

Article

Not peer-reviewed version

Interpolating Strange Attractors via Fractional Brownian Bridges

Sebastian Raubitzek^{*}, Thomas Neubauer, [Jan Friedrich](#), [Andreas Rauber](#)

Posted Date: 18 April 2022

doi: 10.20944/preprints202204.0167.v1

Keywords: time series interpolation; phase space reconstruction; takens' theorem; interpolation; stochastic interpolation; genetic algorithm; time series data; preprocessing; strange attractor; attractor; attractor reconstruction



Preprints.org is a free multidiscipline platform providing preprint service that is dedicated to making early versions of research outputs permanently available and citable. Preprints posted at Preprints.org appear in Web of Science, Crossref, Google Scholar, Scilit, Europe PMC.

Copyright: This is an open access article distributed under the Creative Commons Attribution License which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

Disclaimer/Publisher's Note: The statements, opinions, and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions, or products referred to in the content.

Article

Interpolating Strange Attractors via Fractional Brownian Bridges

Sebastian Raubitzek^{*1}, Thomas Neubauer², Jan Friedrich³ and Andreas Rauber⁴

¹ TU Wien, Information and Software Engineering Group, Favoritenstrasse 9-11/194, 1040 Vienna, Austria; sebastian.raubitzek@gmail.com

² TU Wien, Information and Software Engineering Group, Favoritenstrasse 9-11/194, 1040 Vienna, Austria; thomas.neubauer@tuwien.ac.at

³ ForWind, Institute of Physics, University of Oldenburg, K pkersweg 70, 26129 Oldenburg, Germany; jan.friedrich@uni-oldenburg.de

⁴ TU Wien, Information and Software Engineering Group, Favoritenstrasse 9-11/194, 1040 Vienna, Austria; rauber@ifs.tuwien.ac.at

* Correspondence: sebastian.raubitzek@gmail.com

Abstract: We present a novel method for interpolating univariate time series data. The proposed method combines multi-point fractional Brownian bridges, a genetic algorithm, and Takens' theorem for reconstructing a phase space from univariate time series data. The basic idea is to first generate a population of different stochastically interpolated time series data, and secondly, to use a genetic algorithm to find the pieces in the population which generate the smoothest reconstructed phase space trajectory. A smooth trajectory curve is hereby found to have a low variance of second derivatives along the curve. For simplicity, we refer to the developed method as *PhaSpaSto*-interpolation, which is an abbreviation for **phase-space-trajectory-smoothing stochastic interpolation**. The proposed approach is tested and validated with a univariate time series of the Lorenz system, five non-model data sets and tested against a cubic spline interpolation and a linear linear interpolation. We find that the criterion for smoothness guarantees low errors on known model and non-model data. Finally, we interpolate the discussed non-model data sets, and show the corresponding improved phase space portraits. The proposed method is useful for interpolating low-sampled time series data sets for, e.g., machine learning, regression analysis, or time series prediction approaches.

Keywords: time series interpolation; phase space reconstruction; takens' theorem; interpolation; stochastic interpolation; genetic algorithm; time series data; preprocessing; strange attractor; attractor; attractor reconstruction;

1. Introduction

Many real-life time series data sets originate from complex real-life systems and/or non-linear phenomena. Often these data sets are sparsely sampled as, e.g., long-term temperature, yield, or environmental data sets. The non-linear and stochastic nature of these data sets, in addition to being sparsely sampled, makes predictions and analysis rather difficult. Thus, one tends to employ data augmentation, and, in the case of univariate time-series data, interpolation techniques. There are various interpolation techniques, such as linear, polynomial, fractal, and stochastic interpolation methods. Choosing a suitable interpolation method can be difficult, and one should take into account the characteristics of the data at hand. For fluctuating data with inherent randomness, for example one would choose a stochastic interpolation as discussed in [1]. But when it comes to deterministically chaotic systems where one can reconstruct a phase space based on Taken's theorem [2], the choice is not so obvious. Here, we present a method taking into account the reconstructed phase space properties of time series data. Thereby, we want our reconstructed phase space trajectories to be as smooth as possible. But how can one achieve a smooth phase space trajectory? To solve this

35 issue, we developed a method combining multi-point Brownian bridges [1] and a genetic
36 algorithm. For simplicity, we refer to the developed method as *PhaSpaSto*-interpolation,
37 which is an abbreviation for **phase-space-trajectory-smoothing stochastic** interpolation.
38 We show, that the developed method performs well for the reconstructed phase space of
39 the Lorenz system and several univariate, sparsely sampled time series data. The results
40 show that the presented method effectively can interpolate the Lorenz system and the
41 discussed non-model data sets with comparatively low errors on known data points and
42 convincing phase space portraits.

43 Furthermore, many of today's most employed time series analysis and prediction
44 techniques stem from the domain of machine and/or deep learning. These methods
45 are data-based, i.e., they learn from data; thus, a sufficient amount of data and data of
46 good quality is necessary to, e.g., train a neural network. It is shown that interpolating
47 time series data using a fractal or linear interpolation can improve the accuracy of the
48 algorithm drastically, [3]. We thus suggest the presented method be tested and used for
49 data-based learning algorithms.

50 This article is structured as follows: Section 2 collects publications related to the
51 developed method and discusses them briefly. Section 3 describes the multi-point
52 Brownian bridges [1], the Lorenz system, and the employed genetic algorithm, and
53 further sums up the developed scheme. All results with the corresponding error tables
54 and figures are discussed in Section 4. Section 5 concludes the findings of this article.

55 2. Related Work

56 The presented research is motivated by findings of [3] and [4]. It is further based on
57 the stochastic interpolation method presented in [1]. Thus we will briefly describe the
58 mentioned publications and chronologically, i.e. by their date of publication, list related
59 ideas below.

- 60 • Ref [5]: This publication presents a method to determine if images are blurry.
61 For this purpose, the second derivatives of grey-scale images are taken, and the
62 corresponding variance over all pixels is analyzed. If the variance is below a certain
63 threshold, the image is a blurry image. This concept is used in the presented article.
64 We adapted the idea of variances of second derivatives, which is discussed in
65 Section 3.3.1.
- 66 • Ref [6]: In this research, a combination of inverse distance methods, fuzzy set
67 theory, and a genetic algorithm are applied to interpolate rainfall data. The Genetic
68 algorithm was used to determine the parameters of the corresponding fuzzy mem-
69 bership functions. Thus the idea of improving interpolation techniques is adapted
70 from this publication.
- 71 • Ref [1]: This publication presents a novel interpolation technique where the idea of
72 a Brownian bridge, i.e., a constrained fractional Brownian motion (fBm), is extended
73 to more than two points, i.e., to *multi-point fractional Brownian bridges*. The authors
74 present an explicit construction that operates linearly on the fBm and can thus
75 be interpreted as a Gaussian random process constrained on multiple, prescribed
76 points. Further applications of this method are presented, such as determining
77 optimal Hurst exponents for sparsely sampled time series filled up by multi-point
78 fractional Brownian bridges with varying Hurst exponents. This method is used in
79 the presented research to fuel the genetic algorithm.
- 80 • Ref [3]: In this publication, a fractal interpolation to interpolate univariate time
81 series data is presented. The presented method considers the Hurst exponent of
82 the data under study. The authors show that fractal interpolation can increase the
83 predictability of a given univariate time series data. This research suggests that
84 different interpolation methods for univariate time series data may yield predictions
85 of different qualities. Thus, as presented here, an attractor-based interpolation is an
86 obvious next step in contrast to a fluctuation-based interpolation.

- 87 • Ref [4]: This publication is a continuation of [3]. The fractal interpolation and LSTM
 88 neural network approach is continued as ensembles of predictions. Randomly
 89 parameterized LSTM neural networks are generated from non-, linear-, and fractal-
 90 interpolated data. Afterward, these predictions are filtered based on their signal
 91 complexities. Some of the mentioned complexity measures require a suitable phase
 92 space embedding of the data under study and are related to the presented research
 93 in this article. Further, some of the data sets used here are discussed and predicted.
 94 We expect LSTM neural network predictions of stochastically interpolated data
 95 to outperform other interpolated approaches when considering the reconstructed
 96 phase space.
- 97 • Ref [7] describes a multi-point reconstruction of a given time series. The method
 98 is based on the assumption of Markovianity of the time series and a refinement
 99 algorithm is presented which allows to systematically fill up data points based on
 100 the empirically determined transition probability from one level to the next.

101 3. Methodology

102 The developed method consists of two steps, first generating a population of stochas-
 103 tically interpolated time series data, each with a different Hurst exponent. Firstly, these
 104 stochastically interpolated time series data are generated using multi-point fractional
 105 Brownian bridges, see Section 3.1. Secondly, these multi-point fractional Brownian
 106 bridges are improved using a genetic algorithm, to minimize the variance of second-
 107 order derivatives along the reconstructed phase space trajectory (see Section 3.3). The
 108 whole scheme is depicted in Figure 1. Finally, we briefly discuss the Lorenz system and
 109 its implementation in Section 3.4.

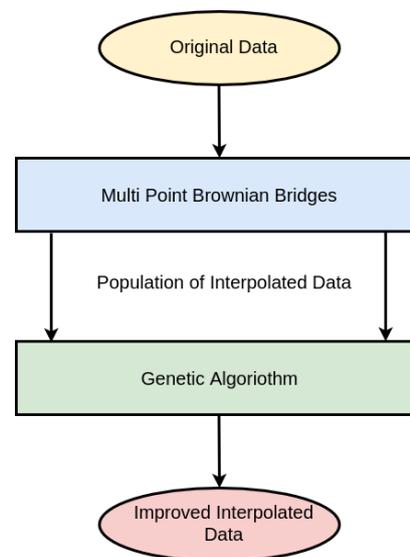


Figure 1. Depiction of the employed scheme.

110 3.1. Multi Point Fractional Brownian Bridges

111 As depicted in Figure 1, the employed genetic algorithm is fueled by a population
 112 of stochastically-interpolated time series data, in our case *multi-point fractional Brown-*
 113 *ian bridges*. To generate these stochastically interpolated time series data, multi-point
 114 fractional Brownian bridges [1] were used. Thus we briefly summarize this approach.

We consider a Gaussian random process $X(t)$ whose covariance is defined as $C(t, t') = \langle X(t)X(t') \rangle$. In the following, we focus on fractional Brownian motion where the covariance is given according to $\langle X(t)X(t') \rangle = \frac{1}{2}(t^{2H} + t'^{2H} - |t - t'|^{2H})$, where H is the Hurst exponent. To elucidate our interpolation scheme, we first define a so-called

fractional Brownian bridge [8],[9], which is a construction of fBm starting from 0 at $t = 0$ and ending at X_1 at $t = t_1$, i.e.,

$$X^B(t) = X(t) - (X(t_1) - X_1) \frac{\langle X(t)X(t_1) \rangle}{\langle X(t_1)^2 \rangle}. \quad (1)$$

This construction ensures that $X^B(t_1) = X_1$, which is also depicted in Figure 10. This single bridge can now be generalized to an arbitrary number of (non-equidistant) prescribed points X_i at t_i by virtue of a multi-point fractional Brownian bridge [1]

$$X^B(t) = X(t) - (X(t_i) - X_i) \sigma_{ij}^{-1} \langle X(t)X(t_j) \rangle, \quad (2)$$

115 where $\sigma_{ij} = \langle X(t_i)X(t_j) \rangle$ denotes the covariance matrix. Furthermore, we imply sum-
 116 mation over identical indices. Latter linear operation on the Gaussian random process
 117 $X(t)$ ensures that the bridge takes on exactly the values X_k at t_k , which can be seen from
 118 $X^B(t_k) = X(t_k) - (X(t_i) - X_i) \sigma_{ij}^{-1} \sigma_{kj} = X(t_k) - (X(t_i) - X_i) \delta_{ik} = X_k$, where δ_{ik} denotes
 119 the Kronecker-delta. Hence, this method allows for the reconstruction of a sparse signal
 120 where small-scale correlations are determined by the choice of the Hurst exponent H .

121 3.2. Phase Space Reconstruction

122 Before describing the fitness of our phase space trajectories, we need to introduce
 123 the concept of reconstructed phase spaces, [10,11].

124 We estimate a phase space embedding for all data under study. To find a suitable
 125 phase space embedding, one has to determine two parameters, the embedding dimension
 126 and the time delay.

127 To estimate the time delay τ , i.e., the delay between two consecutive time steps, we
 128 use the method based on the average information between two signals [12].

129 To estimate the embedding dimension d_E , we use the algorithm of false nearest
 130 neighbors [13]. Also, because the evaluations performed in this paper aim to graphically
 131 depict the embedding space, but with no limitations to the general applicability of our
 132 approach, we only chose data sets with an embedding dimension of three, i.e. $d_E = 3$.

The phase space embedding for a given signal $[x_1, x_2, \dots, x_n]$, thus is:

$$\vec{y}(i) = [x_i, x_{i+\tau}, \dots, x_{i+(d_E-1)\tau}] \quad , \quad (3)$$

and the corresponding three dimensional phase space embedding, thus is

$$\vec{y}(i) = [x_i, x_{i+\tau}, x_{i+2\tau}] \quad . \quad (4)$$

133 3.3. Genetic algorithm

134 We build a simple genetic algorithm to find the best possible interpolation given
 135 the data's phase space reconstruction using Taken's theorem. We want our reconstructed
 136 phase space curve to be as smooth as possible and thus define the trajectory's fitness as
 137 follows.

138 3.3.1. The Fitness of a Trajectory

139 The basic idea is to use a concept from image-processing, i.e., the blurriness of a
 140 picture, and apply it to phase space trajectories. This is because we want our trajectory as
 141 blurry, i.e., as smooth as possible. In image processing the blurriness is determined via
 142 second-order derivatives of grey-scale images at each pixel, [5]. We employ this concept,
 143 but instead of using it at each pixel, we calculate the variance of second-order derivatives
 144 along our phase space trajectories. Similar to the concept from image-processing, where
 145 the low variance of second-order derivatives implies more blurriness, curves with a
 146 low variance of second-order derivatives exhibit comparatively smooth trajectories. The
 147 reason here is intuitively clear, whereas curves with a high variance of second-order

148 derivatives have a range of straight and pointy sections, curves with a low variance
 149 of second-order derivatives have a similar curvature along the trajectory and thus are
 150 smoother. Hence, in order to guarantee smoothness along the trajectory, we want this
 151 variance to be as low as possible, which thus is our loss L . Concluding, our fitness is
 152 maximal when our loss L is minimal.

Again we start with the phase space vector and the corresponding embedding dimension d_E and time delay τ (See Section 3.2) of each signal as

$$\vec{y}(i) = [x_i, x_{i+\tau}, \dots, x_{i+(d_E-1)\cdot\tau}] \quad . \quad (5)$$

Thus we have one component for each dimension of the phase space. Consequently we can write the individual components as:

$$y_j(i) = x_{i+(j-1)\cdot\tau} \quad , \quad (6)$$

where $j = 1, 2, \dots, d_E$. We then take the second-order finite difference central derivative of a discrete function [14]:

$$u_j''(i) = x_{i+(j-1)\cdot\tau+1} - 2x_{i+(j-1)\cdot\tau} + x_{i+(j-1)\cdot\tau-1} \quad , \quad (7)$$

at each point, and for each component. Next we add up all the components as:

$$u''(i) = \sqrt{\sum_{j=1}^{d_E} u_j''(i)^2} \quad . \quad (8)$$

And finally, we use the variance of the absolute values of second derivatives along the phase space curve as our loss L of a phase space trajectory:

$$L = \text{Var}_i[u''(i)] \quad . \quad (9)$$

153 3.3.2. Genetic Algorithm Architecture

154 The employed genetic algorithm consists of the following building blocks:

155 A candidate solution is an interpolated time series using a random Hurst exponent
 156 $H \in]0; 1[$. The corresponding population of candidates is, e.g.m 1000 of these stochas-
 157 tically interpolated time series data with random Hurst exponents. A population of
 158 interpolated time series data is generated using the multi-point Brownian bridges such
 159 that, for each member of the population, a random Hurst exponent with $H \in]0; 1[$ is
 160 chosen, which then defines the interpolation of the member of the population. After
 161 generating the population, all members are sorted with respect to their fitness, i.e., the
 162 lower the loss L , the better an interpolation is. The mating is implemented such that only
 163 the best 50%, with respect to the fitness, can mate to produce new offsprings. The mating
 164 is done such that, for every gene, i.e., each interpolation between two data points, there is
 165 a 50:50 chance to inherit it from either one of the parents. The mutation was implemented
 166 that, in each generation, there is a 20% chance that a randomly chosen interpolated time
 167 series is replaced with a new interpolated time series within a corresponding randomly
 168 chosen new Hurst exponent. Also, we implemented a criterion for aborting the program,
 169 which was fulfilled if the population fitness mean did not change for ten generations.
 170 This described procedure is then performed for 1000 generations. But, actually, the 1000
 171 generations were never reached, as the criterion for abortion always triggered, and the
 172 program was ended, thus yielding the best interpolation with respect to the fitness of
 173 the phase space trajectories before reaching the 1000th generation.

174 3.4. The Lorenz System

175 For this research and to show the applicability of the proposed interpolation method,
 176 we chose the Lorenz system, [15] as a model to illustrate our ideas.

The Lorenz system is a set of three nonlinear equations:

$$\begin{aligned}\frac{dx}{dt} &= 10(-x + y) , \\ \frac{dy}{dt} &= 28x - y - xz , \\ \frac{dz}{dt} &= xy - \frac{8}{3}z .\end{aligned}\tag{10}$$

177 We solved this system using a basic Runge-Kutta 4 approach, [16]. We chose the step
178 size and length of the simulation with respect to the number of interpolation points to
179 test the quality of our interpolation scheme,

$$dt = \frac{0.1}{n_I + 1}, \quad L = 200 \cdot (n_I + 1) ,\tag{11}$$

where dt is the step size and L is the length of the simulation. The initial conditions were chosen to be:

$$x = -8, \quad y = 8, \quad z = 27 .\tag{12}$$

180 Finally, we need a univariate signal for the phase space reconstruction, and to test our
181 method, thus we choose one of the three variables. Accordingly, here we chose $x(t)$.

182 4. Results

183 Here we present the results of the Genetic algorithm for all data sets; first for
184 the Lorenz system, then for five non-model data sets. For both cases we validate the
185 developed method such that we delete data points from the original time series and
186 reconstruct the missing data points using the presented interpolation technique. Further,
187 we tested the presented interpolation technique against the best random interpolation of
188 the population, against a linear interpolation and against a cubic spline interpolation
189 [17]. Both the linear and spline interpolation were performed using the python package
190 `scipy`, `??`. For the validation we put emphasis on the Lorenz system, as the generated
191 model data allows us to test arbitrary settings, i.e. using different numbers of missing
192 data points. Contrary to that, for the non-model data sets we delete every second data
193 point and reconstructed the missing data points using the presented method. For the
194 non-model data sets we also present actual interpolation results, i.e. data sets with
195 smoothed-out phase space trajectories.

196 4.1. Results for the Lorenz System

197 For the actual equations and further details on the Lorenz system, see Section
198 3.4. We perform our interpolation for a number of different interpolation points $N_I =$
199 $\{1, 2, 3, 4, \dots, 20\}$.

200 We develop the following experimental steps to assess the performance of our
201 interpolation scheme:

- 202 1. Obtain a univariate time series from the Lorenz system.
- 203 2. Delete points from the data which will be reconstructed later on.

Given some univariate time series data of the Lorenz system, $[x_1, x_2, \dots, x_n]$, we extract certain data points with respect to the number of interpolation points $n_I \in N_I$ and the interpolated data set such that:

$$\begin{aligned}[x_1, x_{n_I+2}, x_{2n_I+2}, \dots, x_n] & \text{ Original data points to be kept for interpolation,} \\ [x_1, \hat{x}_1, \dots, \hat{x}_{n_I}, x_{n_I+2}, \hat{x}_{n_I+1}, \dots, \hat{x}_{2n_I+1}, x_{2n_I+2}, \hat{x}_{2n_I+2}, \dots, x_n] & \text{ Interpolated data,}\end{aligned}\tag{13}$$

204 where \hat{x}_i are the new found interpolated data points.

- 205 3. Perform the interpolation according to Section 3.3.
- 206 4. Calculate RMSE for the interpolated data points with respect to the previously
207 extracted original data points $[x_2, \dots, x_{n_I+1}, x_{n_I+3}, \dots, x_{2n_I+1}, \dots]$. Do the same
208 for the population mean and each time series of the initial population.

Thus we obtained errors for the mean values of the initial population, for each time series *in* the initial population and the time series that was improved using the presented genetic algorithm. Also, from all randomly generated interpolations we picked the one with the lowest RMSE to test it against the gen. alg. improved ones. Here, the root mean squared error E_{RMSE} , which is applied throughout this article, is given as:

$$E_{RMSE} = \left(\frac{1}{n} \sum_{i=1}^n [\hat{x}_i - x_i]^2 \right)^{\frac{1}{2}}, \quad (14)$$

209 where x_i are the original data points, \hat{x}_i are the predicted (in this case interpolated)
210 values and n is the length of the signal.

Table 1: Errors for the interpolated data on the Lorenz system depending on the number of interpolation points. The errors are shown for the mean interpolation of all populations and for the interpolation that was improved using the presented genetic algorithm. *Lowest RMSE in population* refers to the best randomly interpolated result, i.e. the one interpolation from the population that produced the lowest error by chance. We also featured the results for the linear and spline interpolation. We highlighted the interpolations where the genetic-algorithm-based interpolation outperformed the whole population of interpolations. Further, we give the percentage of how much of the population is outperformed by the genetic algorithm improved interpolation. This table is depicted in Figure 3.

n_I	1	2	3	4	5	6	7	8	9	10
RMSE population mean	0.77419	0.88263	0.91026	0.89442	0.87013	0.86858	0.89120	0.84220	0.90323	0.88777
lowest RMSE in population	0.17939	0.18068	0.16757	0.19206	0.16126	0.18134	0.17782	0.17211	0.18216	0.19371
RMSE linear interpolated	0.42534	0.44752	0.44179	0.44185	0.41574	0.42406	0.42894	0.40883	0.43263	0.43353
RMSE spline interpolated	0.12263	0.11808	0.09968	0.12586	0.09678	0.12195	0.12280	0.10862	0.11554	0.13008
RMSE gen. alg. improved	1.03779	0.24381	0.17517	0.19488	0.16264	0.18182	0.17818	0.17121	0.18239	0.19374
below best %	74.4%	21.6%	4.3%	2.2%	1.4%	1.1%	0.7%	0.1%	0.8%	0.3%

n_I	11	12	13	14	15	16	17	18	19	20
RMSE population mean	0.90844	0.92145	0.91509	0.90686	0.91750	0.90326	0.90789	0.90080	0.88835	0.89651
lowest RMSE in population	0.18789	0.19238	0.18693	0.18632	0.19943	0.19640	0.19208	0.18449	0.18415	0.20291
RMSE linear interpolated	0.43649	0.44211	0.43687	0.43423	0.44142	0.43534	0.43720	0.43170	0.42685	0.43398
RMSE spline interpolated	0.12086	0.12646	0.11530	0.11765	0.12532	0.12873	0.13098	0.12581	0.11912	0.12581
RMSE gen. alg. improved	0.18816	0.19141	0.18670	0.18626	0.19943	0.19663	0.19215	0.18462	0.18441	0.20300
below best %	0.8%	0.1%	0.1%	0.1%	0.1%	0.8%	0.6%	0.6%	0.8%	0.6%

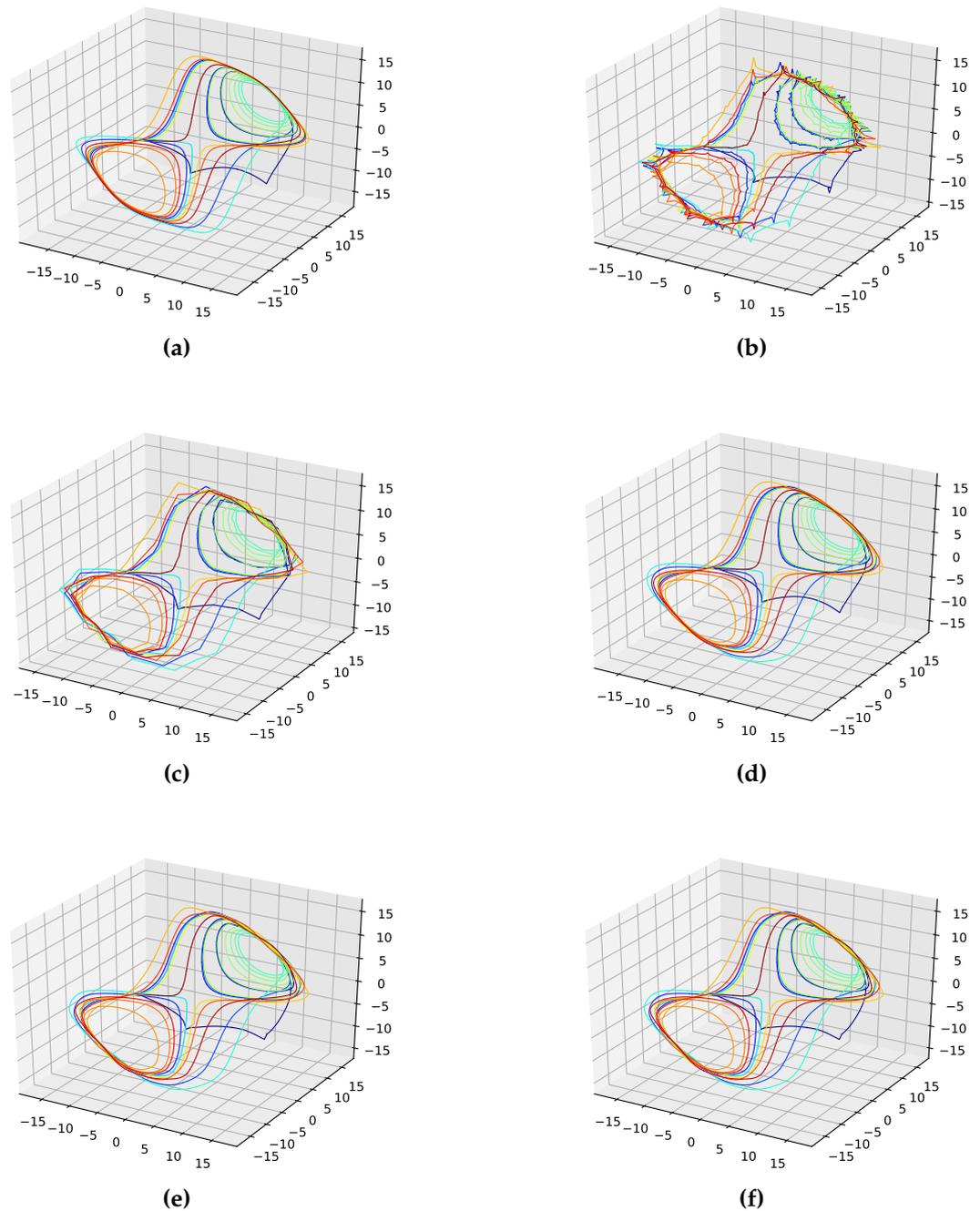


Figure 2. Reconstructed attractors for the interpolated Lorenz system.

(a): Non-interpolated original data (i.e. the one the errors are calculated with);

(b): The average interpolation of the whole population;

(c): Linear interpolated;

(d): Spline interpolated;

(e): The one interpolation of the population that has the lowest RMSE;

(f): Interpolation improved by the presented genetic algorithm approach.

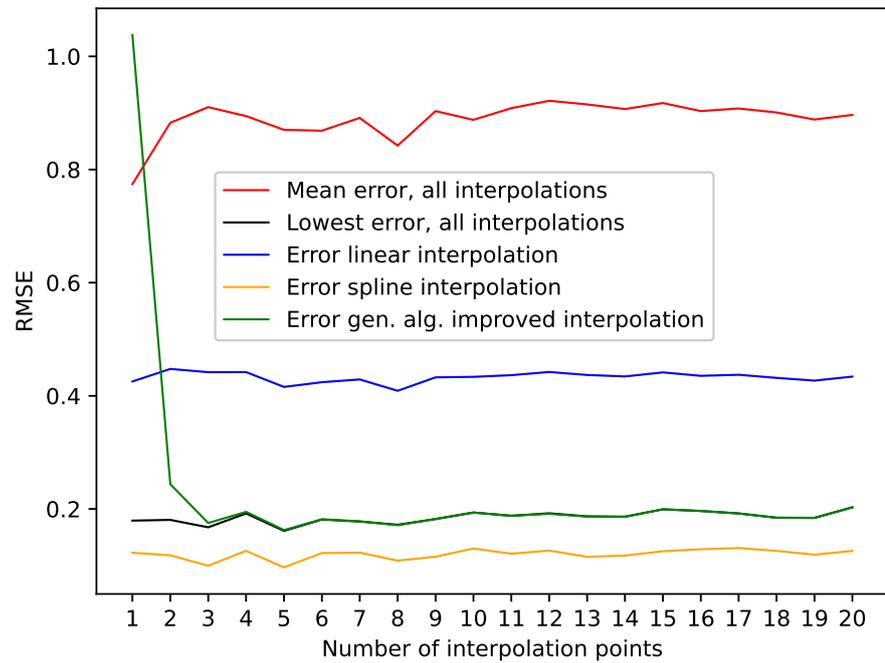


Figure 3. Shown in this Figure are the errors from Table 1 depending on the different numbers of interpolation points.

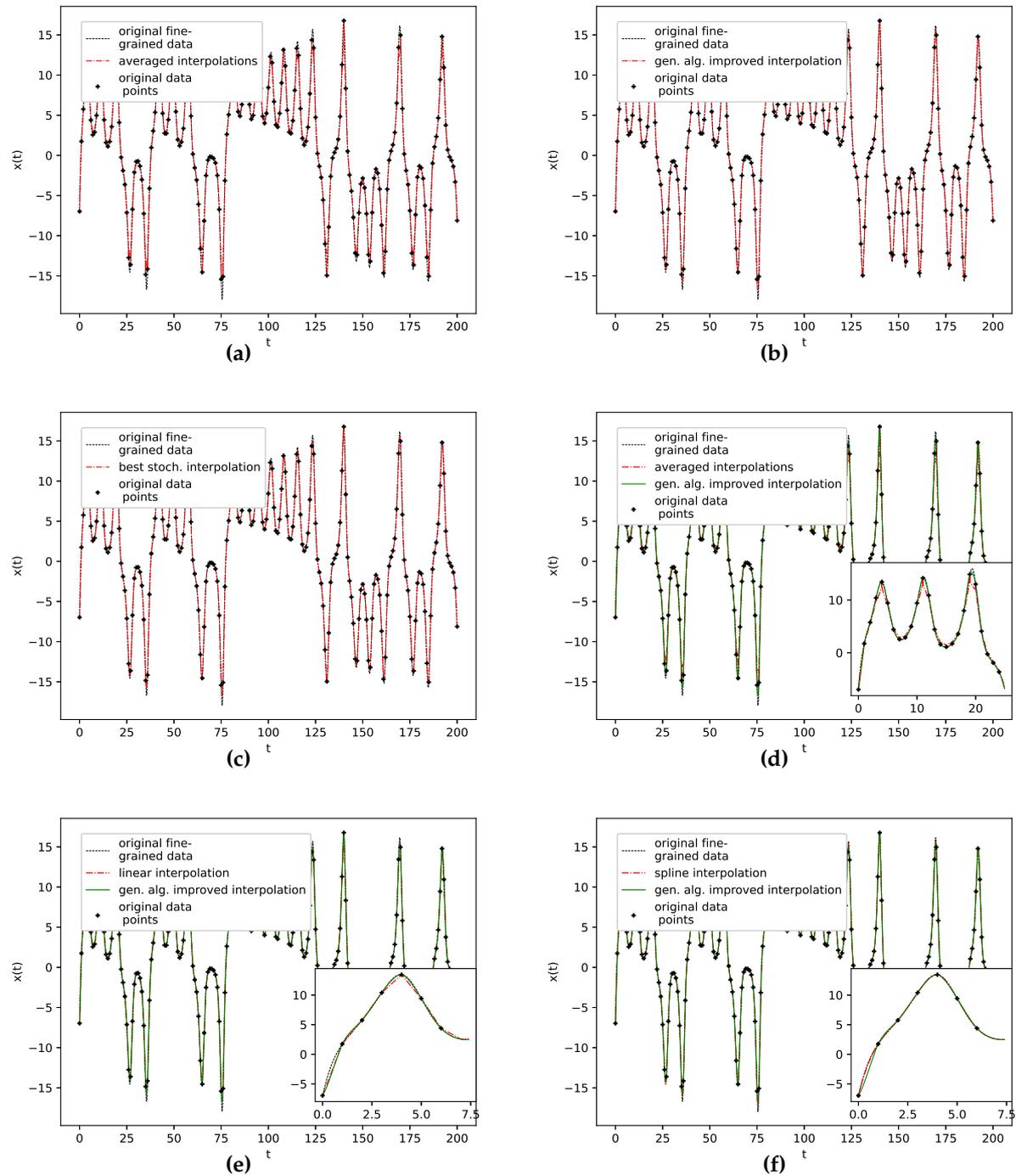


Figure 4. Original vs. interpolated time series data.

- (a): Non-interpolated original data (i.e. the one the error's are calculated with) and population average;
 (b): Genetic-algorithm-improved interpolation;
 (c): The one interpolation of the population that has the lowest RMSE;
 (d): Population average vs. genetic-algorithm-improved interpolation;
 (e): Linear interpolation vs. genetic-algorithm-improved interpolation;
 (f): Spline interpolation vs. genetic-algorithm-improved interpolation.

211 The presented results for the Lorenz system show that the algorithm can identify/
 212 tify/generate the best interpolation in terms of a low RMSE on missing data points

213 out of the given initial population. This can be seen in Table 1, where we highlighted
214 the results where the genetic-algorithm-improved-interpolation outperformed every
215 random interpolation of the population. Still the spline interpolation outperforms the
216 presented approach. This is also depicted in Figure 3, where we plotted the RMSE on
217 missing data points for varying numbers of interpolation points. This graphic shows
218 that the presented approach requires a certain amount of interpolation points, in this
219 case, three, to be close to the best random interpolation of the population. We assume
220 that the reason for this is that the variance of second derivatives along a phase space
221 trajectory requires a certain *density* of phase space points to be able to differ between
222 smooth and edgy phase space trajectories. The spline interpolation on the other hand
223 performs well right from the start.

224 The corresponding reconstructed phase space plots (Figure 2) show that both the
225 best random interpolation (e) and the genetic-algorithm-improved interpolation (f)
226 provide convincing phase space portraits, as both are indeed close to the ground truth
227 (a). On the other hand, the population mean (b) is far off and features many sharp edges
228 and pointy sections. Also, the linear interpolation (c) provides a very edgy phase space
229 portrait, just as one would expect from a linear interpolation. Contrary to that, from
230 all presented phase space portraits, the one for the spline interpolation (d) is most similar
231 to the original one, i.e. even the initial abbreviations caused by the time delay are almost
232 perfectly reconstructed.

233 We further plotted all obtained results for 13 interpolation points as time series in
234 Figure 4. The results show that the population mean (a) is far off the ground truth, and
235 differs drastically at the high and low peaks, as it does not reach the actual data points.
236 Both the genetic-algorithm-improved (b) and the best random interpolation of the initial
237 population (c) capture most of the high and low peaks compared to the population mean.
238 Further, when comparing the genetic-algorithm-improved and the population mean (d),
239 one can see that the improved interpolation provides a smoother curve when depicted
240 as a time series. In contrast, the population mean tends to produce sharp peaks. Finally,
241 we compare the linear interpolation (e) and the spline interpolation (f) to the genetic
242 algorithm improved interpolation. The linear interpolation here is far off, but the spline
243 interpolation reproduces the Lorenz system almost perfectly, and thus outperforms the
244 genetic-algorithm-improved interpolation.

245 4.2. Results for Non-Model Data Sets

246 This section tests our interpolation scheme on real-life data sets that possess only a
247 limited number of sampled data points. But these are actually the focus of the proposed
248 method, i.e., to increase the fine-grainedness of short, sparsely sampled time series data,
249 e.g., environmental or agricultural data sets. We must stress that our method is not
250 restricted to equidistant time series: Due to the general form of the bridge construction (
251 2), non-equidistant time series excerpts can be interpolated as well.

252 For this reason, we chose 5 data sets to demonstrate our method further, i.e., first
253 we validated the interpolation with missing data points and then we present an actual
254 interpolation and the improved phase space trajectories for each time series data. We
255 consider a phase space trajectory to be improved if we achieve smoother trajectories,
256 which exhibit fewer edgy points in a phase space representation. For this reason, we
257 chose only data sets that can be depicted in reconstructed three-dimensional phase space,
258 as an, e.g., four-dimensional phase space makes an intuitive understanding impossible.

259 The validation on these non-model data sets is performed such that every second
260 data point of the original time series is deleted, then all the gaps are interpolated to
261 reconstruct the missing data points. The results are shown for the average prediction of
262 the population, the random interpolation with the lowest RMSE, a linear interpolation, a
263 cubic spline interpolation and the improved interpolation using the presented genetic
264 algorithm. This section features only the validation errors, the corresponding plots are
265 collected in Appendix A to keep the main text focused.

266 4.2.1. NYC Measles outbreaks data Set

267 This is a data set that we obtained from [18], where it is discussed and shown to
268 feature an attractor structure in the embedded phase space. The corresponding original
269 source is [19]. It depicts measles outbreaks in New York City (NYC) from 1928 to 1964,
270 binned every two weeks, with a total of 432 data points.

271 The corresponding phase space embedding, with a time delay $\tau = 1$, was normal-
272 ized such that the range of all data is between $[0, 1]$.

273 The results on how well the presented interpolation can reproduce missing data
274 points of this data set are collected in Table 2 and depicted in Figure 11 (a). These
275 results show that, though the genetic-algorithm-improved interpolation drastically out-
276 performs the average random interpolation, the algorithm did not once outperform the
277 best interpolation of the population. Still, starting with seven interpolation points, the
278 genetic-algorithm-improved interpolation performs well and is very close to the best of
279 1000 randomly interpolated results, i.e., always below or around the best 1% of the pop-
280 ulation. Further, *PhaSpaSto*-interpolation does outperform the cubic spline interpolation
281 starting with five interpolation points. We thus conclude that the presented interpolation
282 technique captures the phase-space properties of this data set and effectively can be used
283 to interpolate this time series data. Also, compared to the cubic and linear interpolation,
284 the proposed method takes at least seven interpolation points to reach peak performance
285 for this data set. All validation plots, are collected in Appendix A.2.

286 An interpolation of the original data set is depicted in Figure 5. Comparing the
287 reconstructed phase space of the original data set, the population mean (c), and the
288 presented interpolation technique (d); we see that the phase space portrait of the latter
289 features a smoothed-out phase space trajectory compared to the original time series
290 (b) and the population mean (c), which are both pointy and have many sharp edges.
291 Further, considering the graph of the actual time series (a), we see that the presented
292 interpolation technique increases the major peaks, thus making extreme events more
293 prominent.

Table 2: Errors for the interpolated data on the NYC measles data set depending on the number of interpolation points. The errors are shown for the mean interpolation of all populations, as well as for the lowest error in the population and for the interpolation that was improved using the presented genetic algorithm. We highlighted the interpolation where the genetic-algorithm-based interpolation performed best. The corresponding plots for the best interpolation are shown in Appendix A.2. Further, we give the percentage of how much of the population is outperformed by the genetic algorithm improved interpolation.

n_I	1	3	5	7	9	11	13	15
RMSE Population Mean	860.56140	860.56165	860.56098	860.56235	860.56124	860.56210	860.56145	860.56220
Lowest RMSE in population	594.27833	594.27832	594.27832	594.27832	594.27748	594.27831	594.27832	594.27833
RMSE linear interpolated	713.61079	713.61089	713.61089	713.61089	713.61089	713.61089	713.61089	713.61089
RMSE spline interpolated	607.03778	607.03778	607.03778	607.03778	607.03778	607.03778	607.03778	607.03778
RMSE gen. alg. improved	1138.28460	621.70136	602.03361	594.36367	594.33054	594.34891	594.34819	594.34132
Below Best %	75.3%	25.8%	13.4%	0.8%	0.8%	0.8%	0.8%	0.8%

n_I	17	19	21	23	25	27	29	31
RMSE Population Mean	860.56196	860.56132	860.56168	860.56090	860.56153	860.56287	860.56138	860.56192
Lowest RMSE in population	594.27901	594.27750	594.27934	594.27834	594.28039	594.27831	594.28069	594.27833
RMSE linear interpolated	713.61089	713.61089	713.61089	713.61089	713.61089	713.61089	713.61089	713.61089
RMSE spline interpolated	607.03778	607.03778	607.03778	607.03778	607.03778	607.03778	607.03778	607.03778
RMSE gen. alg. improved	594.33772	594.33837	594.33400	594.35508	594.31806	594.36050	594.42183	594.39145
Below Best %	0.8%	0.8%	0.8%	0.8%	0.6%	0.8%	1.1%	1.1%

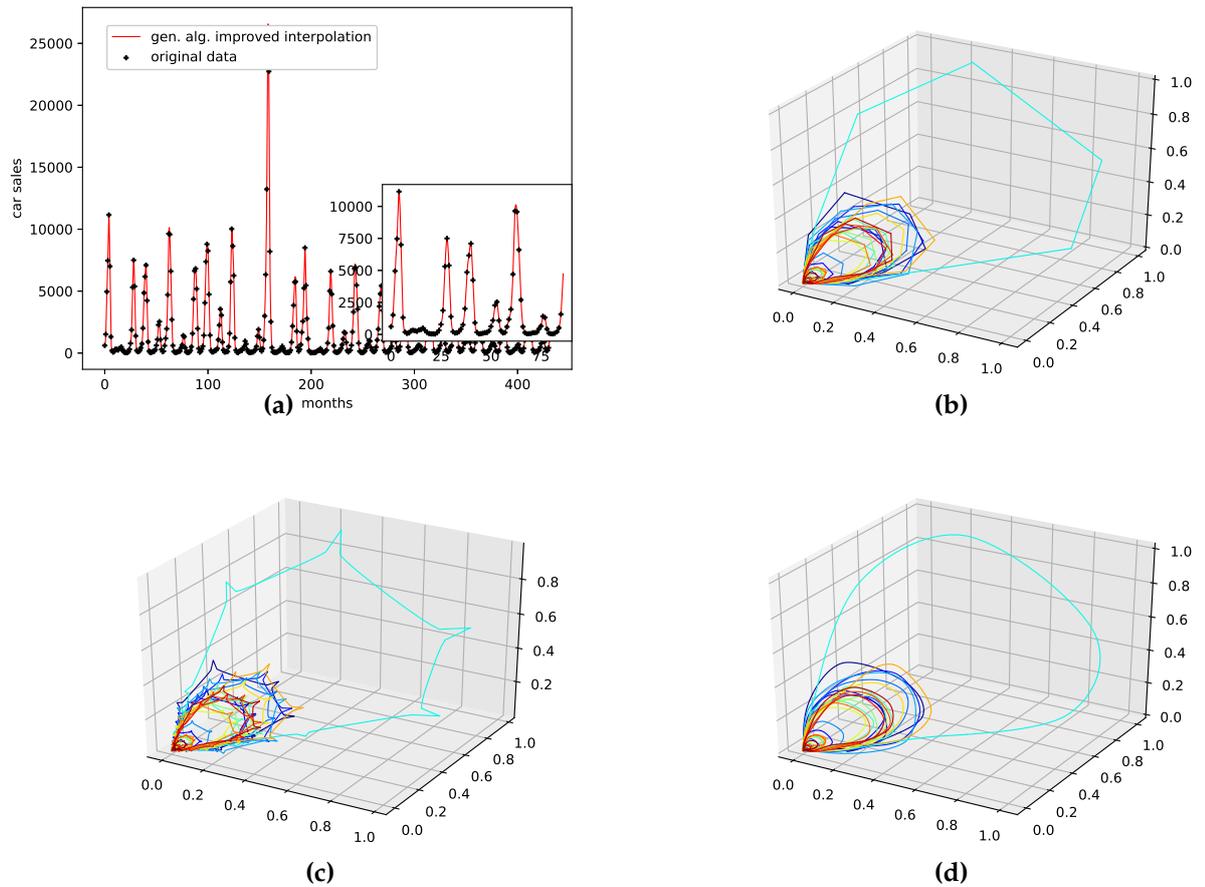


Figure 5. Interpolated data and reconstructed attractors for the NYC measles outbreaks data set.

(a): The original and interpolated time series data;

(b): Phase space reconstruction of the original data;

(c): Phase space reconstruction of the average population data;

(d): Phase space reconstruction of the genetic-algorithm-improved data.

294 4.2.2. Car Sales in Quebec Data Set

295 This is a data set from the Time Series Data Library, [20]. It depicts monthly car
296 sales in Quebec from January 1960 to December 1968, with a total of 108 data points.

297 The corresponding phase space embedding, with a time delay $\tau = 1$, was detrended
298 by subtracting a linear fit from the data and normalized such that the range of all data is
299 between $[0, 1]$.

300 The results on how well the presented interpolation can reproduce missing data
301 points of this data set are collected in Table 3, and depicted in Figure 11 (b). The
302 genetic-algorithm-improved interpolation drastically outperforms the average random
303 interpolation. Further, the *PhaSpaSto*-interpolation always outperforms the cubic spline
304 and linear interpolation. Still, the genetic-algorithm-improved interpolations did not
305 always outperform the best random interpolations, but for one, three, and five inter-
306 polation points. Overall, the genetic-algorithm-improved interpolation performs well
307 and is very close to the best of 1000 randomly interpolated results, i.e., for most cases
308 below or around the best 1% of the population. Thus, we conclude that the presented
309 interpolation technique effectively captures the phase-space properties of this data set
310 and can be used to interpolate this time series data. All additional plots for the validation
311 are collected in Appendix A.3, whereas one can find the reconstructed attractors for all
312 interpolated validation sets and the corresponding time series plots.

313 An interpolation of the original data set is depicted in Figure 6. Here Figures 6
314 (c) and (d) present the population mean and the improved interpolation, respectively.
315 When comparing them, one can see that the genetic algorithm improves the phase space
316 portrait in terms of a smoothed-out phase space trajectory compared to the original time
317 series (b) and the population mean (c), which are both pointy and have many sharp edges.
318 When considering the actual time-series graph (a), the presented interpolation technique
319 increases the major peaks, thus making extreme events more prominent. Further, it
320 provides a rather smooth curve, i.e., no pointy edges, as depicted in the zoomed-in plot
321 in (a).

Table 3: Errors for the interpolated data on the car sales in Quebec data set depending on the number of interpolation points. The errors are shown for the mean interpolation of all populations, the linear interpolation, the cubic spline interpolation, as well as for the lowest error in the population and for the interpolation that was improved using the presented genetic algorithm. We highlighted the interpolation where the genetic-algorithm-based interpolation performed best. The corresponding plots for the best interpolation are shown in Appendix A.3. Further, we give the percentage of how much of the population is outperformed by the genetic algorithm improved interpolation.

n_I	1	3	5	7	9	11	13	15
RMSE Population Mean	2030.11005	2030.11166	2030.11148	2030.11230	2030.11030	2030.11138	2030.11106	2030.11110
Lowest RMSE in population	1954.95010	1954.95013	1954.95016	1954.95013	1954.95005	1954.95009	1954.95020	1954.95015
RMSE linear interpolated	2017.79949	2017.79949	2017.79949	2017.79949	2017.79949	2017.79949	2017.79949	2017.79949
RMSE spline interpolated	1971.23755	1971.23755	1971.23755	1971.23755	1971.23755	1971.23755	1971.23755	1971.23755
RMSE gen. alg. improved	1907.40084	1960.21475	1954.94790	1954.94792	1954.95375	1954.97452	1958.57232	1954.97468
Below Best %	0.1%	17.2%	0.1%	0.1%	0.6%	1.01%	14.6%	1.01%
n_I	17	19	21	23	25	27	29	31
RMSE Population Mean	2030.11260	2030.11057	2030.11226	2030.11047	2030.11078	2030.11105	2030.11171	2030.11013
Lowest RMSE in population	1954.95010	1954.95007	1954.95011	1954.95014	1954.95007	1954.95010	1954.95003	1954.95021
RMSE linear interpolated	2017.79949	2017.79949	2017.79949	2017.79949	2017.79949	2017.79949	2017.79949	2017.79949
RMSE spline interpolated	1971.23755	1971.23755	1971.23755	1971.23755	1971.23755	1971.23755	1971.23755	1971.23755
RMSE gen. alg. improved	1954.97730	1954.99153	1955.00052	1954.99273	1955.02450	1955.02418	1955.01367	1954.98108
Below Best %	1.3%	1.4%	1.4%	1.4%	1.6%	1.6%	1.4%	1.4%

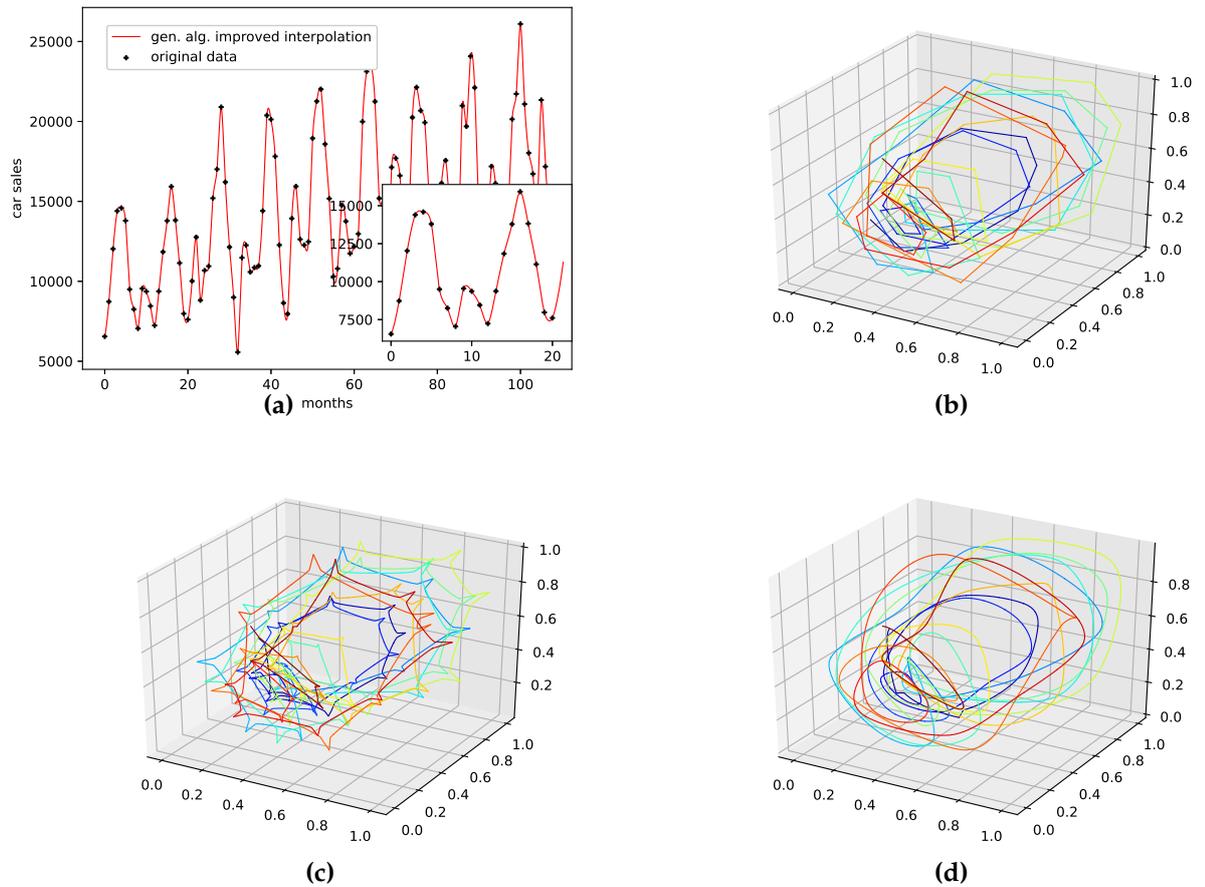


Figure 6. Interpolated data and reconstructed attractors for the car sales in Quebec data set.

(a): The original and interpolated time series data;

(b): Phase space reconstruction of the original data;

(c): Phase space reconstruction of the average population data;

(d): Phase space reconstruction of the genetic-algorithm-improved data.

322 4.2.3. Perrin Freres Champagne Sales Data Set

323 This is a data set from the Time Series Data Library, [20]. It depicts Perrin Freres
324 Champagne sales from January 1964 to September 1972, with a total of 105 data points.

325 The corresponding phase space embedding, with a time delay $\tau = 1$, was normal-
326 ized such that the range of all data is between $[0, 1]$.

327 The results on how well the presented interpolation can reproduce missing data
328 points of this data set are collected in Table 4 and Figure 11 (c).

329 Though the genetic-algorithm-improved interpolation drastically outperforms the
330 average random interpolation, the algorithm did not once outperform the best inter-
331 polation of the population. Still, starting with five interpolation points, the genetic-
332 algorithm-improved interpolation performs well and is very close to the best of 1000
333 randomly interpolated results, i.e., consistently below or around the best 1% of the
334 population. Overall the cubic spline interpolation performed best on this data set. The
335 linear interpolation, though outperforming the population mean, is still far off. We
336 thus conclude that the presented interpolation technique does capture the phase-space
337 properties of this data set from a given population and can be used to interpolate this
338 time series data, but the cubic spline interpolation is the better choice.

339 An interpolation of the original data set is depicted in Figure 5. We again show the
340 population mean (c) and the improved interpolation (d). The presented interpolation
341 technique improves the phase space portrait in terms of a smoothed-out phase space
342 trajectory (d) compared to the original time series (b), and the population mean (c),
343 which are both pointy and have many sharp edges. Further, considering the graph of the
344 actual time series (a), the presented interpolation technique increases the major peaks,
345 thus making extreme events more prominent and provides a thoroughly smooth curve,
346 as depicted in the zoom-in window in (a).

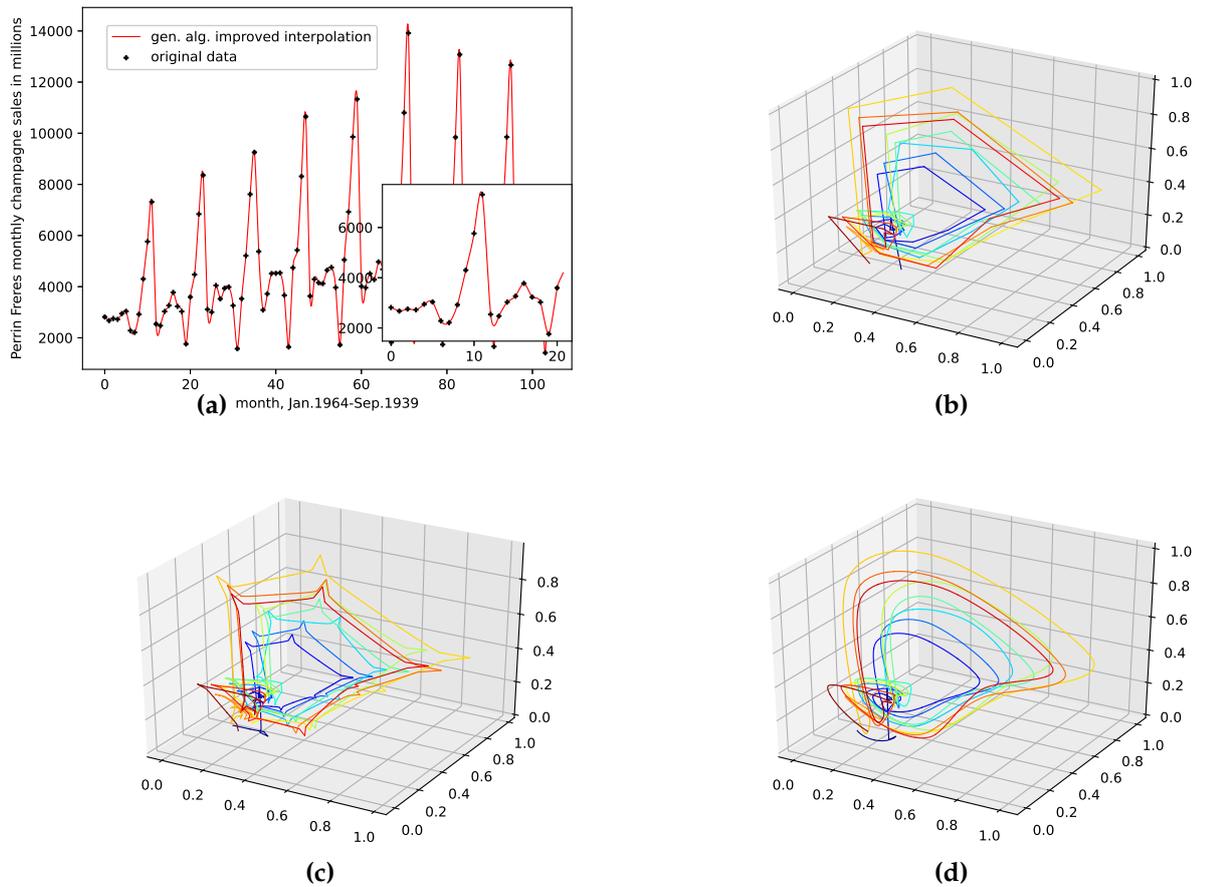


Figure 7. Interpolated data and reconstructed attractors for the Perrin Freres Champagne sales data set.

(a): The original and interpolated time series data;

(b): Phase space reconstruction of the original data;

(c): Phase space reconstruction of the average population data;

(d): Phase space reconstruction of the genetic-algorithm-improved data.

347 4.2.4. Monthly Airline Passengers Data Set

348 This is a data set from the Time Series Data Library, [20]. It depicts monthly
349 international airline passengers from January 1949 to December 1960, with a total of 144
350 data points, given in units of 1000.

351 The corresponding phase space embedding, with a time delay $\tau = 1$, was detrended
352 by subtracting a linear fit from the data and normalized such that the range of all data is
353 between $[0, 1]$.

354 The results on how well the presented interpolation can reproduce missing data
355 points of this data set are collected in Table 5 and depicted in Figure 11 (d). The results
356 show that, though the genetic-algorithm-improved interpolation drastically outperforms
357 the average random interpolation, the algorithm did not once outperform the best inter-
358 polation of the population. Still, starting with three interpolation points, the algorithm
359 did outperform both the linear and the cubic spline interpolation. What's curious,
360 though, is that, for this data set, of all the non-model data sets, the linear interpolation
361 outperforms the cubic spline interpolation.

362 The genetic-algorithm-improved interpolation does not perform that well for this
363 data set compared to a random interpolation of the time series. The improved interpola-
364 tion is only around the best $\approx 40\%$ of the initial population for this data set. We thus
365 conclude that the presented interpolation technique does not capture the phase-space
366 properties of this data set very well. Still, the genetic algorithm does improve the initial
367 population such that the population mean, the linear interpolation, and the cubic spline
368 interpolation are outperformed, starting with three interpolation points. All-time series
369 and reconstructed attractor plots for this data set can be found in Appendix A.5.

370 An actual interpolation of the original data set is depicted in Figure 8. We again
371 show the population mean (c) and the improved interpolation (d). The presented
372 *PhaSpaSto*-interpolation (d) improves the phase space portrait in terms of a smoothed-
373 out phase space trajectory, compared to the original time series (b) and the population
374 mean (c), which are both pointy and have many sharp edges. Further, considering
375 the actual time series (a) graph, we see that the presented interpolation technique
376 slightly increases the major peaks. Also, compared to the other non-model data sets, the
377 improved interpolation does provide a relatively smooth curve, but it looks way sharper
378 than for, e.g., the car sales in Quebec (See Figure 6 (a))data set.

Table 5: Errors for the interpolated data on the monthly airline passengers data set data set depending on the number of interpolation points. The errors are shown for the mean interpolation of all populations, the linear interpolation, the cubic spline interpolation, as well as for the lowest error in the population and for the interpolation that was improved using the presented genetic algorithm. We highlighted the interpolation where the genetic-algorithm-based interpolation performed best. The corresponding plots for the best interpolation are shown in Appendix A.5. Further, we give the percentage of how much of the population is outperformed by the genetic algorithm improved interpolation.

n_I	1	3	5	7	9	11	13	15
RMSE Population Mean	19.93996	19.93999	19.93841	19.94072	19.93976	19.93873	19.94070	19.93889
Lowest RMSE in population	16.55624	16.55779	16.55732	16.55753	16.55558	16.55836	16.55719	16.55776
RMSE linear interpolated	17.39496	17.39496	17.39496	17.39496	17.39496	17.39496	17.39496	17.39496
RMSE spline interpolated	18.33872	18.33872	18.33872	18.33872	18.33872	18.33872	18.33872	18.33872
RMSE gen. alg. improved	18.65257	16.81653	16.84539	17.02728	16.84536	16.84545	16.84540	16.84539
Below Best %	59.4%	35.6%	38.0%	42.20%	38.1%	38.0%	38.1%	38.0%

n_I	17	19	21	23	25	27	29	31
RMSE Population Mean	19.94029	19.94030	19.93985	19.93939	19.93659	19.94172	19.94023	19.93909
Lowest RMSE in population	16.55752	16.55730	16.55810	16.55715	16.55733	16.55603	16.55789	16.55741
RMSE linear interpolated	17.39496	17.39496	17.39496	17.39496	17.39496	17.39496	17.39496	17.39496
RMSE spline interpolated	18.33872	18.33872	18.33872	18.33872	18.33872	18.33872	18.33872	18.33872
RMSE gen. alg. improved	16.84546	16.84545	16.84548	16.84540	16.84535	16.84544	16.84544	16.84546
Below Best %	38.1%	38.1%	38.0%	38.1%	38.1%	38.1%	38.2%	38.1%

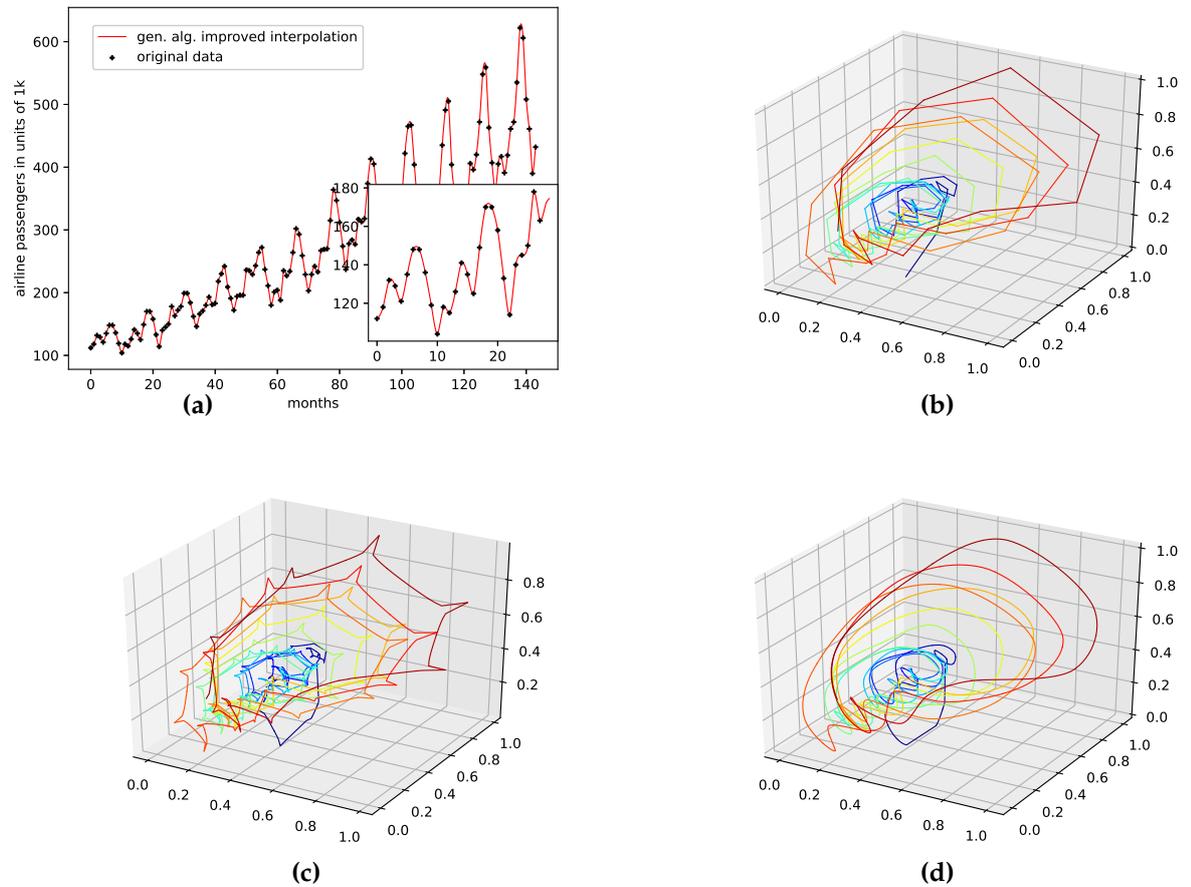


Figure 8. Interpolated data and reconstructed attractors for the monthly international airline passengers data set.

(a): The original and interpolated time series data;

(b): Phase space reconstruction of the original data;

(c): Phase space reconstruction of the average population data;

(d): Phase space reconstruction of the genetic-algorithm-improved data.

379 4.2.5. Monthly Mean Temperature in Nottingham Castle

380 This is a data set from the Time Series Data Library, [20]. It depicts the mean
381 monthly temperature in Nottingham castle from January 1920 to December 1939, given
382 in degrees Fahrenheit, with a total of 240 data points.

383 The corresponding phase space embedding, with a time delay $\tau = 3$, was normal-
384 ized such that the range of all data is between $[0, 1]$, with a time delay $\tau = 1$.

385 The corresponding phase space embedding, with a time delay $\tau = 1$, was detrended
386 by subtracting a linear fit from the data and normalized such that the range of all data is
387 between $[0, 1]$.

388 The results on how well the presented interpolation can reproduce missing data
389 points of this data set are collected in Table 6 and depicted in Figure 11. The results
390 show that, though the genetic-algorithm-improved interpolation drastically outperforms
391 the average random interpolation, the algorithm did not once outperform the best
392 interpolation of the population, although outperforming the linear and the cubic spline
393 interpolation. The genetic-algorithm-improved interpolation does not perform that well
394 for this data set compared to a random interpolation of the time series, as the improved
395 interpolation is only around the best $\approx 34\%$ for this data set. We thus conclude that the
396 presented interpolation technique does not capture the phase-space properties of this
397 data set very well. The corresponding time-series and reconstructed phase space plots
398 are collected in Appendix A.6.

399 An interpolation of the original data set is depicted in Figure 9. We again show the
400 population mean (c) and the improved interpolation (d). The presented interpolation
401 technique improves the phase space portrait (d) in terms of a smoothed-out phase space
402 trajectory compared to the original time series (b) and the population mean (c), which
403 are both pointy and have many sharp edges. Also, given the time-series depiction of the
404 *PhaSpaSto*-interpolation (Figure 9 (a)), we see the same behavior as for all the other data
405 sets; the major peaks are increased.

Table 6: Errors for the interpolated data on the monthly mean temperature in Nottingham castle data set depending on the number of interpolation points. The errors are shown for the mean interpolation of all populations, as well as for the lowest error in the population and for the interpolation that was improved using the presented genetic algorithm. We highlighted the interpolation where the genetic-algorithm-based interpolation performed best. The corresponding plots for the best interpolation are shown in Appendix A.6. Further, we give the percentage of how much of the population is outperformed by the genetic algorithm improved interpolation.

n_I	1	3	5	7	9	11	13	15
RMSE Population Mean	3.09115	3.09170	3.09167	3.09055	3.09088	3.09055	3.09166	3.09165
Lowest RMSE in population	2.47879	2.47858	2.47910	2.47890	2.47886	2.47901	2.47900	2.47875
RMSE linear interpolated	2.61279	2.61279	2.61279	2.61279	2.61279	2.61279	2.61279	2.61279
RMSE spline interpolated	2.59028	2.59028	2.59028	2.59028	2.59028	2.59028	2.59028	2.59028
RMSE gen. alg. improved	2.48413	2.50179	2.50279	2.50406	2.50420	2.50521	2.50512	2.505089
Below Best %	12.6%	31.3%	32.4%	33.5%	33.8%	34.4%	34.4%	34.1%

n_I	17	19	21	23	25	27	29	31
RMSE Population Mean	3.09095	3.09177	3.09115	3.09122	3.09146	3.09143	3.09179	3.09023
Lowest RMSE in population	2.47887	2.47920	2.47925	2.47899	2.47867	2.47941	2.47885	2.47892
RMSE linear interpolated	2.61279	2.61279	2.61279	2.61279	2.61279	2.61279	2.61279	2.61279
RMSE spline interpolated	2.59028	2.59028	2.59028	2.59028	2.59028	2.59028	2.59028	2.59028
RMSE gen. alg. improved	2.50494	2.50541	2.50529	2.50552	2.50505	2.50547	2.50550	2.50533
Below Best %	33.9%	34.6%	34.4%	35%	34.6%	34.6%	34.7%	34.7%

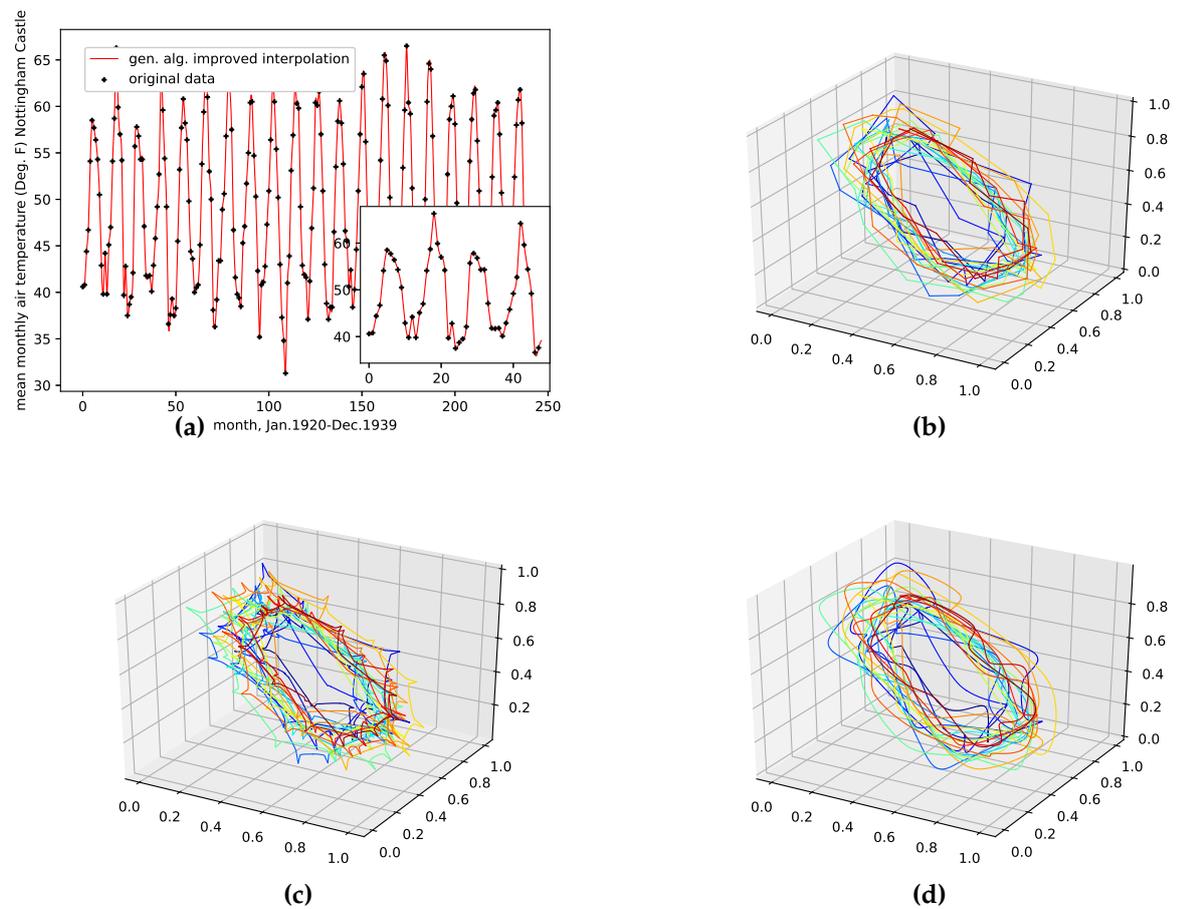


Figure 9. Interpolated data and reconstructed attractors for the monthly mean temperature in Nottingham castle data set.

(a): The original and interpolated time series data;

(b): Phase space reconstruction of the original data;

(c): Phase space reconstruction of the average population data;

(d): Phase space reconstruction of the genetic-algorithm-improved data.

406 4.3. Summary

407 We briefly summarize this research and point out the main findings:

- 408 • We presented a genetic algorithm to improve a stochastic interpolation, i.e., multi-
409 point fractional Brownian bridges, in order to enhance the reconstructed phase
410 space of any given time series data. For simplicity, we named this method *PhaSpaSto*-
411 interpolation.
- 412 • We presented a novel approach to measure the quality of a phase space reconstruc-
413 tion according to Taken's theorem. Here we used an idea from image processing,
414 i.e., to identify blurry images via the variance of second derivatives. These second
415 derivatives are calculated along the reconstructed phase space curve for any given
416 reconstructed phase space. The variance of these second derivatives is used to mea-
417 sure the quality of our phase space reconstruction. Given two interpolated phase
418 space curves of the same time series, the one with the lower variance of second
419 derivatives along the curve is found to be the better phase space reconstruction.
- 420 • We showed that the developed technique performed well in the case of a model
421 data set, i.e., one variable of the Lorenz system. Here we deleted data points
422 from the original time series data and were able to outperform, in some cases, any
423 best random interpolations of this time series data. Also, the presented method
424 outperformed a linear interpolation when it comes to finding the missing data
425 points. Still, the proposed method did not outperform the presented cubic spline
426 interpolation on this task. This was done to validate our method and to show its
427 applicability. Further, the presented reconstructed phase spaces plots show that the
428 interpolated phase space reconstruction is very similar to the original reconstructed
429 phase space. The results for the Lorenz system are collected in Section 4.1.
- 430 • We validated the presented method using five sparsely sampled non-model data
431 sets. The validation was done such that we deleted every second data point from the
432 original time series and reconstructed these missing data points using the developed
433 technique. For three out of five data sets, the developed method effectively can
434 identify the interpolations or parts of it with low errors, i.e., the result is around the
435 best 1% of the population in terms of the RMSE for the reconstructed data points.
436 Also, *PhaSpaSto*-interpolation outperformed the spline interpolation for four of five
437 non-model data sets and the linear interpolation on all non-model data sets. Also,
438 the best random interpolation outperformed the cubic spline interpolation on four
439 of five non-model data sets. For two data sets, the *PhaSpaSto*-interpolation does
440 not perform very well as it is only around the best 30 – 40% of all RMSEs of the
441 population. The interpolation performed well in case of the measles cases in NYC
442 data set (Section 4.2.1), the car sales in Quebec data set (Section 4.2.2) and the Perrin
443 Freres champagne sales data set (Section 4.2.3). The two cases where the presented
444 method did not perform well are the monthly international airline passengers data
445 set (Section 4.2.4) and the monthly mean temperature in Nottingham castle data set
446 (Section 4.2.5).
- 447 • We also used the five non-model data sets to show the applicability of the developed
448 technique as an actual interpolation technique, i.e., no deleted data points. The plots
449 of the reconstructed phase spaces show that it softens the edges and provides are
450 thoroughly smoother and cleaner reconstructed phase space trajectory. Therefore
451 the authors conclude that this technique applies to arbitrary univariate data sets.
452 All of these plots are collected in Section 4.2. We further recommend it when dealing
453 with sparsely sampled time series, e.g., attractor reconstructions of real-life time
454 series, e.g., to improve machine learning time-series predictions.

455 5. Conclusion

456 This article presents a novel approach to interpolate univariate time series data.
457 For simplicity, we named this method *PhaSpaSto*-interpolation. The concept is first to
458 generate a population of, e.g., 1000, different stochastically interpolated time series data.

459 This is done using multi-point Brownian bridges, each assigned with a random Hurst
 460 exponent. Then, as a second step, a genetic algorithm generates one time series out of
 461 the population with a low variance of second-order derivatives along the corresponding
 462 reconstructed phase space trajectory. I.e., we want this curve to be as smooth as possible.
 463 Using the variance of second-order derivatives is adapted from image processing, where
 464 the variance of second-order derivatives is used to differentiate between blurry and
 465 sharp images. To illustrate this idea, we plotted the randomly generated stochastic
 466 interpolations and the improved interpolation in Figure 10 for the monthly international
 467 airline passengers data set. We also tested the discussed approach with different loss
 468 functions that, in the end, did not work. These failed attempts are collected in Appendix
 469 C.

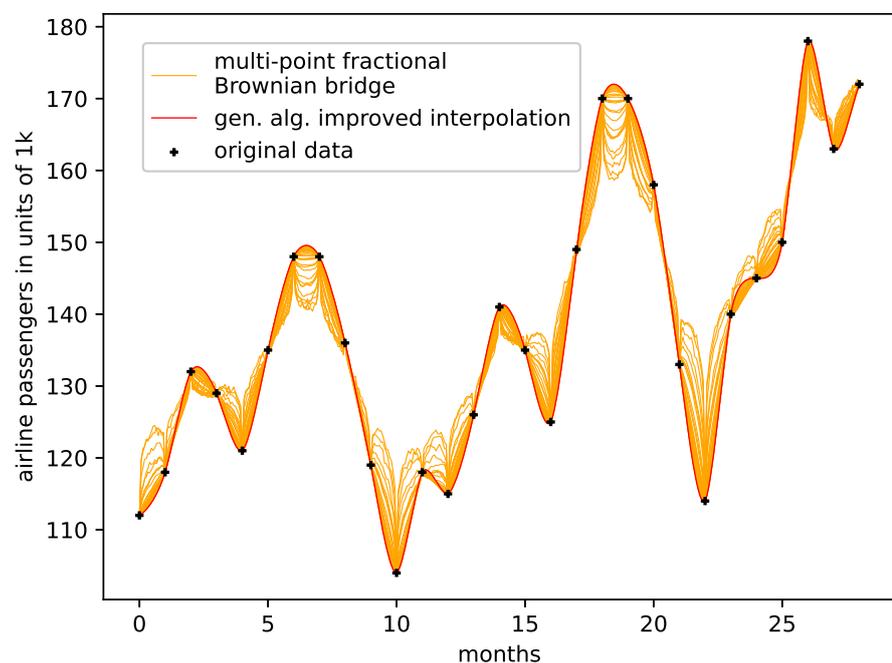


Figure 10. Random and gen. alg. improved interpolations for the monthly international airline passengers data set, see Section 4.2.4. The orange lines are the randomly generated stochastic interpolations, i.e., multi-point fractional Brownian bridges with different Hurst exponents.

470 We then applied the presented interpolation technique to the Lorenz system, or to be
 471 specific, to one of the variables of the Lorenz system, since we're dealing with univariate
 472 data only. We then deleted data points from this time series and interpolated the missing
 473 data points with the presented interpolation technique. We also tested the proposed
 474 approach against a linear and a cubic spline interpolation. The results show that the
 475 presented attractor-based stochastic interpolation can reproduce the Lorenz system, i.e.,
 476 the genetic algorithm can find the best parts of the initial population to reconstruct the
 477 Lorenz system. Still, the spline interpolation outperformed the *PhaSpaSto*-interpolation
 478 for the Lorenz system. Finally, we applied the presented approach to various real-life
 479 and/or benchmark data sets. There's no fine-grained model data available for data sets
 480 like these. We cannot verify the interpolation as we did with the Lorenz system. Instead,
 481 we deleted every second data point of these data sets and reconstructed them using the
 482 developed method, i.e., generated interpolations using a range of interpolation points,
 483 picking the missing data points, and verifying them against the ground truth. *PhaSpaSto*-
 484 interpolation performed well on three of five data sets, as the genetic algorithm can, in

485 fact, identify/build interpolations with low errors for the missing data points. Further,
486 *PhaSpaSto*-interpolation outperformed the spline-interpolation on four of five data sets.
487 Thus, we conclude that the presented method can also be applied to non-model data
488 sets. Lastly, we show actual interpolations on these non-model data sets, i.e., no deleted
489 data points. Given that the reader is familiar with how strange attractors of chaotic
490 systems look, it should be clear from the presented reconstructed phase space portraits
491 that our approach can interpolate real-life data as one would expect a phase space
492 embedding of a strange attractor to look like, see Section 4.2. Future research will also be
493 devoted to generalizations of the bridge process (2) to random processes which exhibit
494 non-Gaussian features [21].

495 We expect the presented research to be useful for predicting and analyzing sparsely
496 sampled time series data, e.g., in agriculture or other fields where fine-grained mea-
497 surements are expensive. We further expect the presented research to be utilized for
498 improving machine and deep learning approaches with insufficient data. We expect that
499 an enhanced phase space structure improves forecasts' accuracy.

500 Future improvements and applications of this technique include the expansion to
501 multi-variate data sets and using the presented loss function, i.e., the variance of second
502 derivatives along phase space trajectories, to find better phase space embeddings. The
503 interested reader is therefore referred to appendix B where we present the loss-surface of
504 the Lorenz system with a varying time delay τ and a varying embedding dimension d_E .

505 And, as previously mentioned, we want to test to what degree improved phase-
506 space embeddings can be beneficial for machine and/or deep learning approaches for
507 learning and predicting time series data.

508 Lastly, the presented interpolation technique code will be available on GitHub from
509 the corresponding author in the future.

510 Acknowledgments

511 The authors acknowledge the funding of the project "DiLaAg – Digitalization and
512 Innovation Laboratory in Agricultural Sciences", by the private foundation "Forum
513 Morgen", the Federal State of Lower Austria, by the FFG; Project AI4Crop, No. 877158
514 and by TU Wien Bibliothek for financial support through its Open Access Funding
515 Program. J.F. acknowledges funding from the Humboldt Foundation within a Feodor-
516 Lynen fellowship.

517 Appendices

518 A. Additional Plots

519 This section provides additional plots for all data sets discussed in Section 4.2. As
520 such we plotted the evolution of errors for the validation depending on the varying
521 number of interpolation points, i.e. the errors from Tables 2, 3, 4, 5 and 6. Further, we
522 added each time series and the corresponding best validation interpolation, and finally,
523 the corresponding phase space plots.

524 A.1. Evolution of Errors

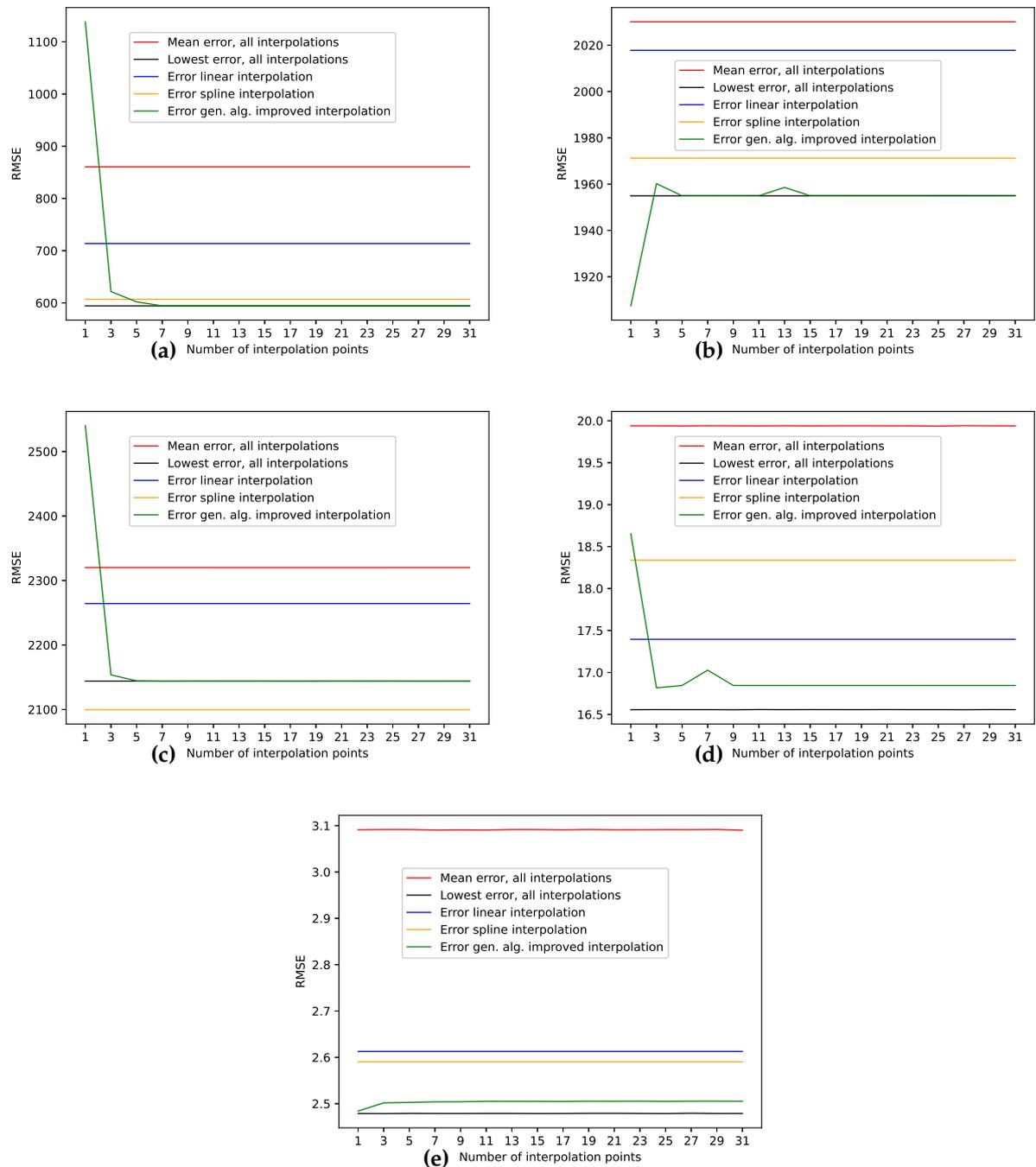


Figure 11. Evolution of errors depending on the number of interpolation points for the non-model data validation.

(a): Measles cases in NYC data set, results from Table 2;

(b): Car Sales in Quebec data set, results from Table 3;

(c): Perrin Freres champagne sales data set, results from Table 4;

(d): Monthly international airline passengers data set, results from Table 5;

(e): Monthly mean temperature in Nottingham castle data set, results from Table 6;

525 A.2. NYC Measles Outbreaks Data Set

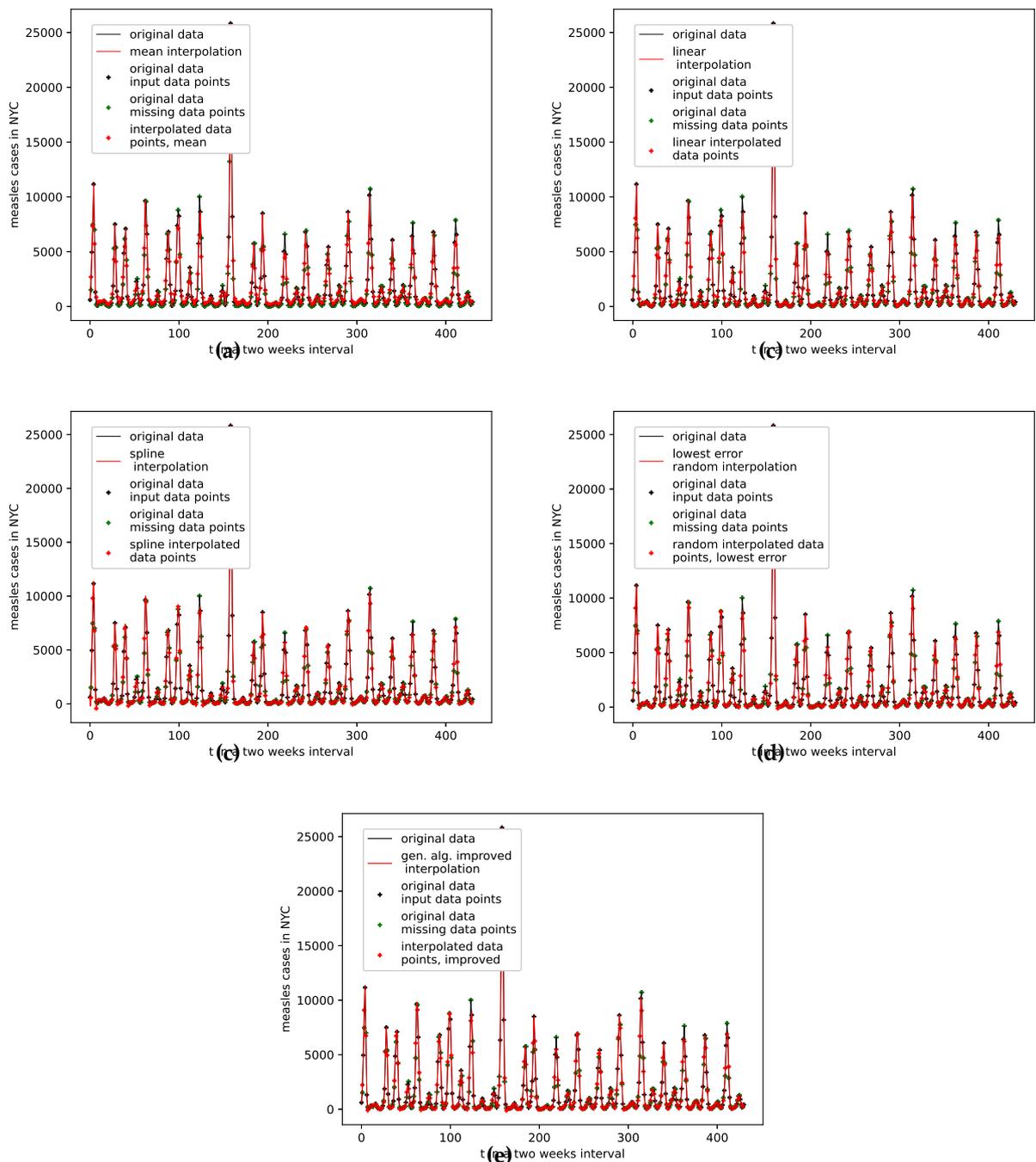


Figure 12. Interpolated validation data (25 interpolation points) for the measles cases in NYC data set.

- (a): Average population validation;
- (b): Validation, linear interpolation;
- (c): Validation, spline interpolation;
- (d): Validation, best random interpolation;
- (e): Validation, gen. alg. improved interpolation;

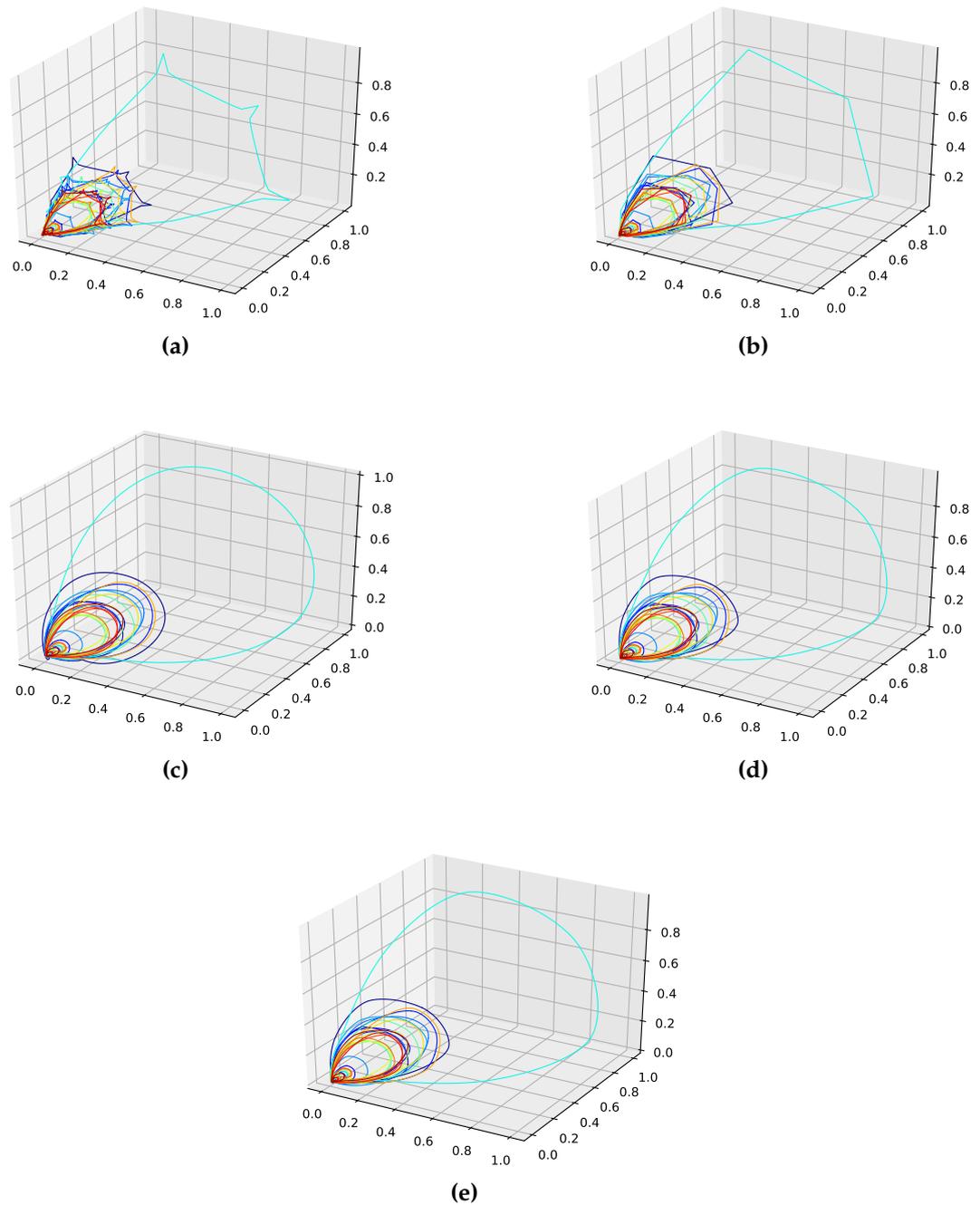


Figure 13. Reconstructed validation attractors (25 interpolation points) for the measles cases in NYC data set.

- (a): Reconstructed attractor, average population validation interpolation;
- (b): Reconstructed attractor, linear interpolation;
- (c): Reconstructed attractor, spline interpolation;
- (d): Reconstructed attractor, best random validation interpolation;
- (e): Reconstructed attractor, gen. alg. improved validation interpolation;

526 A.3. Car Sales in Quebec Data Set

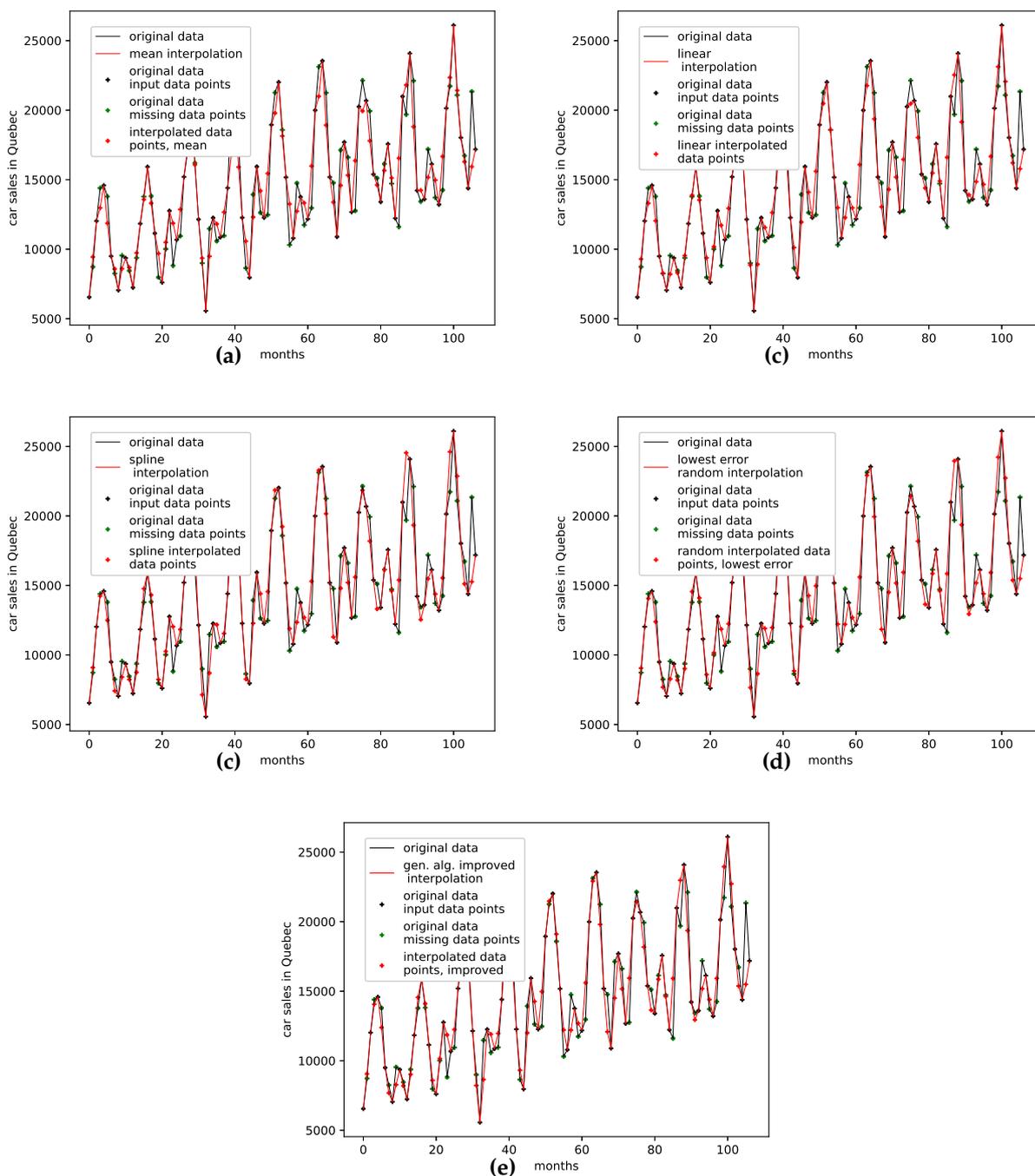


Figure 14. Interpolated validation data (one interpolation point) for the car sales in Quebec data set.

- (a):** Average population validation;
- (b):** Validation, linear interpolation;
- (c):** Validation, spline interpolation;
- (d):** Validation, best random interpolation;
- (e):** Validation, gen. alg. improved interpolation;

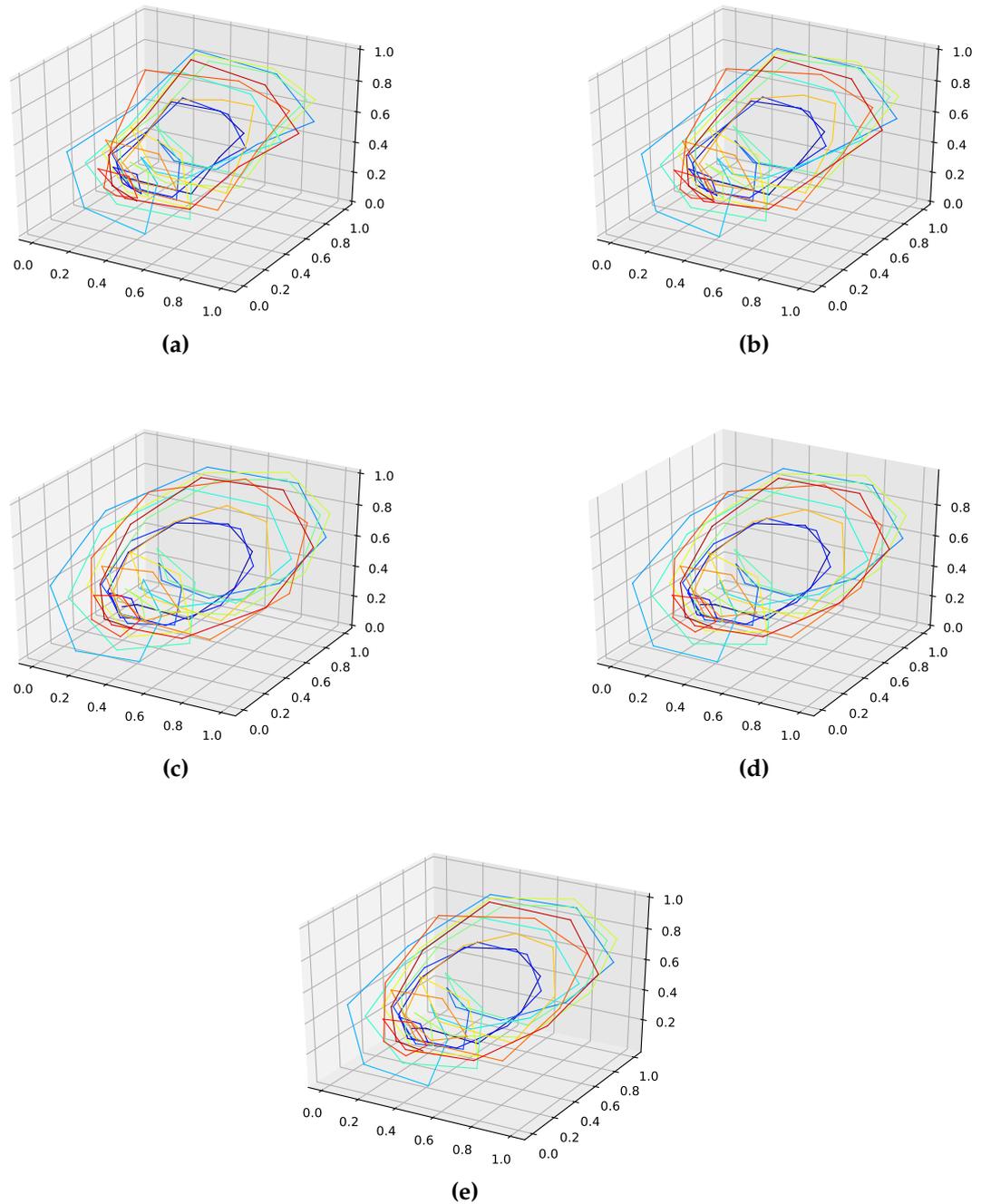


Figure 15. Reconstructed validation attractors (one interpolation point) for the car sales in Quebec data set.

- (a):** Reconstructed attractor, average population validation interpolation;
- (b):** Reconstructed attractor, linear interpolation;
- (c):** Reconstructed attractor, spline interpolation;
- (d):** Reconstructed attractor, best random validation interpolation;
- (e):** Reconstructed attractor, gen. alg. improved validation interpolation;

527 A.4. Perrin Freres Champagne Sales Data Set

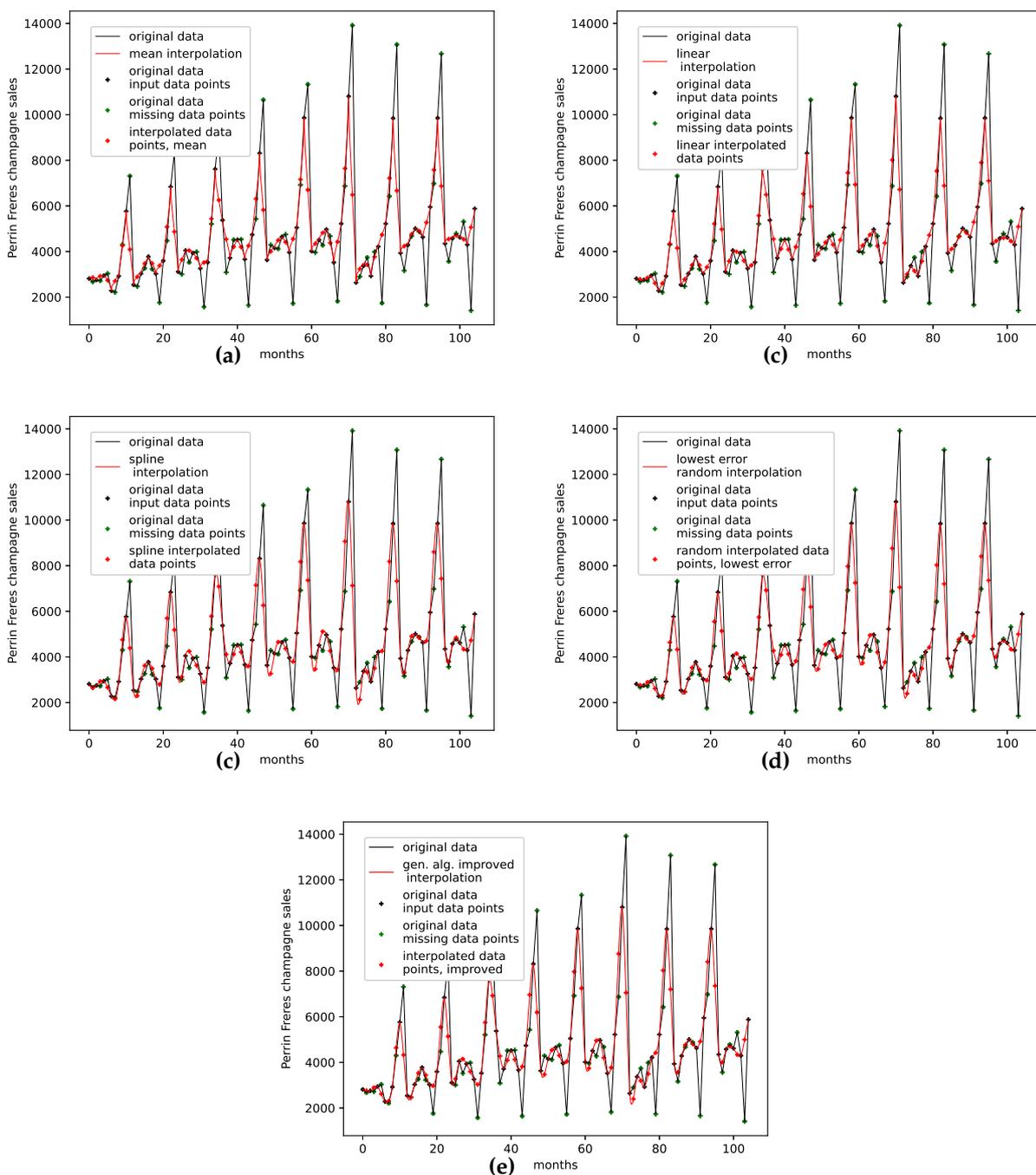


Figure 16. Interpolated validation data (seven interpolation points) for the Perrin Freres champagne sales data set.

- (a): Average population validation;
- (b): Validation, linear interpolation;
- (c): Validation, spline interpolation;
- (d): Validation, best random interpolation;
- (e): Validation, gen. alg. improved interpolation;

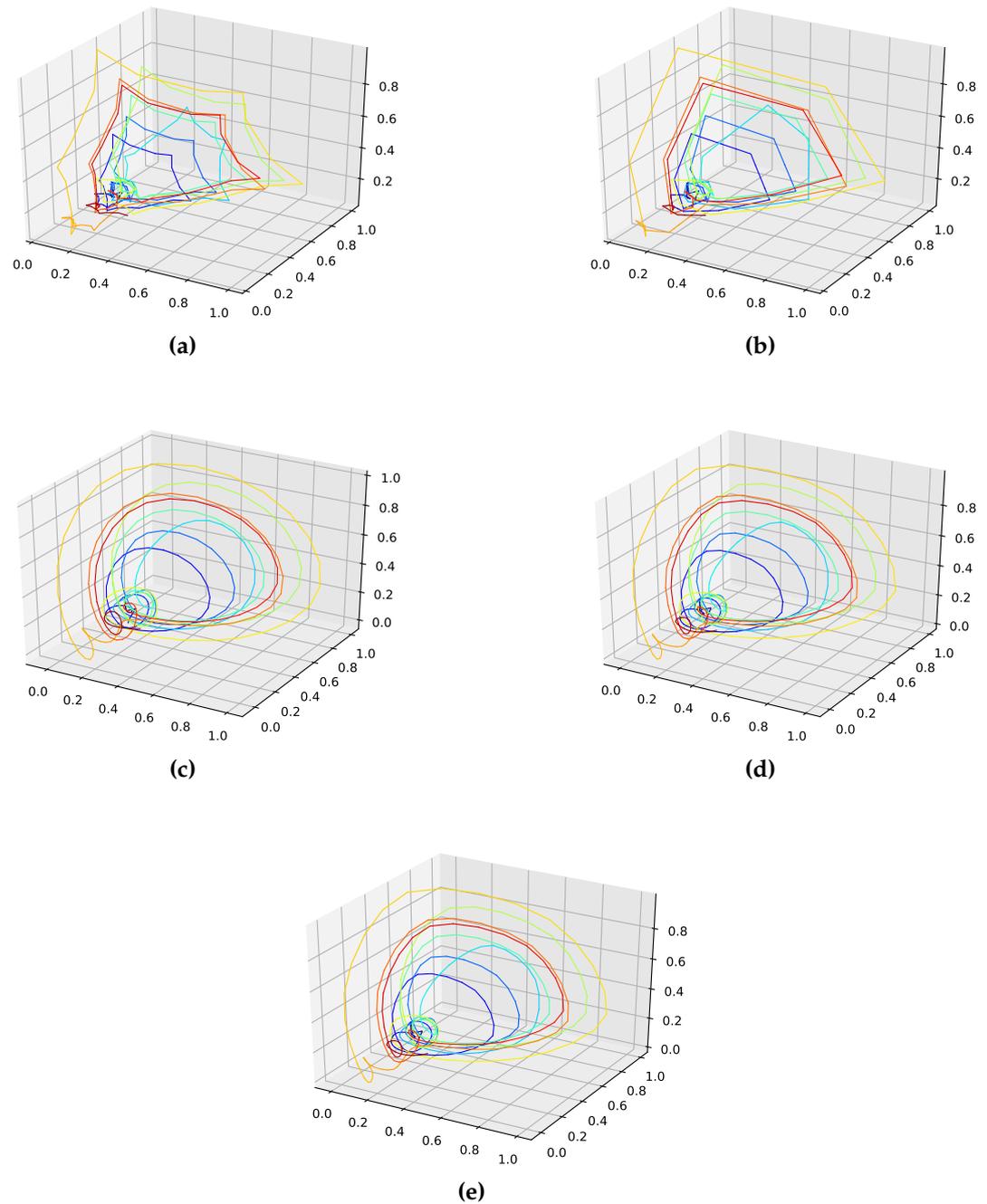


Figure 17. Reconstructed validation attractors (seven interpolation points) for the Perrin Freres champagne sales data set.

- (a):** Reconstructed attractor, average population validation interpolation;
- (b):** Reconstructed attractor, linear interpolation;
- (c):** Reconstructed attractor, spline interpolation;
- (d):** Reconstructed attractor, best random validation interpolation;
- (e):** Reconstructed attractor, gen. alg. improved validation interpolation;

528 A.5. Monthly Airline Passengers Data Set

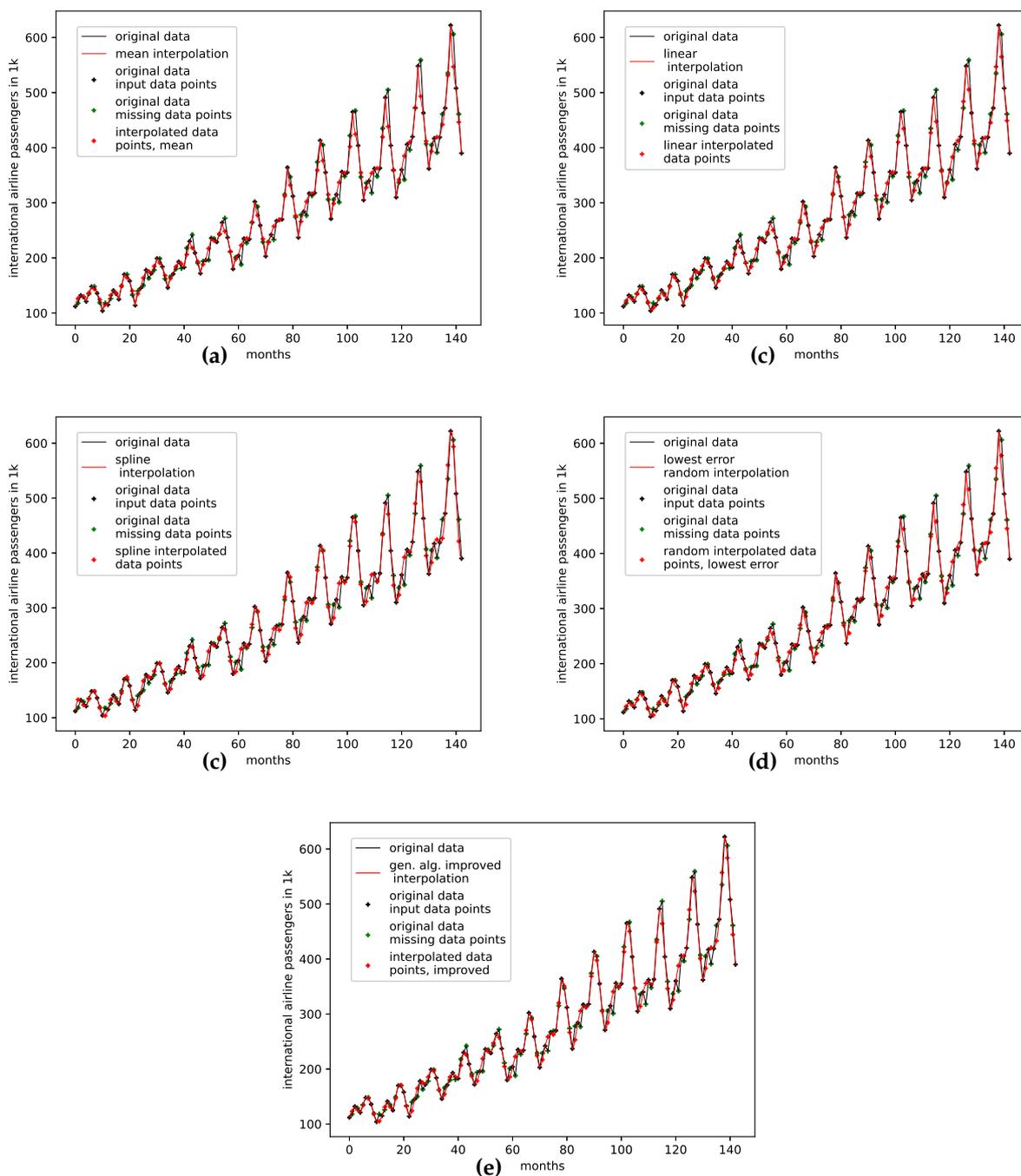


Figure 18. Interpolated validation data (three interpolation points) for the monthly international airline passengers data set.

- (a): Average population validation;
- (b): Validation, linear interpolation;
- (c): Validation, spline interpolation;
- (d): Validation, best random interpolation;
- (e): Validation, gen. alg. improved interpolation;

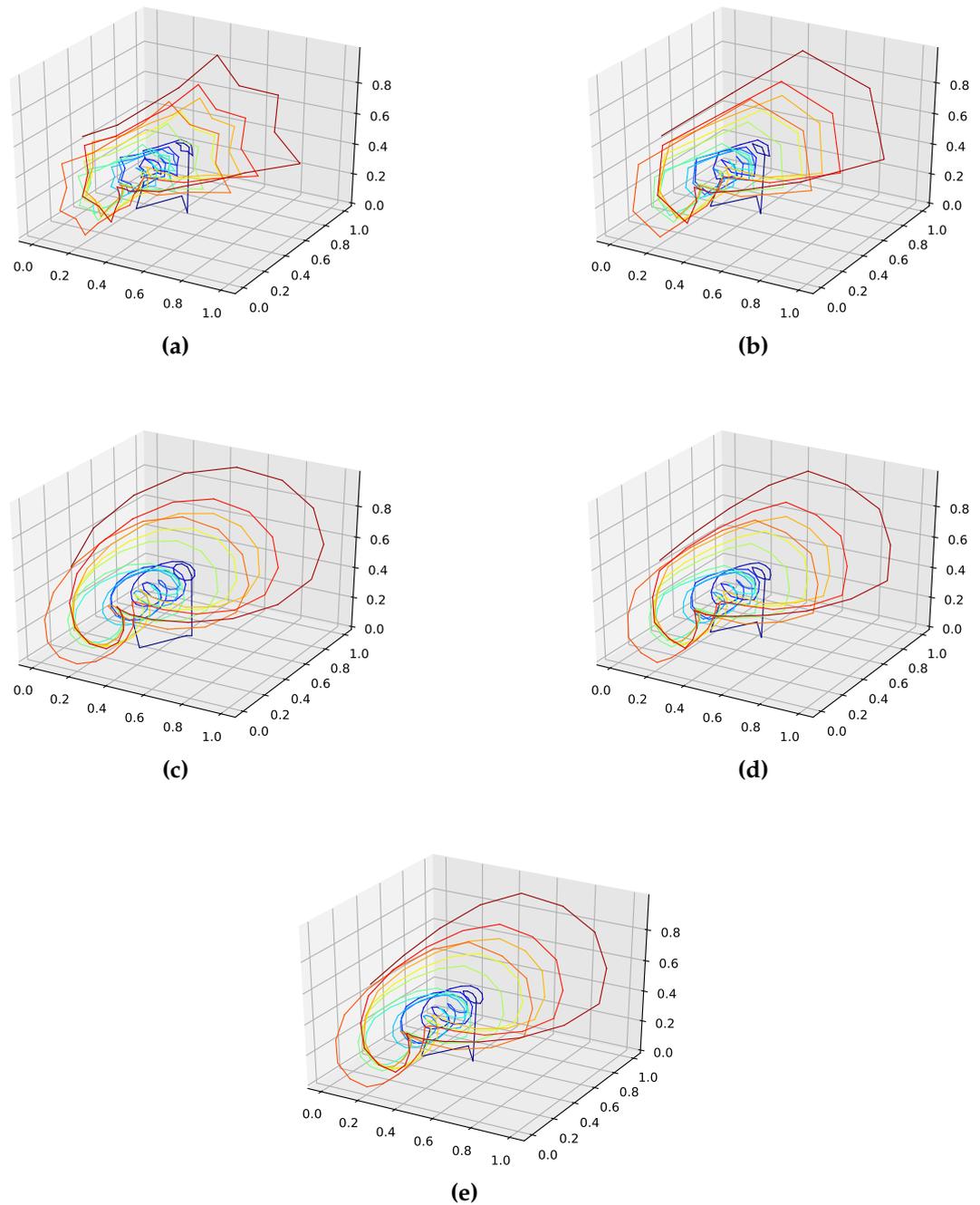


Figure 19. Reconstructed validation attractors (three interpolation points) for the monthly international airline passengers data set.

- (a):** Reconstructed attractor, average population validation interpolation;
- (b):** Reconstructed attractor, linear interpolation;
- (c):** Reconstructed attractor, spline interpolation;
- (d):** Reconstructed attractor, best random validation interpolation;
- (e):** Reconstructed attractor, gen. alg. improved validation interpolation;

529 A.6. Monthly Mean Temperature in Nottingham Castle Data Set

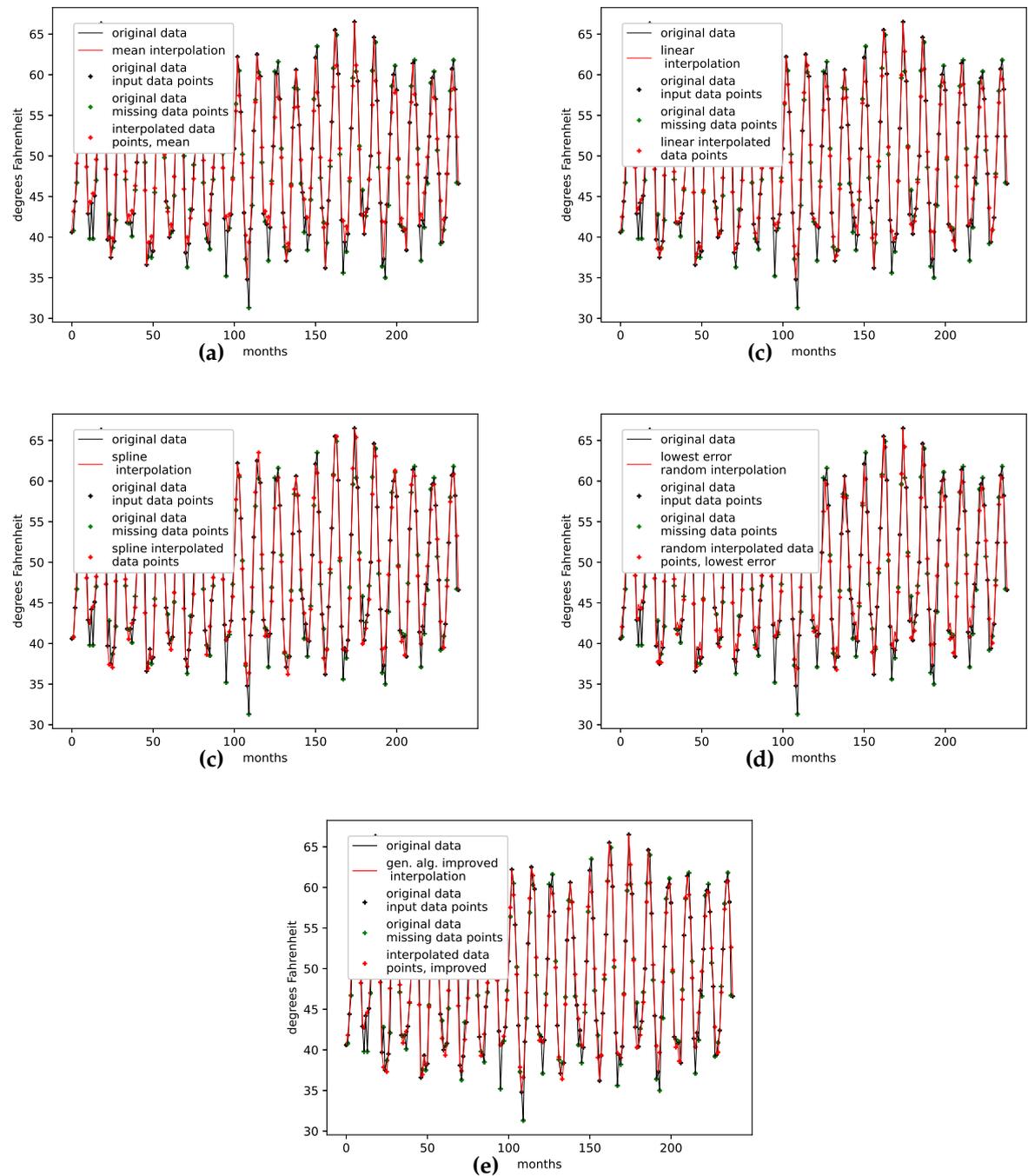


Figure 20. Interpolated validation data (one interpolation point) for the monthly mean temperature in Nottingham castle data set.

- (a): Average population validation;
- (b): Validation, linear interpolation;
- (c): Validation, spline interpolation;
- (d): Validation, best random interpolation;
- (e): Validation, gen. alg. improved interpolation;

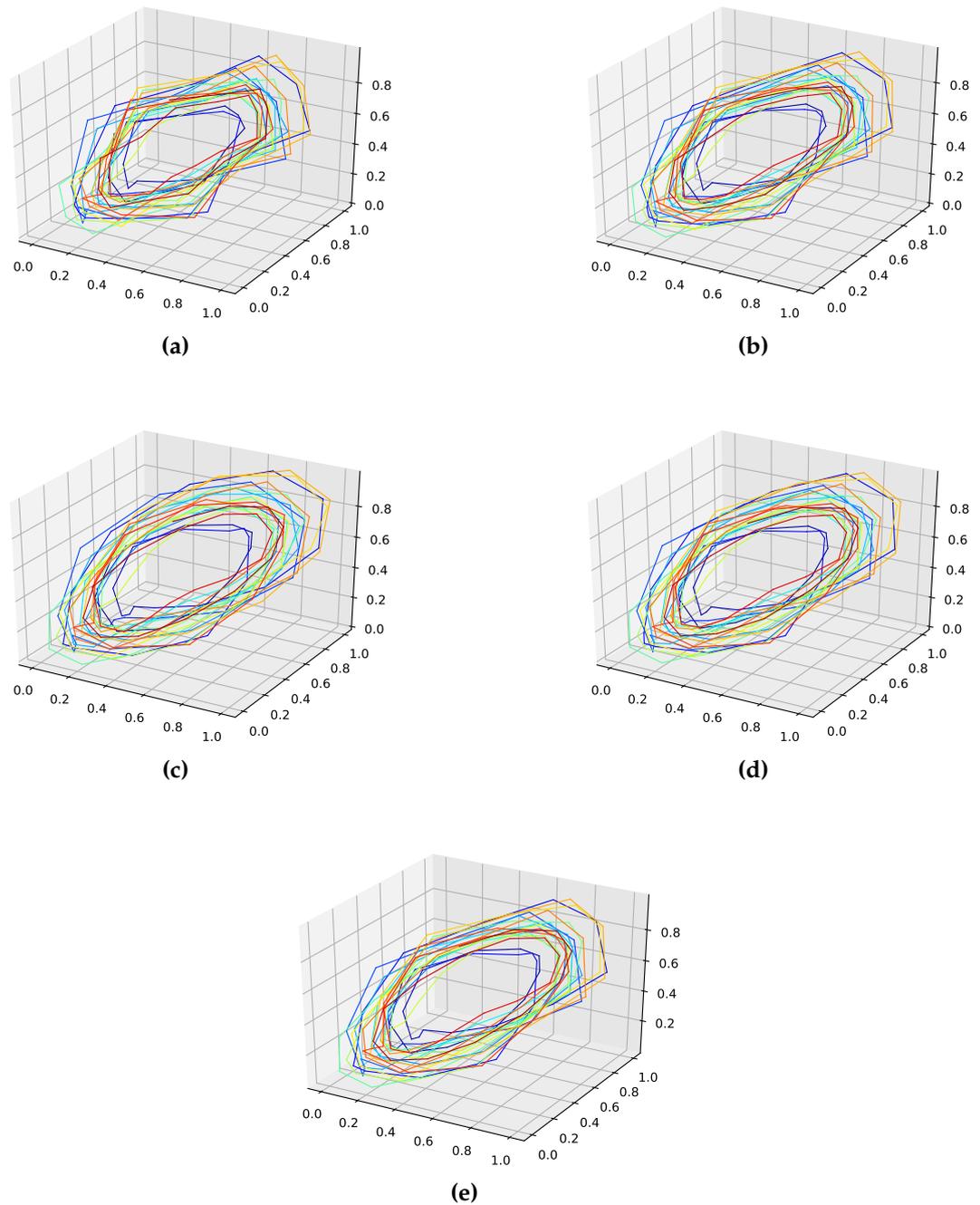


Figure 21. Reconstructed validation attractors (one interpolation point) for the monthly mean temperature in Nottingham castle data set.

- (a): Reconstructed attractor, average population validation interpolation;
- (b): Reconstructed attractor, linear interpolation;
- (c): Reconstructed attractor, spline interpolation;
- (d): Reconstructed attractor, best random validation interpolation;
- (e): Reconstructed attractor, gen. alg. improved validation interpolation;

530 B. Loss Surface

531 We present the loss surface for the Lorenz attractor in Figure 22 from two per-
 532 spectives. The orange dot marks the actual embedding of the Lorenz system. The plot
 533 suggests that the correct phase space embedding is located in an area where the loss

534 surface flattens out. At this point, we did not check for possible ways to locate the correct
 535 phase space embedding in the loss surface. Future approaches might find ways to do so.
 536

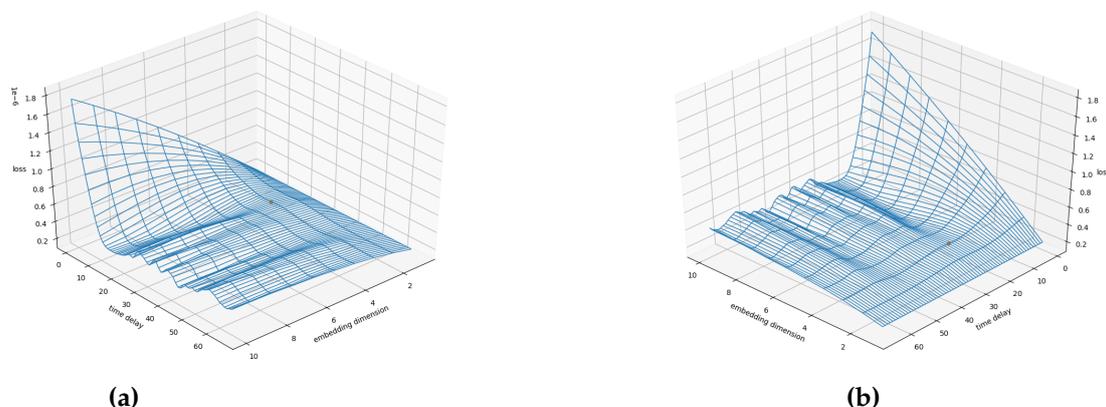


Figure 22. Loss surface for the Lorenz attractor.

(a) and (b) both show the same surface from different angles. This is the employed loss function (Section 3.3.1) depending on a varying embedding dimension and time delay. The orange dot marks the correct embedding dimension and time delay.

537 C. Failed Attempts

538 This section provides additional material for failed attempts to find a smooth
 539 phase space trajectory. For this reason, we provide additional plots (Figure 23) and the
 540 corresponding errors for the Lorenz system in Table 7. These attempts for different loss
 541 functions include:

- 542 • Minimizing the nearest neighbour distance between phase space points.
- 543 • Minimizing the mean of first-order derivatives along the phase space trajectory.
- 544 • Minimizing the variance of first-order derivatives along the phase space trajectory.
- 545 • Minimizing the mean of second order derivatives along the phase space trajectory.

Table 7: Errors for the interpolated data on the Lorenz system for 14 interpolation points and different loss functions. The Errors are shown for the mean interpolation of all populations, the lowest error in the population, and the interpolation that was improved using the presented genetic algorithm. Further, we give the percentage of how much of the population is outperformed by the genetic algorithm improved interpolation. Here, one can see that only methods including the second derivatives performed well. Further, the variance of second-order derivatives along the phase space trajectory performed best.

Loss Function →	Nearest Neighbour Distance	First Derivative Mean	First Derivative Variance	Second Derivative Mean	Second Derivative Variance
RMSE Population Mean	0.90686				
Lowest RMSE in population	0.18632				
RMSE gen. alg. improved	1.13779	0.67649	0.54291	0.19274	0.18626
Below Best %	73.9%	62.9%	55.5	4.2%	0.1%

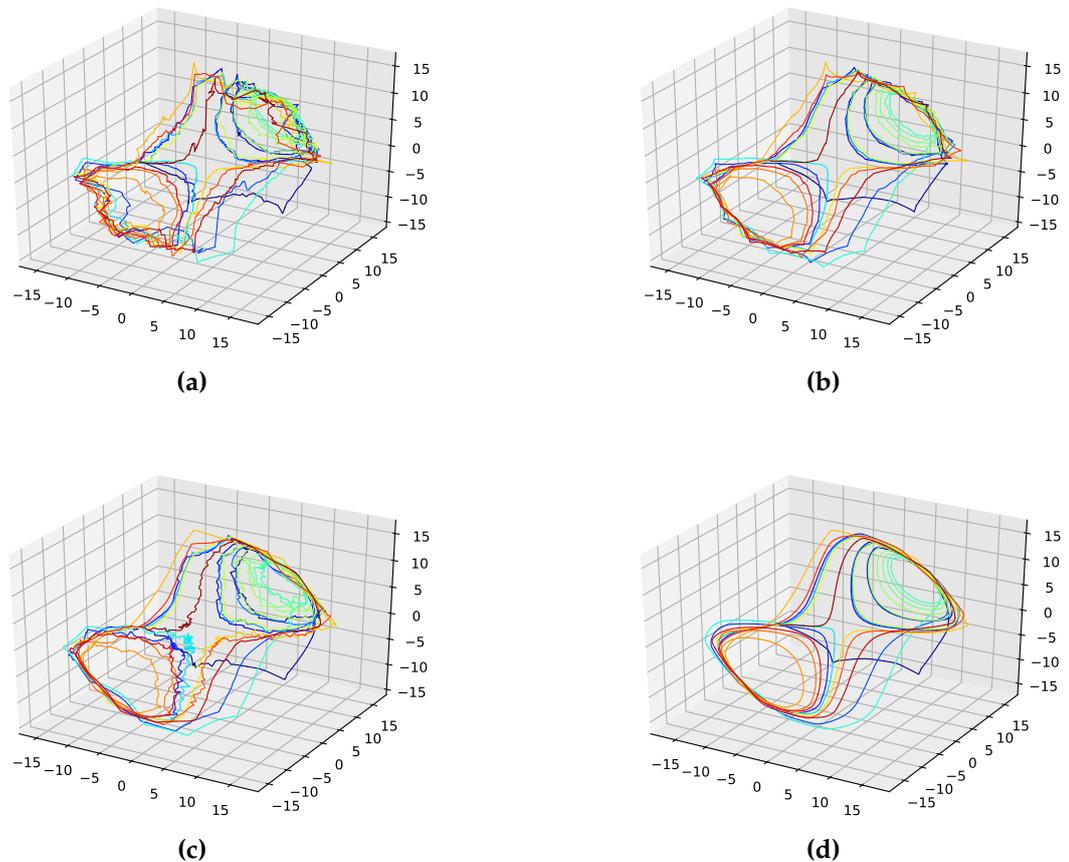


Figure 23. Reconstructed attractors for the interpolated Lorenz system for different loss functions.
(a): Nearest neighbour distance loss function;
(b): First derivative mean loss function;
(c): First derivative variance loss function;
(d): second derivative mean loss function;

References

1. Friedrich, J.; Gallon, S.; Pumir, A.; Grauer, R. Stochastic Interpolation of Sparsely Sampled Time Series via Multipoint Fractional Brownian Bridges. *Phys. Rev. Lett.* **2020**, *125*, 170602. doi:10.1103/PhysRevLett.125.170602.
2. Stark, J.; Broomhead, D.S.; Davies, M.E.; Huke, J. Takens embedding theorems for forced and stochastic systems. *Nonlinear Analysis: Theory, Methods & Applications* **1997**, *30*, 5303 – 5314. doi:https://doi.org/10.1016/S0362-546X(96)00149-6.
3. Raubitsek, S.; Neubauer, T. A fractal interpolation approach to improve neural network predictions for difficult time series data. *Expert Systems with Applications* **2021**, *169*, 114474. doi:https://doi.org/10.1016/j.eswa.2020.114474.
4. Raubitsek, S.; Neubauer, T. Taming the Chaos in Neural Network Time Series Predictions. *Entropy* **2021**, *23*. doi:10.3390/e23111424.
5. Pech-Pacheco, J.; Cristobal, G.; Chamorro-Martinez, J.; Fernandez-Valdivia, J. Diatom autofocusing in brightfield microscopy: a comparative study. Proceedings 15th International Conference on Pattern Recognition. ICPR-2000, 2000, Vol. 3, pp. 314–317 vol.3. doi:10.1109/ICPR.2000.903548.
6. Chang, C.; Lo, S.; Yu, S. Applying fuzzy theory and genetic algorithm to interpolate precipitation. *Journal of Hydrology* **2005**, *314*, 92–104. doi:10.1016/j.jhydrol.2005.03.034.
7. Sinhuber, M.; Friedrich, J.; Grauer, R.; Wilczek, M. Multi-level stochastic refinement for complex time series and fields: A data-driven approach. *New J. Phys.* **2021**.
8. Delorme, M.; Wiese, K.J. Extreme-value statistics of fractional Brownian motion bridges. *Phys. Rev. E* **2016**, *94*, 052105. doi:10.1103/PhysRevE.94.052105.
9. Sottinen, T.; Yazigi, A. Generalized Gaussian bridges. *Stochastic Processes and their Applications* **2014**, *124*, 3084–3105. doi:https://doi.org/10.1016/j.spa.2014.04.002.
10. Takens, F. Detecting strange attractors in turbulence. In *Dynamical Systems and Turbulence*, Warwick 1980, *Lecture Notes in Mathematics*; Rand, D.; Young, L.S., Eds.; Springer-Verlag: Berlin Heidelberg, 1981; Vol. 898, pp. 366–381.

11. Packard, N.H.; Crutchfield, J.P.; Farmer, J.D.; Shaw, R.S. Geometry from a Time Series. *Phys. Rev. Lett.* **1980**, *45*, 712–716. doi:10.1103/PhysRevLett.45.712.
12. Fraser, A.M.; Swinney, H.L. Independent coordinates for strange attractors from mutual information. *Phys. Rev. A* **1986**, *33*, 1134–1140. Publisher: American Physical Society, doi:10.1103/PhysRevA.33.1134.
13. Rhodes, C.; Morari, M. The false nearest neighbors algorithm: An overview. *Computers & Chemical Engineering* **1997**, *21*, S1149–S1154. Supplement to Computers and Chemical Engineering, doi:https://doi.org/10.1016/S0098-1354(97)87657-0.
14. Quarteroni, A.; Sacco, R.; Saleri, F. *Numerical Mathematics*; Vol. 37, 2007. doi:10.1007/b98885.
15. Lorenz, E.N. Deterministic nonperiodic flow. *Journal of atmospheric sciences* **1963**, *20*, 130–141.
16. Albrecht, P. The Runge–Kutta theory in a nutshell. *SIAM Journal on Numerical Analysis* **1996**, *33*, 1712–1735.
17. Hall, C.A.; Meyer, W. Optimal error bounds for cubic spline interpolation. *Journal of Approximation Theory* **1976**, *16*, 105–122. doi:https://doi.org/10.1016/0021-9045(76)90040-X.
18. Brunton, S.; Brunton, B.; Proctor, J.; Kaiser, E.; Kutz, J. Chaos as an Intermittently Forced Linear System. *Nature Communications* **2016**, *8*. doi:10.1038/s41467-017-00030-8.
19. LONDON, W.P.; YORKE, J.A. RECURRENT OUTBREAKS OF MEASLES, CHICKENPOX AND MUMPS: I. SEASONAL VARIATION IN CONTACT RATES¹. *American Journal of Epidemiology* **1973**, *98*, 453–468, [<https://academic.oup.com/aje/article-pdf/98/6/453/269755/98-6-453.pdf>]. doi:10.1093/oxfordjournals.aje.a121575.
20. Hyndman, R.; Yang, Y. Time Series Data Library v0.1.0. pkg.yangzhuoranyang.com/tsdl, 2018.
21. Friedrich, J.; Peinke, J.; Pumir, A.; Grauer, R. Explicit construction of joint multipoint statistics in complex systems. *J. Phys. Complexity* **2021**, *2*, 045006.