixed-point argument to obtain the existence of a solution "nearby" the approximate one as well as an error bound (in the Sobolev space $H(\Omega)$).

Finally, if our computer-assisted proof provides the existence of a solution $u^*$ to the weak formulation for the velocity we additionally prove existence of a corresponding pressure $p^*$ such that the pair $(u^*, p^*)$ is a weak solution to the transformed Navier-Stokes equations. The existence result is obtained by purely analytical techniques, nevertheless, for a given approximate solution to the pressure our methods provide an error bound (in a dual norm) as well.

# A numerical verification method on time-global solutions of autonomous systems of complex functions

Nobito Yamamoto[1] and Koki Nitta[1]

[1] The University of Electro-Communications, Tokyo, Japan
`yamamoto@im.uec.ac.jp`

We propose a numerical verification method to prove existence of time-global solutions. Our problems are autonomous ODEs on a complex plane with analytic functions in the right-hand sides:

$$\dot{z} = f(z),$$
$$z(0) = z_0 \in D \subset \mathbb{C}.$$

The method is based on the features of analytic functions and can be applied to the cases where Lyapunov functions are hardly constructed. Moreover we mention a numerical verification method to prove existence of closed orbits distributing continuously. An example problem derived from nonlinear Schrödinger equation is treated:

$$\begin{cases} \dot{a}_0 &=& i(a_0^2 + a_1^2), \\ \dot{a}_1 &=& i(-1 + 2a_0)a_1, \end{cases}$$
$$a_0(0), a_1(0) \in D \subset \mathbb{C}.$$

We apply our method to the above system and prove that it has a time-global solution within a bounded region $D$ and that there are periodic solutions which distribute continuously in a subset of $D$.

## References

[1] R. SVERDLOVE: *Vector fields defined by complex functions*, JDE, 34 (1979), 427–439

[2] M. HUKUHARA, T. KIMULWA, & K. MATUDA: *Equations Differentielles Ordinaires du Premier Ordre dans le Champ Complexe*, Mathematical Society of Japan, Tokyo, (1961)

# Verified algorithm for high-order partial derivatives using nilpotent matrix

Naoya Yamanaka[1], Takeo Uramoto[2]

[1] Meisei University, Tokyo, Japan
[2] Kyushu University, Tokyo, Japan
naoya.yamanaka@meisei-u.ac.jp

A verified algorithm for higher order derivatives using several features of nilpotent matrix is proposed. Since the calculations of derivatives are used in many applications, there are many previous studies. For example, the automatic differentiation (AD) algorithm is well known as an accurate and efficient algorithm for lower order derivatives, however, for arbitrarily higher order derivatives the AD algorithm quickly grow very complicated.

Recently, hyper-dual numbers are proposed by Fike and Alonso. The numbers are an extension of dual numbers, which are based on the non-real term

$$\epsilon^2 = 0 \quad \text{where} \quad \epsilon \neq 0.$$

Similar to this features, hyper-dual numbers are based on the non-real term

$$\epsilon_i^2 = 0 \ (1 \leqq i \leqq n) \quad \text{where} \quad \epsilon_i \neq 0, \epsilon_i \neq \epsilon_j, \epsilon_i \epsilon_j = \epsilon_j \epsilon_i \neq 0 \ (i \neq j).$$

These numbers enable us to get accurate results of second or higher derivatives. In addition to that, the matrix representation of hyper-dual numbers, which consists of the nilpotent matrix, is very useful for calculation of high order derivatives. In the AD algorithm, every operation needs to be made for calculation, however, the matrix representation of hyper-dual numbers enable us to avoid it.

In this talk, we propose verification algorithm using features of the nilpotent matrix. Numerical examples are presented for illustrating effectiveness of the proposed algorithm.

## References

[1] J. A. FIKE, J. J. ALONSO: The Development of Hyper-Dual Numbers for Exact Second Derivative Calculations, *AIAA paper*, 2011-886, 49th AIAA Aerospace Sciences Meeting, January 4-7, 2011.

[2] J. A. FIKE: Multi-objective optimization using hyper-dual numbers: *Ph.D. thesis, Stanford University*, 2013.

[3] Y. IMOTO, N. YAMANAKA, T. URAMOTO, M. TANAKA, M. FUJIKAWA, N. MITSUME : Fundamental theorem of matrix representations of hyper-dual numbers for computing higher-order derivatives, *JSIAM Letters*, 12 (2020), 29–32.

# Verification of artificial neural networks via Taylor models of INTLAB

Dániel Zombori[1], Tamás Szabó[1], János Horváth[1], Attila Szász[1], Tibor Csendes[1], and Balázs Bánhelyi[1]

[1] University of Szeged, Szeged, Hungary
`banhelyi@inf.u-szeged.hu`

Currently, neural networks are used in many fields. Whether we just think of the face recognition in mobile phones or self-driving cars, they have become more and more a part of our everyday life. A lot has happened in the last decades, newer and newer methods have been invented. Their robustness is usually checked by random testing, i.e. by looking for examples that are close to known examples and for which the neural network gives no longer the desired result. These cases are called adversarial examples. A more sophisticated method than random testing is currently not widely used according to the literature. Most of the verification techniques are based either on optimization problems [3], on simple linear [1] or affine [2] approximations.

Our goal is to develop a method to improve the procedures based on this linear approach and control the robustness of different artificial neural networks with greater efficiency. This system is based on Taylor polynomials. Another drawback of the current systems is that while they attempt to provide a guarantee of neural network robustness, the systems themselves are not reliable against adversarial examples. Such examples can now be found in the literature. There are many cases in the literature where a neural network that is considered to be good can easily be "cheated". For example, for number recognition systems (MNIST), the neural network can be completely diverted from the correct answer [4]. This behavior is highly undesirable in such systems.

Currently, neural network robustness control systems are basically optimized for fast execution, which leads to rounding errors that can mask potentially hostile examples. The approach we have taken has not really appeared in this field before, so it may be an interesting direction and it can provide novel results if implemented successfully. Moreover, a proper implementation can address the reliability gaps of previous systems. It is expected that runtime will increase, but considering that this algorithm only needs to be run once on the final neural network and since it can be well parallelized, to some degree this is tolerable, for a possibly more reliable response. Especially for systems where this is desirable at a higher level.

A MATLAB/INTLAB implementation of the Taylor model was used [3]. The variables of the Taylor models are the inputs of the neural network. The initial Taylor models were defined in the first layer of the neural network. The Taylor model is used in this paper for the Tanh, Sigmoid, and ReLu activation functions. In the last case, the adoption is not trivial because the derivative of the ReLu function is not continuous. In this case, we applied some tricks, for example, we used the soft-max function for inclusion.

In the presentation, our MATLAB implementation is shown, and we measure the speed of the MATLAB version of the Taylor model. We also present the advantages of our method with some examples.

## References

[1] DAVID L. APPLEGATE, WILLIAM COOK, SANJEEB DASH AMD DANIEL G. ESPINOZA: Exact solutions to linear programming problems. *Operations Research Letters*, 35:6, 693–699, 2007.

[2] WANG LIN, ZHENGFENG YANG, XIN CHEN, QINGYE ZHAO, XIANGKUN LI, ZHIMING LIU AND JIFENG HE: Robustness Verification of Classification Deep Neural Networks via Linear Programming. *Conference on Computer Vision and Pattern Recognition*, 2019.

[3] SIEGFRIED M. RUMP: INTLAB – INTerval LABoratory. In Tibor Csendes, editor, *Developments in Reliable Computing*, pages 77-104. Kluwer Academic Publishers, Dordrecht, 1999.

[4] VINCENT TJENG, KAI XIAO AND RUSS TEDRAKE: Evaluating Robustness of Neural Networks with Mixed Integer Programming. *The 7th International Conference on Learning Representations (ICLR)*, 2019, https://openreview.net/pdf?id=HyGIdiRqtm.

[5] DÁNIEL ZOMBORI, BALÁZS BÁNHELYI, TIBOR CSENDES, ISTVÁN MEGYERI AND MÁRK JELASITY: Fooling a Complete Neural Network Verifier. *The 9th International Conference on Learning Representations (ICLR)*, 2021, https://openreview.net/pdf?id=4IwieFS44l.

# Verification of artificial neural networks via MIPVerify and SCIP

Dániel Zombori[1], Tamás Szabó[1], János Horváth[1], Attila Szász[1], Tibor Csendes[1], and
Balázs Bánhelyi[1]

[1] University of Szeged, Szeged, Hungary
`banhelyi@inf.u-szeged.hu`

We have introduced the problem of adversarial examples in our other two talks: "Adversarial Example Free Zones for Specific Inputs and Neural Networks" and "Verification of artificial neural networks via Taylor models of INTLAB". Obviously, the verification of artificial neural networks is a challenging new field for the application of reliable computation techniques. The present talk summarizes the results of our computational experiments with the neural network verification method MIPVerify [3] when used together with the rational arithmetic based mixed integer optimizer SCIP [2]. The use of the latter one was implied by the fact that MIPVerify proved to be unreliable with the suggested solver Gurobi [1], see our recent paper [4].

In the presentation, our new implementation of MIPVerify and SCIP is shown, and we present the measured speed of the procedure. Ten handwritten numbers were applied from the MNIST database, and some artificial neural networks from the ERAN (ETH Robustness Analyzer for Neural Networks) set were applied. We also present the advantages of the new method with some examples.

## References

[1] Gurobi guidelines for numerical issues, 2017. URL http://files.gurobi.com/Numerics.pdf.

[2] SCIP, Solving Constraint Integer Programs. https://www.scipopt.org/

[3] VINCENT TJENG, KAI XIAO AND RUSS TEDRAKE: Evaluating Robustness of Neural Networks with Mixed Integer Programming. *The 7th International Conference on Learning Representations (ICLR)*, 2019, https://openreview.net/pdf?id=HyGIdiRqtm

[4] DÁNIEL ZOMBORI, BALÁZS BÁNHELYI, TIBOR CSENDES, ISTVÁN MEGYERI AND MÁRK JELASITY: Fooling a Complete Neural Network Verifier. *The 9th International Conference on Learning Representations (ICLR)*, 2021, https://openreview.net/pdf?id=4IwieFS44l.

# LIST OF AUTHORS

**Asai, Taisei**:  Waseda University, Japan

**Auer, Ekaterina**:  University of Applied Sciences Wismar, Germany

**Balogh, Nándor**:  Redink Ltd., Hungary

**Bánhelyi, Balázs**:  University of Szeged, Hungary

**Bartha, Ferenc Agoston**:  Bolyai Institute, University of Szeged, Hungary

**Bourgois, Auguste**:  FORSSEA, France

**Ceberio, Martine**:  University of Texas at El Paso, United States

**Chaabouni, Amine**:  Ecole polytechnique, France

**Chausova, Elena**:  Tomsk State University, Russia

**Churavy, Valentin**:  Department of Electrical Engineering and Computer Science, Massachusetts Institute of Technology, United States

**Contreras, Jonatan**:  The University of Texas at El Paso, United States

**Csendes, Tibor**:  University of Szeged, Hungary

**De Angelis, Marco**:  University of Liverpool, United Kingdom

**Dhouib, Mounir**:  ENAU, Tunisia

**Dózsa, Tamás**:  Eötvös Loránd University, Hungary

**Dudás, János**:  University of Szeged, Szeged, Hungary, Hungary

**Fasi, Massimiliano**:  Örebro University, Sweden

**Fedorchenko, Lyudmila**:  Saint Petersburg Research Center of the Russian Academy of Sciences, Russia

**Fekete, Imre**:  Eötvös Loránd University & MTA-ELTE NUMNET Research Group, Hungary, Hungary

**G. Casado, Leocadio**:  University of Almeria, Spain

**G.-Tóth, Boglárka**:  University of Szeged, Hungary

**Geyda, Alexander**:  St. Petersburg Federal Research Center of the Russian Academy of Sciences, Russia

**Gierzkiewicz, Anna**:  University of Agriculture in Krakow, Poland

**Gillner, Lorenz**:  University of Applied Sciences Wismar, Germany

**Hendrix, E.M.T.**:  Universidad de Málaga, Spain

**Horváth, János**:  University of Szeged, Hungary

**Jaulin, Luc**:  Lab-STICC, ENSTA-Bretagne, France

**Jézéquel, Fabienne**:  LIP6, Sorbonne Université, France

**Kearfott, Ralph Baker**:  University of Louisiana at Lafayette, United States

**Kinoshita, Takehiko**:  Department of Information Science, Saga University, Japan

**Kosheleva, Olga**:  University of Texas at El Paso, United States

**Kreinovich, Vladik**:  University of Texas at El Paso, United States

**Krisztin, Tibor**:  Bolyai Institute, University of Szeged, Hungary

**Krisztin, Tibor**:  University of Szeged, Szeged, Hungary, Hungary

**Kubica, Bartlomiej**:  Department of Applied Informatics, Warsaw University of Life Sciences, Poland

**Kumkov, Sergey**:  Institue of Mathematics and Mechanics Ural Branch Russian Academy of Sciences, Russia

**Lange, Marko**:  Hamburg University of Technology, Germany

**Lessard, Jean-Philippe**:  McGill University, Canada

**Liu, Dun**:  University of Louisiana at Lafayette, United States

**Liu, Xuefeng**:  Niigata University, Japan

**Lorenc, Matyáš**:  Charles University, Faculty of Mathematics and Physics, Czechia

**Luther, Wolfram**:  University of Duisburg-Essen, Germany

**Malti, Rachid**:  IMS Laboratory, University of Bordeaux, France

**Markót, Mihály Csaba**:  University of Vienna, Austria

**Matsue, Kaname**:  Institute of Mathematics for Industry / International Institute for Carbon-Neutral Energy Research, Kyushu University, Japan

**Megyeri, István**:  University of Szeged, Hungary

**Messine, Frédéric**:  ENSEEIHT-Toulouse INP, France

**Mikaitis, Mantas**:  University of Manchester, United Kingdom

**Mireles James, Jason**:  Florida Atlantic University, United States

**Miyajima, Shinya**:  Iwate University, Japan

**Morozov, Alexander**:  FRC CSC RAS and MAI, Russia

**Nakao, Mitsuhiro T.**:  Faculty of Science and Engineering, Waseda University, Japan

**Nitta, Koki**:  The University of Electro-Communications, Japan

**Ogawa, Shota**:  Systemers Inc., Japan

**Oishi, Shin'ichi**:  Waseda University, Japan

**Okayama, Tomoaki**:  Hiroshima City University, Japan

**Ozaki, Katsuhisa**:  Department of Mathematical Sciences, Shibaura Institute of Technology, Japan

**Rauh, Andreas**:  ENSTA-Bretagne, Lab-STICC, France

**Reviznikov, Dmitry**:  MAI and FRC CSC RAS, Russia

**Revol, Nathalie**:  Inria, France

**Rump, Siegfried M.**:  Hamburg University of Technology, Germany

**Saito, Yuuki**:  Waseda University, Japan

**Sanders, David**:  Facultad de Ciencias, Universidad Nacional Autónoma de México & Department of Mathematics, MIT, United States

**Sekine, Kouta**:  Toyo University, Japan

**Sekine, Kouta**:  Toyo University, Japan

**Selivanov, Victor**:  A. P. Ershov Institute of Informatics Systems, Novosibirsk, Russia

**Shary, Sergey**:  Institute of computational technologies SD RAS, Russia

**Soussi, Imene**:  ENAU, United Arab Emirates

**Szabó, Tamás**:  University of Szeged, Hungary

**Szász, Attila**:  University of Szeged, Hungary

**Takamatsu, Naoki**:  Waseda University, Japan

**Takayasu, Akitoshi**:  University of Tsukuba, Japan

**Tanaka, Kazuaki**:  Waseda University, Japan

**Tóth, Richárd**:  University of Szeged, Hungary

**Uramoto, Takeo**:  Kyushu University, Japan

**Vígh, Alexandra**:  University of Szeged, Hungary

**Watanabe, Yoshitaka**:  Research Institute for Information Technology, Kyushu University, Japan

**Wunderlich, Jonathan**:  Karlsruhe Institute of Technology, Department of Mathematics, Germany

**Yamamoto, Nobito**:  The University of Electro-Communications, Japan

**Yamanaka, Naoya**:  Meisei University, Japan

**Zgliczynski, Piotr**:  Jagiellonian University, Poland

**Zombori, Dániel**:  University of Szeged, Hungary

NOTES

**SCAN2020**