

Error correction of a predictive ocean wave model using local model approximation

Vladan Babovic^a, S.A. Sannasiraj^{b,*}, Eng Soon Chan^c

^a*Tectrasys AG, Sihleggstrasse 23, 8832 Wollerau, Switzerland*

^b*Department of Ocean Engineering, Indian Institute of Technology Madras, Chennai, 600 036, India*

^c*Tropical Marine Science Institute, National University of Singapore, Singapore*

Received 3 June 2002; accepted 24 May 2004

Available online 19 September 2004

Abstract

Constructing models from time series with nontrivial dynamics is a difficult problem. The classical approach is to build a model from first principles and use it to forecast on the basis of the initial conditions. Unfortunately, this is not always possible. For example, in fluid dynamics, a perfect model in the form of the Navier–Stokes equations exists, but initial conditions and accurate forcing terms are difficult to obtain. In other cases, a good model may not exist. In either case, alternative approaches should be examined. This paper describes an alternative approach of combining observations and numerical model results in order to produce an accurate forecast. The approach is based on application of a method inspired by chaos theory for building nonlinear models from data called Local Models. Embedding theorem based on the time lagged embedded vectors is the basis for the local model. This technique is used for analysis and updating of numerical model output variables to forecast and correct the errors created by numerical model. The local model approximation is a powerful tool in the forecasting of chaotic time series and has been employed for wave prediction in a forecasting horizon from a few hours to 24 h. The efficacy of the local model as an error correction tool (by combining the model predictions with the observations) compared with the predictions of linear auto regressive models has been brought up. In the present study, the parameters driving the local model are optimized using evolutionary algorithms.

© 2004 Elsevier B.V. All rights reserved.

Keywords: Wave forecasting; Chaos theory; Optimization; Embedding theorem; Genetic algorithm; Local model

1. Introduction

Numerical models are far from being perfect. A numerical model is indeed only a model of reality. It employs a number of simplifying assumptions, such as depth averaging of velocities in vertically integrated two-dimensional models, which inevitably

* Corresponding author. Fax: +91 44 22578625.

E-mail addresses: vb62@bluewin.ch (S.A. Babovic),
sasraj@iitm.ac.in (V. Sannasiraj), tmsdir@nus.edu.sg (E.S. Chan).

produce inaccuracies. In a numerical model, one also discretises the domain and is therefore not able to resolve numerous subgrid-scale phenomena. Errors in the model parameterisation (mainly because most model parameters cannot be directly measured) may contribute significantly to the overall error in a numerical model. It is also impossible to precisely define initial conditions and forcing terms in the entire computational domain. All of these inaccuracies and uncertainties could accumulate to produce poor model results, despite our ‘perfect’ knowledge of the governing laws.

To combat the inevitable presence of such model errors, a number of approaches for correcting the model results are often employed. Data assimilation is a methodology that utilizes information from observations, and combines it with (or assimilate it into) numerical models.

A number of different data assimilation procedures can be adopted. These are designed to either improve description of initial conditions at the time of forecast or provide correction of model predictions during a forecast period. The data assimilation procedures may be classified according to the variables modified during the updating process. In WMO Report (WMO, 1992; Refsgaard, 1997), four

different methodologies have been defined (Fig. 1) as follows:

- Updating of input parameters: This is the classical method justified by the fact that input uncertainties may be the dominant error source in operational forecasting.
- Updating of state variables: The correction of the state variables can be done in different ways. The theoretically most comprehensive methodology is based on Kalman filtering (Gelb, 1974). Kalman filtering is the optimal updating procedure for linear systems, but can, with some modifications, also provide an approximate solution for nonlinear hydrodynamic systems.
- Updating of model parameters: The continuous adaptation of model parameters is a matter of continuous debate. The prevailing view seems to be that for hydrodynamic models of nontrivial complexity, recalibration of the model parameters at every time step has no real advantages, as the operation of any hydrodynamic system cannot significantly change over the short interval of time.
- Updating of output variables (error prediction): The deviations between the model simulation and the observed variables such as wave height are

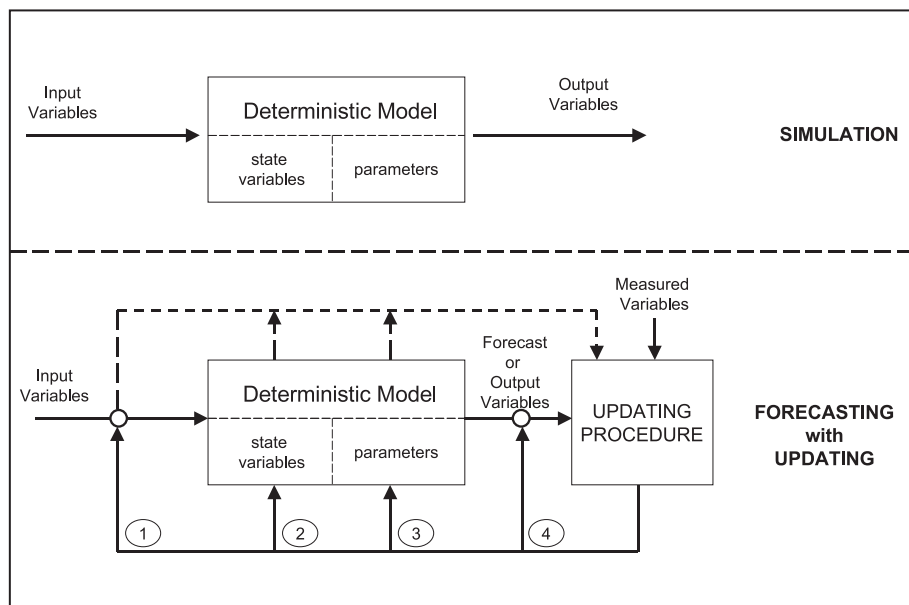


Fig. 1. Schematic diagram of simulation and forecasting with emphasis on the four different updating methodologies.

model errors. Possibility of forecasting these errors and superimposing them to the simulation model forecasts usually give good performance. This method is most often referred to as error prediction and is the method employed in the present study.

Alternatively, if forecasting interest is limited to only a few variables at some specific locations with a high degree of accuracy and for a considerably long lead-time forecast, statistical forecasting models based on updating of output variables (error prediction) may be the most suitable approach. These models can be either linear such as autoregressive integrated moving average (ARIMA) models or nonlinear models based on chaos. In model error prediction, other techniques such as artificial neural networks (Babovic, 1996; Henk et al., 1996; Minns, 1998), genetic programming (Khu et al., 2001; Babovic and Keijzer, 2003) or an approach based on chaos theory (Babovic and Keijzer, 1999) have demonstrated good forecast skill. By using these techniques, one can combine the forecast of the numerical model (model output) at the point of interest with the latest observed data in order to obtain an improved forecast. Another advantage of such an approach is that it allows the combination of different variables (for example, atmospheric data such as wind speed and oceanography consequent such as wave height) to improve the accuracy. This cannot be done in conventional data assimilation methods where the data has to be introduced in the model state in order to be assimilated.

In this paper, the advantages of both the numerical and chaotic methods have been explored. Bringing the underlying dynamics of a chaotic time series into a deterministic predictive wave model would improve the prognostic capability with good accuracy. The prediction of nonlinear dynamics is based on the embedding theorem proposed by Takens (1980) and on the works of Abarbanel (1996). The numerical model is based on the formulation of the third-generation wave model, WAM (The WAMDI group, 1988). The merger would be defined as an error correction tool at individual stations, where sufficient length of past observations is available. The efficiency of the local

model forecast has shown significant improvements over the prediction from a globally linear autoregressive model.

2. Ocean wave modelling

The importance of prediction of the offshore wave environment accelerated in the sixties, which laid foundations for the first generation wave model. The challenging task then was the formulation of the momentum transfer from wind to sea surface. The works of Miles (1957) and Philips (1957) were among the early advancements of the wave modelling through better definition of wind–wave relation. The formulation, however, overestimated the wind input. The first-generation wave models were also incomplete without nonlinear wave–wave interactions and the high-frequency spectral tail was prescribed (Philips, 1958). The second-generation models evolved about a decade later with the inclusion of nonlinear wave–wave interaction terms, but not in a full sense. It was due to the lack of computational power in early eighties as well as the difficult task in the integration of nonlinear interaction terms. The second-generation models were unable to predict complex wind–seas generated by rapidly changing wind fields due to hurricanes, intense, small-scale cyclones or fronts. The models also encountered basic difficulties in treating the transition between wind–sea and swell. The shortcomings in the second-generation models were eradicated in the third-generation models. The two-dimensional wave spectrum was allowed to grow without any imposed conditions. The wind–wave momentum transfer was defined in a better manner (Snyder et al., 1981; Janssen, 1989). Hasselmann and Hasselmann (1985) suggested discrete integral approximation scheme to evaluate the complex integration of nonlinear wave–wave interaction terms. The integration of the nonlinear interaction terms is still the time expensive part of the wave prediction. The emergence of third-generation wave model in the late 1980s is still playing a vital role in operational wave predictions. Presently, WAM (The WAMDI group, 1988) and WaveWatch III (Tolman, 1999) are widely used operational third-generation

wave models in most of the meteorological stations worldwide.

2.1. The wave model

In the present paper, a third-generation wave model, WAM (The WAMDI group, 1988; Komen et al., 1994) was considered. WAM estimates the evolution of the energy spectrum for ocean waves by solving the wave transport equation explicitly without any presumptions on the shape of the wave spectrum.

$$\frac{\partial F(f, \theta; x, t)}{\partial t} + v \nabla_x F(f, \theta; x, t) = S \quad (1)$$

where $F(f, \theta; x, t)$ is the wave energy spectrum in terms of frequency f and propagation direction θ at the position vector x and at time t ; v is the group velocity. The second term on the left-hand side is the divergence of the convective energy flux ($v \nabla_x F$). The net source function S takes into account all physical processes, which contribute to the evolution of the wave spectrum. The source function is represented as superposition of source terms due to wind input, nonlinear wave–wave interaction, dissipation due to wave breaking and bottom friction.

$$S = S_{in} + S_{nl} + S_{ds} + S_{bot} \quad (2)$$

The wind input source function, S_{in} was adopted from Snyder et al. (1981) and Komen et al. (1994). The nonlinear source function, S_{nl} is represented by the discrete interaction operator parameterisation proposed by Hasselmann and Hasselmann (1985). The dissipation source function, S_{ds} , is proposed by Komen et al. (1984) and modified by The WAMDI group (1988) to enhance the stability of the implicit numerical scheme. The additional dissipation term, S_{bot} , representing the energy loss due to bottom friction and percolation in shallower waters, is included following the definitions of Hasselmann et al. (1973).

The synthesis of these source terms signifies the current state of understanding of the physical processes of wind waves, namely, the inputs from the processes of wind field, nonlinear interaction, white capping and dissipation due to bottom friction balance each other to form self similar spectral shapes corresponding to the measured wind wave spectra. Except for the nonlinear source term, all the other

source terms are individually parameterised to be proportional to the action density spectrum, F . The nonlinear source uses the discrete interaction approximation that simulates a nonlinear transfer process formulated by the four-wave resonant interaction Boltzmann equation and characterises the third-generation model.

On solving the wave transport Eq. (1), the directional wave spectrum in each spatial grid location at every time step has been estimated. The spectral parameters such as significant wave height H_s , mean wave period T_m and mean wave direction θ_m can be derived. These parameters are the first guess model values from the numerical model for the given wind field.

The accuracy of wave forecasts thus depends on the wave models and the accuracy of prescribed driving force, wind. To improve the accuracy of forecasting, hence, one can update the model definition and driving coefficients as well as prescribe accurate surface wind. Being wave models reached a stagnant stage for more than a decade now; it is hard to improve except few model parameters such as dissipation and refractive terms. Even then, the improvement is insignificant in most of the cases. The other factor, the prescription of accurate driving force, would give the correct forecast. The wind is again forecasted by meteorological models characterized by many atmospheric variables which question the accuracy of prediction. The wave forecasts by many operational weather centres are limited by the quality of wind (Christopoulos, 1997; Bidlot and Holt, 1999). The seasonal variation of root-mean-square error (RMSE) of significant wave height was reported to be 0.3 to 0.6 m during December 95 to August 98 (Bidlot and Holt, 1999). In this situation, time series forecasting models and data assimilation enter into the system for steering the model predictions towards observations. The following section describes an efficient time series forecasting technique.

3. Characterisation of model errors

Many of the important aspects of analyzing dynamical systems are carried out by the study of observable variables of the system (such as signifi-

cant wave height, mean wave period and mean wave direction) as a function of time. However, it may be argued that not enough attention is given to the representation of the data obtained from field or laboratory experiments. As a result, observations, which appear randomly under time series representation, are typically discarded as noise in most cases.

Recent developments in nonlinear dynamics have demonstrated that irregular or random behavior in natural systems may arise from purely deterministic dynamics with unstable trajectories. Although some observations might appear randomly, beneath their random behavior may lay an order or pattern. Such types of nonlinear dynamical systems, which are also highly sensitive to initial conditions, are popularly known as *chaotic systems*. Major advances in the analysis of chaotic systems have helped to develop forecast skill that can also be applied for prediction of nonlinear systems. Inspired by this line of thought, the present paper employs some of the techniques developed in the study of chaotic systems to characterise and forecast errors of deterministic wave models.

The globally linear (such as auto regressive or moving average) models have dominated the field of time series analysis and forecasting for more than half a century. These models transform the signal into a small number of coefficients plus residual white noise. However, such appealing simplicity of linear models can entirely be misleading even when weak nonlinearities occur. It has been shown (e.g., Gershenfeld and Weigend, 1993) that the power spectra of such globally linear models and related autocorrelation coefficients contain the same information about a system driven by uncorrelated white noise. Thus, if and only if the power spectrum fully characterises the relevant features of a time series, a linear model would be the appropriate description of the system. At the same time, it is widely known that two time series can have similar broadband spectra but can be generated from systems with different properties. Such is a linear system that is driven stochastically by external noise and a deterministic (noise-free) nonlinear system with a small number of degrees of freedom. Therefore, in order to forecast a nonlinear system, a different class of models is called upon.

3.1. Embedding theorem

Let a real time process result in a time series $x(t)$ ($=\{x(t_0+i\Delta t)\}$, $i=0,1,\dots,n$) sampled at intervals Δt and initiated at t_0 . The process $x(t)$ represents either significant wave height H_s , mean wave period T_m or mean wave direction θ_m . The purpose of forecasting is to predict the state of the system $x(t)$ at a time horizon T in the future, say $x(t+T)$. An auto-regressive moving average (ARMA) model can achieve forecasting using the immediately preceding values which can be written as,

$$x_t = ax_{t-i} + be_t, \quad i = 1..M \quad (3)$$

where, $b=\{b_0, b_1, \dots, b_N\}$, $e=\{e_0, e_1, \dots, e_N\}$, $a=\{a_0, a_1, \dots, a_M\}$ and $x=\{x_0, x_1, \dots, x_M\}$ denote lag vectors and also referred to as tapped delay lines. Here, $\{a\}$ are auto-regressive coefficients of the order M and $\{b\}$ are moving average coefficients of the order N .

Packard et al. (1980) explored such time-lagged vectors to realize the underlying dynamics and proposed the *Time-Delay Embedding theorem*. The theorem was made impeccable by Takens (1980) and later strengthened by Sauer and Yorke (1991). The *Time-Delay Embedding theorem* can be stated as follows:

Consider a dynamical system with a δ -dimensional space and an evolving solution $g(t)$. Let x_t be some observation $x(g(t))$ and τ_i be specified time lags. The lag vector can be defined as $\mathbf{x}(t) \equiv \{x_t, x_{t-\tau_1}, x_{t-\tau_2}, x_{t-\tau_3}, \dots, x_{t-\tau_{\delta-1}}\}$. In general, the space of vectors $\mathbf{x}(t)$ generated by the dynamics contains all of the information of the space of solution vectors $\mathbf{g}(t)$. The mapping between them is smooth and invertible. This property is referred to as an embedding. Thus, the study of the time series $x(t)$ is also the study of the solutions of the underlying dynamical system $g(t)$ via a particular coordinate system given by the observable x .

The embedding theorem establishes that, given a scalar time series from a dynamical system, it is possible to reconstruct a phase space from this single variable that is, in theory, an embedded space with dimensions consisting of various time lags of the variable itself. The embedded space can also be

created from many dynamic variables. For example, the measured wave height record could be thought of as a projection of its own and as time histories of other variables. The additional variable in the environment for a wave height prediction may be wind speed, characteristic wave frequency and/or wave direction. This projected wave record contains information about all the other contributing phenomena, although it is measured at only one point in the geographic space. According to the embedding theorem, the underlying structure cannot be seen in the space of the original scalar time series, but rather only when unfolded into an embedded (or phase) space. Time series can correspondingly be forecasted based on this structure in the phase space.

3.2. Local model

Local models are particularly well suited for the forecasting of chaotic time series because these share many fundamental ideas with the time-delay embedding theorem. A rather effective method of simulating the evolution of a dynamical system is by means of a local approximation, using only the most similar trajectories from the past to make predictions of the future. The following steps describe the local modelling in the context of wave forecasting using single-variable H_s .

3.2.1. Step 1: Embedding the time series into a phase (embedded) space

Let $x(t)$ be the time series of wave height, H_s with the records sampled at a uniform interval of Δt hours. For an embedding dimension δ , with a time lag τ , a time-delay embedding space can be constructed. The coordinates of the points in such an embedded space, $H(\tau, t)$ are represented at any given time, t as,

$$\begin{aligned} H_i(\tau, t) &= [x(t - (\delta - 1)\tau), x(t - (\delta - 2)\tau), \dots, \\ &\quad x(t - \tau), x(t)], i = 1, \dots, n - (\delta - 1); \\ t &= (i - 1)\Delta t + (\delta - 1)\tau \end{aligned} \quad (4)$$

Let n be the length of the historical wave height record and then the number of projected points in the embedded space is $(n - (\delta - 1))$. The above equation represents the embedded space using a uniformly

distributed time lag. The embedding can also be achieved with different time lags in each dimensional space ($\tau_i, i=1, \dots, \delta$). The historical data set, $x(t)$ is also referred to as training series. The historical records (training series) need not be extended up to the present time horizon.

3.2.2. Step 2: Search for neighbours in the embedded space

Let $t=t_n$ indicates the time step from which the forecast is to be evaluated. The most similar points of the starting forecast point, $H(\tau, t_n)$ within the phase space, $H(\tau, t)$ are chosen. Let $H'()$ are the selected neighbourhood points from the phase space, $H'() \in H()$. The neighbours can be chosen either within a circle of constant radius (distance) from the forecast point or by specifying the number of nearest neighbour, k .

3.2.3. Step 3: Constructing the “expected value” vector

Having constructed the phase space and pooled the most similar points corresponding to the present time horizon, t_n , the expected (forecast) value vector, $X(t)$ is formed for each point in the neighbourhood domain, $H'()$ as

$$X_j(t) = x_i(t_n + T) \quad j = 1, \dots, k; \quad i = f(j) \quad (5)$$

where, $T (=T_f - t_n)$ is the forecast lead-time with reference to the present time t_n and the forecasting time period T_f .

3.2.4. Step 4: Performing a regression on the local neighbourhood to obtain a forecast

The regression has been performed using the neighbourhood coordinates ($H'()$) as inputs, and their corresponding expected values (X) as outputs. The problem then can be cast as,

$$X_j(t) = H'_j(\tau, t)\beta \quad (6)$$

in which β 's are the regression coefficients. In this study, the regression order is chosen as linear (polynomial degree one) and hence, called as local linear model. Although, a local linear model makes use of a linear approximation for each separate prediction, the resulting overall model can be highly nonlinear, as each of these linear approximations are

made for each separate neighbourhood (Babovic and Keijzer, 1999). However, the local approximation can be a polynomial of any degree, which should be tested for different cases according to the dynamics of the time series.

3.2.5. Step 5: Deriving the forecast wave height at the time period T_f

Using the regression coefficients, β obtained from the regression analysis of Eq. (6), the wave height forecast, $X(T_f = t_n + T)$ could be evaluated as,

$$X(T_f) = H(\tau, t_n)\beta \quad (7)$$

The first step in making a local model forecast is to embed the time series into a phase space. This typically involves the selection of a time lag, τ and an embedding dimension, δ . An approximate assumption would give a feeling on the efficiency of the *local modelling*. There are, however, methods available for the optimal selection of these parameters. Methods for selecting prescription values for the necessary parameters using average mutual information (AMI) and false nearest neighbour (FNN) analyses are typically recommended in the literature (Abarbanel, 1996). However, these methods have been shown to generally be suboptimal selections (Babovic et al., 2000) and an alternate strategy using genetic algorithms (GA) was suggested. This approach is employed in this study as it has been shown to demonstrate significant improvements in the delineation of prescription values.

3.3. Choice of optimal time delays

The optimal embedding could be ensured from the proper selection of optimal time delays $\{\tau\}$. The underlying fact of such selection should be that each component of the vector provides new information about the signal source in a given time. The dynamical difference between the components is achieved by the evolution of the signal source over a time $\{\tau\}$. These times must be large enough so that all dynamical degrees of freedom coupled to the variables, $x(t)$, have the opportunity to influence the value of $x(t)$. At the same time, τ must be small enough so that inherent instabilities in nonlinear systems do not contaminate

the measurements at a later time, $t+T$. The embedding theorem provides assurance that, when an infinite amount of infinitely accurate data is available, any τ would work. Because the amount of data is usually finite with a finite precision, the optimal selection process has to be addressed. Smaller values of τ would not add significant new information about the dynamics and larger values of τ create uncorrelated elements in $H(\tau, t)$. The chosen τ values thus should be balanced.

A typical simplifying assumption is the creation of an embedding time-delay vector as integer multiplications of a certain elementary embedding time τ , i.e., $\tau_1=\tau$, $\tau_2=2\tau$, $\tau_3=3\tau$, etc. The optimised value of τ is the one at which either value of autocorrelation goes through first zero, or the value of average mutual information takes the first minimum. Such recipes provide robust, but in principle, suboptimal choices of embedding parameters, thus resulting in a suboptimal embedding properties as well as a suboptimal forecast skill. The concept of evolutionary embedding provides the solution of optimal embedding parameters.

4. Evolutionary embedding

The ultimate purpose of time series analysis and characterization is prediction. In this context, the establishment of optimal embedding parameters is to provide the best possible forecast skill. Good forecast skill implies good embedding properties. It is therefore prudent to investigate a scheme in which embedding parameters (δ and τ) as well as the parameter of local model (neighbourhood size k) is object of the search procedure. The selection of these parameters at a most generic level is an optimization problem. For large embedding dimension δ , the optimal choice of embedding parameters is consequently a very difficult one. In the sequel, an approach in which genetic algorithms are used to select the optimal parameters is described.

4.1. Genetic algorithms

Evolutionary algorithms are engines simulating simplified processes occurring in nature and implemented in artificial media such as a computer. The

fundamental idea is that of emulating the Darwinian theory of evolution. According to Darwin, evolution is best depicted as the process of the adaptation of species to their environment as the one of *natural selection*. Perceived in this way, all species inhabiting our planet are actually results of this process of adaptation. Evolutionary algorithms provide a similar approach for problem solving in which solutions to the problem are evolved rather than the problems being solved directly. The family of evolutionary algorithms is divided into four main streams: Evolution Strategies (Schwefel, 1981), Evolutionary Programming (Fogel et al., 1966), Genetic Algorithms (Holland, 1975) and Genetic Programming (Koza, 1992).

In the present study, Genetic Algorithm (GA) was used for optimizing local model parameters. GA is a nongradient optimization algorithm used for the search of local extremes (minimum or maximum) of functions with many variables and functional extremes. In

principle, an initial population of individuals (an individual represents a set of local model parameters which need to be optimized) is created in a computer and is allowed to evolve using the principles of inheritance (so that children resemble parents), variability (the process of child creation is not perfect due to mutations) and selection (more fit individuals are allowed to reproduce more often and vice versa so that their *genealogical* trees disappear in time). The process of mating and child creation is continued until an entirely new population is generated with the hope that strong parents will create a fitter generation of children. In practice, the average fitness of the population tends to increase with each new generation. The fitness of each of the children is determined and the process of selection/crossover/mutation is repeated. Successive generations are created until very fit individuals are obtained. The flow diagram for a simple GA is shown in Fig. 2. The main blocks of the procedure are explained as follows (Fig. 3).

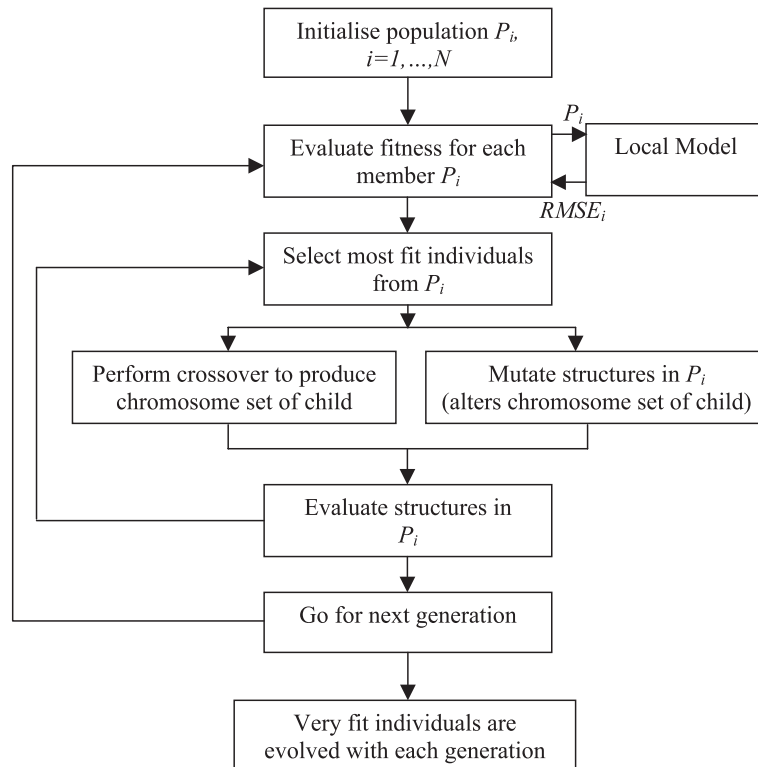


Fig. 2. Flow diagram for a typical genetic algorithm. Fitness criterion is the root-mean-square error (RMSE) of the variable say significant wave height H_s .

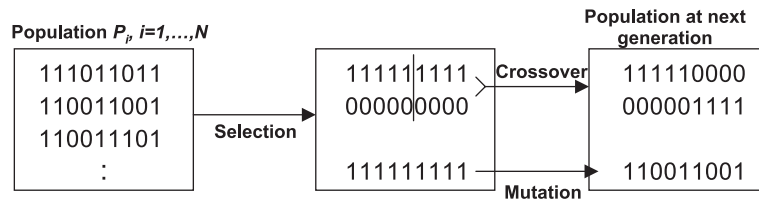


Fig. 3. Schematic illustration of a genetic algorithm. The parental genotypes are indicated as either all 1s or all 0s, for the sake of clarity.

4.1.1. Initialisation

A number of populations are randomly generated within the specified ranges of parameters. Each individual in the population, $P_i = \{\delta, \tau, k\}$, represents a set of local model parameters from the whole search space. A subcondition item (parameter represented like a gene) is represented as a bit-vector. The number of bits depends on the total number of possible variations in all the parameters. Each parameters, δ , τ and k , are prescribed within selected ranges of the order of 2. For example, if δ varies between 1 and 8 (range= 2^3) and τ varies between 1 and 32 h (range= 2^5) and k varies between 1 and 256 (range= 2^8), then the total number of bit vectors is 16.

4.1.2. Mutation

One individual is chosen from the population and it is transformed to a new individual by changing its conditions or its action. Mutation simply flips the randomly selected bit. The operation mutation may bring some new genes that do not exist in the current population.

4.1.3. Crossover

Two individuals that are not muted are selected from the population and their subconditions or actions are exchanged according to some constraints in order to generate two new sets of parameters. The newly generated individuals (children) are kept with their parents. This operation corresponds to sexual reproduction. Either one or two children in each crossover could be reproduced.

4.1.4. Evaluation

For each strategy, evaluate its benefit by applying it to a period of data.

4.1.5. Sort and selection

During mutation, the parents are replaced with new genes. However, in the process of crossover, a

new set of individuals is generated. Thus, the population will become crowded at the reproduction stage as it consists of both parents and children. Thus, all the individuals are sorted according to the fitness and the best individuals are selected from the entire population to carry forward to next generation. If the interest is of wave height forecasting, the fitness criterion may be the root-mean-square error between the observed and the predicted wave heights.

4.1.6. Termination

A possible termination criterion is the number of generations or user-defined threshold for the fitness. For every individual, the local model is executed to evaluate the fitness criterion. The member of the population set, $P_i = \{\delta, \tau, k\}$, which satisfies the prescribed threshold fitness value, is the best set of local model parameters.

One of the main advantages of GA is their domain independence. GAs can evolve almost anything, given an appropriate representation of evolving structures. Several applications of GAs in the field of water resources are described by one of the authors: Babovic et al. (1994, 1999); Babovic and Abbott (1997a,b); Babovic and Keijzer (2000); Keijzer and Babovic (1999).

5. Results and discussion

The capability of the *local model* has been demonstrated for forecasting errors in a wave prediction model applied to South China Sea domain (Fig. 4). An arbitrary observation station was chosen in the domain and the wave observation was simulated by WAM driving with real time wind records. The model prediction, so-called *first guess* value, was simulated by WAM driving with error-induced wind records. A long historical record of the observations and first

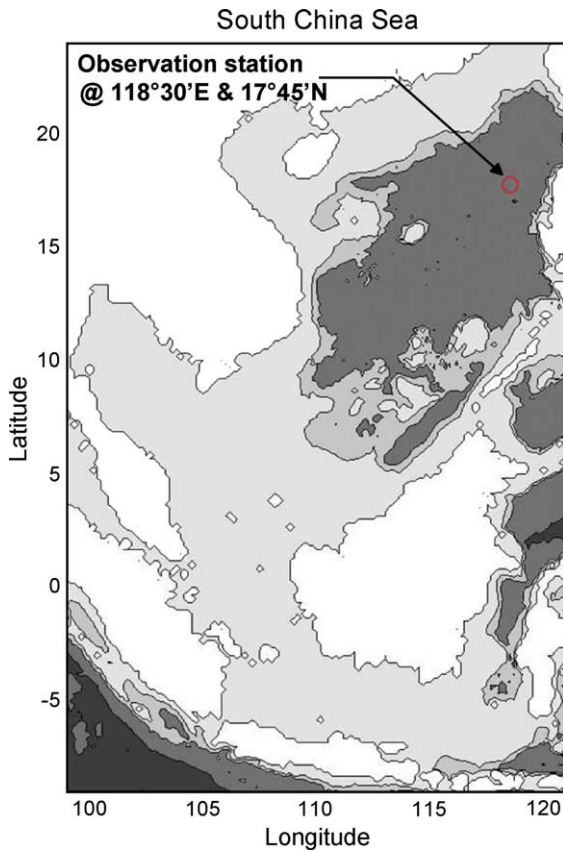


Fig. 4. South China Sea model domain.

guess values was generated initially. The local linear model was built on this historical data set to evaluate errors in a forecast horizon.

WAM was setup in a domain covering longitudes from 99°E to 121°E and latitudes from 9°S to 24°N in a grid spacing of 15' on either direction. The time period was chosen from 1st March 2001 to 31st May 2001. The analysed hindcast wind records were used to drive WAM and the resulting output wave parameters were taken as observations. These were best estimates of observations for the real time wind. The ECMWF wind data have been validated relative to TOPEX wind speed. The wave observation station was positioned at a geographical location of 118°30' E and 17°45' N and the forecasting would be carried out for this station. In principle, correction by local models can be applied to any measurement point in the domain irrespective of its physical location and its proximity to the boundary or shore.

The wave parameters such as significant wave height, mean wave period and mean wave direction were calculated at 1-h interval ($\Delta t = 1$ h) throughout the simulation period. A total of 2200 records were obtained. The time interval chosen in the present study was arbitrary to create a large data bank. A longer historical data set would be better to predict the odd events in the time series such as severe storm condition. Initial 1600 time steps ($n = 1600$ and $t_n = 1600$ h) were taken as training series and the remaining time steps as testing series. The testing series has information in the forecasting time horizon, which was used for validation purposes. The forecasting period starts from 1700 h of 6th May 2001.

To simulate model first guess values over the same period, WAM was again setup with error-induced wind records. These wind vectors were obtained from the real time wind records by multiplying wind speed at each geographical grid with a uniformly distributing factor in the range 0.75 to 1.25 and, varying the wind direction, θ_w in the band of $\pm 30^\circ$. In the domain with the given number of grid locations, the multiplication factor was uniformly distributed and the wind direction was systematically varied in the band of $\theta_w - 30^\circ$ to $\theta_w + 30^\circ$ to avoid creation of arbitrary wind directions between two adjacent grid points. The first guess wave parameters were also estimated at 1-h intervals.

The identified variables, which need to be forecasted were significant wave height H_s (or total wave energy), mean wave period T_m and mean wave direction θ_m . The errors ($\varepsilon_{H_s} = H_{so} - H_{sf}$; $\varepsilon_{T_m} = T_{mo} - T_{mf}$ and $\varepsilon_{\theta_m} = \theta_{mo} - \theta_{mf}$) were forecasted rather than the physical value. The subscript f denotes model first guess values and o denotes observations. The error forecasting led a way to effective long-term predictions. These variables (ε_{H_s} , ε_{T_m} and ε_{θ_m}) were embedded together to predict future of any individual variable. For each variable prediction at each forecast time horizon, the local model parameters (δ , τ and k) were optimized using genetic algorithm.

Following the idea of evolutionary embedding, a steady-state GA has been applied for the optimization of local model parameters, in which the evolving individuals represent embedding vector $x(t)$ ($= \varepsilon_{H_s} : \varepsilon_{T_m} : \varepsilon_{\theta_m}$). The specified ranges of parameters are presented in Table 1. In this application, three variables (ε_{H_s} , ε_{T_m} and ε_{θ_m}) were embedded simulta-

Table 1

Local modelling parameters (δ , τ and k) setting for the selection of population (P_i , $i=1 \dots N$)

Parameter	Embedding dimension, δ	Time lag, τ (hours)	Number of neighbours, k
Range	2–9	1–32	50–305
Number of possibilities (the number of bit vectors)	8	32	256

neously, and hence, the number of combinations of parameters (seven parameters in this case: 3δ , 3τ , k) would cost more than 10^5 possible permutations to find an optimized set. The genetic algorithm was thus used in this study. An initial population N was set to 100 and the number of child per parent was set to 2. The fitness criterion of each individual was the root mean square error of the variable. GA was set to propagate up to 50 generations to establish the best individual.

The optimized parameters embedding dimension (δ), the time lags (τ) and the number of nearest neighbours (k) for the variables ε_{H_s} , ε_{T_m} and ε_{θ_m} are presented in Tables 2–4, respectively. Corresponding forecast skill and improvement by the local models constructed on the basis of embedding prescription are given in Tables 5–7, as well as in Fig. 5. The nearest neighbours are neighbouring points for the forecasting point in the embedded domain. If there are only few neighbours exhibiting similar forecasting information, then, for that particular forecast, only a small number of nearest neighbours exist. The number of neighbours is thus independent of forecast lead-time. Similarly, the embedding structure (with a dimension δ and a time lag τ) is unique for every variable and also the structure is different for various forecast lead-

time. In the present application, the optimised time lag (τ) for 1-h forecast was obtained as 1-h which correspond to auto regressive model structure. Hence, the forecasting efficiency was of the same order of auto-regressive model prediction.

The time series forecasting has also been carried out using globally linear auto regressive model. The order of auto-regressive model was taken as 50 and the number of tuning points as 250. The further increase in the order of auto regressive model could not significantly increase the efficiency of the prediction. The prediction skill of the auto regressive model is presented in Tables 8–10 for H_s , T_m and θ_m , respectively. In general, it can be seen that the efficiency of the local model was better than the linear auto-regressive model. For 1-h forecast, the linear model prediction was as good as local model estimation. But as the forecast lead-time increases, the accuracy of linear model prediction deteriorates and fails to predict nonlinearity in the time series. The 24-h prediction was inferior than first guess values due to negative error corrections. However, the local model keeps better accuracy up to 24-h forecast. RMSE for significant wave height for 24-h forecast was only 50% of first guess values and the correlation with the observation was as high as 99.7% for the local model forecast.

Figs. 6–8 depict times series of observed and forecasted wave parameters. The differences in values of observed and first guess wave parameters occurred due to the systematic variation of wind direction. In general, if the change in the wind direction (for the case of first guess wind vectors) enhances the fetch availability, the first guess values overpredict the observations. The shift in the wind direction of nearly $\pm 30^\circ$ would make considerable differences in fetch lengths.

Table 2

Embedding characteristics for various forecast lead times for the prediction of significant wave height, ε_{H_s}

Forecast horizon (T ; h)	1			4			8			12			24		
Neighbours (k)	275			250			275			192			110		
Variable	H_s	T_m	θ_m	H_s	T_m	θ_m	H_s	T_m	θ_m	H_s	T_m	θ_m	H_s	T_m	θ_m
Embedding dimension (δ)	2	5	5	6	9	8	8	8	2	9	6	9	9	2	9
Time lag (τ ; h)	1	22	24	25	21	15	22	7	5	21	14	6	16	17	6

Table 3

Embedding characteristics for various forecast lead times for the prediction of mean wave period, ε_{T_m}

Forecast horizon (T ; h)	1			4			8			12			24		
Neighbours (k)	258			250			250			94			258		
Variable	H_s	T_m	θ_m	H_s	T_m	θ_m	H_s	T_m	θ_m	H_s	T_m	θ_m	H_s	T_m	θ_m
Embedding dimension (δ)	2	4	3	3	3	4	6	3	4	4	6	3	2	3	4
Time lag (τ ; h)	1	1	12	1	8	12	4	6	6	17	9	20	1	21	15

Table 4

Embedding characteristics for various forecast lead times for the prediction of mean wave direction, ε_{θ_m}

Forecast horizon (T ; h)	1			4			8			12			24		
Neighbours (k)	266			192			135			209			168		
Variable	H_s	T_m	θ_m	H_s	T_m	θ_m	H_s	T_m	θ_m	H_s	T_m	θ_m	H_s	T_m	θ_m
Embedding dimension (δ)	8	2	8	5	6	3	6	3	3	6	4	6	6	9	4
Time lag (τ ; h)	32	2	1	31	28	1	5	9	14	13	6	5	5	27	5

Table 5

Mean absolute error (MAE), root-mean-square error (RMSE) as well as correlation coefficient (γ) for significant wave height (H_s) for different forecast horizons based on local model

Forecast horizon (T ; h)	MAE (m)	RMSE (m)	γ
Model first guess	0.1363	0.1695	0.9797
1	0.0124	0.0276	0.9993
4	0.0255	0.0395	0.9987
8	0.0339	0.0646	0.9962
12	0.0435	0.0696	0.9957
24	0.0611	0.0845	0.9960

Table 6

Mean absolute error (MAE), root-mean-square error (RMSE) as well as correlation coefficient (γ) for mean wave period (T_m) for different forecast horizons based on local model

Forecast horizon (T ; h)	MAE (s)	RMSE (s)	γ
Model first guess	0.1414	0.1719	0.9940
1	0.0112	0.0195	0.9999
4	0.0172	0.0267	0.9998
8	0.0266	0.0363	0.9997
12	0.0712	0.0907	0.9972
24	0.0688	0.0932	0.9969

One-hour forecast successfully resolves rising and falling tendencies of observed variables. The observations and forecasts are fairly close to each other. This is further confirmed by the corresponding high value of correlation ($\gamma > 0.99$) for 1-h forecast. It should be noted that as the forecasting horizon increases, the mean absolute error (MAE) increases, as anticipated. The correlation coefficient is found to be marginally reduced for longer forecast horizons, but the local model improves

Table 7

Mean absolute error (MAE), root-mean-square error (RMSE) as well as correlation coefficient (γ) for mean wave direction (θ_m) for different forecast horizons based on local model

Forecast horizon (T ; h)	MAE	RMSE	γ
Model first guess	7.83°	12.35°	0.9819
1	1.07°	1.75°	0.9995
4	1.52°	2.21°	0.9993
8	1.94°	3.24°	0.9984
12	2.91°	5.59°	0.9951
24	4.27°	6.20°	0.9942

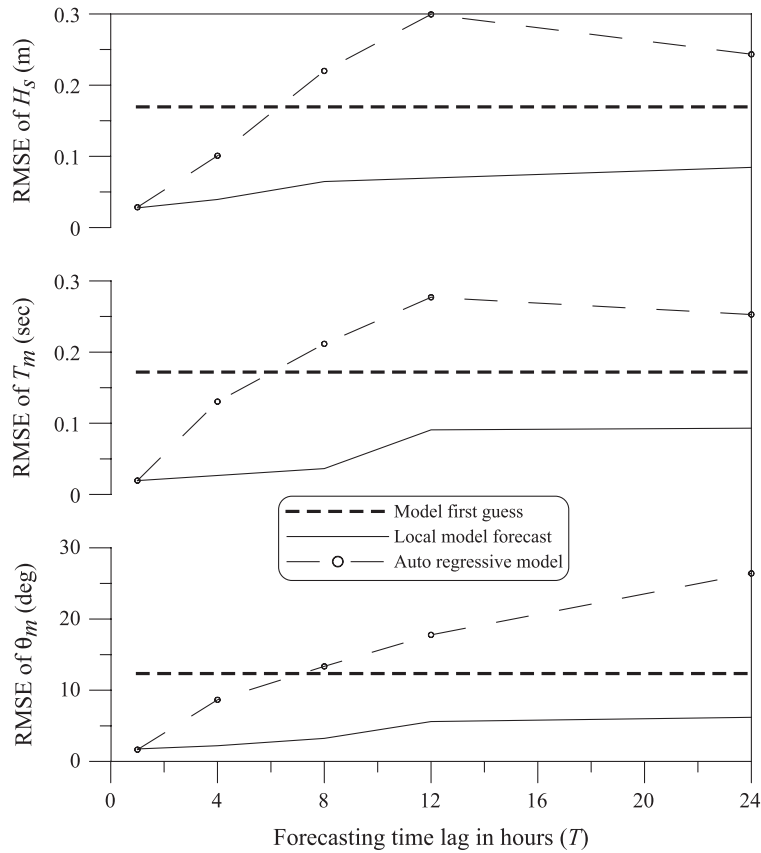


Fig. 5. Evolution of forecast error as a function of forecast horizon for the first guess numerical model and for the one with error correction based on local model and auto regressive model.

the first guess results by reducing errors more than 50% for 24-h forecast.

The linear auto-regressive model failed for forecast lead-time more than say 6-h. The prediction for

more than 6 h deteriorates the first guess values, and hence, trivial for the value addition into the system. The time series correlation is also poor compared to the first guess values for longer forecast periods.

Table 8

Mean absolute error (MAE), root-mean-square error (RMSE) as well as correlation coefficient (γ) for significant wave height (H_s) for different forecast horizons based on auto-regressive model

Forecast horizon (T ; h)	MAE (m)	RMSE (m)	γ
Model first guess	0.1363	0.1695	0.9797
1	0.0132	0.0285	0.9992
4	0.0693	0.1008	0.9909
8	0.1829	0.2200	0.9666
12	0.2569	0.2994	0.9533
24	0.2018	0.2434	0.9472

Table 9

Mean absolute error (MAE), root-mean-square error (RMSE) as well as correlation coefficient (γ) for mean wave period (T_m) for different forecast horizons based on auto-regressive model

Forecast horizon (T ; h)	MAE (s)	RMSE (s)	γ
Model first guess	0.1414	0.1719	0.9940
1	0.0129	0.0195	0.9999
4	0.1043	0.1304	0.9976
8	0.1752	0.2116	0.9927
12	0.2308	0.2770	0.9908
24	0.2040	0.2528	0.9851

Table 10

Mean absolute error (MAE), root-mean-square error (RMSE) as well as correlation coefficient (γ) for mean wave direction (θ_m) for different forecast horizons based on auto-regressive model

Forecast horizon (T ; h)	MAE	RMSE	γ
Model first guess	7.83°	12.35°	0.9819
1	0.85°	1.65°	0.9996
4	5.37°	8.65°	0.9885
8	8.22°	13.35°	0.9775
12	11.96°	17.78°	0.9611
24	18.88°	26.41°	0.9101

These clearly show the requirement of nonlinear models such as local model for the time series forecasting of wave parameters. In the present application, measurement error is ignored. However,

the fitness criterion in the selection of local model parameters can be tailored for the known measurement uncertainties.

6. Conclusions

The present paper describes an efficient error correction algorithm and its application in ocean wave prediction. The historical wave data were set-up to develop a local model, which was then used to forecast the errors of deterministic model. The results demonstrate significant increase in accuracy of a resulting hybrid model (combining a deterministic model with the stochastic local model). Such hybrid model provides the combination of the better of the two worlds: high-quality forecast skill based on utilisation of data, which can also be extended

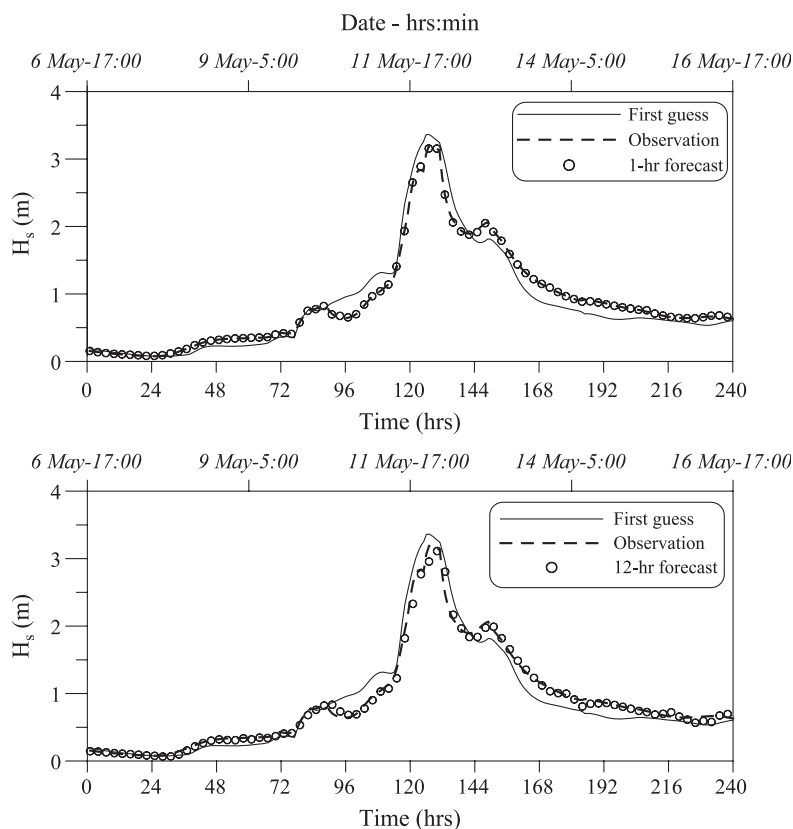


Fig. 6. Time series of first guess and observed significant wave heights, as well as the wave heights after the correction by local model. Forecast horizon: 1 and 12 h.

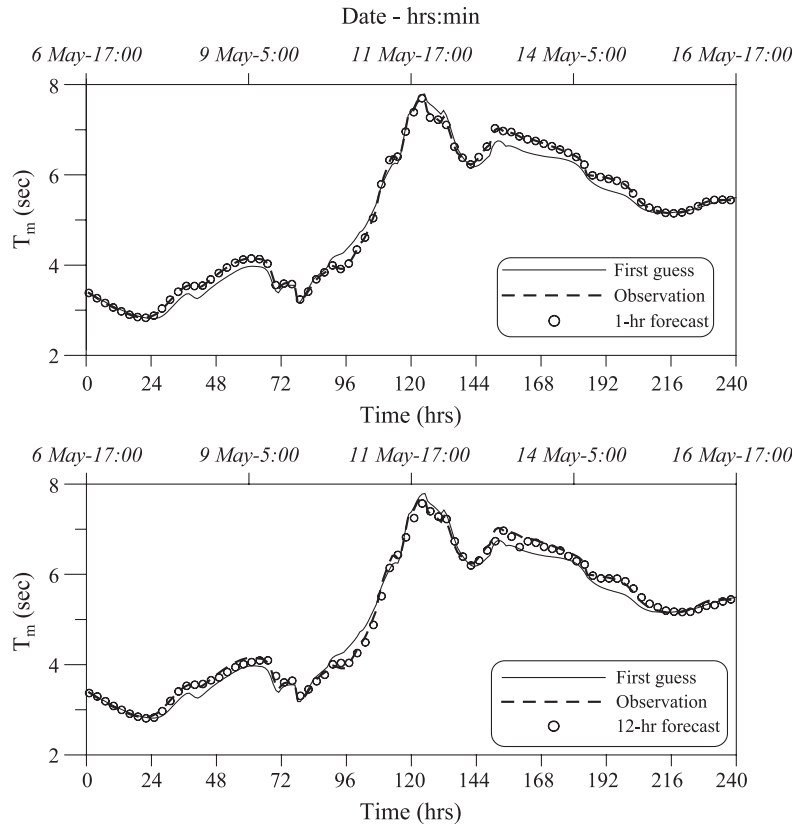


Fig. 7. Time series of first guess and observed mean wave periods as well as the periods after the correction by local model. Forecast horizon: 1 and 12 h.

for longer lead-time due to the use of the deterministic model. The described method is based on computationally inexpensive approaches which are fast to execute and yet with superior accuracies. The local model is built over the concept of embedding theorem and the embedding parameters are optimized using genetic algorithm.

The forecasting of mean wave parameters has been carried out using the local model and using a linear auto-regressive model. For shorter forecast lead-time, up to 6 h, both the time series prediction models performed well. However, for longer forecasting period, the linear model failed to improve the first guess values and the nonlinear models such as local model performed well up to 24-h forecast. The improvement over first guess values was shown up to 50% in the reduction of RMSE for significant wave height.

List of symbols

Δt	time step (h)
δ	embedding dimension
ε_{xi}	difference of observation and first guess value of the variable x_i
γ	correlation coefficient
θ_m	mean wave direction ($^\circ$)
τ	time lag (h)
H_s	significant wave height (m)
k	number of nearest neighbours
n	number of records in the historical data set
P	population ($=\{\delta_i, \tau_i, k\}$)
T	forecast lead-time (h)
T_f	time ($t=T_f$) at which forecast is required (h)
T_m	mean wave period (s)
$x(t)$	time series of the variable
MAE	mean absolute error
RMSE	root-mean-square error

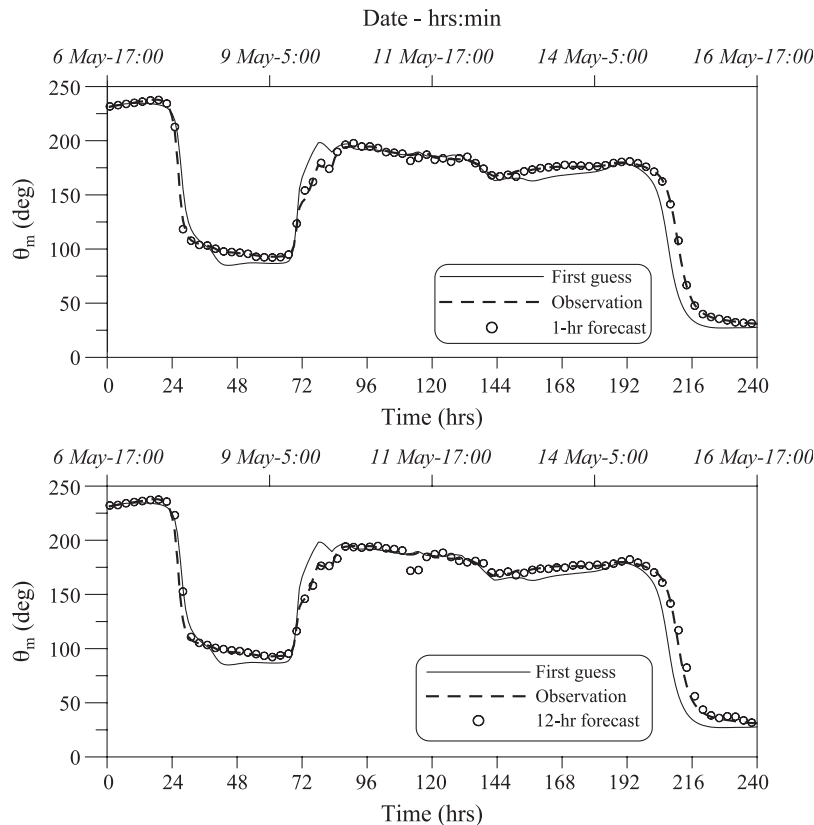


Fig. 8. Time series of first guess and observed mean wave directions as well as the direction after the correction by local model. Forecast horizon: 1 and 12 h.

References

- Abarbanel, H.D.I., 1996. Analysis of Observed Chaotic Data. Springer-Verlag, New York.
- Babovic, V., 1996. Emergence, Evolution, Intelligence: Hydroinformatics. Balkema, Rotterdam.
- Babovic, V., Abbott, M.B., 1997a. Evolution of equation from hydraulic data. Part I: theory. *Journal of Hydraulic Research* 35 (3), 1–14.
- Babovic, V., Abbott, M.B., 1997b. Evolution of equation from hydraulic data. Part II: application. *Journal of Hydraulic Research* 35 (3), 15–34.
- Babovic, V., Keijzer, M., 1999. Forecasting of river discharges in the presence of chaos and noise. In: Marsalek, Jiri (Ed.), *Coping with Floods: Lessons Learned from Recent Experiences*, NATO ARW Series. Kluwer, Dordrecht, pp. 405–420.
- Babovic, V., Keijzer, M., 2000. Genetic programming as a model induction engine. *Journal of Hydroinformatics* 2 (1), 35–60.
- Babovic, V., Keijzer, M., 2003. Rainfall runoff modeling based on genetic programming. *Nordic Hydrology* 34, 1.
- Babovic, V., Keijzer, M., Stefansson, M., 2000. Optimal embedding using evolutionary algorithms. *Proceedings of the Fourth International Conference on Hydroinformatics*, Iowa City, Iowa, USA, 8 pp.
- Babovic, V., Larsen, L.C., Wu, Z., 1994. Calibrating hydrodynamic models by means of simulated evolution. *Proceedings of the First International Conference on Hydroinformatics*. Balkema, Rotterdam, pp. 193–200.
- Babovic, V., Keijzer, M., Mahbub, R., 1999. Analysis and prediction of chaotic time series. D2K Technical report 0399-2, Danish Hydraulic institute, <http://www.d2k.dk>.
- Bidlot, J.-R., Holt, M.W., 1999. Numerical wave modeling at operational weather centres. *Coastal Engineering* 17, 409–429.
- Christopoulos, S., 1997. Wind-wave modeling aspects within complicate topography. *Annales Geophysicae* 15, 1340–1353.
- Fogel, L.J., Owens, A.J., Walsh, M.J., 1966. Artificial Intelligence through Simulated Evolution. Ginn, Needham Height.
- Gelb, A., 1974. Applied Optimal Estimation. MIT press, Cambridge.
- Gershenfeld, N.A., Weigend, A.S., 1993. In *Time Series Prediction: Forecasting the Future and Understanding the Past*. Addison Wesley, Reading, MA.
- Hasselmann, S., Hasselmann, K., 1985. Computations and parametrizations of the nonlinear energy transfer in a gravity-wave spectrum. Part I: a new method for efficient computations of the

- exact nonlinear transfer integral. *Journal of Physical Oceanography* 15, 1369–1377.
- Hasselmann, K., Barnett, T.P., Bouws, E., Carlson, H., Cartwright, D.E., Enke, K., Ewing, J.A., Gienapp, H., Hasselmann, D.E., Kruseman, P., Meerburg, A., Muller, P., Olbers, D.J., Richter, K., Sell, W., Walden, H., 1973. Measurements of windwave growth and swell decay during the Joint North Sea Wave Project (JONSWAP). *Deutsche Hydrographische Zeitschrift. Ergänzungsheft. Reihe A* 8 (12), 95 pp.
- Henk, P.C., Boogard, Van den, Kruisbrinck, A.C.H., 1996. Hybrid modeling by integrating neural networks and numerical models. *Proceedings of the Second International Conference on Hydroinformatics vol. 2*. Balkema, Rotterdam, pp. 471–478.
- Holland, J.H., 1975. *Adaptation in Natural and Artificial Systems*. University of Michigan, Ann Arbor.
- Janssen, P.A.E.M., 1989. Wave induced stress and the drag of air over sea waves. *Journal of Physical Oceanography* 19, 745–754.
- Keijzer, M., Babovic, V., 1999. Error correction of a deterministic model in Venice lagoon by local linear models. *Proc. 'Modeli Complessie Metodi Computazionali Intensivi per la Stima e la Previsione' Conference*, Venice.
- Khu, S.T., Liong, S.Y., Babovic, V., Madsen, H., Muttill, N., 2001. Genetic programming and its application in real-time runoff forecasting. *Journal of the American Water Resources Association* 37 (2), 439–451.
- Komen, G.J., Hasselmann, S., Hasselmann, K., 1984. On the existence of a fully developed windsea spectrum. *Journal of Physical Oceanography* 14, 1271–1285.
- Komen, G.J., Cavaleri, L., Donelan, M., Hasselmann, K., Hasselmann, S., Janssen, P.A.E.M., 1994. *Dynamics and Modeling of Ocean Waves*. Cambridge Univ. Press, New York.
- Koza, J.R., 1992. *Genetic Programming: On the Programming of Computers by Means of Natural Selection*. MIT press, Cambridge, MA, USA.
- Miles, J.W., 1957. On the generation of surface waves by shear flows. *Journal of Fluid Mechanics* 3, 185–204.
- Minns, A.W., 1998. *Artificial neural networks on subsymbolic process descriptors*. PhD thesis, Balkema, Rotterdam.
- Packard, N.H., Crutchfield, J.P., Farmer, J.D., Shaw, R.S., 1980. Geometry from a time series. *Physical Review Letters* 45, 712–716.
- Philips, O.M., 1957. On the generation of waves by turbulent wind. *Journal of Fluid Mechanics* 2, 417–445.
- Philips, O.M., 1958. The equilibrium range in the spectrum of wind-generated water waves. *Journal of Fluid Mechanics* 4, 426–434.
- Refsgaard, J.C., 1997. Validation and intercomparison of different updating procedures for real-time forecasting. *Nordic Hydrology* 28, 65–84.
- Sauer, T., Yorke, J.A., 1991. Rigorous verification of trajectories for the computer simulation of dynamical systems. *Nonlinearity* 4, 961–979.
- Schwefel, H.-P., 1981. *Numerical Optimization of Computer Models*. Wiley, Chichester.
- Snyder, R.J., Dobson, F.W., Elliott, J.A., Long, R.B., 1981. Array measurements of atmospheric pressure fluctuations above surface gravity waves. *Journal of Fluid Mechanics* 102, 1–59.
- Takens, F., 1980. Detecting strange attractors in turbulence. In: Rand, D.A., Young, L.-S. *Lecture Notes in Mathematics: Dynamical Systems and Turbulence*, vol. 898. Springer-Verlag, Berlin, pp. 366–381.
- The WAMDI group, 1988. The WAM model—a third generation ocean wave prediction model. *Journal of Physical Oceanography* 18, 1775–1810.
- Tolman, H.L., 1999. User manual and system documentation of wavewatch-III version 1.18. Technical Note, vol. 16. National Oceanic and Atmospheric Administration, Washington, DC.
- WMO, 1992. Simulated real-time intercomparison of hydrological models. WMO operational hydrology report 38-WMO No. 779, World Meteorological Organisation, Geneva.