# WINNING ENTRY OF THE K. U. LEUVEN TIME-SERIES PREDICTION COMPETITION[*]

J. McNAMES[†]

*Information Systems Laboratory,*
*Stanford University, Stanford, CA 94305-9510, USA*

J. A. K. SUYKENS[‡] and J. VANDEWALLE[§]
*Katholieke Universiteit Leuven,*
*Department of Electrical Engineering,*
*ESAT–SISTA, Kardinaal Mercierlaan 94,*
*B-3001 Leuven (Heverlee), Belgium*

In this paper we describe the winning entry of the time-series prediction competition which was part of the International Workshop on *Advanced Black-Box Techniques for Nonlinear Modeling*, held at K. U. Leuven, Belgium on July 8–10, 1998. We also describe the source of the data set, a nonlinear transform of a 5-scroll generalized Chua's circuit. Participants were given 2000 data points and were asked to predict the next 200 points in the series.

The winning entry exploited symmetry that was discovered during exploratory data analysis and a method of local modeling designed specifically for the prediction of chaotic time-series. This method includes an exponentially weighted metric, a nearest trajectory algorithm, integrated local averaging, and a novel multistep ahead cross-validation estimation of model error for the purpose of parameter optimization.

## 1. Introduction

The prediction of chaotic time-series is a challenging problem for nonlinear modeling techniques. Although local instability on attractors prohibits accurate long-term predictions, short term predictions can be made with varying degrees of accuracy depending upon the technique that is used.

A time-series prediction competition sponsored by the Santa Fe institute was held in 1991. The purpose of the competition was twofold: first, to bring together researchers from diverse disciplines that study time-series (physics, biology, economics, etc.) and quantitatively compare different approaches within a common domain of problems. Second, the competition was held to incite collective discussion and research on common problems within the field of time-series analysis.

The K. U. Leuven competition was held to take a current assessment of where this discussion

has been led to and to also, give a fresh quantitative evaluation of leading time-series prediction techniques.

The competition data set consisted of 2000 points[1] generated by a nonlinear transform of a 5-scroll generalized Chua's circuit. The data set is shown in Fig. 2. The data set was made available in November 1997 and the deadline for entries was in April 1998. No information was given about the source of the data. Participants were asked to predict the next 200 points in the series. Seventeen entries were submitted.

The winning entry in the competition was generated by a local modeling method that is described in this paper. The method was inspired by the success of the second place entry in the Santa Fe competition [Sauer, 1994] and the success of local models as reported by other researchers [Casdagli *et al.*, 1992; Farmer & Sidorowich, 1988; Kugiumtzis *et al.*, 1998]. Other entries in the K. U. Leuven competition that incorporated local models also performed very well [Suykens & Vandewalle, 1998a, 1998b].

This paper is organized as follows. Section 2 describes the origin of the competition data set. Section 3 describes a general approach to time-series prediction including assumptions about the source of the time-series, issues that arise when applying a local model to this type of problem, and some innovations that can improve the accuracy of local models for time-series prediction. Section 4 describes the exploratory data analysis, user specified decisions and the parameter optimization that generated the winning entry.

## 2. Origin of Data

Consider the following generalized Chua's circuit

$$\dot{x}_1 = \alpha[x_2 - h(x_1)]$$
$$\dot{x}_2 = x_1 - x_2 + x_3 \qquad (1)$$
$$\dot{x}_3 = -\beta x_2$$

with piecewise-linear characteristic

$$h(x_1) = m_{2q-1}x_1$$
$$+ \frac{1}{2}\sum_{i=1}^{2q-1}(m_{i-1}-m_i)(|x_1+c_i|-|x_1-c_i|) \qquad (2)$$

where $q$ denotes a natural number. A complete family of $n$-scroll attractors, allowing an even and odd number of scrolls, can be obtained from this circuit (see [Suykens *et al.*, 1997]). For the parameters $\alpha = 9$, $\beta = 14.286$ and the vectors $m = [m_0\,m_1\cdots m_{2q-1}]$, $c = [c_1\,c_2\cdots c_{2q-1}]$ one obtains

- $q = 1$: double scroll (2-scroll) [Chua *et al.*, 1986; Chua, 1994; Madan, 1993]

$$m = [-1/7 + 2/7]$$
$$c = 1 \qquad (3)$$

- $q = 3$: 5-scroll [Suykens *et al.*, 1997]

$$m = [0.9/7\ -3/7\ 3.5/7\ -2.7/7\ 4/7\ -2.4/7]$$
$$c = [1\ \ 2.15\ \ 3.6\ \ 6.2\ \ 9]. \qquad (4)$$

In order to generate the competition data we have selected the 5-scroll attractor (see Fig. 1).

The generalized Chua's circuit was simulated with the initial state vector $[0.1\ \ -0.2\ \ 0.3]^T$ using a Runge–Kutta integration rule (ode23 in Matlab). The time-series was generated by taking a nonlinear combination of the three state variables,

$$y = W\tanh(Vx), \qquad (5)$$

where $x$ is the three-dimensional state vector and the nonlinearity is a multilayer perceptron with
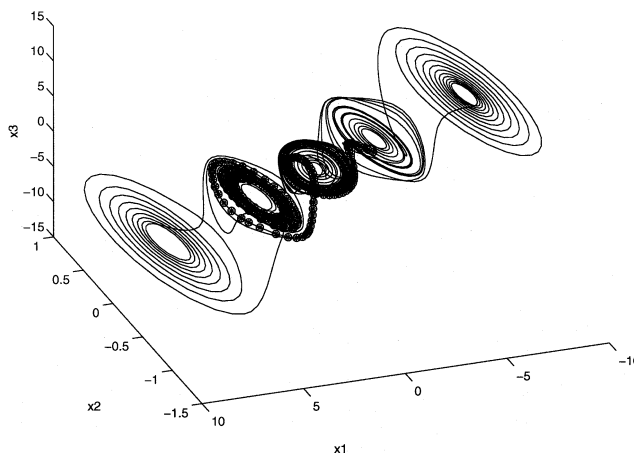


Fig. 1. The 5-scroll attractor [Suykens *et al.*, 1997] from which the competition data have been generated. The marked part of the trajectory corresponds to the data which have to be predicted.

[1]The data set was available at http://www.esat.kuleuven.ac.be/sista/workshop/
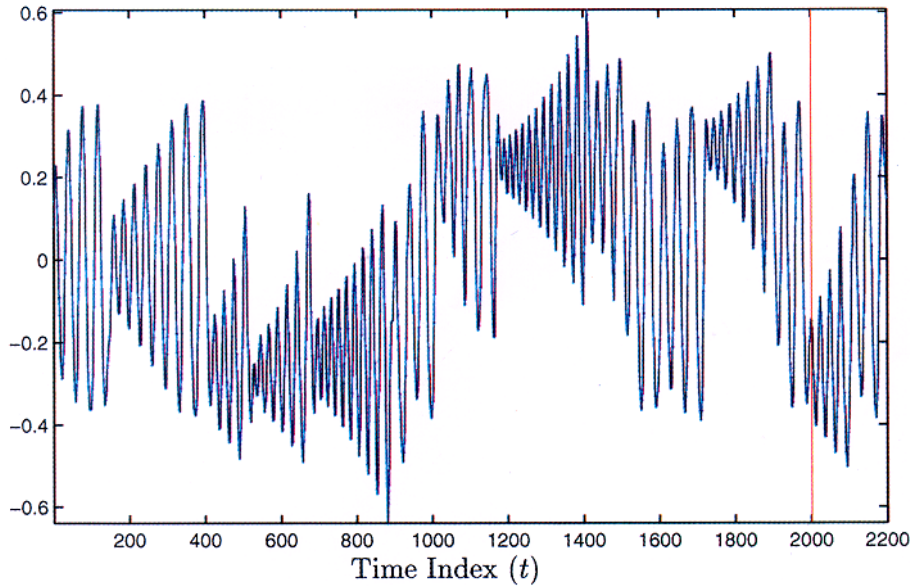
Fig. 2. Competition data. The first 2000 points (data before the red vertical line) were given to participants. The next 200 points (after the red vertical line) were to be predicted and judged according to the mean squared error (MSE).

three hidden units, interconnection matrices

$$W = [-0.0124 \quad 0.3267 \quad 1.2288] \,,$$

$$V = \begin{bmatrix} -0.1004 & -0.1102 & -0.2784 \\ 0.0009 & 0.5792 & 0.6892 \\ 0.1063 & -0.0042 & 0.0943 \end{bmatrix}$$

and a zero bias vector. This multilayer perceptron hides the underlying structure of the attractor.

The resulting time-series $y(t)$ is shown in Fig. 2. 2000 data points were given (data before the vertical line) and the aim was to predict the next 200 data points (between the two vertical lines) so as to minimize the mean squared error (MSE),

$$\text{MSE} \triangleq \frac{1}{200} \sum_{i=2001}^{2200} (y_i - \hat{y}_i)^2 \tag{6}$$

In Fig. 1 we have marked with a star the data points that were to be predicted. One can observe that the 2000 data points cover the whole attractor, i.e. the 5 scrolls. Within the 200 data points to be predicted the time-series makes a transition between scrolls that is difficult to predict.

A physical implementation of a related family of $n$-double scroll attractors has been given in [Arena *et al.*, 1996a, 1996b]. Further extensions to hyperchaotic $n$-double scroll hypercubes in cellular neural networks (one-dimensional array) have been

proposed in [Suykens & Chua, 1997]. The $n$-scroll attractors have been used for secure communication applications exploiting chaos, such as the method of nonlinear $H_\infty$ synchronization of chaotic Lur'e systems in [Suykens *et al.*, 1997; Suykens *et al.*, 1997]. Steps towards unmasking chaotic communication schemes have been taken in [Short, 1994, 1996; Yang, 1995], which are related to the nonlinear system identification. The latter is a motivation for understanding the limits of performance of time-series prediction methods.

## 3. Local Modeling

Local models generate predictions by finding segments of the time-series that closely resemble the segment of the points immediately preceding the point to be predicted. The prediction is an estimate of the average change that occurred immediately after these similar segments of points.

Previous studies have shown that forecasting methods based on local models produce predictions that are better than or comparable to competing models and they have a number of favorable properties not shared by other methods [Casdagli *et al.*, 1992; Farmer & Sidorowich, 1987, 1988; Liu *et al.*, 1997; Sauer, 1994].

To use local models for time-series prediction there are many decisions that one must make. For example, how should *local* be defined

mathematically, how should the input vector be constructed, what should the model be designed to predict, what type of local model should be used, and how values for model parameters should be chosen. This section describes how these decisions were made for the method of local modeling that generated the winning entry.

### 3.1.  *Assumptions*

The local modeling method described here was designed to predict time-series generated by nonlinear dynamic systems that can be described by a set of nonlinear ordinary differential equations:

$$\dot{w}_t = f_c(w_t) \tag{7}$$

$$y_t = g_c(w_t) \tag{8}$$

where $w_t \in \Re^{n_s}$ is the state of the system, $n_s$ is the order of the system, and $y_t \in \Re^1$ is the time-series to be predicted. An equivalent discrete-time system also exists,

$$w_t = f_d(w_{t-1}) \tag{9}$$

$$y_t = g_d(w_t) \tag{10}$$

where we have assumed that the sampling period is one for ease of presentation. This system can be thought of as a nonlinear oscillator as illustrated in Fig. 3.

We assume that $f_d(w_t)$ and $g_d(w_t)$ are smooth in the sense that they have continuous and bounded derivatives. We also assume that the system is time invariant, that $f_d(w_t)$ has a finite number of equilibria, and that the order of the system is small ($n_s < 10$).

Many equivalent mathematical representations exist for any given set of ordinary differential equations. For the type of nonlinear dynamic system discussed here Takens has shown that an equivalent representation exists where the equivalent state is defined as a finite window of past values of the time-series [Takens, 1981]. This is illustrated in Fig. 4. This window is called a time delay embedding and is defined as

$$x_t \triangleq [y_t,\, y_{t-\delta},\, y_{t-2\delta},\, \dots,\, y_{t-(m-1)\delta}] \tag{11}$$

where $m$ is the embedding dimension and $\delta$ is the embedding delay, both of which are natural numbers. Takens' work was later generalized and shown
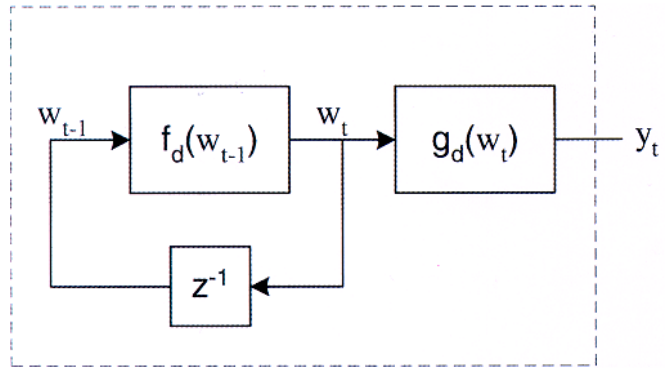


Fig. 3.   Nonlinear oscillator. Chaotic time-series generated by nonlinear dynamic systems can be thought of as the output of a nonlinear oscillator as illustrated in this diagram.
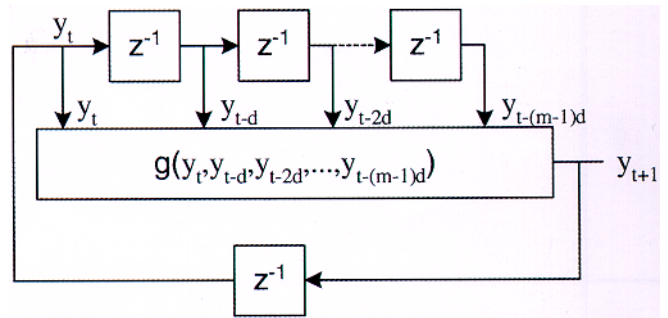


Fig. 4.   Equivalent oscillator. Takens has shown that the nonlinear oscillator shown in Fig. 3 has an equivalent representation shown here.

to apply to a broader class of systems by Sauer *et al.* [1991].

This theorem is important because it implies that an exact one-step ahead prediction is possible given a finite window of previous values:

$$y_{t+1} = g(x_t) \tag{12}$$

$$= g(y_t,\, y_{t-\delta},\, y_{t-2\delta},\, \dots,\, y_{t-(m-1)\delta}) \tag{13}$$

for some unknown function $g(x_t)$. This effectively reduces the one-step ahead prediction problem to a standard nonlinear modeling problem.

Multistep ahead predictions can be generated by making further iterations. For example, if the embedding delay is 1 then the following equations illustrate how iterative prediction can be used to predict more than one-step ahead

$$\begin{aligned}
\hat{y}_{t+1} &= \hat{g}(y_t,\, y_{t-1},\, y_{t-2},\, \dots,\, y_{t+1-m}) \\
\hat{y}_{t+2} &= \hat{g}(\hat{y}_{t+1},\, y_t,\, y_{t-1},\, \dots,\, y_{t+2-m}) \\
\hat{y}_{t+3} &= \hat{g}(\hat{y}_{t+2},\, \hat{y}_{t+1},\, y_t,\, \dots,\, y_{t+3-m})
\end{aligned} \tag{14}$$

$$\vdots$$

The generalization to other values of $\delta$ is straightforward.

Any of the standard nonlinear models, such as neural networks, radial basis functions, local models, or support vector machines could be used to create an approximation for $g(x_t)$. However, a more accurate model can be built if the aforementioned assumptions are incorporated into the model structure. Several techniques that tailor local models specifically for time-series prediction are described in the following section.

## 3.2. *Introduction to Local Modeling*

Generating a prediction with a local model can generally be divided into two steps: first, the $k$ nearest neighbors in the data set of a given input query vector $x_\tau$ [Eq. (11)] are found and, second, a simple model is constructed using only the $k$ neighboring points to generate the prediction $\hat{y}_{\tau+1}$.

### 3.2.1. *Local Model Types*

The second step, building a model given $k$ data points, is an unusual function approximation problem because the dimension of the input vector ($m$) is often large compared to the number of data points ($k$) that are used to construct the model. Most nonlinear modeling methods are not viable under these conditions because they require that the number of data points be larger than the dimension of the input vector. Models that are computationally expensive to construct are also not viable since a different model must be constructed for every query. Instead a relatively simple model must be used that can be cheaply constructed and evaluated with just a few data points.

For this reason, the two most common types of local models are local averaging models and local linear models. Local linear models generally produce more accurate predictions but are very sensitive to the method of regularization [Kugiumtzis *et al.*, 1998] and generally require more neighbors than local averaging models. Local averaging models can be used with smaller neighborhoods, are more stable, and are often more accurate for very short data sets. Since only 2000 points were available in the competition data set and because of its simplicity a local averaging model was used to generate the winning entry.

### 3.2.2. *Discontinuities*

One commonly mentioned disadvantage of local models is that they generate discontinuous predictions [Casdagli *et al.*, 1992; Lillekjendlie & Christophersen, 1994; Liu *et al.*, 1997] For example, consider a local averaging model:

$$\hat{y}_\tau = \frac{1}{k} \sum_{i=1}^{k} y_{c_i} \qquad (15)$$

where $c_i$ is the data set index of the $i$th nearest neighbor. In some regions of the input space where two data points are nearly equidistant to the input vector, an arbitrarily small perturbation can result in a different $k$th nearest neighbor. This can cause a discontinuous change in the summation above, which in turn causes a discontinuous change in $\hat{y}_\tau$.

This problem can be remedied by decreasing the influence of neighbors that are relatively far away. For example, consider the weighted average

$$\hat{y}_\tau = \frac{\displaystyle\sum_{i=1}^{k} y_{c_i} w_i}{\displaystyle\sum_{i=1}^{k} w_i} \ . \qquad (16)$$

where $w_i$ are non-negative weights that are a monotonically decreasing function of the distance to the $i$th nearest neighbor.

It is generally accepted that model accuracy is insensitive to the type of weighting function used [Atkeson *et al.*, 1997]. A good choice is the biweight function,

$$w_i = \left(1 - \frac{d_i^2}{d_{k+1}^2}\right)^2 \qquad (17)$$

where $d_i$ is the distance to the $i$th nearest neighbor. Since the biweight function is smooth in the sense that it has a continuous first derivative, the weighted average produces an estimate that is also a smooth function of the input vector.

### 3.2.3. *Fast nearest neighbor algorithms*

Another common criticism of local models is that they are computationally expensive. For most local modeling techniques the majority of the computation is used finding the $k$ nearest neighbors in the data set. For time-series prediction problems the data sets can be very large and the computational cost of the naïve brute force approach can be prohibitive.

To solve this problem a fast nearest neighbor algorithm can be used. Many of these algorithms fall

into two categories: axis-partitioning algorithms and triangle inequality-based algorithms. Axis-partitioning algorithms divide the $m$-dimensional input space into hypercubes and establish a lower bound on the distance from the query point to all points contained in each hypercube. If the lower bound is greater than the distance to the $k$th nearest neighbor found so far, all of the points contained by the hypercube can be eliminated without explicitly calculating the distance to each point [Friedman *et al.*, 1977; Kim & Park, 1986; Nene & Nayar, 1997; Niemann & Goppert, 1988; Tai *et al.*, 1996].

Triangle inequality-based algorithms use the triangle inequality to find lower bounds on the distance to points. If the lower bounds is greater than the distance to the $k$th nearest neighbor found so far, the point can be eliminated from consideration without explicitly calculating the distance to that point. Subsets of points can also be eliminated by finding a lower bound on the distance to an enclosing hypersphere [Fukunaga & Narendra, 1975; Kalantari & McDonald, 1983; Kamgar-Parsi & Kanal, 1985; McNames, 1998; Micó *et al.*, 1996; Vidal, 1986].

If the data set variables are independently distributed and there are relatively few of them (i.e. a low-dimensional space) then both types of nearest neighbor algorithms perform drastically better than the brute force approach. In low dimensions most of these algorithms achieve $O(\log n_c)$ search time and require only $O(n_c \log n_c)$ preprocessing time and storage, where $n_c$ is the number of points in the data set.

In high-dimensional spaces, say $m > 15$, the bounding techniques of both types of algorithms are ineffectual and the performance is much worse. In fact, the computational cost of these algorithms can be significantly higher than the brute force approach due to overhead imposed for ordering the search of the input space and for calculating the lower bounds. Generally, the search time increases exponentially with $m$ up to a limit where the distance is calculated for nearly all of the points in the data set.

If all the points in the data set lie in a low-dimensional subspace of $\Re^m$, the triangle inequality algorithms have an advantage over axis-partitioning algorithms — they can achieve search times comparable to low-dimensional data sets. If the dimension of the subspace is held constant, the search time grows roughly linearly with $m$ rather than exponentially. This is especially relevant when the

data points are taken from a time-series generated by a nonlinear dynamic system because the points will often lie on a manifold of low dimension that is independent of $m$ [Gershenfeld & Weigend, 1994; Sauer *et al.*, 1991; Takens, 1981].

Figure 5 illustrates the advantage of using a fast triangle inequality-based algorithm when the points lie on a low-dimensional manifold, which is always the case for time-series generated by the type of nonlinear dynamic systems discussed here (Sec. 3.1). This figure shows the average query time for four different nearest neighbor algorithms applied to the competition data set. The trace labeled *Brute* represents the brute force algorithm. The trace labeled *KD Tree* represents an axis-partitioning algorithm proposed by Friedman *et al.* [1977]. The traces labeled *FN* and *FNM* represent triangle inequality algorithms proposed by Fukunaga and Narendra [1975] and McNames [1998] respectively.

### 3.3. *Nearest neighbor metrics*

Local models are constructed using only the nearest points on the $k$ nearest neighboring trajectory segments, as described in the next section. Choosing an appropriate measure of nearness, or metric, is an important decision that can strongly affect accuracy [García *et al.*, 1996; Murray, 1993; Tanaka *et al.*, 1995]. The most common metric is the square of the Euclidean distance between the query vector $x_\tau$ and a vector taken from the data set, $x_t$ [Eq. (11)]

$$\mathrm{D}_E(x_\tau, x_t) \triangleq \sum_{i=0}^{m-1} (y_{\tau-i\delta} - y_{t-i\delta})^2 \qquad (18)$$

In this case the only parameters are the embedding dimension ($m$) and the embedding delay ($\delta$). If the time-series is generated from a continuous-time dynamic system then the Euclidean distance is an estimate of the integrated squared error,

$$\mathrm{ISE}(x_\tau, x_t) \triangleq \frac{1}{\delta T} \int_0^{m\delta T} (y_{\tau T-\upsilon} - y_{tT-\upsilon})^2 d\upsilon \quad (19)$$

$$\approx \sum_{i=0}^{m-1} (y_{\tau-i\delta} - y_{t-i\delta})^2 \qquad (20)$$

$$= \mathrm{D}_E(x_\tau, x_t) \qquad (21)$$

where $T$ is the sampling period. If the window length, $m\delta T$, is fixed then smaller values of $\delta$ produce more accurate estimates of the ISE. However,
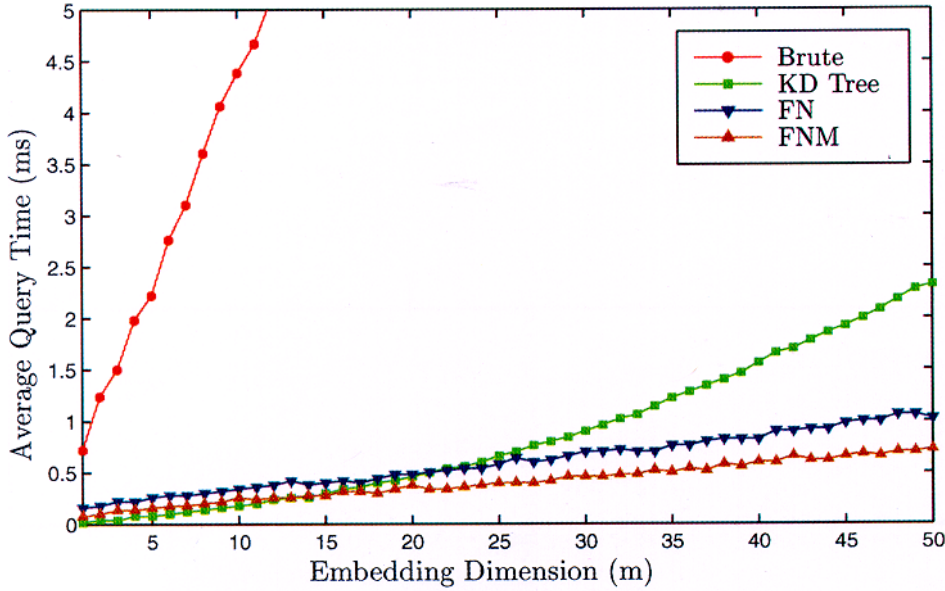
Fig. 5. Average query time to find 5 nearest neighbors in the Leuven time-series with 1900 points. The query vectors were taken from the time-series continuation that was released after the competition deadline. *Brute* is the brute force algorithm. *KD Tree* is an axis partitioning algorithm. *FN* is a triangle inequality algorithm proposed by Fukunaga and Narendra, [1975] and *FNM* is a modified version of the same algorithm.

smaller values of $\delta$ also require more computation so $\delta$ is treated here as a user-selected parameter that should be chosen as small as possible within the limits of available computational resources.

Choosing a more general measure of nearness with more parameters can greatly increase model accuracy. For example, a diagonally weighted Euclidean distance,

$$D_{WE}(x_\tau, x_t) \triangleq \sum_{i=0}^{m-1} \lambda_i (y_{\tau-i\delta} - y_{t-i\delta})^2 \qquad (22)$$

could be used where $\lambda_i \geq 0 \, \forall i$. However, short time-series are often too short to accurately estimate the best $m$ parameters ($\lambda_i$) and optimizing over so many parameters is computationally impractical for longer time-series.

Murray has proposed using an exponentially weighted diagonal metric similar to the one described by the following equation,

$$D_\lambda(x_\tau, x_t) = \sum_{i=0}^{m-1} \lambda^{i-1} (y_{\tau-i\delta} - y_{t-i\delta})^2 \qquad (23)$$

where $0 < \lambda \leq 1$ [Murray, 1993]. Although this metric is not optimal [Tanaka *et al.*, 1995] it is intuitively appealing because the variables closest in time to the prediction are given exponentially more

weight. It is especially appropriate for chaotic systems when $m$ is large because neighboring states are known to diverge exponentially with time.

The metric $D_\lambda$ is defined by the parameters $m$, $\delta$ and $\lambda$. Assuming $\delta$ is fixed, if $m$ is sufficiently large and $\lambda$ is sufficiently small then the influence of the later elements of the sum multiplied by $\lambda^{i-1}$ (for large $i$) will have a negligible influence. Thus, if a model is constructed with $m$ sufficiently large, only $\lambda$ will affect model accuracy. This can be implemented practically by defining a cutoff value, $\lambda_{min} \triangleq \lambda^{m-1}$, for which larger values of $i$ in Eq. (23) will have a negligible effect on $D_\lambda$:

$$\lambda = \lambda_{min}^{\frac{1}{(m-1)}} \qquad (24)$$

$$D_\lambda(x_\tau, x_t) \triangleq \sum_{i=1}^{m} \lambda^{i-1} (y_{\tau-i\delta} - y_{t-i\delta})^2 \qquad (25)$$

$$\approx \sum_{i=1}^{\infty} \lambda^{i-1} (y_{\tau-i\delta} - y_{t-i\delta})^2 . \qquad (26)$$

This modified metric has one free parameter ($m$) and two user-selected parameters ($\delta$ and $\lambda_{min}$). Figure 6 shows the weight $\lambda^{i-1}$ of each element in the sum for different values of $m$. Since the model accuracy is sensitive to the value of $m$ it is usually
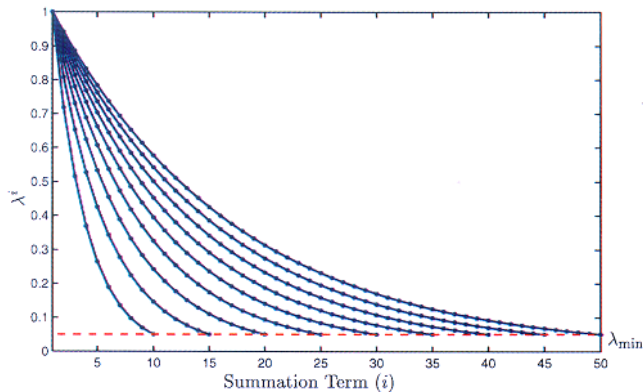
Fig. 6. The weight, $\lambda^{i-1}$, of each element in the metric summation given in Eq. (23) for various values of the embedding dimension ($m$).

worth the computational cost to find the value that maximizes model accuracy.

## 3.4. *Nearest trajectories*

One of the differences between a regression problem and a time-series prediction problem is that it is possible to increase the number of data points by upsampling, or interpolating, the time-series; more data can be generated without adding any new information. This is problematic because a signal sampled at a high enough rate will have all $k$ of its nearest neighbors adjacent to each other in the time-series. For example, in Fig. 7 many of the points on the nearest trajectory segment are closer than the points on the third nearest trajectory segment.

This problem can be solved by finding the *nearest trajectory segments* instead of the nearest neighbors [Farmer & Sidorowich, 1988]. Fortunately, it is possible to modify existing nearest neighbor algorithms to find the nearest trajectory segments. This is advantageous because, as described previously, much research has gone into the development of efficient nearest neighbor algorithms.

The nearest trajectory modification described by McNames [1998] is repeated here. Suppose a nearest neighbor algorithm has found a point $x_i$ that is closer than the $k$ nearest points found by the algorithm so far.

1. Calculate the distance to the points preceding $x_i$, $\{x_{i-1}, x_{i-2}, \ldots\}$, until the nearest local minimum is found. Repeat this procedure for the points succeeding $x_i$. The local minimum found by this procedure, $x_{\min}$, is the closest point in the trajectory segment.
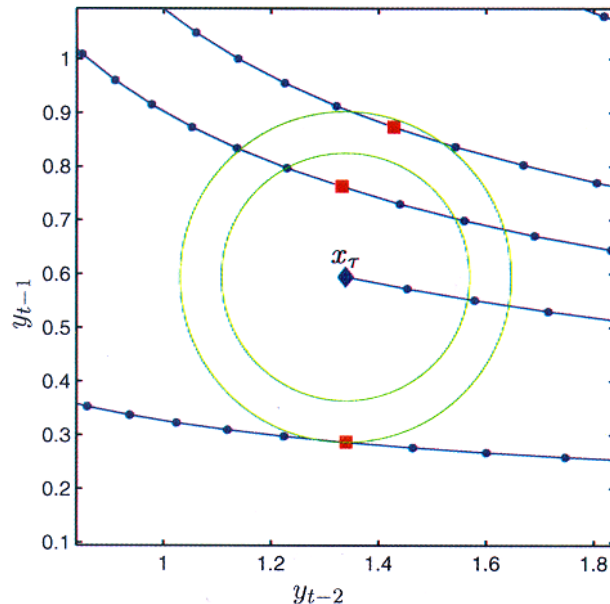


Fig. 7. Trajectory segments from the Lorenz time-series in a two-dimensional embedding space. The query point, $x_\tau$, is shown by a diamond. The three nearest points on each of the three nearest trajectory segments are shown by red squares. The inner green circle shows the distance to the third nearest neighbor and the outer green circle shows the distance to the nearest point on the third nearest trajectory segment. Note that the three nearest neighbors all lie on the nearest trajectory segment.

2. Calculate the distance to the points preceding $x_{\min}$ until either a maximum is found or the distance becomes greater than the distance to the $k$th nearest neighbor found so far. Call this point $x_{\max}$. Eliminate all points between $x_{\min}$ and $x_{\max}$ from consideration by the nearest neighbor algorithm.
3. Repeat the previous step for the points succeeding $x_{\min}$.
4. Replace the $k$th nearest neighbor found so far with $x_{\min}$ and continue with the nearest neighbor algorithm.

## 3.5. *Local averaging models*

Although local averaging models are reasonably accurate at interpolation they are generally poor at extrapolation. Since the local model described previously is a weighted average, the output will never be greater than or less than any of the points that make up the average. In high-dimensional spaces this can be a significant disadvantage since nearly every point is an outlier [Friedman, 1991]. For example, Fig. 8 shows a prediction generated by averaging the next point of the three nearest
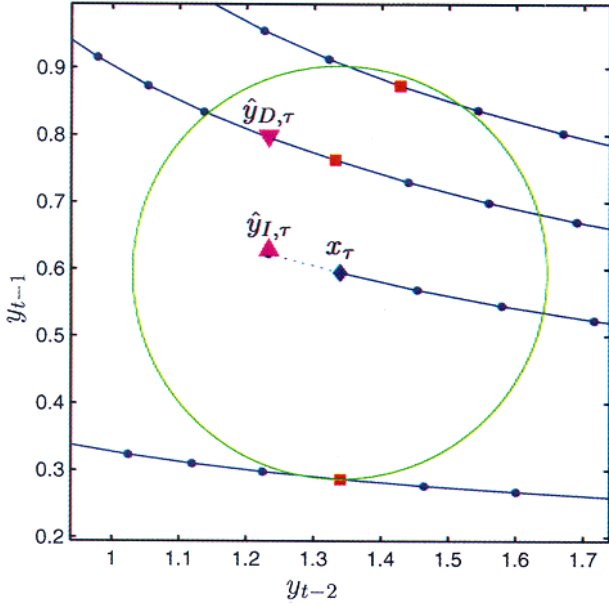
Fig. 8. Trajectory segments from the Lorenz time-series in a two-dimensional embedding space. The query point, $x_\tau$, is shown by ◆. The prediction generated by direct averaging is shown by ▼. The prediction generated by integrated averaging is shown by ▲.

trajectories. Because two of the trajectories are above the query point the predicted next point is biased upward.

One way to fix this problem is to build a model that predicts the change, or residual, $y_{t+1} - y_t$, rather than $y_{t+1}$ directly. The final prediction can then be found by integrating from the current state. The equations for direct local averaging ($\hat{y}_{D,\tau}$) and integrated local averaging ($\hat{y}_{I,\tau}$) are given below:

$$\hat{y}_{D,\tau} \triangleq \frac{\sum\limits_{i=1}^{k} w_i y_{c_i}}{\sum\limits_{i=1}^{k} w_i} \tag{27}$$

$$\hat{y}_{I,\tau} \triangleq \hat{y}_{I,\tau-1} + \frac{\sum\limits_{i=1}^{k} w_i (y_{c_i} - y_{c_i-1})}{\sum\limits_{i=1}^{k} w_i} \tag{28}$$

where $w_i$ are the weights defined in Eq. (17) and $c_i$ is the data set index of the $i$th nearest neighbor.

Figure 8 illustrates why integrated local averaging may be better than direct local averaging. In this case the integrated local average is much closer to the actual $y_\tau$ than the direct local average.

For long-term predictions these methods have very similar performance. However, if the time-series is a smoothly varying function of time then integrated local averaging is usually more accurate for very short-term predictions.

### 3.6. *Assessing model accuracy*

The type of local averaging discussed so far has four parameters: the embedding dimension ($m$), the embedding delay ($\delta$), the minimum metric weight ($\lambda_{\min}$), and the number of neighbors ($k$). Reasonable values of $\lambda_{\min}$ and $\delta$ can be chosen by the user but there is no obvious means of picking the best values of $m$ and $k$, which can strongly affect the model accuracy [Casdagli *et al.*, 1992; Casdagli & Weigend, 1994; Gershenfeld & Weigend, 1994; Smith, 1994]. However, these two parameters can be optimized globally if an accurate estimate of the model accuracy can be calculated cheaply.

One advantage that local models have over global models is that it is possible to calculate the leave-one-out cross-validation error very cheaply. This method begins by constructing a model to predict one step ahead with all but one of the $n_c$ points in the data set. The prediction error for the omitted point is calculated and the process is repeated for all $n_c$ points. The average leave-one-out error is used as a measure of the model accuracy.

Any of a variety of measures of error could be used. Absolute error, squared error, and the coefficient of correlation are common choices. For example, if squared error is used the cross-validation error is defined as

$$\text{MSE}_1 \triangleq \frac{1}{n_c} \sum_{i=1}^{n_c} (y_{i+1} - f_{i+1}^{-(i+1)}(x_i))^2 \tag{29}$$

$$\approx \text{E}[(y_{t+1} - f_{t+1}^{-(t+1)}(x_t))^2] \tag{30}$$

$$\approx \text{E}[(y_{t+1} - f_{t+1}(x_t))^2] \tag{31}$$

where $f_{i+1}^{-(i+1)}(x_i)$ is the model constructed to predict one step ahead with the $(i+1)$th point left out of the data set and $\text{E}[\,]$ denotes expectation.

A disadvantage of using one-step ahead cross-validation error (OSCV) for parameter optimization is that it does not take into account the model sensitivity to errors in the input vector that occur with iterated prediction [Eq. (14)]. Hence, the parameter values that minimize the OSCV error are generally not the same values that minimize the $p$-steps ahead cross-validation error.

It is better to choose an error measure that represents the cost of making poor predictions in the application at hand. In many cases an average (possibly weighted) model accuracy over $p$-steps ahead is appropriate,

$$\text{MSE}_{1,p} \triangleq \frac{1}{p\,n_p} \sum_{j=1}^{p} \sum_{i=1}^{n_p} (y_{i+j} - f_{i+j}^{-(i+1,\,i+p)}(x_i))^2$$

(32)

where $n_p = n_c - p + 1$ and $f_{i+j}^{-(i+1,\,i+p)}(x_i)$ is the model output $j$-steps ahead constructed with the points $(i+1)$ to $(i+p)$ left out. Kantz and Jaeger [1997] have argued that multi step ahead cross-validation (MSCV) also reduces bias.

To reduce the computation required to calculate $\text{MSE}_{1,p}$, fewer than $n_p$ terms could be used to approximate the expected error. For example, windows of $p$-points could be used that are spaced $s$ points apart.

$$\text{MSE}_{1,p} \approx \frac{1}{p\,n_v} \sum_{j=1}^{p} \sum_{i=1}^{n_v} (y_{v_i+j} - f_{v_i+j}^{-(v_i+1,\,v_i+p)}(x_{v_i}))^2$$

(33)

where $v_1 = 1$, $v_2 = s$, $v_3 = 2s$, and so on. Here $n_v$ is the number of windows, or validation sets.

For large values of $p$ and short time-series, removing $p$ values from the data set may significantly affect the model that is constructed. In this case, the cross-validation error will not be representative of the expected error using all $n_c$ points. One approach to reduce this effect is to omit only the trajectory segment that surrounds the point being predicted, $y_{v_i+\rho}$, instead of all $p$ points.

# 4.  The Winning Entry

The winning competition entry was generated using the method of local modeling described in the previous section. The exploratory data analysis, selection of user-specified parameters, and the model analysis that was performed prior to the competition deadline is described in this section.

## 4.1.  *Measure of accuracy*

Entrants were told that the predictions would be compared using mean squared error [Eq. (6)]. Hence the cross-validation MSE of Eq. (33) was chosen for the purpose of parameter optimization and accuracy estimation. This is equivalent to minimizing the average normalized mean squared error,

which is defined as

$$\text{NMSE}_\rho \triangleq \frac{\dfrac{1}{n_v} \sum_{i=1}^{n_v} (y_{v_i+\rho} - f_{v_i+\rho}^{-(v_i+1,\,v_i+200)}(x_{v_i}))^2}{\dfrac{1}{n} \sum_{i=1}^{n} (y_i - \bar{y})^2}$$

(34)

where $y_{v_i+\rho}$ is the $(v_i + \rho)$th point in the time-series, $v_i$ is the index of the first point in the $i$th validation set, and $\bar{y}$ is the average value of $y_t$. $f_{v_i+\rho}^{-(v_i+1,\,v_I+200)}(x_{v_i})$ is the prediction at time $v_i + \rho$ of a model that has been iterated $\rho$ times with the points $y_{v_i+1}, y_{v_i+2}, \ldots, y_{v_i+200}$ left out of the data set.

Similarly, the normalized root mean squared error is defined as

$$\text{NRMSE}_\rho \triangleq \sqrt{\text{NMSE}_\rho}.$$

(35)

$\text{NRMSE}_\rho$ is a convenient measure of error because it is independent of the scale of the time-series and $\text{NRMSE}_\rho = 1$ can be interpreted as meaning the model prediction error is no better than predicting the sample mean $\bar{y}$, on average.

## 4.2.  *Exploratory data analysis*

By visual inspection, the data set appeared to be noise-free (Fig. 2). A plot of the power spectrum (Fig. 9) indicated that the series had been sampled at approximately five times the Nyquist frequency and, therefore, little could be gained by upsampling the time-series. This is also confirmed by considering that more than 99.3% of the estimated signal power is at frequencies less than one-tenth the Nyquist frequency.

Further inspection indicated that the series was probably chaotic with three unstable equilibria at approximately $-0.25$, 0 and 0.25, as shown in Fig. 10. The period and growth of the oscillations about the upper and lower equilibria visually appeared to be the same. This manifested the possibility that the source of the time-series might be symmetrical; that is, it may just as likely have generated the time-series reflected about the horizontal axis (i.e. multiplied by $-1$), as shown in Fig. 11.

If the system was symmetrical a more accurate model could be built using both the original and reflected time-series. To determine if this was reasonable for the purpose of prediction, the cross-validation error *on the original time-series* was compared for two models: one built with the original
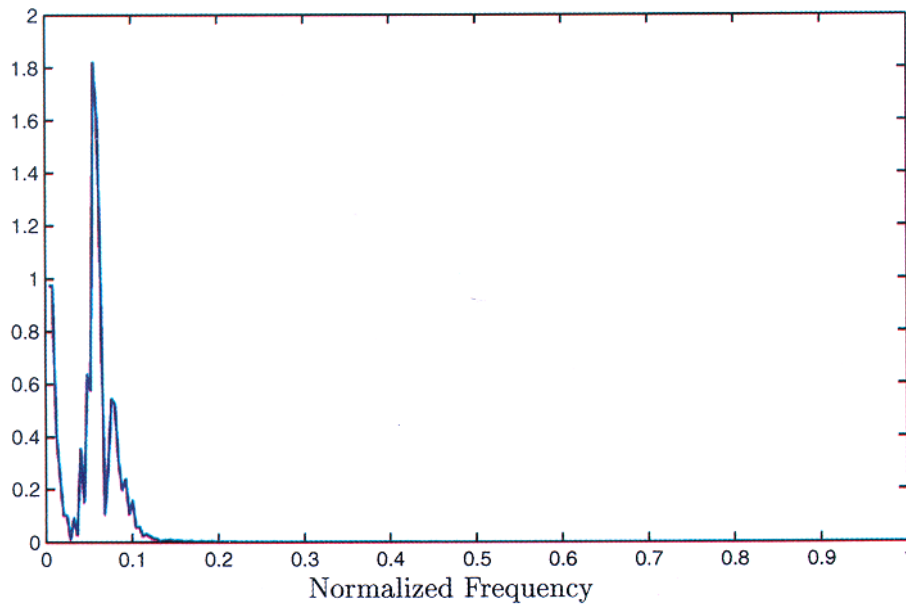
Fig. 9. The estimated power spectrum of the workshop time-series. The spectrum was estimated using Bartlett's method of periodogram averaging with four nonoverlapping windows. The $x$-axis has been scaled so that 1 represents half the Nyquist frequency.
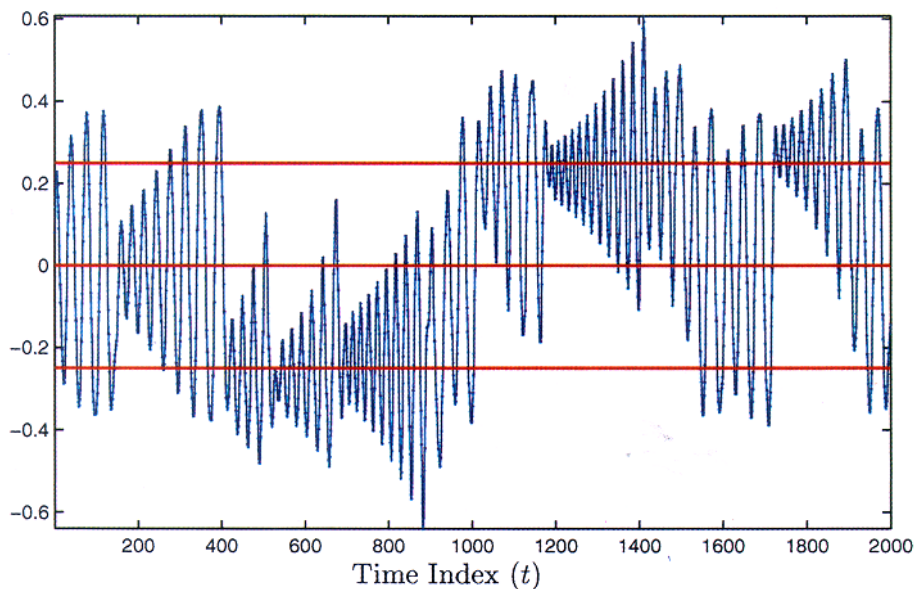


Fig. 10. The competition time-series. The three red horizontal lines show the approximate location of the three unstable equilibria in the series.

time-series and the other built with both time-series. Figure 12 illustrates that using both time-series results in significantly less cross-validation error than using only the original time-series. Consequently, both time-series were used to build the model that generated the winning entry. It has since been shown that the actual system that generated the time-series is noise-free and symmetric (see Sec. 2).

Figure 12 also illustrates that the average prediction error becomes worse than predicting the sample mean $(\bar{y})$ after approximately 80 steps ahead. However, this is only an *average* measure of error and does not represent the expected error for this particular prediction interval and data set. The accuracy of any single prediction depends upon many factors such as how densely the attractor is covered by the data set in the prediction region and
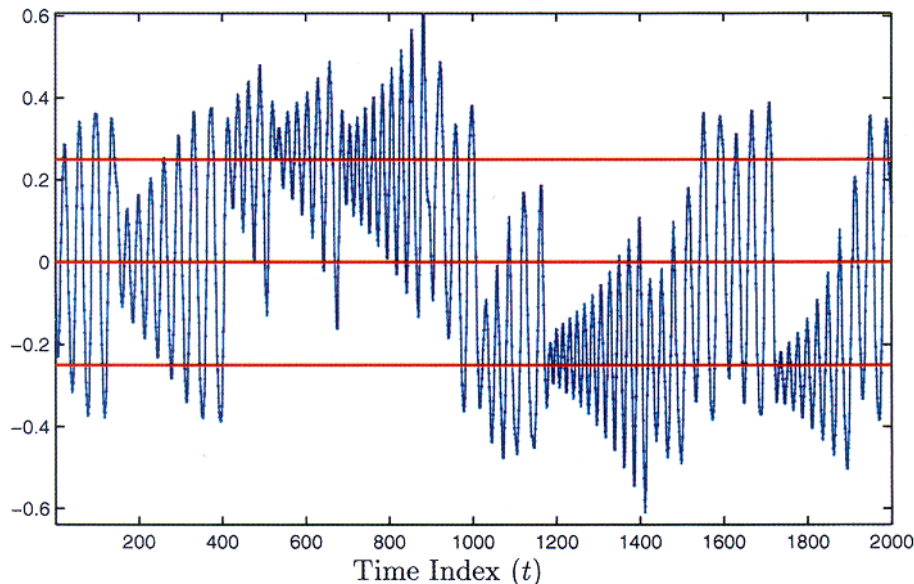
Fig. 11.   The competition time-series reflected about the horizontal axis. The three red horizontal lines show the approximate location of the three unstable equilibria in the series.
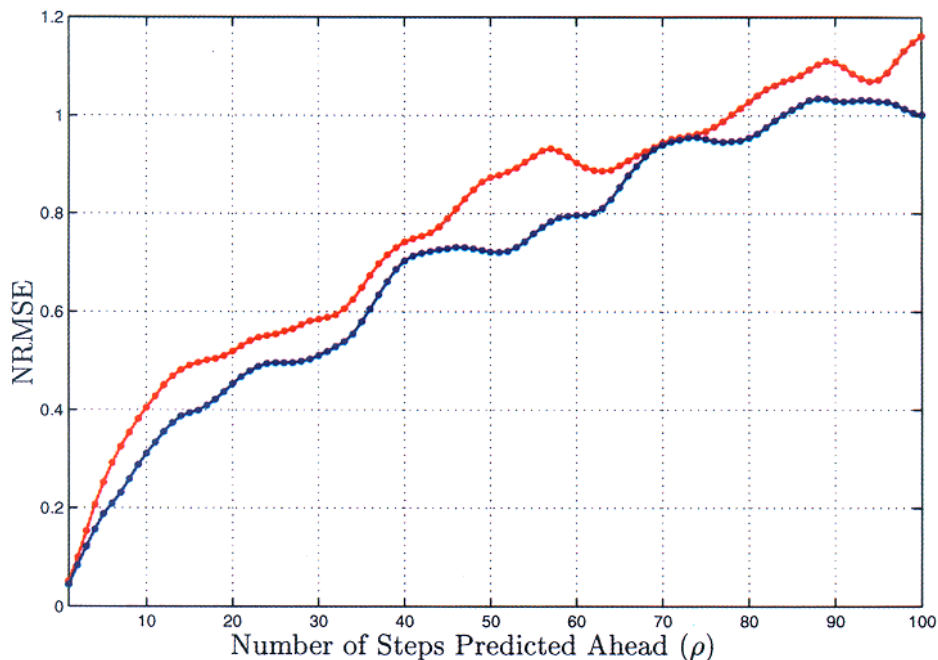


Fig. 12.   Plot of the NRMSE$_\rho$ as a function of $\rho$, the number of steps predicted ahead. The upper red trace shows the error using only the original time-series. The lower blue trace shows the error using both the original time-series and the reflected time-series.

the value of the largest Lyapunov exponent in the prediction region.

It is now known that the prediction region contained a single transition between scrolls (Sec. 2). Near the transition local trajectories diverge very rapidly which makes prediction after this point especially difficult.

### 4.3.   *Parameter selection and optimization*

The user-selected parameters that had to be chosen are the minimum metric weight ($\lambda_{\min}$), the number of points used to calculate the MSE ($n_v$), and the embedding delay ($\delta$). All three of these parameters

Table 1.  Summary of local modeling method that generated the winning entry.

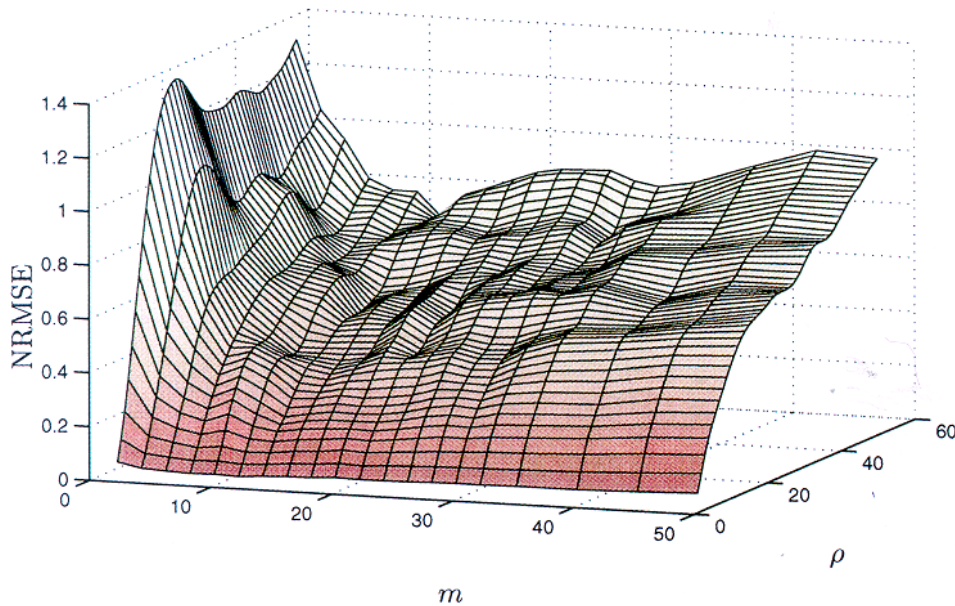| | |
|---|---|
| Preprocessing: | Both the original data set and the reflected data set were used to construct the local model. |
| Nearest Neighbor Algorithm: | ANNA as described by McNames [1998].  However, any nearest neighbor algorithm that incorporates the nearest trajectory modification described in Sec. 3.4 could have been used. |
| Nearest Neighbor Metric: | Exponentially weighted Euclidean [Eq. (25)] with user-selected parameters $\lambda_{\min} = 0.05$ and $\delta = 2$. |
| Local Model Type: | Integrated local averaging [Eq. (28)] |
| Local Weighting: | Biweight function [Eq. (17)] |
| Parameter Optimization: | Global minimization of multistep ahead cross-validation mean squared error [Eq. (33) with $n_v = 33$]. |
| Optimized Parameters: | Embedding dimension: $m = 16$. Number of nearest neighbors: $k = 1$ |



Fig. 13.  $\mathrm{NRMSE}_\rho$ as a function of the number of steps predicted ahead ($\rho$) and the embedding dimension ($m$).  $m = 4$ minimized the one-step ahead cross-validation error and $m = 16$ minimized the multistep ahead cross-validation error. This plot was generated using 200 cross-validation points.

were constrained only by limited computational resources. The values that were selected are given in Table 1. The choice of these parameters was fortuitous; slight changes to these parameters result in worse predictions after the transition point, though the predictions up to that point are much less sensitive.

Both the embedding dimension ($m$) and the number of neighbors ($k$) were chosen so as to minimize MSE.[2] Figure 13 shows the $\mathrm{MSE}_\rho$ for different values of $m$ and $\rho$. This figure illustrates the importance of minimizing the average multistep ahead cross-validation (MSCV) error rather than the one-step ahead cross-validation (OSCV) error. If OSCV had been used for parameter optimization a much smaller value of $m$ would have been chosen ($m = 22$) which would have had a larger MSCV and generated a worse prediction.

Similarly, Fig. 14 shows the $\mathrm{NRMSE}_\rho$ for different values of $k$ and $\rho$. The same problem with

---

[2]An equivalent approach was actually used to generate the winning entry where the embedding dimension ($m$) was held constant at 50 and $\lambda$ was optimized [Eq. (23)]. However, the method described here produces the same prediction and has the advantages of being slightly cheaper computationally and the user-selected parameter is more intuitive ($\lambda_{\min}$ instead of $m$). See [McNames, 1998] for more details.
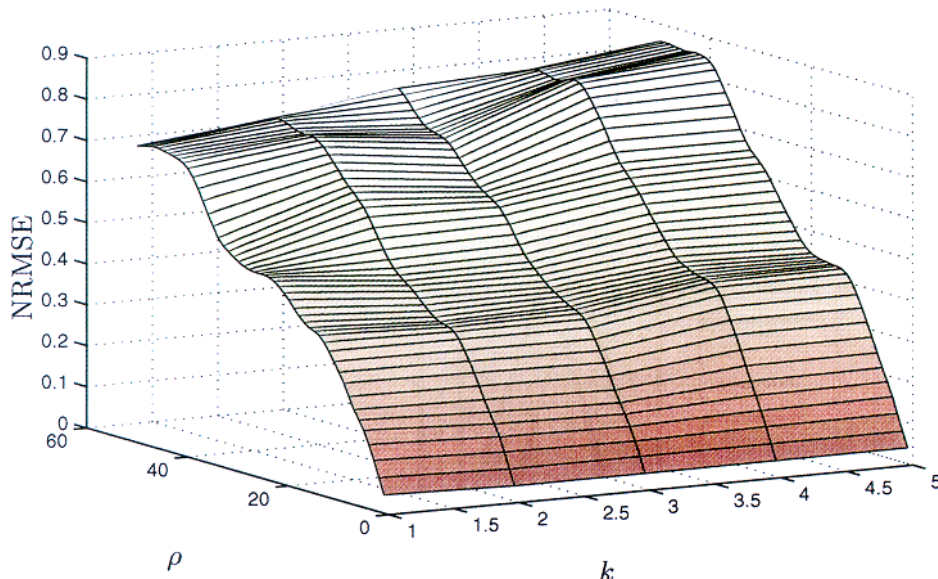
Fig. 14. NRMSE$_\rho$ as a function of the number of steps predicted ahead ($\rho$) and the number of neighbors ($k$). $k = 5$ minimized the one-step ahead cross-validation error and $k = 1$ minimized the multistep ahead cross-validation error. This plot was generated using 200 cross-validation points.
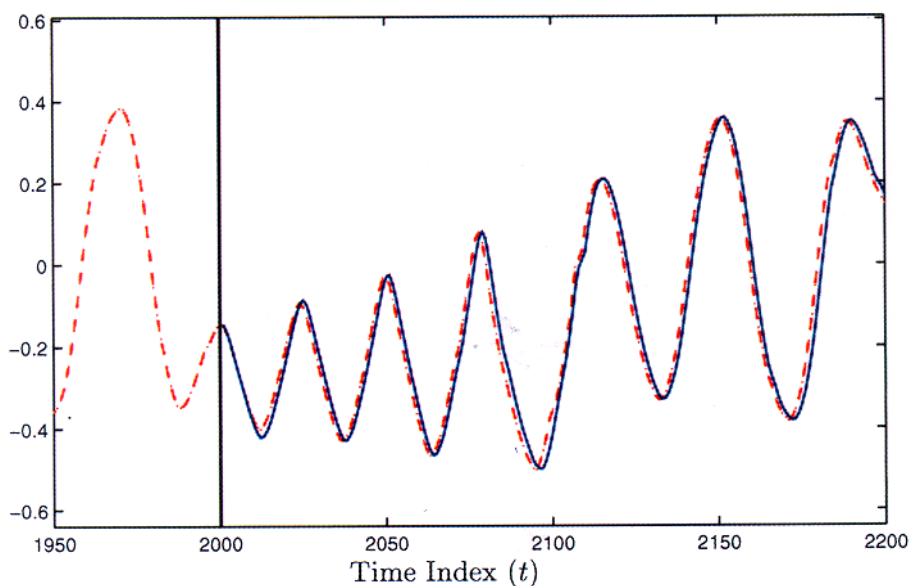


Fig. 15. The dashed trace shows the winning competition entry and the full trace shows the true time-series. The method used to generate the winning entry is summarized in Table 1.

OSCV is illustrated by this figure; the value of $k$ that minimized the average OSCV ($k = 5$) would not have minimized the MSCV and would have generated a worse prediction.

The values of $m$ and $k$ that minimized the cross-validation error are given in Table 1. The winning prediction generated by the model using these values is shown in Figs. 15 and 16. A summary of the method that generated the winning entry is also given in the table.

## 5. Conclusions

Takens' theorem provides a sound theoretic basis for using general nonlinear modeling methods for time-series prediction. However, if the time-series is generated by a nonlinear dynamic system and exhibits chaotic behavior, a nonlinear model that incorporates this knowledge will generally perform better.
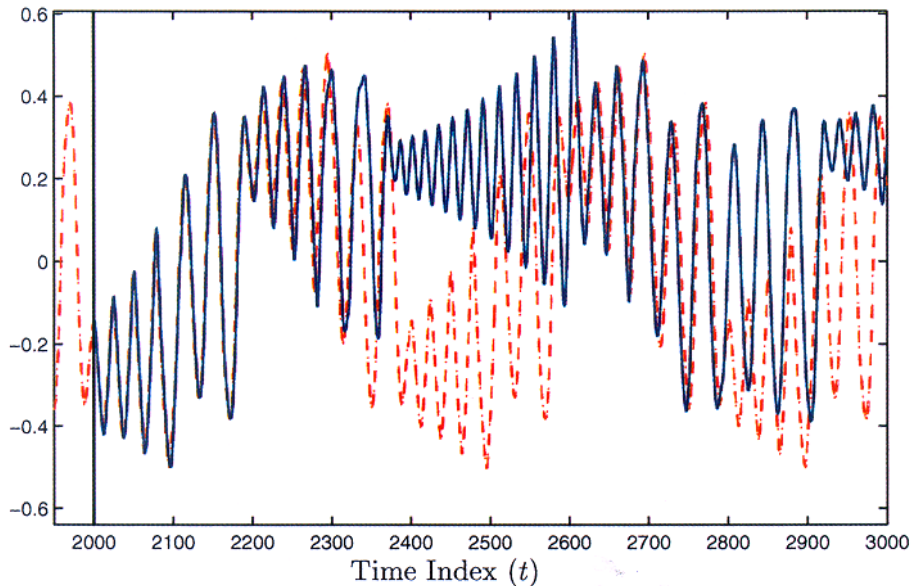
Fig. 16. A continuation of the predicted data for the winning entry. An accurate prediction is obtained over a time horizon of about 300 points. The dashed red trace shows the true continuation and the full blue trace shows the 1000-step ahead prediction.

Several such modifications to local modeling have been described in this paper including the use of integrated local averaging, an exponentially weighted Euclidean metric, a fast triangle inequality-based nearest neighbor algorithm, a *nearest trajectory* algorithm, and a novel multistep ahead cross-validation method of accuracy estimation for parameter optimization.

The prediction interval contained a single transition between scrolls. All five of the competition entries that employed local modeling generated accurate predictions up to this transition (about 80 steps ahead) [Suykens & Vandewalle, 1998a, 1998b]. A synergy of exploiting the data set symmetry, a fortuitous choice of user-selected parameters, and a method of local modeling tailored specifically to chaotic time-series prediction produced the winning entry.

## References

Arena, P., Baglio, S., Fortuna, L. & Manganaro, G. [1996a] "Generation of *n*-double scrolls via cellular neural networks," *Int. J. Circuit Th. Appl.* **24**, 241–252.

Arena, P., Baglio, S., Fortuna, L. & Manganaro, G. [1996b] "State controlled CNN: A new strategy for generating high complex dynamics," *IEICE Trans. Fundamentals* **E79-A**, 1647–1657.

Atkeson, C. G., Moore, A. W. & Schaal, S. [1997] "Locally weighted learning," *Artif. Intell. Rev.* **11**(1–5), 11–73.

Bakamidis, S. G. [1993] "An exact fast nearest neighbor identification technique," in *1993 IEEE Int. Conf. Acoustics, Speech, and Signal Processing*, Vol. 5, pp. 658–661.

Casdagli, M. C., Des Jardins, D., Eubank, S., Farmer, J. D., Gibson, J. & Theiler, J. [1992] "Nonlinear modeling of chaotic time-series: Theory and applications," in *Applied Chaos*, eds. Kim, J. H. & Stringer, J. (John Wiley & Sons), pp. 335–380.

Casdagli, M. C. & Weigend, A. S. [1994] "Exploring the continuum between deterministic and stochastic modeling," in *Time-Series Prediction*, eds. Weigend, A. S. & Gershenfeld, N. A., Santa Fe Institute Studies in the Sciences of Complexity (Addison-Wesley), pp. 347–366.

Chua, L. [1994] "Chua's circuit: An overview ten years later," *J. Circuits Syst. Comput.* **4**(2), 117–159.

Chua, L., Komuro, M. & Matsumoto, T. [1986] "The double scroll family," *IEEE Trans. Circuits Syst. I* **33**(11), 1072–1118.

Farmer, J. D. & Sidorowich, J. J. [1987] "Predicting chaotic time-series," *Phys. Rev. Lett.* **59**(8), 845–848.

Farmer, J. D. & Sidorowich, J. J. [1988] "Exploiting chaos to predict the future and reduce noise," in *Evolution, Learning and Cognition*, ed. Lee, Y. C. (World Scientific), pp. 277–330.

Friedman, J. H. [1991] "An overview of predictive learning and function approximation," in *From Statistics to Neural Networks*, eds. Cherkassky, V., Friedman, J. H. & Wechsler, H., Computer and Systems Sciences (Springer-Verlag), Vol. 13, pp. 1–61.

Friedman, J. H., Bentley, J. L. & Finkel, R. A. [1977] "An algorithm for finding best matches in

logarithmic expected time," *ACM Trans. Math. Software* **3**(3), 209–226.

Fukunaga, K. & Narendra, P. M. [1975] "A branch and bound algorithm for computing *k*-nearest neighbors," *IEEE Trans. Comput.* **C-24**, 750–753.

García, P., Jiménez, J., Marcano, A. & Moleiro, F. [1996] "Local optimal metrics and nonlinear modeling of chaotic time-series," *Phys. Rev. Lett.* **76**(9), 1449–1452.

Gershenfeld, N. A. & Weigend, A. S. [1994] "The future of time-series: Learning and understanding," in *Time-Series Prediction*, eds. Weigend, A. S. & Gershenfeld, N. A., Santa Fe Institue Studies in the Sciences of Complexity (Addison-Wesley), pp. 1–70.

Kalantari, I. & McDonald, G. [1983] "A data structure and an algorithm for the nearest point problem," *IEEE Trans. Software Engin.* **SE-9**(5), 631–634.

Kamgar-Parsi, B. & Kanal, L. N. [1985] "An improved branch and bound algorithm for computing *k*-nearest neighbors," *Patt. Recogn. Lett.* **3**, 7–12.

Kants, H. & Jaeger, L. [1997] "Improved cost functions for modeling of noisy chaotic time-series," *Physica* **D109**, 59–69.

Kim, B. S. & Park, S. B. [1986] "A fast *k*-nearest neighbor finding algorithm based on the ordered partition," *IEEE Trans. Patt. Anal. Mach. Intell.* **8**(6), 761–766.

Kugiumtzis, D., Lingjærde, N. & Christophersen, N. [1998] "Regularized local linear prediction of chaotic time-series," *Physica* **D112**, 344–360.

Lillekjendlie, D. K. & Christophersen, N. [1994] "Chaotic time-series. Part II. System identification and prediction," *Modeling, Identification and Control* **15**(4), 225–243.

Liu, Z., Ren, X. & Zhu, Z. [1997] "Equivalence between different local prediction methods of chaotic time-series," *Phys. Lett.* **A227**, 37–40.

Madan, R. [1993] *Chua's Circuit: A Paradigm for Chaos* (World Scientific, Singapore).

McNames, J. [1998] "A nearest trajectory strategy for time-series prediction," in *Proc. Int. Workshop on Advanced Black-Box Techniques for Nonlinear Modeling*, Katholieke Universiteit Leuven, Belgium, pp. 112–128.

Micó, L., Oncina, J. & Carrasco, R. C. [1996] "A fast branch & bound nearest neighbor classifier in metric spaces," *Patt. Recogn. Lett.* **17**, 731–739.

Murray, D. B. [1993] "Forecasting a chaotic time-series using an improved metric for embedding space," *Physica* **D68**, 318–325.

Nene, S. A. & Nayar, S. K. [1997] "A simple algorithm for nearest neighbor search in high dimensions," *IEEE Trans. Patt. Anal. Mach. Intell.* **19**(9), 989–1003.

Niemann, H. & Goppert, R. [1988] "An efficient branch-and-bound nearest neighbor classifier," *Patt. Recogn. Lett.* **7**, 67–72.

Sauer, T. [1994] "Time-series prediction by using delay coordinate embedding," in *Time-Series Prediction* eds. Weigend, A. S. & Gershenfeld, N. A., Santa Fe Institue Studies in the Sciences of Complexity (Addison-Wesley), pp. 175–193.

Sauer, T., Yorke, J. A. & Casdagli, M. [1991] "Embedology," *J. Stat. Phys.* **65**(3), 579–616.

Short, K. M. [1994] "Steps toward unmasking secure communications," *Int. J. Bifurcation and Chaos* **4**(4), 957–977.

Short, K. M. [1996] "Unmasking a modulated chaotic communications scheme," *Int. J. Bifurcation and Chaos* **6**(2), 367–375.

Smith, L. A. [1994] "Local optimal prediction: Exploiting strangeness and the variation of sensitivity to initial condition," *Philos. Trans. R. Soc.* **A348**, 371–381.

Suykens, J. A. K. & Chua, L. O. [1997] "*n*-double scroll hypercubes in 1D-CNNS," *Int. J. Bifurcation and Chaos* **7**(8), 1873–1885.

Suykens, J. A. K., Curran, P., Vandewalle, J. & Chua, L. O. [1997] "Robust nonlinear $h_\infty$ synchronization of chaotic lur'e systems," *IEEE Trans. Circuits Syst. I, Special Issue on Chaos Synchronization, Control and Applications* **44**(10), 891–904.

Suykens, J. A. K., Huang, A. & Chua, L. O. [1997] "A family of *n*-scroll attractors from a generalized Chua's circuit," *Archiv fur Elektronik und Ubertragungstechnik, Int. J. Electron. Commun.* **51**(3), 131–138.

Suykens, J. A. K. & Vandewalle, J. (eds.) [1998a] *Proc. Int. Workshop on Advanced Black-Box Techniques for Nonlinear Modelings*, Katholieke Universiteit Leuven, Belgium.

Suykens, J. A. K., Vandewalle, J. & Chua, L. O. [1997] "Nonlinear $h_\infty$ synchronization of chaotic lur'e systems," *Int. J. Bifurcation and Chaos* **7**(6), 1323–1335.

Suykens, J. A. K. & Vandewalle, J. [1998b] *Nonlinear Modeling Advanced Black-Box Techniques* (Kluwer Academic Publishers).

Tai, S. C., Lai, C. C. & Lin, Y. C. [1996] "Two fast nearest neighbor searching algorithms for image vector quantization," *IEEE Trans. Commun.* **44**(12), 1623–1628.

Takens, F. [1981] "Detecting strange attractors in turbulence," in *Dynamical Systems and Turbulence*, eds. Rand, D. A. & Young, L. S., Lecture Notes in Mathematics **898** (Springer-Verlag), pp. 336–381.

Tanaka, N., Okamoto, H. & Naito, M. [1995] "An optimal metric for predicting chaotic time-series," *Jpn. J. Appl. Phys.* **34**(1), 388–394.

Vidal, E. [1986] "An algorithm for finding nearest neighbors in (approximately) constant average time," *Patt. Recogn. Lett.* **4**, 145–157.

Yang, T. [1995] "Recovery of digital signal from chaotic switching," *Int J. Circuit Th. Appl.* **33**(6), 611–615.