

Differential radial basis function network for sequence modelling

Kojo Sarfo Gyamfi^{1,*}, James Brusey, Elena Gaura

*School of Computing, Electronics and Mathematics, Coventry University, CV1 5FB,
Coventry, United Kingdom*

Abstract

We propose a differential radial basis function (RBF) network termed RBF-DiffNet—whose hidden layer blocks are partial differential equations (PDEs) linear in terms of the RBF—to make the baseline RBF network robust to noise in sequential data. Assuming that the sequential data derives from the discretisation of the solution to an underlying PDE, the differential RBF network learns constant linear coefficients of the PDE, consequently regularising the RBF network by following modified backward-Euler updates. We experimentally validate the differential RBF network on the logistic map chaotic timeseries as well as on 30 real-world timeseries provided by Walmart in the M5 forecasting competition. The proposed model is compared with the normalised and unnormalised RBF networks, ARIMA, and ensembles of multilayer perceptrons (MLPs) and recurrent networks with long short-term memory (LSTM) blocks. From the experimental results, RBF-DiffNet consistently shows a marked reduction over the baseline RBF network in terms of the prediction error (e.g., 26% reduction in the root mean squared scaled error on the M5 dataset); RBF-DiffNet also shows a comparable performance to the LSTM ensemble at less than one-sixteenth the LSTM computational time. Our proposed network consequently enables more accurate predic-

*Corresponding author

Email addresses: `kojo.gyamfi@loblaw.ca` (Kojo Sarfo Gyamfi),
`james.brusey@coventry.ac.uk` (James Brusey), `elena.gaura@coventry.ac.uk`
(Elena Gaura)

¹Present address:

Data Insights and Analytics, Loblaw Companies Limited, 1 President's Choice Circle,
Brampton, ON L6Y 5S5, Canada

tions—in the presence of observational noise—in sequence modelling tasks such as timeseries forecasting that leverage the model interpretability, fast training, and function approximation properties of the RBF network.

Keywords: Radial basis function, neural network, sequence modelling

1. Introduction

The radial basis function network (RBFN) is an artificial neural network first introduced by Broomhead and Lowe in the 1980s [1] but still very much in vogue now [2, 3, 4, 5] due to its robustness as a universal function approximator. Architecturally, it is a three-layer network having one input layer, one hidden layer, and one output layer, as shown in Figure 1, with activation units in the hidden layer made up of radial basis functions (RBFs).

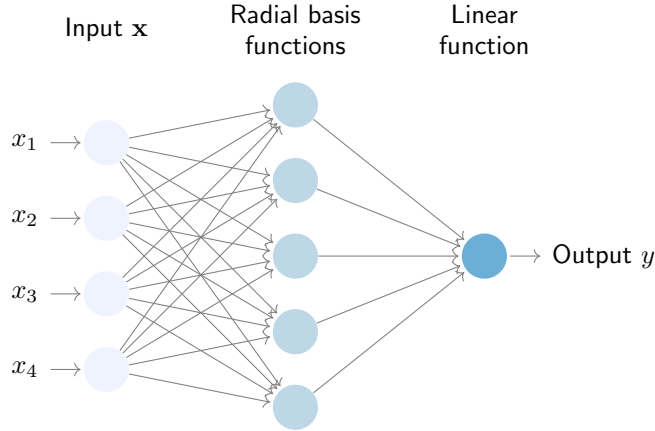


Figure 1: Radial basis function network (RBFN)

The input layer is often connected to the hidden layer via direct connections whose weights are frozen at unity and thus not trainable, while the output layer is a simple linear layer. Though simple, this architecture is particularly efficient as a universal function approximator [6, 7, 8, 9]; one reason for this is that the hidden layer of the RBF network performs a similar kernel transformation to an infinite-dimensional inner product space employed by other kernel machines such as the support vector machine [10, 2]. The RBF network has therefore seen many uses especially in timeseries forecasting, control and classification problems that occur in many real-world applications

such as fraud detection, speech recognition, manufacturing, medical diagnosis, and face recognition [11, 12]. One other area in which the radial basis function network has seen increasing adoption is in the linear approximation of the value function in reinforcement learning in terms of the state-action variables [13, 14, 15].

One reason for the ubiquity of RBF networks in many machine learning tasks is the interpretability of their outputs, since the hidden layer essentially performs a fuzzy nearest-neighbour association of an input vector to a set of well-defined exemplars in the training data; thus, for a given input, the influence of different features on the output can be estimated from the relative importance of the features in these exemplars. Another reason for the sustained use of the RBF network is the speed in training the network, since only the final linear layer is often trained; for the least-squares error (with ridge regularisation), there is, in fact, a closed-form solution for the network weights. This comes at the expense of a usually unsupervised step of selecting the number, centres and widths of the radial basis function activation units. Thus, the eventual performance of the RBF network is highly susceptible to the widths and centre locations of the RBF units [16, 17, 18, 19], which in turn can be heavily influenced by the presence of noise in the data [4, 3]. Other neural network architectures may not be so exceptionally sensitive to noise. Several other approaches mainly using forward selection [20, 21, 22] and sophisticated clustering methodologies or the self-organising map [11, 23, 24, 16] have thus been proposed to better locate the RBF centres.

Crucially however, for sequence modelling tasks such as timeseries forecasting, the sequence data has to be reshaped such that the RBF network takes as inputs the l lagged values of the sequence (in addition to any exogenous inputs) and outputs some next points in the sequence; therefore, if there is a single noisy observation in the sequence, this observation likely gets replicated l times in the reshaped data that is input to the RBF network. In the extreme case, the sequence data may be so corrupted that the data contains no more meaningful information for prediction [3]. This profoundly degrades the optimal placement of the RBF centres and consequently the performance of the network. It is worth noting that this problem of noise propagation, as described, is not prevalent in the application of the RBF network to other regression or classification tasks where there are no temporal correlations in the data, in which case each noisy observation occurs only once in training. For example, in reinforcement learning settings, since

the state transition in the sequential decision process is typically considered Markovian, there is functional dependence on only the last state vector, i.e. $l = 1$, and thus a single noisy observation does not replicate itself in the reshaped training data.

In this paper, we propose a novel neural network architecture, termed the differential RBF network (RBF-DiffNet), that is designed to learn a representation of the sequence data that ignores signal noise. By utilising activation blocks made up of partial differential equations (PDEs) linear in terms of the radial basis function—with the constant linear coefficients of the PDE being trainable—we introduce a regularisation based on backward Euler updates that makes the network robust to noise. The intuition behind this contribution stems from the observation that many real-world sequence data derive from phenomena (such as in biology, economics or physics) whose underlying dynamics, in the absence of noise or control inputs, may be modelled by a set of partial differential equations, however complex. For example, the sequence of air temperature data observed in a car cabin may very well be described by a set of heat balance equations [25], which are PDEs.

Our main contribution in this paper is therefore as follows: we introduce the differential RBF network and analyse its mathematical properties that enable the network to show robustness to noisy perturbations to sequential data in applications such as timeseries forecasting, where the baseline RBF network would be rather susceptible to the noisy training inputs.

This proposed network is detailed in Section 3. Section 4 presents an experimental validation of the proposed architecture on the logistic map chaotic timeseries [26, 27] and 30 different real-world retail timeseries from the M5 competition [28]; this section also includes performance comparisons with the baseline RBF network, autoregressive integrated moving average (ARIMA) model, ensembles of MLPs and recurrent neural networks with long-short-term memory (LSTM) blocks [28]. Conclusions are given in Section 5, while the problem statement and related work are presented in the next section.

2. Problem statement and related work

We begin by considering a set of input-output pairs $\{\mathbf{x}_n, y_n\}_{1:N}$ from which we wish to train a function $f : \mathcal{X} \rightarrow \mathcal{Y}$, where $\mathbf{x}_n \in \mathcal{X}$ is a d -dimensional input vector and $y_n \in \mathcal{Y}$ is a scalar-valued output. For a sequence data $\{s_\pi\}_{1:T}$ of length T , we assume that this is reshaped into

input-output pairs $\{\mathbf{x}_n, y_n\}_{1:N}$ as before, where \mathbf{x}_n is given by the last l realisations of the sequence prior to timestep $t_{\pi+1}$, i.e., $\mathbf{x}_n^\top = [s_{\pi-l+1}, \dots, s_\pi]$, and $s_{\pi+1} = f(\mathbf{x}_n)$ where y_n is the true output. If there are other exogenous inputs \mathbf{e} on which the sequence $\{s_\pi\}_{1:T}$ has dependence, we assume this is again captured in \mathbf{x}_n as $\mathbf{x}_n^\top = [\mathbf{e}^\top, s_{\pi-l+1}, \dots, s_\pi]$. Note that if there are exogenous inputs, then $l < d$; otherwise $l = d$.

The radial basis function network (RBFN) as shown in fig. 1, can be written concisely as:

$$f(\mathbf{x}_n) = \sum_{j=1}^c w_j \phi_j(\mathbf{x}_n) + w_0, \quad (1)$$

where c is the number of RBF centres, corresponding to the number of neurons in the hidden layer, and

$$\phi_j(\mathbf{x}_n) = e^{-\beta_j \|\mathbf{x}_n - \boldsymbol{\mu}_j\|^2} \quad (2)$$

is the Gaussian radial basis function, which is the RBF we consider in this paper; the choice of non-linearity in other radial basis functions has been found not to be too crucial to the performance of the network [29]. In (2), $\boldsymbol{\mu}_j$ are exemplars in the training data; β_j , as defined in (10), is inversely proportional to the width of the j th RBF, and w_j, w_0 are the RBF weights and bias to be optimised.

From (1) and (2), the hidden layer of the RBF network essentially performs a fuzzy nearest-neighbour association of the input vector \mathbf{x}_n to the set of exemplars $\boldsymbol{\mu}_j$, based on a Mahalanobis distance criterion with spherical covariance in terms of β_j . Thus, for any given input vector, the influence of different features on the output can be estimated from the relative importance of the features in each of the c exemplars, each weighted by its corresponding RBF weight w_j ; this property makes the output of the RBF network interpretable in terms of its inputs.

For the sake of brevity, we drop the subscript n in \mathbf{x}_n in the following. In matrix form, (1) is equivalent to:

$$f(\mathbf{x}) = \boldsymbol{\phi}(\mathbf{x})^\top \tilde{\mathbf{w}}, \quad (3)$$

where,

$$\boldsymbol{\phi}(\mathbf{x}) = [1, \phi_1(\mathbf{x}), \dots, \phi_c(\mathbf{x})]^\top, \quad (4)$$

and

$$\tilde{\mathbf{w}} = [w_0, w_1, \dots, w_c]^\top \quad (5)$$

If one considers the entire training set, then we have the following system of equations:

$$\mathbf{f} = \Phi \tilde{\mathbf{w}}, \quad (6)$$

where

$$\mathbf{f} = [f(\mathbf{x}_1), f(\mathbf{x}_2), \dots, f(\mathbf{x}_N)]^\top, \quad (7)$$

and

$$\Phi = \begin{bmatrix} 1 & \phi_1(\mathbf{x}_1) & \cdots & \phi_c(\mathbf{x}_1) \\ \vdots & \ddots & \ddots & \vdots \\ 1 & \phi_1(\mathbf{x}_N) & \cdots & \phi_c(\mathbf{x}_N) \end{bmatrix} \quad (8)$$

For the least-squares error and low to moderate number of RBF nodes, $\tilde{\mathbf{w}}$ can be optimised in closed-form as:

$$\tilde{\mathbf{w}}^* = (\Phi^\top \Phi + \gamma \mathbf{I})^{-1} \Phi \mathbf{y}, \quad (9)$$

where γ is a regularisation coefficient, \mathbf{I} is the identity matrix of size $c+1$, and $\mathbf{y}^\top = [y_1, \dots, y_n]$. In general, the weights $\tilde{\mathbf{w}}$ can be trained for any arbitrary loss function using different optimisation routines.

Separate from the RBF network weights are the hyperparameters c , $\boldsymbol{\mu}_j$ and β_j that require tuning. Most commonly, these are obtained from an unsupervised preprocessing step; the number of RBF centres is fixed and the exemplars $\boldsymbol{\mu}_j$ are determined as the cluster centres obtained from some variants of K-Means clustering [16, 30, 2, 3], while β_j is derived from the compactness of the individual clusters. Specifically, β_j is defined as:

$$\beta_j = \frac{1}{2\sigma_j^2}, \quad (10)$$

where popular choices of σ_j [31, 32, 33, 34] include:

$$\sigma_j := \frac{d_{max}}{\sqrt{2c}}, \quad \forall j \in \{1, \dots, c\} \quad (11)$$

$$\sigma_j := \frac{1}{n_j} \sum_{\mathbf{x} \in \mathcal{C}_j} \|\mathbf{x} - \boldsymbol{\mu}_j\|, \quad \forall j \in \{1, \dots, c\} \quad (12)$$

$$\sigma_j := \frac{1}{r} \sum_{\mathbf{x} \in \mathcal{R}_j} \|\mathbf{x} - \boldsymbol{\mu}_j\|, \quad \forall j \in \{1, \dots, c\} \quad (13)$$

$$\sigma_j := \frac{1}{cr} \sum_{j=1}^c \sum_{\mathbf{x} \in \mathcal{R}_j} \|\mathbf{x} - \boldsymbol{\mu}_j\|, \quad \forall j \in \{1, \dots, c\}, \quad (14)$$

where d_{max} is the maximum distance between the centres of any two clusters, \mathcal{C}_j is the set of all points in the j th cluster, n_j is the number of points in the j th cluster, and \mathcal{R}_j is the set of the r closest points to the j th cluster centre. An alternative to K-Means clustering employed for RBF networks is the self-organising map which allows for a more intuitive determination of the number of cluster centres c when the data is projected onto two or three dimensions [35, 24].

Since the RBF network parameters affect the network’s performance quite significantly, other approaches utilising other performance metrics have been proposed for the selection of the hyperparameters. For example, the RBF nodes may be incrementally added in order to maximise the Fisher class-separability ratio and the leave-one-out cross validation RMSE for classification and regression tasks respectively [20, 21]. Nevertheless, the presence of observational noise tends to influence the optimal placement of centres [4] or any metrics computed based on them, such as the Fisher class-separability ratio. The prevalence of observational noise tends to have an even more profound effect for sequence data [3], where upon reshaping into input-output pairs, noisy observations get replicated throughout the data if the lookback window l is much greater than 1.

One way to make the RBF network robust is to normalise the network to achieve the so-called *partition of unity* property where the c normalised radial basis functions sum up to one for every point in the input space \mathcal{X} , making the RBF network less susceptible to noisy observations in arbitrary regions of the input space [36, 37]. Normalisation thus results in the RBF network losing its local characteristics, but often results in improved generalisability

[38]. The normalised RBF network f_{norm} is given mathematically as:

$$f_{norm}(\mathbf{x}_n) = \frac{\sum_{j=1}^c w_j \phi_j(\mathbf{x}_n)}{\sum_{j=1}^c \phi_j(\mathbf{x}_n)} + w_0, \quad (15)$$

Since normalising the RBF network has several known side effects such as shifts in the maxima of the RBFs [36, 37], in this paper, we take a different perspective towards making the RBF network robust specifically for its application to modelling sequential data. This involves utilising activation blocks made up of partial differential equations linear in terms of the radial basis function and based on backward Euler discretisation of the sequence data.

It is worth mentioning that our proposed differential RBF network architecture, although it utilises a similar Euler discretisation as the neural ordinary differential equation (ODE-Net) [39], differs from the ODE-Net as follows: while the ODE-Net parameterises the derivatives of the hidden state of a residual network with another neural network in the limit of a large number of hidden layers when the discrete step size approaches 0, the differential RBF network parameterises the solution to the differential equation with an RBF network and directly evaluates the derivatives of this solution according to the backward Euler updates given in Section 3 in such a way as to regularise the original RBF network.

Furthermore, our work differs from other works that have employed the RBF network in solving ordinary or partial differential equations [40, 41, 42, 43, 44] in that, instead of a well-defined set of differential equations to solve, we have only some discrete realisations from phenomena that are assumed to be described by unknown differential equations. Thus, we start with a solution to an unknown differential equation in the form of the RBF network and work backwards to find an optimal differential equation with constant coefficients that best describes the sequence data, based on backward Euler updates.

3. Proposed differential RBF network

3.1. Intuition

We first consider the univariate series $\{s_\pi\}_{1:T}$, and assume that this sequence data is generated by some underlying differential equation given by:

$$\frac{dz}{dt} = g(t), \quad (16)$$

that is, the sequence $\{s_t\}_{1:T}$ comprises discrete realisations of this underlying process. These discrete realisations may be approximated by backward Euler updates [45] of the form:

$$s_{\pi+1} = s_\pi + hg(t_{\pi+1}, s_{\pi+1}), \quad (17)$$

where h is some small interval between the occurrences of s_π and $s_{\pi+1}$, i.e., $h = t_{\pi+1} - t_\pi$. Note that in a noiseless system, $s_\pi = z(t_\pi)$. The backward Euler update in (17) can be thought of as a first-order Taylor's approximation, and thus to improve this approximation and reduce the truncation errors, we may consider a higher-order Taylor's expansion up to degree ν as follows:

$$s_{\pi+1} = s_\pi + h z'(t_{\pi+1}, s_{\pi+1}) + \dots + \frac{h^\nu}{\nu!} z^{(\nu)}(t_{\pi+1}, s_{\pi+1}). \quad (18)$$

However, since we wish to predict the next state of the sequence based on its prior values, we replace the functional dependence of z on time $t_{\pi+1}$ and $s_{\pi+1}$ with the last l realisations of the series, obtaining an analogous form, using multi-index notation as:

$$\begin{aligned} s_{\pi+1} &= \boldsymbol{\lambda}^\top \mathbf{s}_{\pi-l+1:\pi} + \mathbf{h} D z(s_{\pi-l+1}, \dots, s_\pi) + \frac{\mathbf{h}^2}{2} D^2 z(s_{\pi-l+1}, \dots, s_\pi) + \dots \\ &+ \frac{\mathbf{h}^\nu}{\nu!} D^{(\nu)} z(s_{\pi-l+1}, \dots, s_\pi), \end{aligned} \quad (19)$$

where \mathbf{h} is now a small vector interval, $D^i z$ is an i th-order tensor, with Dz and $D^2 z$ being the gradient and Hessian respectively of z , and $\boldsymbol{\lambda}^\top \mathbf{s}_{\pi-l+1:\pi}$ is the weighted sum of the last l realisations of the sequence, with $\boldsymbol{\lambda}$ being the vector of weights.

We may then generalise (19) to consider a multivariate series as follows:

$$s_{\pi+1} = \boldsymbol{\lambda}^\top \mathbf{s}_{\pi-l+1:\pi} + \mathbf{h} D z(\mathbf{x}) + \frac{\mathbf{h}^2}{2} D^2 z(\mathbf{x}) + \dots + \frac{\mathbf{h}^\nu}{\nu!} D^{(\nu)} z(\mathbf{x}). \quad (20)$$

Here, as mentioned previously in Section 2, \mathbf{x} encapsulates the lagged inputs of the series as well as exogenous inputs, and $s_{\pi+1} = z(\mathbf{x}) \in \mathcal{Y}$. For the special case where $\nu = 2$, since Dz and D^2z are respectively the gradient and Hessian of z , we have from (20) that,

$$s_{\pi+1} = \boldsymbol{\lambda}^\top \mathbf{s}_{\pi-l+1:\pi} + \sum_{i=1}^d h_i \frac{\partial z(\mathbf{x})}{\partial x_i} + \sum_{i=1}^d \sum_{p=1}^d h_i h_p \frac{\partial^2 z(\mathbf{x})}{\partial x_i \partial x_p}, \quad (21)$$

where h_i, h_p are respectively the i th and p th components of \mathbf{h} .

To keep the computation in (20) tractable, we ignore all mixed derivatives (such as $\frac{\partial^2 z(\mathbf{x})}{\partial x_i \partial x_p}$ in the summation in (21) involving the second-order tensor, where $i \neq p$) thus approximating (20) with the following partial differential equation:

$$s_{\pi+1} \approx \boldsymbol{\lambda}^\top \mathbf{s}_{\pi-l+1:\pi} + \sum_{k=1}^{\nu} \sum_{i=1}^d a_{k,i} \frac{\partial^{(k)} z(\mathbf{x})}{\partial x_i^k}, \quad (22)$$

where the coefficients $a_{k,i}$ (which replace \mathbf{h} in (20)) now have to be optimised to maximise the fit to the data.

Since $z : \mathcal{X} \rightarrow \mathcal{Y}$, where $\mathbf{x} \in \mathcal{X}$ and $s_{\pi+1} = z(\mathbf{x}) \in \mathcal{Y}$, a good choice for z is the radial basis function network f defined in (1), i.e.,

$$z(\mathbf{x}) := f(\mathbf{x}) = \sum_{j=1}^c w_j \phi_j(\mathbf{x}) + w_0. \quad (23)$$

With z defined as above, we have from (22) that:

$$f(\mathbf{x}) = \boldsymbol{\lambda}^\top \mathbf{s}_{\pi-l+1:\pi} + \sum_{k=1}^{\nu} \sum_{i=1}^d a_{k,i} \frac{\partial^{(k)} f(\mathbf{x})}{\partial x_i^k}. \quad (24)$$

The utility in defining the function z as a radial basis function network f is in its parameter efficiency, i.e., we are able to obtain closed-form expressions for its higher-order derivatives without any increase in the number of variables

that parameterise z or these higher derivatives.

The relationship in (24) then represents a regularisation of the parameters of f in (1). Specifically, while the radial basis function network f is ordinarily given by (1), we are now subjecting it to the additional constraint that the next observation $s_{\pi+1} = f(\mathbf{x})$ is such that it satisfies the multivariate, higher-order extension to the backward Euler updates given by (22), assuming that the discrete realisations of the sequence are driven by an underlying differential equation; this can be expressed by the following optimisation problem:

$$\begin{aligned} \min \sum_n L(f(\mathbf{x}_n), y_n), \quad \text{where} \quad f(\mathbf{x}_n) &= \sum_{j=1}^c w_j \phi_j(\mathbf{x}_n) + w_0 \\ \text{subject to :} \\ s_{\pi+1} = f(\mathbf{x}_n) &= \boldsymbol{\lambda}^\top \mathbf{s}_{\pi-l+1:\pi} + \sum_{k=1}^{\nu} \sum_{i=1}^d a_{k,i} \frac{\partial^{(k)} f(\mathbf{x}_n)}{\partial x_i^k}, \end{aligned} \quad (25)$$

where L is an arbitrary loss function. This regularisation thus makes the proposed network relatively more robust to noisy observations which affect the choice of RBF centres and widths, causing the vanilla RBF network to underperform in the presence of noise. Consequently, while fitting the data to the differential equation given by (24), we are implicitly only learning a regularised version of the parameters w_j , i.e., regularised by the equality constraint in (25).

3.2. Higher-order derivatives of the RBF

The fundamental idea in our proposed approach is fitting the data to the backward-Euler-regularised RBF network given by the differential equation in (24), rather than training the RBF network to approximate f directly as a function of \mathbf{x} as in (1). The form of (22) thus requires us to obtain the expressions for the first and higher-order derivatives of the radial basis function network.

Deriving from (23), the first-order partial derivative of f is given by:

$$\frac{\partial f(\mathbf{x})}{\partial x_i} = \sum_{j=1}^c w_j \frac{\partial \phi_j(\mathbf{x})}{\partial x_i}, \quad (26)$$

where

$$\frac{\partial \phi_j(\mathbf{x})}{\partial x_i} = -2\beta_j(x_i - \mu_{j,i})e^{-\beta_j\|\mathbf{x}-\boldsymbol{\mu}_j\|^2} \quad (27)$$

Accordingly, the second-order partial derivative of f (ignoring mixed derivatives) is given by:

$$\frac{\partial^2 f(\mathbf{x})}{\partial x_i^2} = \sum_{j=1}^c w_j \frac{\partial^2 \phi_j(\mathbf{x})}{\partial x_i^2}, \quad (28)$$

where

$$\frac{\partial^2 \phi_j(\mathbf{x})}{\partial x_i^2} = 2\beta_j[2\beta_j(x_i - \mu_{j,i})^2 - 1]e^{-\beta_j\|\mathbf{x}-\boldsymbol{\mu}_j\|^2}. \quad (29)$$

This generalises to a higher-order ν as follows:

$$\frac{\partial^{(\nu)} z(\mathbf{x})}{\partial x_i^\nu} = \sum_{j=1}^c w_j \frac{\partial^{(\nu)} \phi_j(\mathbf{x})}{\partial x_i^\nu}, \quad (30)$$

Similar results have been obtained for the Gaussian RBF [46] and the multiquadric radial basis functions [42].

Due to the form of the Gaussian radial basis function in (2), $\phi_j(\mathbf{x})$ is infinitely differentiable. In the following, we derive the formula for finding the k th derivative of $\phi_j(\mathbf{x})$ given by (2). First, we consider the generalised Leibniz rule for finding the k th derivative of the product uv given as:

$$(uv)^{(k)} = \sum_{m=0}^k \binom{k}{m} u^{(k-m)} v^{(m)}. \quad (31)$$

If we start from the first-order partial derivative of $\phi_j(\mathbf{x})$ in (27), we can express this as a product uv , where:

$$u = -2\beta_j(x_i - \mu_{j,i}), \quad (32)$$

and

$$v = e^{-\beta_j\|\mathbf{x}-\boldsymbol{\mu}_j\|^2} = \phi_j(\mathbf{x}). \quad (33)$$

From Leibniz rule, we can then derive the k th order derivative of the RBF

as:

$$\frac{\partial^{(k)}\phi_j(\mathbf{x})}{\partial x_i^k} = \begin{cases} \phi_j(\mathbf{x}), & k = 0 \\ \sum_{m=0}^{k-1} \binom{k-1}{m} u^{(k-m-1)} \frac{\partial^{(m)}\phi_j(\mathbf{x})}{\partial x_i^{(m)}}, & k \geq 1 \end{cases}, \quad (34)$$

which computes the k th partial derivative recursively.

If we now substitute the radial basis function network f and its higher-order partial derivatives into the PDE in (22), we obtain:

$$f = \boldsymbol{\lambda}^\top \mathbf{s}_{\pi-l+1:\pi} + \sum_{j=1}^c w_j \sum_{k=1}^{\nu} \sum_{i=1}^d a_{k,i} \frac{\partial^{(k)}\phi_j(\mathbf{x})}{\partial x_i^k}, \quad (35)$$

whose network architecture is given in Fig. 2

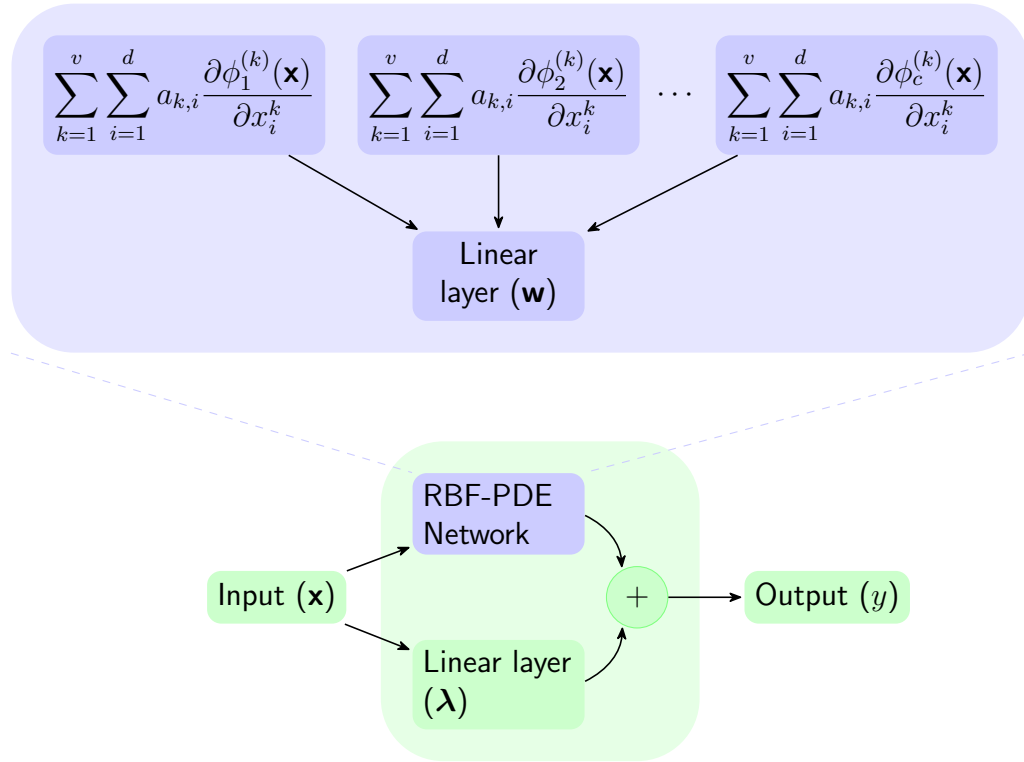


Figure 2: Proposed differential RBF network architecture

In matrix form, (35) is equivalent to:

$$f = \boldsymbol{\lambda}^\top \mathbf{s}_{\pi-l+1:\pi} + \mathbf{w}^\top \boldsymbol{\Theta}(\mathbf{x}) \mathbf{a}, \quad (36)$$

where,

$$\mathbf{w} = [w_1, \dots, w_c]^\top, \quad (37)$$

$$\mathbf{a} = [a_{1,1}, \dots, a_{1,d}, \dots, a_{\nu,1}, \dots, a_{\nu,d}]^\top, \quad (38)$$

and

$$\Theta(\mathbf{x}) = \begin{bmatrix} \frac{\partial^{(1)}\phi_1(\mathbf{x})}{\partial x_1} & \dots & \frac{\partial^{(1)}\phi_1(\mathbf{x})}{\partial x_d} & \dots & \frac{\partial^{(\nu)}\phi_1(\mathbf{x})}{\partial x_1^\nu} & \dots & \frac{\partial^{(\nu)}\phi_1(\mathbf{x})}{\partial x_d^\nu} \\ \vdots & \ddots & \vdots & \ddots & \vdots & \ddots & \vdots \\ \frac{\partial^{(1)}\phi_c(\mathbf{x})}{\partial x_1} & \dots & \frac{\partial^{(1)}\phi_c(\mathbf{x})}{\partial x_d} & \dots & \frac{\partial^{(\nu)}\phi_c(\mathbf{x})}{\partial x_1^\nu} & \dots & \frac{\partial^{(\nu)}\phi_c(\mathbf{x})}{\partial x_d^\nu} \end{bmatrix} \quad (39)$$

Using (36), we may then optimise \mathbf{w} , $\boldsymbol{\lambda}$ and \mathbf{a} jointly on the training data for arbitrary objective functions. For example, for the mean-squared error ϵ given as:

$$\epsilon = \frac{1}{N} \sum_{n=1}^N (y_n - \mathbf{w}^\top \boldsymbol{\phi}(\mathbf{x}_n) - w_0 - \mathbf{w}^\top \Theta(\mathbf{x}_n) \mathbf{a})^2, \quad (40)$$

we can optimise the network weights \mathbf{w} , $\boldsymbol{\lambda}$ and \mathbf{a} to an arbitrary accuracy with a given choice of an optimiser such as stochastic gradient descent (SGD).

4. Experimental Validation

In this section, we test the proposed differential RBF network on the logistic map chaotic timeseries [26, 27], as well as on 30 different timeseries from the M5 competition [28]; the reasons for the choice of these datasets are given in Sections 4.1 and 4.2 respectively.

4.1. Logistic map

We consider the logistic map given as:

$$s_{\pi+1} = 4s_\pi[1 - s_\pi], \quad (41)$$

and we generate 1000 observations in the forward pass as our timeseries using an initial value of 0.1; the first 200 observations are shown in Figure 3: We split the timeseries into training and testing sets, with the last 100 observations as the test set. We have selected the logistic map because, although it is an archetypal chaotic series, it arises from simple polynomial mappings that can be easily modelled by the RBF network. Furthermore, to simulate observational noise, we add some Gaussian noise with variance ω to

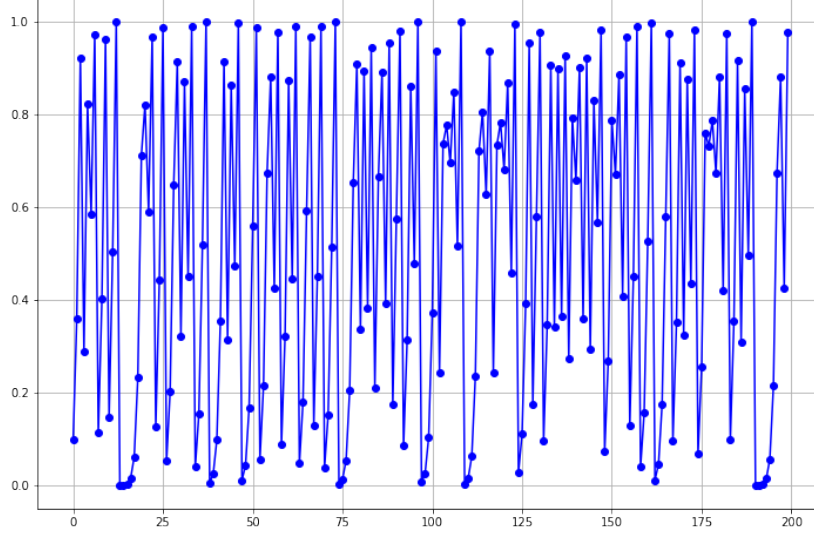


Figure 3: Logistic map chaotic timeseries

$\{0, 0.02, 0.04, 0.08, 0.12\}$ to the timeseries (whose original chaotic behaviour is not due to noise) in order to demonstrate the noise susceptibility of the RBF network. We have selected ω so that it does not exceed the variance of the original timeseries which is 0.12, in order not to corrupt the information contained in the series. We perform mean-normalisation, and then reshape the training data into input-output pairs using a lookback window of $l \in \{1, 2, 4, 8, 16\}$ with which we train the unnormalised and normalised RBF networks and the proposed differential RBF network. Although $l = 1$ suffices for the logistic map because it depends only its previous input, as shown in (41), we experiment with exponentially increasing l in order to demonstrate the replication of noisy observations as inputs to the RBF network. The number of RBF centres used is $c = \max(5, 2l)$, which is what showed the best performance after experimenting with different configurations. On the test set, we perform one-step prediction and report the mean absolute error (MAE). We defer multi-step prediction to Section 4.2.

For all the RBF networks (i.e., unnormalised, normalised and differential variants), we determine the exemplars μ_j via K-Means clustering while the RBF parameters β_j are defined according to (12). We minimise the least-squares error with lasso regularisation, with the regularisation coefficient determined via cross-validation. Furthermore, the order ν of the PDE in the

differential RBF network is set empirically to $\nu = 2$ based on cross-validation results, and we initialise the differential RBF network weights \mathbf{w} , $\boldsymbol{\lambda}$ and \mathbf{a} with the following settings:

$$\boldsymbol{\lambda} = \frac{1}{l}, \quad (42)$$

$$a_{i,k} = \frac{0.001^k}{k!}, \quad (43)$$

$$\mathbf{w} = [\tilde{w}_1^*, \dots, \tilde{w}_c^*]^\top, \quad (44)$$

where i, k index the k th-order partial derivative of the RBF with respect to the component x_i of the input \mathbf{x} , and $\tilde{w}_1^*, \dots, \tilde{w}_c^*$ are obtained from (9) with lasso regularisation, given that $\tilde{\mathbf{w}}^* = [\tilde{w}_0^*, \tilde{w}_1^*, \dots, \tilde{w}_c^*]^\top$. The above choices of weight initialisations for the proposed network are justified as follows:

1. By setting $\boldsymbol{\lambda}$ as in (42), equal weights are assigned to the last l values of the sequence.
2. By setting $a_{i,k}$ as in (43), the PDE weights are defined as Taylor expansion coefficients with 0.001 being the interval between any two realisations of the series.
3. By setting \mathbf{w} as in (44), we utilise the RBF weights obtained from the unnormalised RBF network.

We subsequently optimise the network weights in Figure 2 using the Broyden, Fletcher, Goldfarb, and Shanno (BFGS) method [47].

The MAE results from these experiments are given in Tables 1 to 5, while the model predictions are given in the appendix in Figures A.4 to A.8.

l	RBF-DiffNet	u-RBFN	n-RBFN
1	0.277	0.146	0.0994
2	0.133	0.16	0.166
4	0.0879	0.277	0.311
8	0.0717	0.223	0.313
16	0.0542	0.236	0.314

Table 1: Noise variance $\omega = 0$

l	RBF-DiffNet	u-RBFN	n-RBFN
1	0.286	0.152	0.103
2	0.124	0.162	0.167
4	0.082	0.283	0.312
8	0.0648	0.199	0.308
16	0.0517	0.221	0.315

Table 2: Noise variance $\omega = 0.02$

l	RBF-DiffNet	u-RBFN	n-RBFN
1	0.286	0.14	0.0987
2	0.143	0.161	0.157
4	0.0582	0.291	0.313
8	0.063	0.216	0.315
16	0.0631	0.244	0.314

Table 3: Noise variance $\omega = 0.04$

l	RBF-DiffNet	u-RBFN	n-RBFN
1	0.3	0.109	0.0916
2	0.151	0.139	0.132
4	0.0659	0.301	0.309
8	0.107	0.248	0.312
16	0.11	0.293	0.315

Table 4: Noise variance $\omega = 0.08$

l	RBF-DiffNet	u-RBFN	n-RBFN
1	0.295	0.138	0.106
2	0.18	0.15	0.144
4	0.0852	0.299	0.31
8	0.141	0.275	0.309
16	0.162	0.293	0.314

Table 5: Noise variance $\omega = 0.12$

Mean absolute error (MAE) performance on the logistic map at different noise levels. Best values are in boldface. l is the lookback window; RBF-DiffNet is the proposed differential RBF network; u-RBFN and n-RBFN are the unnormalised and normalised RBF networks respectively.

4.2. M5 dataset

The M5 competition is the latest in a series of M-forecasting competitions organised by the University of Nicosia that has drawn the interest of top forecasting practitioners including Google [48] and Uber [49]. The M5 competition comprises 42840 individual timeseries data made available by Walmart. This dataset consists of unit sales of products across the United States at different levels of aggregation. This dataset has been chosen because, like many real-world datasets, it exhibits a non-deterministic (noisy) behaviour as shown in Figures B.9 to B.13. Furthermore, by not including exogenous inputs such as prices that can affect product sales, we impose further *deterministic noise* [10] on the series. Thus, the dataset is able to sufficiently illustrate the susceptibility of the baseline RBF network to noisy inputs in sequential data.

Due to computational constraints, we work at the level 8 aggregation that consists of 30 different timeseries data of unit sales of products aggregated for each of 10 stores (across three US states) and 3 product categories (i.e., foods, household and hobbies). Each of the 30 different timeseries has 1941 daily observations (spanning the period between 29 January 2011 and 19 June 2016), with the last 28 days set for validation.

Different benchmark algorithms (categorised under statistical, machine learning and combinations of both) have been provided by the M5 competition organisers, and in this section, we limit our comparison to one statistical benchmark: autoregressive integrated moving average (ARIMA), and one machine learning benchmark: the multi-layer perceptron (MLP). The benchmark MLP architecture specified in the M5 competition is given as follows [28]:

1. Number of input layer nodes = 14, corresponding to the last 14 days of sales data
2. Number of hidden layer nodes = 28,
3. Number of output nodes = 1
4. Hidden layer activation function: logistic; output layer activation function: linear,
5. Number of iterations = 500,
6. Ensemble size = 10, i.e, 10 different MLPs are trained and their predictions aggregated in terms of the median operator to account for poor weight initialisations.

Using the above network parameters, we also implement an ensemble of 10 recurrent networks with LSTM blocks which use tanh activation for further comparison.

To make the RBF networks comparable to the MLP architecture, we set d (the RBF network input size) to 14 and c (the number of RBF nodes) to 28. The same settings of the RBF networks used in Section 4.1 are used here; the order ν of the PDE in the differential RBF network is however set to empirically to $\nu = 1$, based on cross-validation results.

For each timeseries, we set the last 28 observations as the test set, and the remaining as the training set. Since preprocessing for timeseries forecasting typically involves stationarising the series to make the series easier to predict [50, 51], we first perform a stationarity test on the training set using the augmented Dickey-Fuller test [52]; if the timeseries does not exhibit stationarity according to the test’s null hypothesis, we proceed to difference the timeseries incrementally until it shows stationarity, up to a differencing order of 10. We then perform mean-normalisation and then reshape the data into input-output pairs using a lookback window of $l = 14$.

On the test set, we perform reverse-differencing and normalisation operations to recover the original series and range, and we perform multi-step prediction over the forecast horizon of 28 days. We report the root mean squared scaled error (RMSSE), which is the error metric used in the M5 competition. The RMSSE is defined as:

$$\text{RMSSE} = \frac{(n_{\text{train}} - 1)}{n_{\text{test}}} \frac{\sum_{\pi=n_{\text{train}}+1}^{n_{\text{train}}+n_{\text{test}}} (s_{\pi} - \hat{s}_{\pi})^2}{\sum_{\pi=2}^{n_{\text{train}}} (s_{\pi} - s_{\pi-1})^2}, \quad (45)$$

where \hat{s}_{π} is a model’s estimate of the timeseries at time t_{π} , $n_{\text{train}} = 1913$ is the length of training series, and $n_{\text{test}} = 28$ is the length of the forecast horizon.

We do not investigate cross-learning, where information from other timeseries are used in training a single global model to predict each timeseries [53]; this is out of the scope of this paper.

4.3. Results and discussion

4.3.1. Logistic map data

From the results in Tables 1 to 5, there is an increase in the MAE values for the unnormalised RBF network (u-RBFN) as the lookback window l increases. This is especially observed for high noise levels ω , since the noisy

observations gets replicated many times in the reshaped data that is input to the RBF network, illustrating the susceptibility of the unnormalised RBF network to noisy observations.

The normalised RBF network (n-RBFN), on the other hand, shows robustness in its predictions across different noise levels, with the MAE values remaining relatively constant. However, as the lookback window length l increases, its performance also deteriorates for each of the noise levels. This performance degradation is mainly due to the side effects of normalisation [37, 36] whereby the maxima of the RBF shifts when the RBFs have unequally spaced centres or different widths. As Shorten [36] notes, this side effect becomes more pronounced with increase in the input dimension. This causes points far away from the RBF centre to have high functional values, and thus, rather than a few RBFs getting activated, many more RBFs tend to contribute uniformly to the output for any given input. This is manifested in the n-RBFN output approaching the mean of the timeseries for large values of l as shown in Figures A.4 to A.8.

The differential RBF network (RBF-DiffNet), however, shows robustness across different noise variance ω and different values of l , significantly outperforming the benchmark RBF networks. In particular, for low to medium noise levels, RBF-DiffNet is able to utilise more information in the lagged values of the series as l increases to improve the prediction accuracy, without any noise enhancement or side effects of normalisation, unlike for u-RBFN and n-RBFN. While the MAE performance can be seen to degrade as the noise variance increases, this is to be expected as the information in the timeseries increasingly gets corrupted. However, RBF-DiffNet noticeably underperforms for $l = 1$ across all noise levels; this is likely due to the fact that the regularisation due to the backward-Euler updates causes the RBF-DiffNet to be overly smooth resulting in underfitting, even though the reshaped data input to the RBF network may not contain many noisy observations for $l = 1$. Thus, RBF-DiffNet is most suited to timeseries forecasting tasks for which the series to be predicted is noisy and influenced by more than one previous realisations of the series: an example of this is demand forecasting of grocery items, where future sales may be autocorrelated with historical sales up to $l = 14$ days prior; other examples include weather forecasting or sentence completion in natural language processing.

In terms of computational performance, u-RBFN and n-RBFN having closed-form solutions given in (9) outperform the proposed RBF-DiffNet which requires an iterative optimisation routine. For a large number of RBF

nodes, the solution given by (9) may be computationally intractable and u-RBFN and n-RBFN may require iterative optimisation techniques such as SGD. Since u-RBFN and n-RBFN each has $c + 1$ optimising variables as compared to $c + d\nu + l + 1$ for the differential RBF network, the existing RBF networks may still have a computational edge over RBF-DiffNet in an iterative optimisation technique.

4.3.2. M5 datasets

Although the different models are able to maintain relatively constant prediction accuracy across the forecast horizon as shown in Figures B.9 to B.13 due to the seasonalities in the timeseries, the results in Table 6 show that the differential RBF network (RBF-DiffNet) achieves superior predictive accuracy over the unnormalised RBF network (u-RBFN) and normalised RBF network (n-RBFN). Over the 30 different timeseries, RBF-DiffNet achieves a 26% reduction over u-RBFN and an 18% reduction over n-RBFN in terms of the mean RMSSE. This performance improvement is shown to be statistically significant at the 0.95 confidence level, as seen from the p-values in Table 7.

The RBF-DiffNet is also shown to yield a 4% reduction over ARIMA (significant at the 0.95 confidence level), 0.2% over the MLP (statistically insignificant), and 2.2% over the LSTM (statistically insignificant) in terms of the RMSSE.

The reason for the performance gain of RBF-DiffNet over u-RBFN and n-RBFN is the relative robustness of RBF-DiffNet to noise as described in Section 4.3.1 due to its backward-Euler regularisation, while the performance gain of RBF-DiffNet over ARIMA is because the RBF network generally learns complex non-linear mappings as compared to the simple linear mapping in the ARIMA model.

Despite their comparable performance, it is worth noting that RBF-DiffNet uses only one weight initialisation as given in (42), (43) and (44), whereas the LSTM and MLP ensembles respectively use 10 different weight initialisations and the median operator to mitigate the effects of poor weight initialisations. The single weight initialisation is possible in RBF-DiffNet because its network weights have more intuitive meanings which allow for an informed choice of initialisation in such a way that multiple random restarts yield only minor performance improvements. Consequently, for the comparable performance given in Table 6, the RBF network requires at least 16 times less computational time than the LSTM ensemble, as shown in Table

8. Moreover, RBF-DiffNet has the added advantage of being interpretable in terms of the influence of different input features on the output, since the hidden layer performs a fuzzy nearest-neighbour association of a given input to the RBF centres.

5. Conclusion

The radial basis function (RBF) network has good function approximation properties, fast training time and is easily interpretable in terms of the influence of different input features on the output. However, it is especially sensitive to noisy inputs due to the usually unsupervised step in selecting the RBF centres and widths. This paper extends the RBF network for use on noisy sequential data. For this purpose, we have introduced a novel differential RBF network (RBF-DiffNet) whose hidden layer blocks are higher-order constant-coefficient partial differential equations in terms of the RBF. RBF-DiffNet exhibits robustness to noisy perturbations to the input sequence by regularising the network following backward-Euler updates, assuming the sequence data originates from an underlying differential equation. Notably, RBF-DiffNet differs from the neural ordinary differential equation (ODE-Net) in that RBF-DiffNet parameterises the solution to the underlying differential equation with an RBF network and directly evaluates the derivatives of the network based on the backward-Euler discretisation, while the ODE-Net parameterises the derivatives of the hidden state of a residual network with another neural network.

The proposed differential RBF network is experimentally validated on the logistic map and 30 other real-world and noisy timeseries data from the M5 competition (level 8 aggregation). Experimental results show RBF-DiffNet significantly outperforms the normalised and unnormalised RBF networks at different noise levels on the logistic map. On the M5 dataset, RBF-DiffNet outperforms the normalised and unnormalised RBF networks by 16% and 28% respectively; furthermore, because its network weights have intuitive meanings, RBF-DiffNet allows for an informed choice of initialisation for the network weights, thereby achieving comparable performance to an ensemble of 10 LSTM networks (built from 10 different weight initialisations) at less than one-sixteenth the LSTM computational time.

Our proposed network therefore enables more accurate predictions—in spite of the presence of observational noise—in sequence modelling tasks

Series ID	RBF-DiffNet	u-RBFN	n-RBFN	ARIMA	MLP	LSTM
1	0.5	<i>0.461</i>	0.53	0.513	0.433	0.46
2	<i>0.693</i>	0.78	0.701	0.686	0.76	0.749
3	0.422	0.621	0.532	0.512	0.467	0.489
4	0.761	1.85	<i>0.695</i>	0.867	0.801	0.683
5	<i>0.702</i>	0.992	0.972	0.884	0.833	0.581
6	0.54	0.745	0.564	0.749	0.664	0.598
7	<i>0.625</i>	0.662	1.35	0.492	0.438	0.653
8	<i>0.587</i>	1.47	0.786	0.582	0.552	0.554
9	0.609	1.33	<i>0.538</i>	0.556	0.528	0.616
10	<i>0.751</i>	0.771	1.45	0.691	0.743	0.663
11	<i>0.945</i>	1.14	1.09	0.978	0.935	0.928
12	<i>0.895</i>	2.25	1.76	0.875	1.12	1.16
13	<i>0.859</i>	0.921	0.936	0.862	0.697	1
14	<i>0.746</i>	0.945	0.775	0.923	0.759	0.725
15	0.847	<i>0.836</i>	1.02	0.83	0.837	1.03
16	<i>0.724</i>	0.9	1.38	0.55	0.521	0.571
17	<i>0.663</i>	0.796	0.704	0.667	0.686	0.637
18	0.526	0.967	0.643	0.536	0.578	0.659
19	1.55	0.85	1.37	1.11	1.38	1.12
20	0.783	1.48	0.977	0.916	0.877	0.817
21	0.858	1.07	0.661	0.809	0.75	1.33
22	0.627	0.778	0.719	0.753	0.704	0.634
23	0.466	0.489	0.472	0.522	0.648	0.485
24	0.637	0.625	0.674	0.687	0.78	0.788
25	1.45	1.42	1.52	1.87	1.46	1.44
26	0.68	0.841	0.973	0.741	0.731	0.73
27	<i>1.02</i>	1.95	1.54	1.12	0.988	0.975
28	0.816	0.971	<i>0.793</i>	1.03	0.705	0.761
29	0.808	<i>0.797</i>	0.858	0.821	0.778	0.81
30	<i>0.676</i>	1.04	0.849	0.697	0.67	0.639
Mean	0.759	1.02	0.927	0.794	0.76	0.776
Median	<i>0.713</i>	0.911	0.821	0.751	0.737	0.704

Table 6: Root mean squared scaled error (RMSSE) results on the M5 dataset (level 8 aggregation). RBF-DiffNet is the proposed differential RBF network; u-RBFN and n-RBFN are the unnormalised and normalised RBF networks respectively. Best values across all algorithms are in boldface. Best values across the three RBF networks (RBF-DiffNet, u-RBFN, n-RBFN) are in italics.

Algorithm	p-value
u-RBFN	3.473×10^{-5}
n-RBFN	2.098×10^{-4}
ARIMA	2.619×10^{-2}
MLP	4.300×10^{-1}
LSTM	2.783×10^{-1}

Table 7: p-values from Wilcoxon signed rank test for differences in median between proposed differential RBF network (RBF-DiffNet) and other algorithms. u-RBFN and n-RBFN are the unnormalised and normalised RBF networks respectively.

Algorithm	Training time (seconds)
RBF-DiffNet	38.97 ± 0.39
u-RBFN	1.27 ± 0.19
n-RBFN	1.20 ± 0.16
ARIMA	33.34 ± 10.74
MLP	307.23 ± 4.30
LSTM	633.61 ± 6.84

Table 8: Model training time on on the M5 datasets (level 8 aggregation) given in the form: mean \pm s.d. obtained on a 32GB-RAM Tesla P100-PCIE 128GB GPU processor with 8 Intel E5-2650 CPU cores @2.2GHz. RBF-DiffNet is the proposed differential RBF network; u-RBFN and n-RBFN are the unnormalised and normalised RBF networks respectively.

such as timeseries forecasting that leverage the fast training and function approximation properties of the RBF network, where model interpretability is desired.

Nevertheless, the current development of RBF-DiffNet is limited to scalar outputs, and on-going work is investigating its extension to vector outputs. The feasibility of the proposed approach to cross-learning, where information from multiple series is used in training a single global model, is also an area we are currently researching.

References

- [1] D. S. Broomhead, D. Lowe, Radial basis functions, multi-variable functional interpolation and adaptive networks, Tech. rep., Royal Signals and Radar Establishment Malvern (United Kingdom) (1988).
- [2] Q. Que, M. Belkin, Back to the future: radial basis function network revisited, *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2019).
- [3] M. Masnadi-Shirazi, S. Subramaniam, Attractor ranked radial basis function network: A nonparametric forecasting approach for chaotic dynamic systems, *Scientific reports* 10 (1) (2020) 1–10.
- [4] P. Dey, M. Gopal, P. Pradhan, T. Pal, On robustness of radial basis function network with input perturbation, *Neural Computing and Applications* 31 (2) (2019) 523–537.
- [5] P. Teng, Machine-learning quantum mechanics: Solving quantum mechanics problems using radial basis function networks, *Physical Review E* 98 (3) (2018) 033305.
- [6] J. Park, I. W. Sandberg, Universal approximation using radial-basis-function networks, *Neural computation* 3 (2) (1991) 246–257.
- [7] F. B. Rizaner, A. Rizaner, Approximate solutions of initial value problems for ordinary differential equations using radial basis function networks, *Neural Processing Letters* 48 (2) (2018) 1063–1071.
- [8] F. Scarselli, A. C. Tsoi, Universal approximation using feedforward neural networks: A survey of some existing methods, and some new results, *Neural networks* 11 (1) (1998) 15–37.

- [9] F. Girosi, T. Poggio, Networks and the best approximation property, *Biological cybernetics* 63 (3) (1990) 169–176.
- [10] Y. S. Abu-Mostafa, M. Magdon-Ismail, H.-T. Lin, *Learning from data*, Vol. 4, AMLBook New York, NY, USA:, 2012.
- [11] C. S. K. Dash, A. K. Behera, S. Dehuri, S.-B. Cho, Radial basis function neural networks: a topical state-of-the-art survey, *Open Computer Science* 1 (open-issue) (2016).
- [12] Y. W. Wong, K. P. Seng, L.-M. Ang, Radial basis function neural network with incremental learning for face recognition, *IEEE Transactions on Systems, Man, and Cybernetics, Part B (Cybernetics)* 41 (4) (2011) 940–949.
- [13] R. S. Sutton, A. G. Barto, *Reinforcement learning: An introduction*, MIT press, 2018.
- [14] A. d. M. S. Barreto, C. W. Anderson, Restricted gradient-descent algorithm for value-function approximation in reinforcement learning, *Artificial Intelligence* 172 (4-5) (2008) 454–482.
- [15] R. M. Kretchmar, C. W. Anderson, Comparison of cmacs and radial basis functions for local function approximators in reinforcement learning, in: *Proceedings of International Conference on Neural Networks (ICNN’97)*, Vol. 2, IEEE, 1997, pp. 834–837.
- [16] E. A. Lim, W. H. Tan, A. K. Junoh, Distance weighted k-means algorithm for center selection in training radial basis function networks, *IAES International Journal of Artificial Intelligence* 8 (1) (2019) 54.
- [17] A. Iske, *Optimal distribution of centers for radial basis function methods* (2000).
- [18] M. J. Orr, Regularization in the selection of radial basis function centers, *Neural computation* 7 (3) (1995) 606–623.
- [19] F. Scheibel, N. C. Steele, R. Low, Centre and variance selection for gaussian radial basis function artificial neural networks, in: *Artificial Neural Nets and Genetic Algorithms*, Springer, 1999, pp. 141–147.

- [20] S. Chen, X. Wang, X. Hong, C. J. Harris, Kernel classifier construction using orthogonal forward selection and boosting with fisher ratio class separability measure, *IEEE transactions on neural networks* 17 (6) (2006) 1652–1656.
- [21] S. Chen, X. Hong, B. L. Luk, C. J. Harris, Construction of tunable radial basis function networks using orthogonal forward selection, *IEEE Transactions on Systems, Man, and Cybernetics, Part B (Cybernetics)* 39 (2) (2008) 457–466.
- [22] J. B. Gomm, D. L. Yu, Selecting radial basis function network centers with recursive orthogonal least squares training, *IEEE Transactions on Neural networks* 11 (2) (2000) 306–314.
- [23] J. Huilan, G. Ying, L. Dongwei, X. Jianqiang, Self-adaptive clustering algorithm based rbf neural network and its application in the fault diagnosis of power systems, in: *2005 IEEE/PES Transmission & Distribution Conference & Exposition: Asia and Pacific*, IEEE, 2005, pp. 1–6.
- [24] A. S. Kamalabady, K. Salahshoor, New siso and miso adaptive nonlinear predictive controllers based on self organizing rbf neural networks, in: *2008 3rd International Symposium on Communications, Control and Signal Processing*, IEEE, 2008, pp. 703–708.
- [25] J. Brusey, D. Hintea, E. Gaura, N. Beloe, Reinforcement learning-based thermal comfort control for vehicle cabins, *Mechatronics* 50 (2018) 413–421.
- [26] H. Maathuis, L. Boulogne, M. Wiering, A. Sterk, Predicting chaotic time series using machine learning techniques, in: *Preproceedings of the 29th Benelux Conference on Artificial Intelligence (BNAIC 2017)*. Groningen, 2017, pp. 326–40.
- [27] J. D. Farmer, J. J. Sidorowich, Predicting chaotic time series, *Physical review letters* 59 (8) (1987) 845.
- [28] S. Makridakis, *The m5 competition* (2020).

- [29] S. Chen, C. F. Cowan, P. M. Grant, Orthogonal least squares learning algorithm for radial basis function networks, *IEEE Transactions on neural networks* 2 (2) (1991) 302–309.
- [30] B. Scholkopf, K.-K. Sung, C. J. Burges, F. Girosi, P. Niyogi, T. Poggio, V. Vapnik, Comparing support vector machines with gaussian kernels to radial basis function classifiers, *IEEE transactions on Signal Processing* 45 (11) (1997) 2758–2765.
- [31] C. McCormick, Radial basis function network (rbfn) tutorial (2013).
- [32] Y. Wu, H. Wang, B. Zhang, K.-L. Du, Using radial basis function networks for function approximation and classification, *ISRN Applied Mathematics* 2012 (2012).
- [33] N. Benoudjit, M. Verleysen, On the kernel widths in radial-basis function networks, *Neural Processing Letters* 18 (2) (2003) 139–154.
- [34] J. Moody, C. J. Darken, Fast learning in networks of locally-tuned processing units, *Neural computation* 1 (2) (1989) 281–294.
- [35] G.-F. Lin, L.-H. Chen, Time series forecasting by combining the radial basis function network and the self-organizing map, *Hydrological Processes: An International Journal* 19 (10) (2005) 1925–1937.
- [36] R. Shorten, R. Murray-Smith, On normalising radial basis function networks, in: *Proceedings of the fourth Irish Conference on Neural Networks, INNC*, Vol. 94, 1994, pp. 213–217.
- [37] R. Shorten, R. Murray-Smith, Side effects of normalising radial basis function networks, *International Journal of Neural Systems* 7 (02) (1996) 167–179.
- [38] G. Bugmann, Normalized gaussian radial basis function networks, *Neurocomputing* 20 (1-3) (1998) 97–110.
- [39] R. T. Chen, Y. Rubanova, J. Bettencourt, D. K. Duvenaud, Neural ordinary differential equations, in: *Advances in neural information processing systems*, 2018, pp. 6571–6583.

- [40] Y. Chen, S. Gottlieb, A. Heryudono, A. Narayan, A reduced radial basis function method for partial differential equations on irregular domains, *Journal of Scientific Computing* 66 (1) (2016) 67–90.
- [41] E. Larsson, V. Shcherbakov, A. Heryudono, A least squares radial basis function partition of unity method for solving pdes, *SIAM Journal on Scientific Computing* 39 (6) (2017) A2538–A2563.
- [42] N. Mai-Duy, T. Tran-Cong, Numerical solution of differential equations using multiquadric radial basis function networks, *Neural networks* 14 (2) (2001) 185–199.
- [43] I. E. Lagaris, A. Likas, D. I. Fotiadis, Artificial neural networks for solving ordinary and partial differential equations, *IEEE transactions on neural networks* 9 (5) (1998) 987–1000.
- [44] R. Schaback, H. Wendland, Using compactly supported radial basis functions to solve partial differential equations, *WIT Transactions on Modelling and Simulation* 23 (1970).
- [45] K. Atkinson, W. Han, D. E. Stewart, Numerical solution of ordinary differential equations, Vol. 108, John Wiley & Sons, 2011.
- [46] N. Mai-Duy, T. Tran-Cong, Approximation of function and its derivatives using radial basis function networks, *Applied Mathematical Modelling* 27 (3) (2003) 197–220.
- [47] J. Nocedal, S. Wright, Numerical optimization, Springer Science & Business Media, 2006.
- [48] C. Fry, M. Brundage, The m4 forecasting competition-a practitioner’s view (2019).
- [49] S. Smyl, A hybrid method of exponential smoothing and recurrent neural networks for time series forecasting, *International Journal of Forecasting* 36 (1) (2020) 75–85.
- [50] I. A. Gheyas, L. S. Smith, A neural network approach to time series forecasting, in: *Proceedings of the World Congress on Engineering*, Vol. 2, 2009, pp. 1–3.

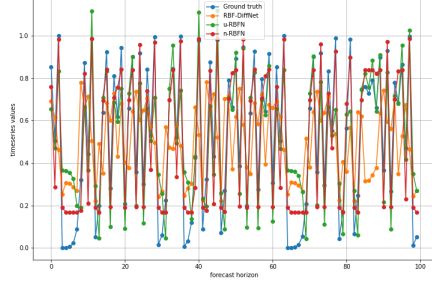
- [51] A. Pal, P. Prakash, Practical Time Series Analysis: Master Time Series Data Processing, Visualization, and Modeling using Python, Packt Publishing Ltd, 2017.
- [52] R. Mushtaq, Augmented dickey fuller test (2011).
- [53] S. Makridakis, E. Spiliotis, V. Assimakopoulos, The m4 competition: 100,000 time series and 61 forecasting methods, International Journal of Forecasting 36 (1) (2020) 54–74.

Appendix A. Logistic map results

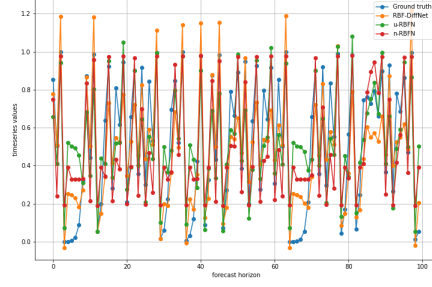
This section contains graphical results from model predictions on the logistic map at different noise levels as described in Section 4.1.

Appendix B. M5 results

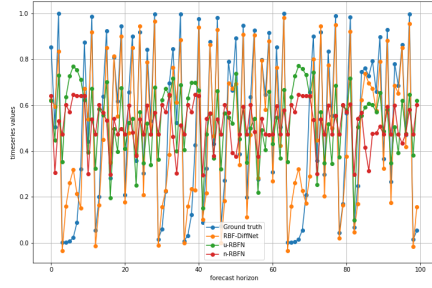
This section contains graphical results from model predictions on the M5 dataset (level 8 aggregation) as described in Section 4.2.



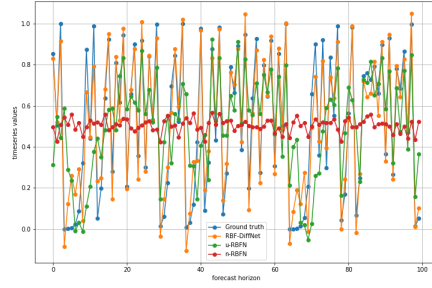
(a) $\omega = 0, l = 1$



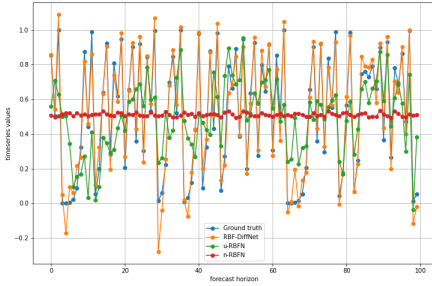
(b) $\omega = 0, l = 2$



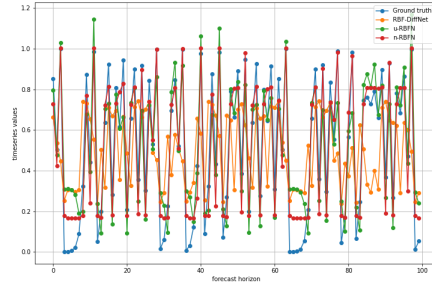
(c) $\omega = 0, l = 4$



(d) $\omega = 0, l = 8$

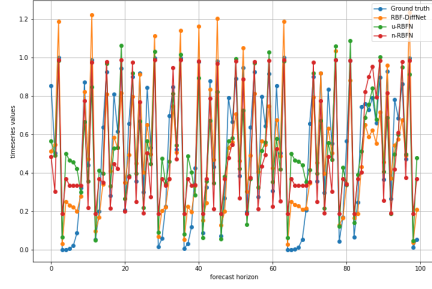


(e) $\omega = 0, l = 16$

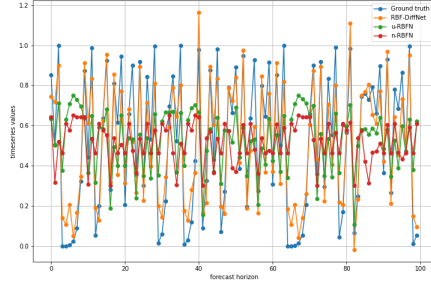


(f) $\omega = 0.02, l = 1$

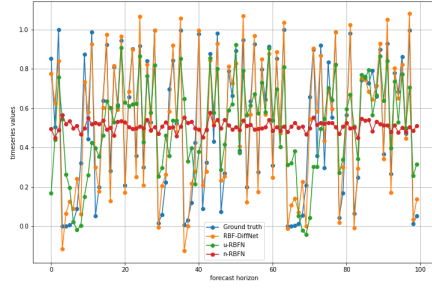
Figure A.4: Model predictions on the logistic map at different noise levels. RBF-DiffNet is the proposed differential RBF network; u-RBFN and n-RBFN are the unnormalised and normalised RBF networks respectively. ω is the noise variance; l is the lookback window length.



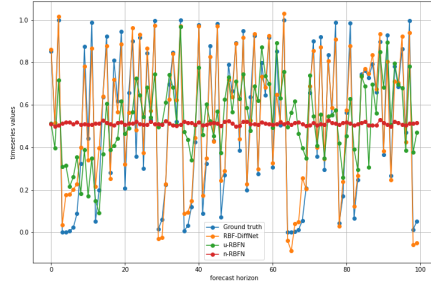
(a) $\omega = 0.02, l = 2$



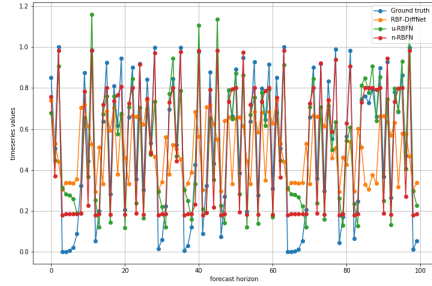
(b) $\omega = 0.02, l = 4$



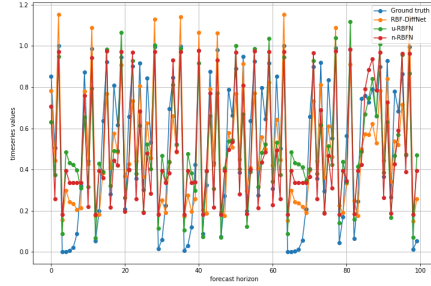
(c) $\omega = 0.02, l = 8$



(d) $\omega = 0.02, l = 16$

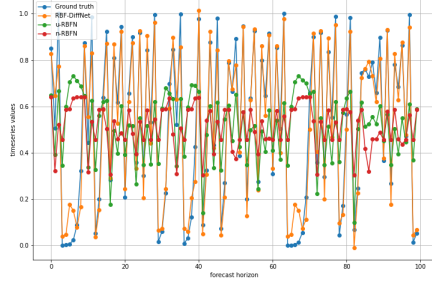


(e) $\omega = 0.04, l = 1$

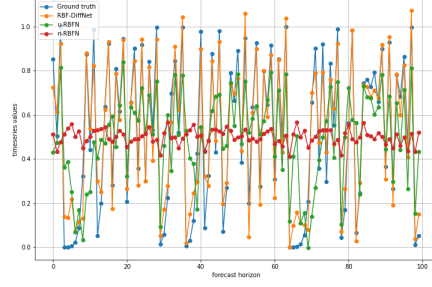


(f) $\omega = 0.04, l = 2$

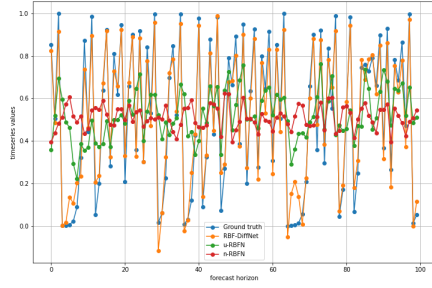
Figure A.5: Model predictions on the logistic map at different noise levels. RBF-DiffNet is the proposed differential RBF network; u-RBFN and n-RBFN are the unnormalised and normalised RBF networks respectively. ω is the noise variance; l is the lookback window length.



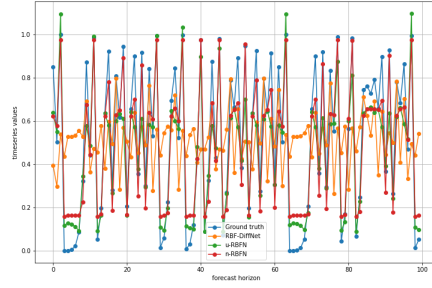
(a) $\omega = 0.04, l = 4$



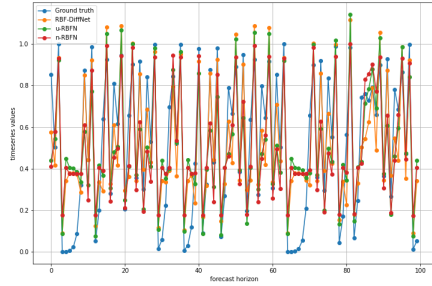
(b) $\omega = 0.04, l = 8$



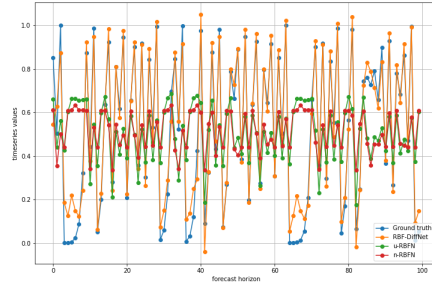
(c) $\omega = 0.04, l = 16$



(d) $\omega = 0.08, l = 1$

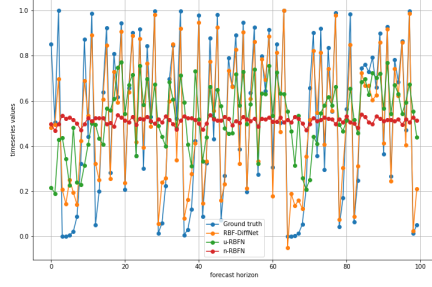


(e) $\omega = 0.08, l = 2$

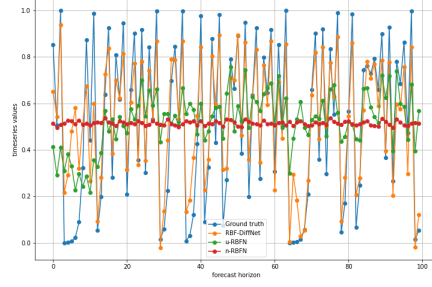


(f) $\omega = 0.08, l = 4$

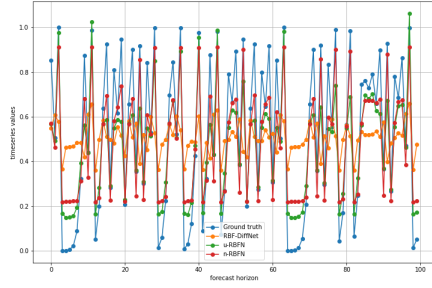
Figure A.6: Model predictions on the logistic map at different noise levels. RBF-DiffNet is the proposed differential RBF network; u-RBFN and n-RBFN are the unnormalised and normalised RBF networks respectively. ω is the noise variance; l is the lookback window length.



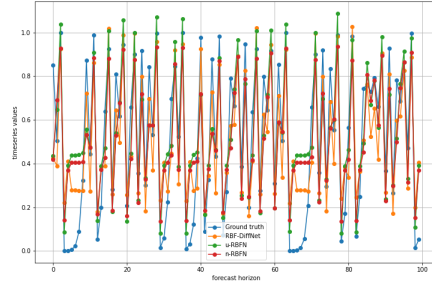
(a) $\omega = 0.08, l = 8$



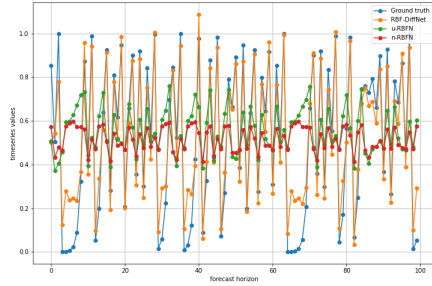
(b) $\omega = 0.08, l = 16$



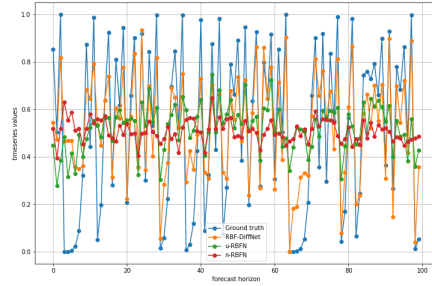
(c) $\omega = 0.12, l = 1$



(d) $\omega = 0.12, l = 2$

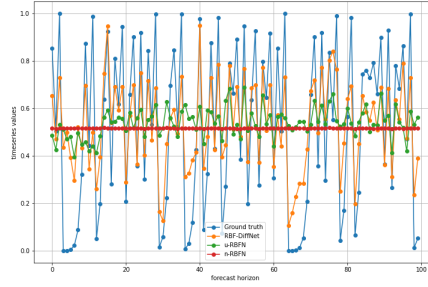


(e) $\omega = 0.12, l = 4$



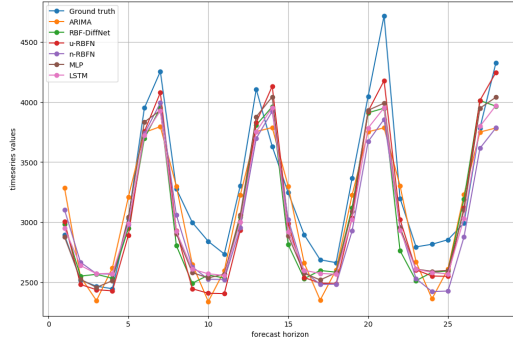
(f) $\omega = 0.12, l = 8$

Figure A.7: Model predictions on the logistic map at different noise levels. RBF-DiffNet is the proposed differential RBF network; u-RBFN and n-RBFN are the unnormalised and normalised RBF networks respectively. ω is the noise variance; l is the lookback window length.

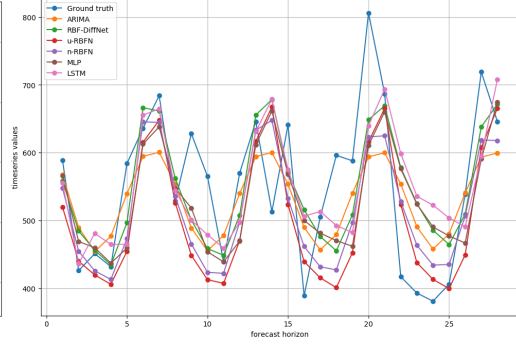


(a) $\omega = 0.12$, $l = 16$

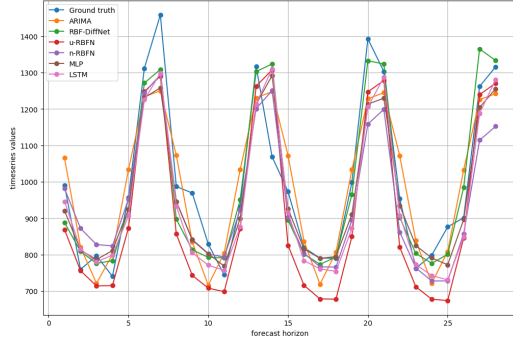
Figure A.8: Model predictions on the logistic map at different noise levels. RBF-DiffNet is the proposed differential RBF network; u-RBFN and n-RBFN are the unnormalised and normalised RBF networks respectively. ω is the noise variance; l is the lookback window length.



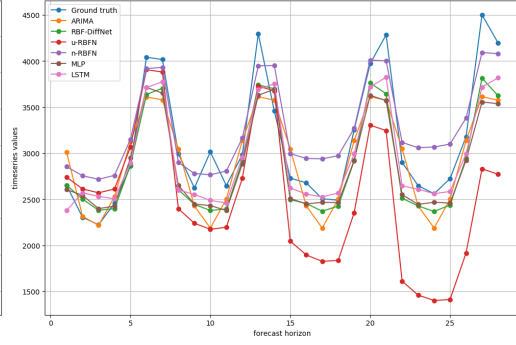
(a) Series 1



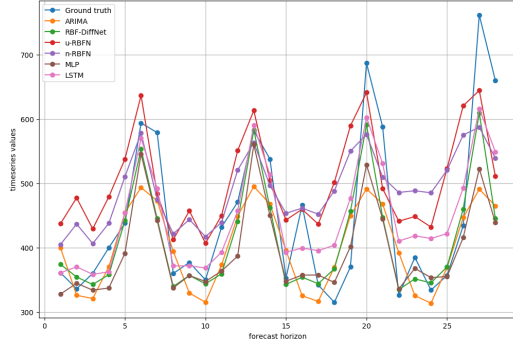
(b) Series 2



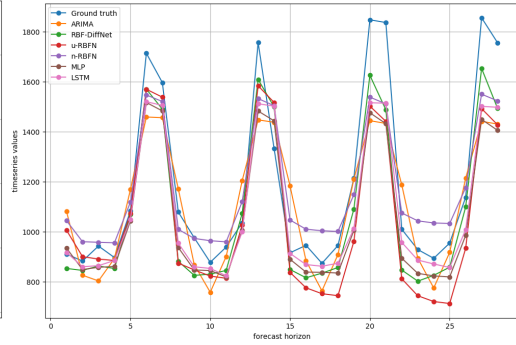
(c) Series 3



(d) Series 4

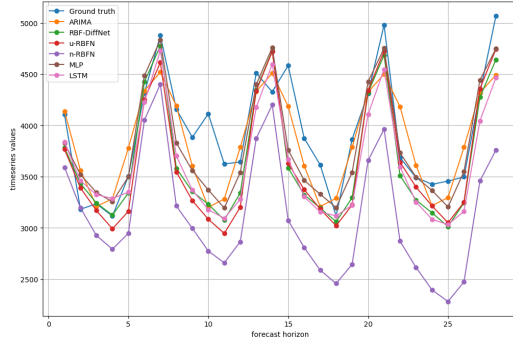


(e) Series 5

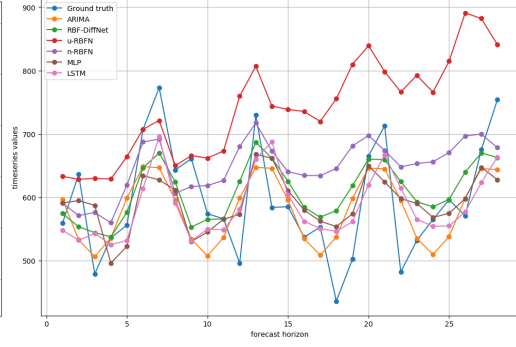


(f) Series 6

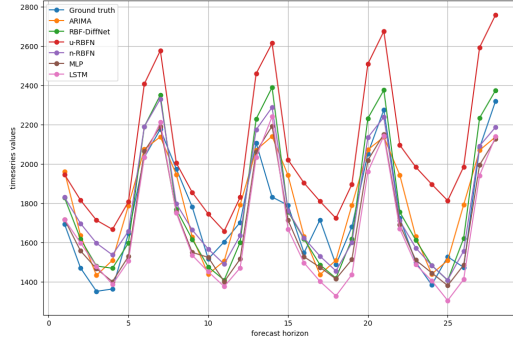
Figure B.9: Model predictions on M5 dataset (level 8 aggregation). RBF-DiffNet is the proposed differential RBF network; u-RBFN and n-RBFN are the unnormalised and normalised RBF networks respectively.



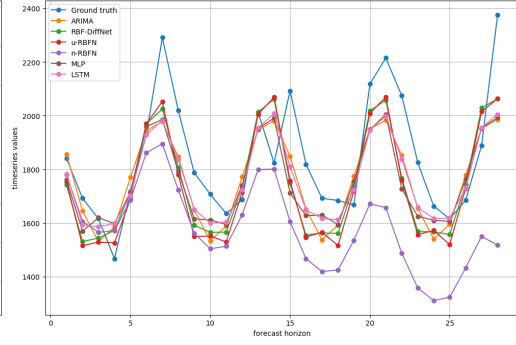
(a) Series 7



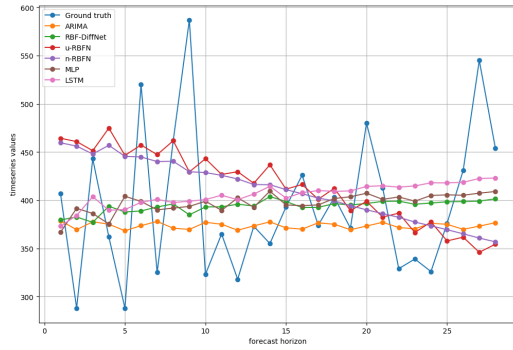
(b) Series 8



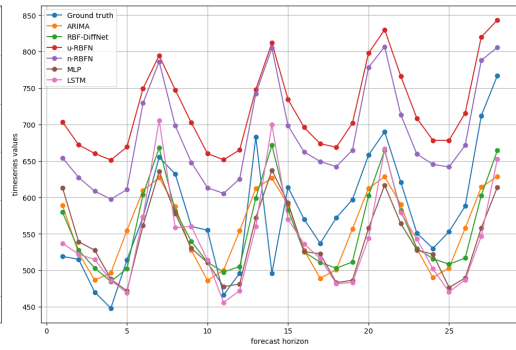
(c) Series 9



(d) Series 10

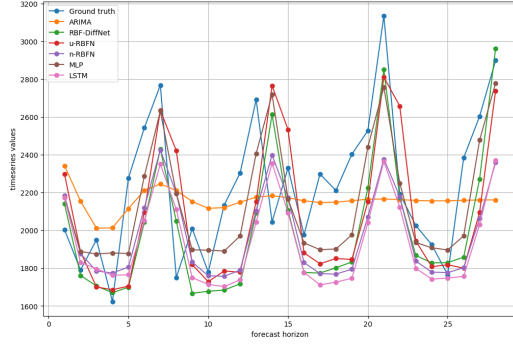


(e) Series 11

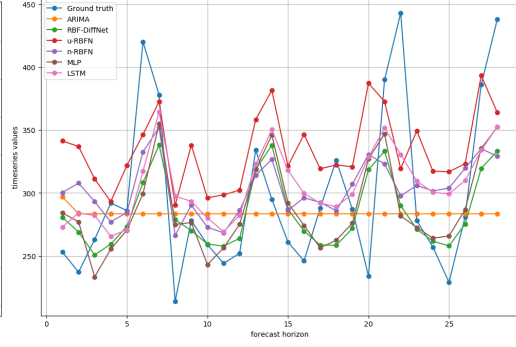


(f) Series 12

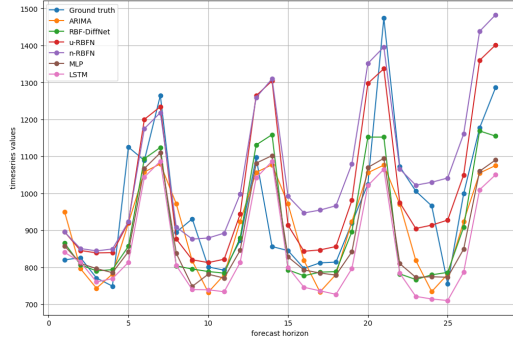
Figure B.10: Model predictions on M5 dataset (level 8 aggregation). RBF-DiffNet is the proposed differential RBF network; u-RBFN and n-RBFN are the unnormalised and normalised RBF networks respectively.



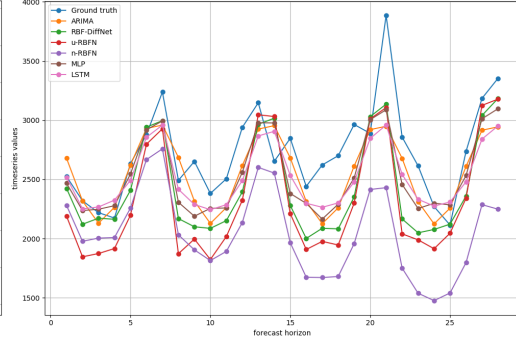
(a) Series 13



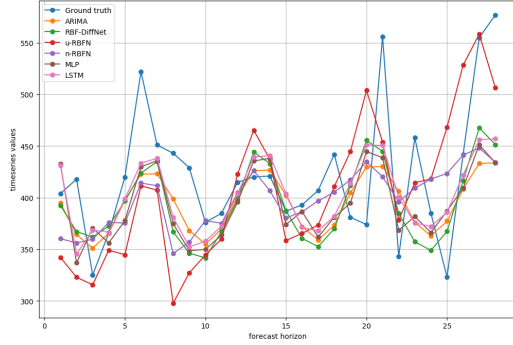
(b) Series 14



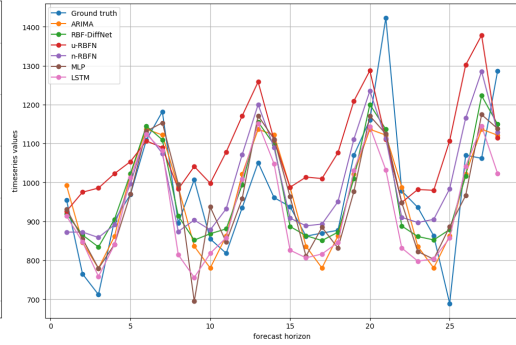
(c) Series 15



(d) Series 16

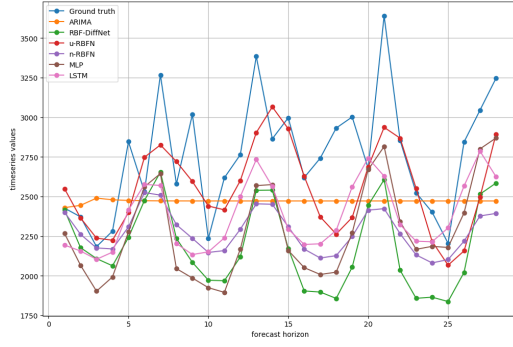


(e) Series 17

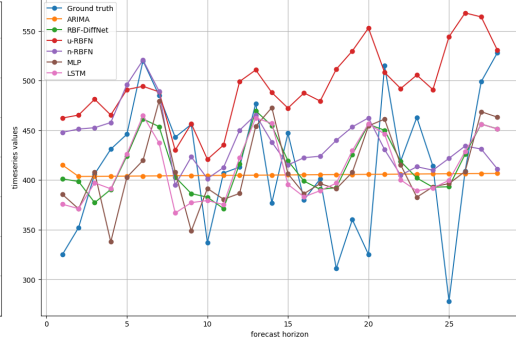


(f) Series 18

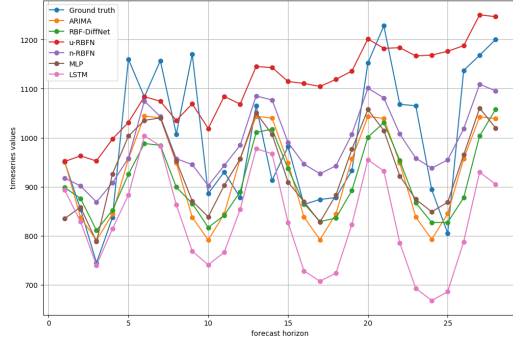
Figure B.11: Model predictions on M5 dataset (level 8 aggregation). RBF-DiffNet is the proposed differential RBF network; u-RBFN and n-RBFN are the unnormalised and normalised RBF networks respectively.



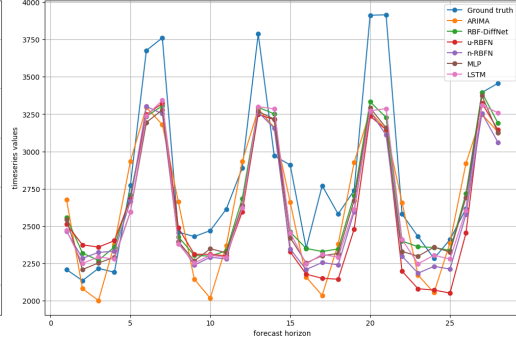
(a) Series 19



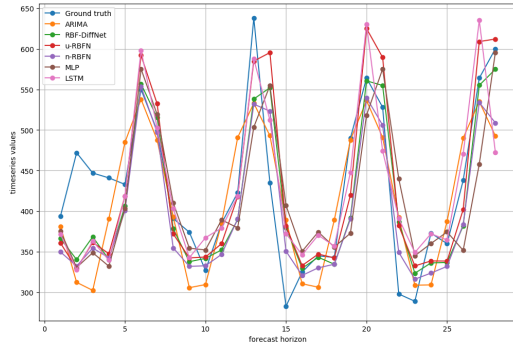
(b) Series 20



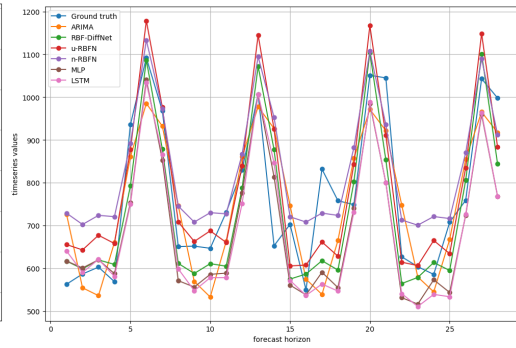
(c) Series 21



(d) Series 22

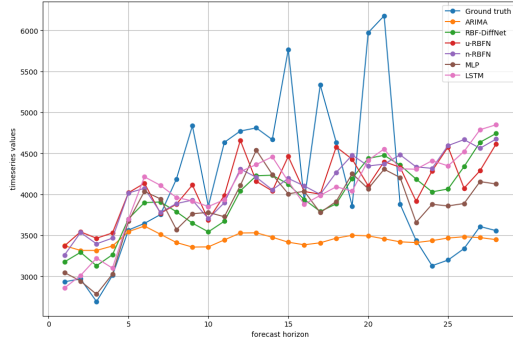


(e) Series 23

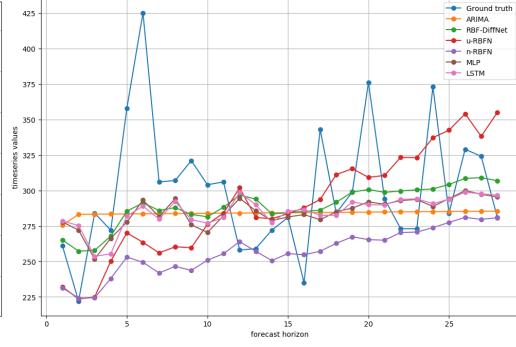


(f) Series 24

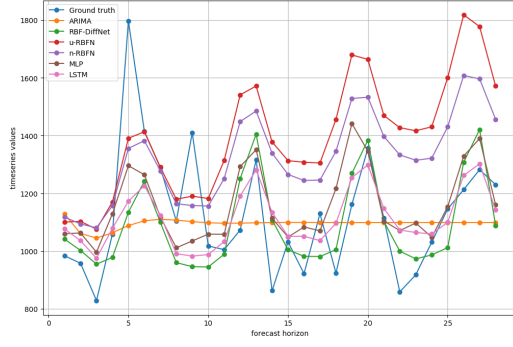
Figure B.12: Model predictions on M5 dataset (level 8 aggregation). RBF-DiffNet is the proposed differential RBF network; u-RBFN and n-RBFN are the unnormalised and normalised RBF networks respectively.



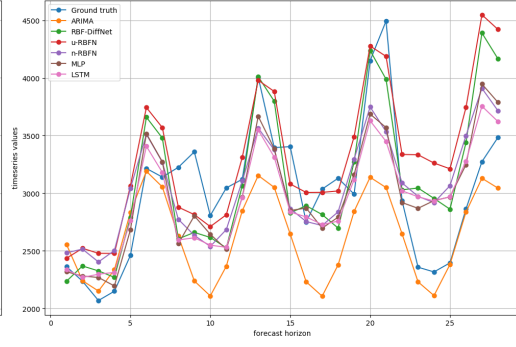
(a) Series 25



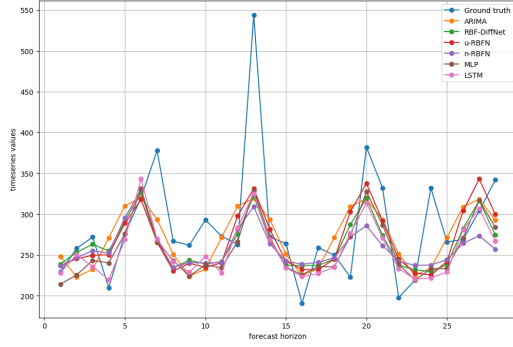
(b) Series 26



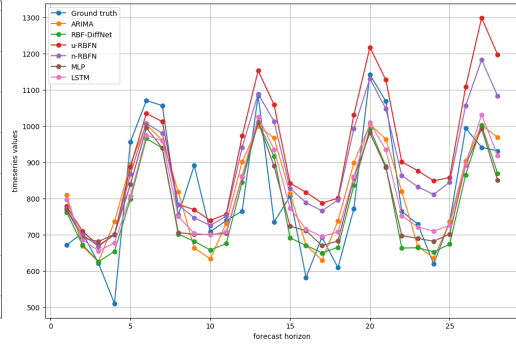
(c) Series 27



(d) Series 28



(e) Series 29



(f) Series 30

Figure B.13: Model predictions on M5 dataset (level 8 aggregation). RBF-DiffNet is the proposed differential RBF network; u-RBFN and n-RBFN are the unnormalised and normalised RBF networks respectively.