

Brief Papers

An Assessment of Qualitative Performance of Machine Learning Architectures: Modular Feedback Networks

Mo Chen, Temujin Gautama, and Danilo P. Mandic

Abstract—A framework for the assessment of qualitative performance of machine learning architectures is proposed. For generality, the analysis is provided for the modular nonlinear pipelined recurrent neural network (PRNN) architecture. This is supported by a sensitivity analysis, which is achieved based upon the prediction performance with respect to changes in the nature of the processed signal and by utilizing the recently introduced delay vector variance (DVV) method for phase space signal characterization. Comprehensive simulations combining the quantitative and qualitative analysis on both linear and nonlinear signals suggest that better quantitative prediction performance may need to be traded in order to preserve the nature of the processed signal, especially where the signal nature is of primary importance (biomedical applications).

Index Terms—Delay vector variance, nonlinearity, pipelined recurrent neural networks (PRNNs), qualitative performance, sensitivity.

I. INTRODUCTION

Most real-world signals contain both linear and nonlinear, as well as deterministic and stochastic components. It is, therefore, essential to characterize the signal nature before the actual processing approach is applied. In addition, a change in the linear, nonlinear, deterministic, or stochastic nature of a signal can convey important information¹ about the underlying signal generation mechanism. In such applications, it is important not only to obtain a good quantitative performance for the signal under study, but also that the signal nature is retained.

To establish a general framework for the quantitative performance analysis of machine learning algorithms and architectures, for generality, we will focus on a modular network, with feedback and nonlinear processing elements. One such architecture is the pipelined recurrent neural network (PRNN), which consists of a number of small-scale recurrent neural networks (RNNs), and maintains a relatively low computational complexity considering the entire number of neurons in its architecture [2]. In a previous study, we have analytically described the core features of the PRNN for the prediction application [3]; based on those quantitative results, we will use the PRNN as a convenient computational model for the analysis of the qualitative performance.

This way, we will combine our recent results on signal modality characterization [4], [5] and machine learning [6] to provide insights into the changes of the nature of the processed signals during online learning.

Manuscript received November 29, 2006; revised May 6, 2007; accepted May 12, 2007.

M. Chen and D. P. Mandic are with the Department of Electrical and Electronic Engineering, Communication and Signal Processing, Imperial College London, London SW7 2BT, U.K. (e-mail: mo.chen02@imperial.ac.uk; d.mandic@imperial.ac.uk).

T. Gautama is with the Philips Leuven, Leuven B-3001, Belgium (e-mail: temujin.gautama@philips.com).

Color versions of one or more of the figures in this paper are available online at <http://ieeexplore.ieee.org>.

Digital Object Identifier 10.1109/TNN.2007.902728

¹In the electrocardiogram (ECG) and heart rate variability (HRV) analysis, where the change in the signal nature from the linear stochastic to nonlinear deterministic, provides an indication of health hazard [1].

II. BACKGROUND

A. Nature of a Signal and Surrogate Data

By the signal “nature,” we adhere to a number of signal properties described in [4]: linear, nonlinear, deterministic, and stochastic signal behavior. A linear signal is generated by a linear time-invariant system, driven by white Gaussian noise, measured by a static, monotonic, and possibly nonlinear observation function; otherwise, the signal is considered to be nonlinear. A signal is considered deterministic if it can be precisely described by a set of equations; otherwise, it is considered as stochastic.

Surrogate time series, or “surrogate” for short, is a nonparametric randomized linear version of the original data which preserves the linear properties of the original data. In our experiments, we choose iterative amplitude adjust Fourier transform (iAAFT) method to generate surrogates, as it preserves the amplitude distribution of the original data and yields reliable results [7].

B. Quality Assessment Tool—Delay Vector Variance Method

Several methods for detecting nonlinear nature of a signal have been proposed over the past few years [8], [9]. For our purpose, it is desirable to have a method which is straightforward to visualize, and which facilitates the analysis of predictability, which is a core notion in online learning. One such method is our recently proposed delay vector variance (DVV) method [4], based upon examining the predictability of a signal in the phase space, and examining simultaneously the determinism and nonlinearity within a signal. This method can be summarized as follows. For an optimal² embedding dimension m , we have the following.

- Generate delay vectors (DVs): $\mathbf{x}(k) = [x_{k-m}, \dots, x_{k-1}]^T$ and the corresponding target x_k .
- The mean μ_d and standard deviation σ_d are computed over all pairwise Euclidean distances between delay vectors (DVs), $\|\mathbf{x}(i) - \mathbf{x}(j)\|$ ($i \neq j$).
- The sets $\Omega_k(r_d)$ are generated such that $\Omega_k(r_d) = \{\mathbf{x}(i) \mid \|\mathbf{x}(k) - \mathbf{x}(i)\| \leq r_d\}$, i.e., sets which consist of all DVs that lie closer to $\mathbf{x}(k)$ than a certain distance r_d , taken from the interval $[\max\{0, \mu_d - n_d\sigma_d\}; \mu_d + n_d\sigma_d]$, where n_d is a parameter controlling the span over which to perform the DVV analysis.
- For every set $\Omega_k(r_d)$, the variance of the corresponding targets $\sigma_k^2(r_d)$ is computed. The average over all sets $\Omega_k(r_d)$, normalized by the variance of the time series σ_x^2 yields the “target variance” $\sigma^{*2}(r_d) = (1/K) \sum_{k=1}^K \sigma_k^2(r_d) / \sigma_x^2$, where K is the total number of the sets $\Omega_k(r_d)$.

To illustrate the meaning of “signal nature,” consider a linear benchmark signal [AR(4)] [6]

$$x(k) = 1.79x(k-1) - 1.85x(k-2) + 1.27x(k-3) - 0.41x(k-4) + n(k) \quad (1)$$

²We adopt Cao’s method [10] for choosing the optimal embedding dimension in all of our simulations, which yields four for the linear benchmark signal (1) and two for the nonlinear benchmark signal (2).

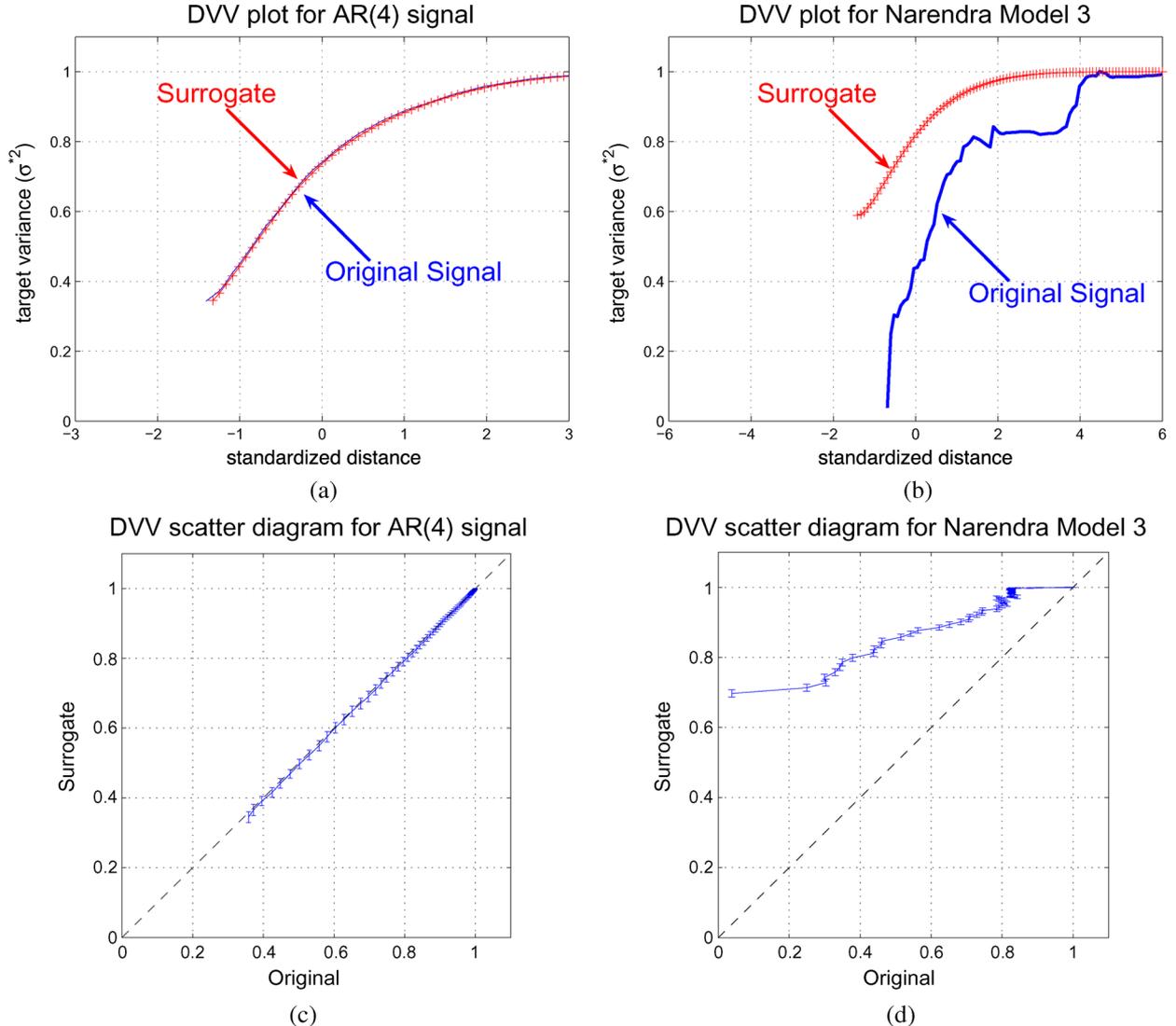


Fig. 1. Nonlinear and deterministic nature of signals. Diagrams (a) and (b) are DVV plots for the AR(4) signal and Narendra model 3, where the light line with crosses denotes the DVV plot for the average of 25 iAAFT-based surrogates while the dark solid line denotes that for the original signal. Diagrams (c) and (d) denote the DVV scatter diagrams for those two signals, where error bars denote the standard deviation of the target variances of surrogates.

and a nonlinear benchmark signal (a Narendra model 3 realization), given by [6]

$$z(k) = \frac{z(k-1)}{1+z^2(k-1)} + x^3(k) \quad (2)$$

where $x\{k\}$ denotes the AR(4) signal (1) and $\{n(k)\}$ is white Gaussian noise $n(k) \in \mathcal{N}(0, 1)$.

In the following step, the linear or nonlinear nature of the time series is examined by performing DVV analyses on both the original and 25 of surrogate time series. The DVV plot (target variance, $\sigma^{*2}(r_d)$) is a function of the standardized distance $(r_d - \mu_d)/\sigma_d$. At the extreme right, the DVV plots smoothly converge to unity, since for maximum spans, all DVs belong to the same set, and the variance of the targets is equal to the variance of the time series, illustrated in Fig. 1(a) and (d). The minimal target variance, e.g., the lowest point of the curve, is a measure for the amount of noise which is present in the time series. DVV scatter diagram is constructed in the way where the horizontal axis corresponds to the DVV plot of the original time series, and the vertical to that of the surrogate time series. A linear signal will have

similar DVV plot as its surrogate, resulting in the DVV scatter diagram coinciding with the bisector line as illustrated in Fig. 1(c), whereas a nonlinear signal will result in a deviation of DVV scatter diagram from the bisector line as illustrated in Fig. 1(d). This provides a very convenient tool for the qualitative analysis in machine learning since the deviation from bisector line in the DVV scatter diagram can be used to indicate the changes in the signal nature before and after processing.

III. PRNN ARCHITECTURE

The PRNN is a modular neural network that consists of M nested recurrent neural networks (RNNs) [11] as its modules, with each module consisting of N neurons. All the modules operate using the same weight matrix \mathbf{W} , as shown in Fig. 2. Besides its modularity, it is not immediately obvious that the PRNN performs *nesting* [12] of its constituting modules, and at the same time *data-reusing* [13].

In our analysis, the PRNN is utilized to perform one-step-forward prediction on two benchmark signals. In the experiments, the real-time recurrent learning (RTRL) algorithm [2] was used to train RNNs within the PRNN, and the activation function of a neuron was chosen to be the

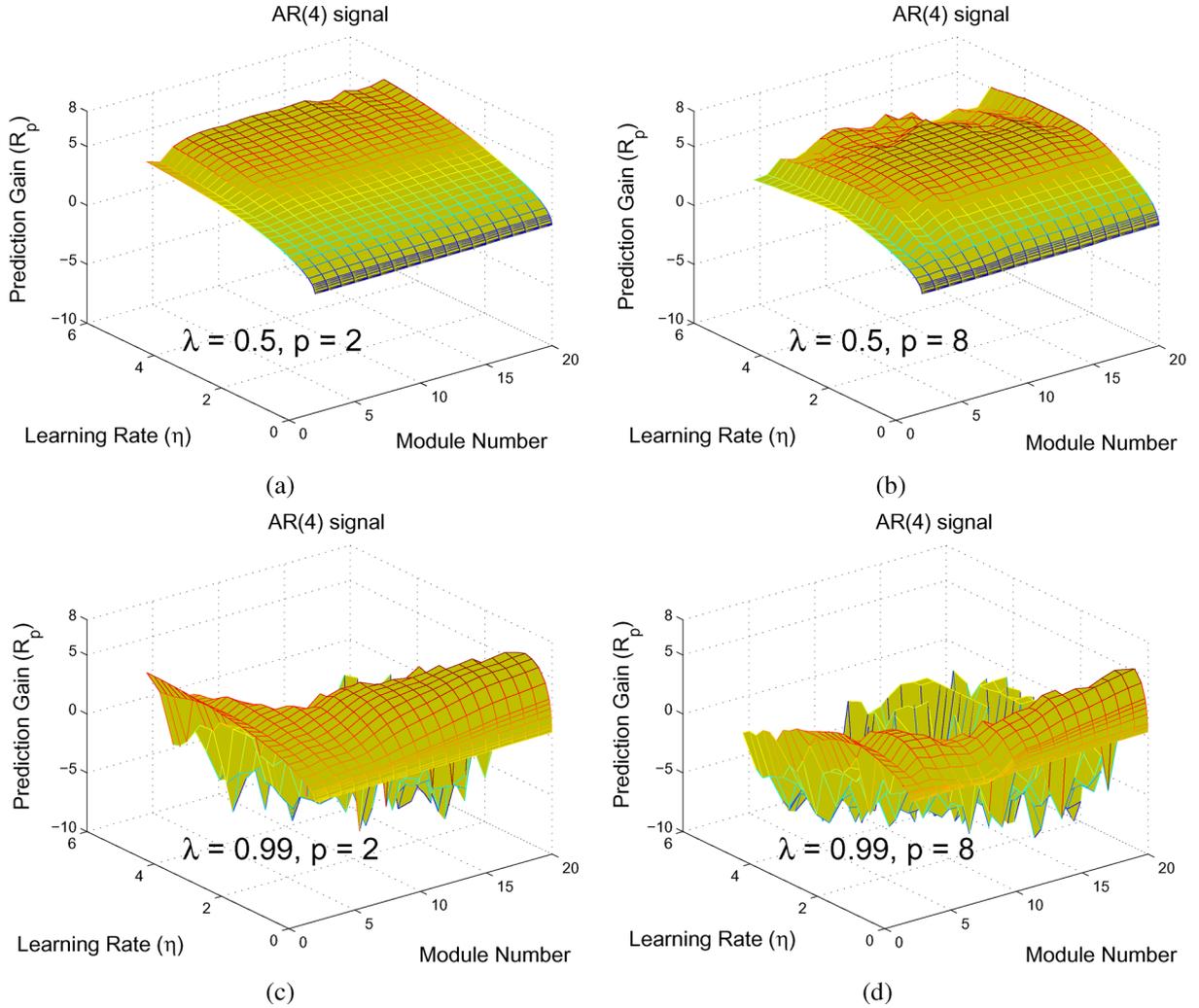


Fig. 3. Prediction performance of PRNN on a linear benchmark AR(4) signal.

TABLE I
PARAMETER SETTINGS OF THE PRNN AT THE MAXIMAL PREDICTION GAIN FOR AR(4) AND NARENDRA MODEL 3 SIGNAL

p	λ	AR(4)			Narendra		Model 3
		M	η	R_p^{max} (dB)	M	η	R_p^{max} (dB)
2	0.50	14	5.0	5.83297	20	3.6	1.31530
2	0.99	19	1.0	6.08877	20	0.4	1.34499
8	0.50	14	3.4	5.71395	20	2.0	1.45836
8	0.99	20	0.4	5.20809	20	0.2	1.00176

the AR(4) signal (R_p^{max}) with a larger number of modules at smaller η . When λ was fixed, with the increase of p , PRNN achieved its best performance at smaller η . These findings are consistent with the previous analysis and discussion.

We next conducted a similar set of experiments on a nonlinear benchmark signal (2). From Table I, when p was fixed, with the increase of λ , the PRNN achieved its best performance at smaller η . When λ was fixed, with the increase of p , PRNN achieved its best performance at smaller η . At all the settings, PRNN always had its highest R_p for a larger number of modules.

B. Assessment of the Qualitative Performance of the PRNN

Next, in order to ascertain whether high quantitative performance yields good qualitative performance, we examined the possible change

in the nature of a signal processed by the PRNN. The experimental setting was the same as the one used previously.

The left two diagrams in Fig. 4 illustrate the DVV scatter diagram for the PRNN consisting of five modules, evaluated on the prediction on the Narendra model 3 signal for different forgetting factors (λ). From Fig. 4, R_p decreased with the increase in λ whereas the qualitative performance improved as indicated by the decrease in ε (the DVV scatter diagrams of the original and predicted signal being closer). This can be explained in the following context: with fewer modules, the PRNN cannot capture the full dynamics of the Narendra model 3, but as the forgetting factor λ increases, the remote modules start to contribute more significantly to the weight update. Thus, the PRNN will learn more about the dynamics of the signal being predicted, which contributes to the improvement in preserving the signal nature. As for the

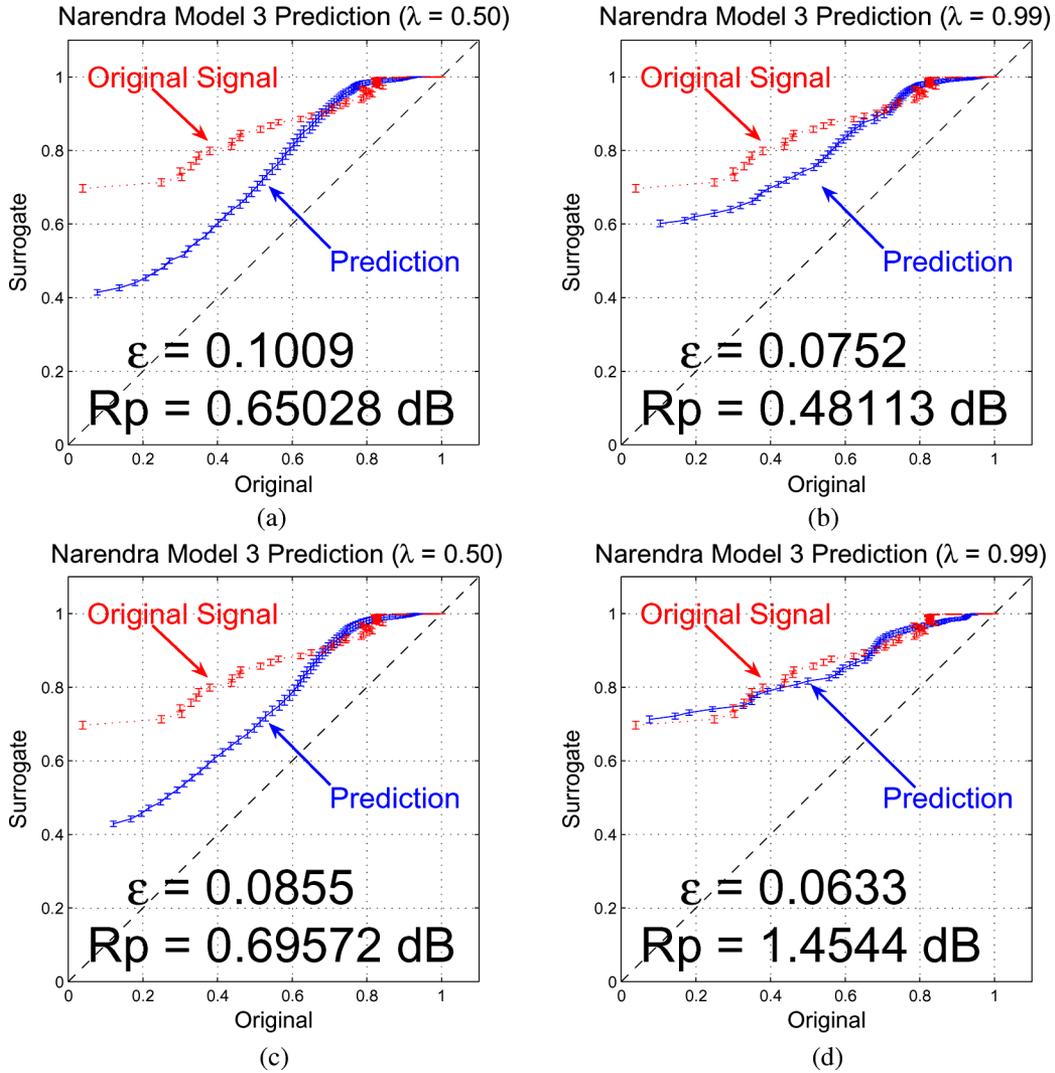


Fig. 4. Qualitative and quantitative comparison of the PRNN with five modules for the one-step-ahead prediction of the nonlinear benchmark signal (Narendra model 3) with different λ and M . (a) $M = 5$ and $\lambda = 0.5$. (b) $M = 5$ and $\lambda = 0.99$. (c) $M = 10$ and $\lambda = 0.5$. (d) $M = 10$ and $\lambda = 0.99$.

TABLE II
PREDICTION GAIN (R_p) OF INDIVIDUAL MODULES FOR THE PRNN WITH FIVE MODULES ON PREDICTING THE NARENDR MODEL 3 SIGNAL, λ VARIES FROM 0.5 TO 0.99. THE 1ST MODULE IS THE OUTPUT OF THE PRNN

λ	1 st	2 nd	3 rd	4 th	5 th	(Module)
0.5	0.6503 dB	2.946 dB	4.749 dB	4.266 dB	3.227 dB	
0.99	0.4811 dB	2.688 dB	4.462 dB	5.561 dB	5.085 dB	

decrease in the prediction gain, this is mainly due to the fact that since using five modules is not enough to capture the nonlinear nature of the signal, emphasizing the contribution of remote modules in fact reduces the portion of the contribution of the first module, which is the main drive for the weight update. That is why for the Narendra model 3, the PRNN with five modules shows an increasingly better nature-preserving capability while it has worse prediction gain with the increase in λ , as illustrated by a decrease in ϵ . This also demonstrates that a high quantitative performance is not necessarily followed by a high qualitative performance. As the number of modules M increased from five to ten, the PRNN obtained much more information about the signal in hand, which not only caused the prediction gain to increase, but also the nature of the processed signal was better preserved, as illustrated in Fig. 4(c) and (d), where the DVV scatter diagrams become very close.

In the next experiment, we investigated the prediction performance of individual modules. Table II illustrates the prediction performance for individual modules of a PRNN with five modules on predicting Narendra model 3 signal. Observe that when $\lambda = 0.5$, the third module yields the highest R_p ; when $\lambda = 0.99$, it was the fourth module. In our previous findings [13], the most remote module, in this case, the fifth module, will exhibit most pronounced “data-reusing” effect, and is supposed to have the highest R_p . In fact, it is the competition of the “forgetting” and “data-reusing” effect that determines the prediction performance of the individual module.

We will now finally analyze the qualitative performance for PRNN on prediction of the linear benchmark signal (1). From Fig. 5, the PRNN was able to preserve the nature of the processed signal, as illustrated by the fact that all the DVV scatter diagrams coincided with the bisector

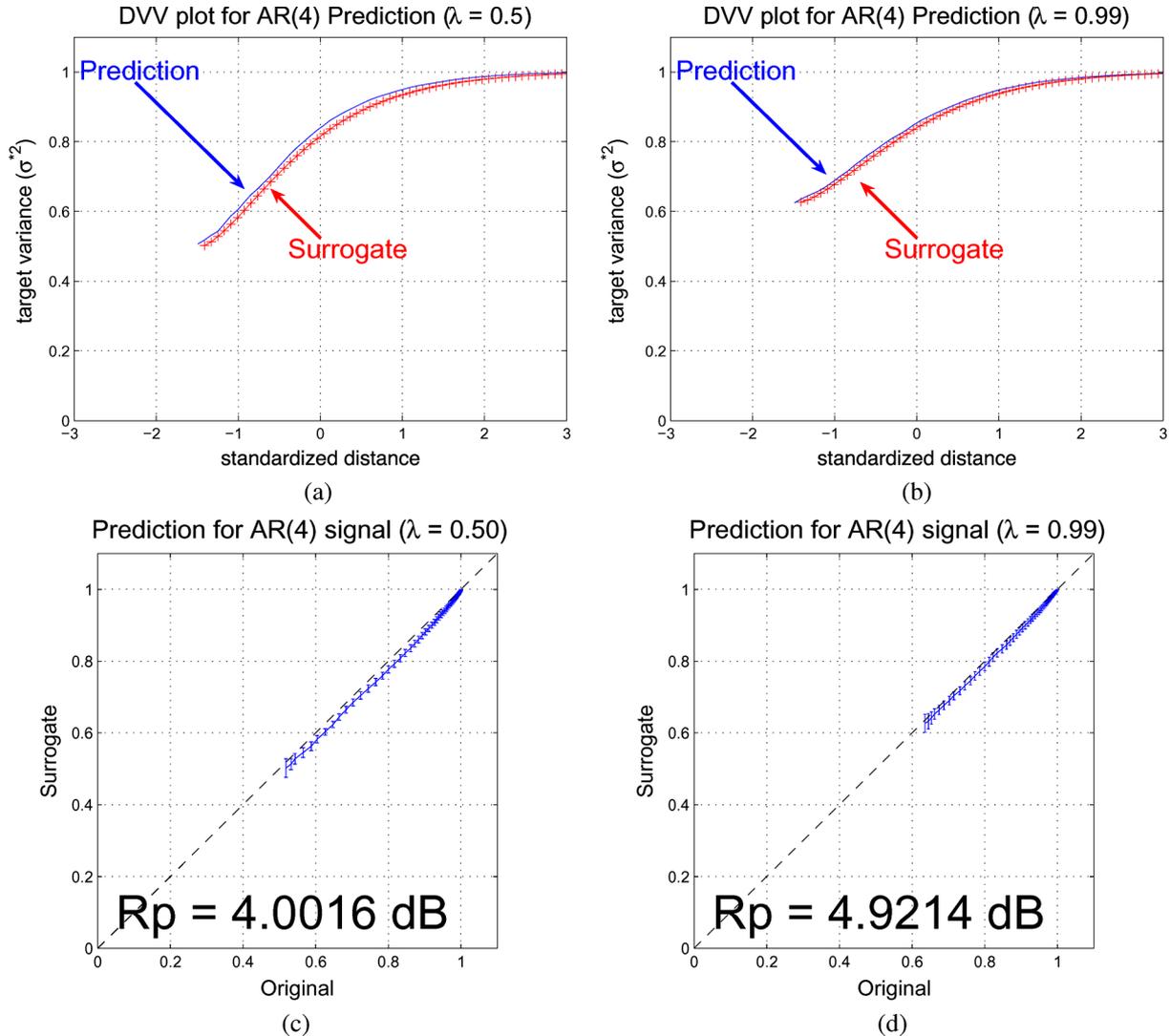


Fig. 5. Qualitative and quantitative comparison of the PRNN for the prediction of the linear benchmark AR(4) signal. λ varies from 0.5 to 0.99.

line, which indicates the preservation of the linear nature of the original signal. However, the noise level increases with the increase in the forgetting factor, illustrated by the leftmost point (minimal target variance, σ^{*2}) in the DVV plots starting further right. This is again a demonstration that a large forgetting factor makes the weights jitter, which, in turn, introduces stochastic noise into the predicted signal.

V. CONCLUSION

In this letter, we have provided new insight into the performance of online machine learning architectures. For generality, this is achieved for the PRNN, and the analysis is conducted based upon examining one-step-forward prediction on both linear and nonlinear benchmark signals and for a range of parameter settings. The qualitative performance assessment is achieved by examining whether there is a change in the nature of the processed signal, based upon the recently proposed “delay vector variance” (DVV) method for signal modality characterization. The qualitative performance analysis of simpler architectures follows naturally, by using only one module, or by canceling feedback, or by using linear neurons. It has been shown that there is a need for a tradeoff between the qualitative and quantitative performance index especially when the nature of a signal conveys some important information, for instance, health hazards in medical engineering.

REFERENCES

- [1] C. S. Poon and C. K. Merrill, “Decrease of cardiac chaos in congestive heart failure,” *Nature*, vol. 389, pp. 492–495, 1997.
- [2] S. Haykin and L. Li, “Nonlinear adaptive prediction of nonstationary signals,” *IEEE Trans. Signal Process.*, vol. 43, no. 2, pp. 526–535, Feb. 1995.
- [3] D. P. Mandic and J. A. Chambers, “On the choice of parameters of the cost function in nested modular RNNs,” *IEEE Trans. Neural Netw.*, vol. 11, no. 2, pp. 315–322, Mar. 2000.
- [4] T. Gautama, D. P. Mandic, and M. M. Van Hulle, “The delay vector variance method for detecting determinism and nonlinearity in time series,” *Physica D*, vol. 190, no. 3–4, pp. 167–176, 2004.
- [5] T. Gautama, D. P. Mandic, and M. M. Van Hulle, “Indications of nonlinear structures in brain electrical activity,” *Phys. Rev. E, Stat. Phys. Plasmas Fluids Relat. Interdiscip. Top.*, vol. 67, pp. 046204-1–046204-5, 2003.
- [6] D. P. Mandic and J. A. Chambers, *Recurrent Neural Networks for Prediction: Learning Algorithms, Architectures and Stability*. New York: Wiley, 2001.
- [7] T. Schreiber and A. Schmitz, “Improved surrogate data for nonlinearity tests,” *Phys. Rev. Lett.*, pp. 635–638, 1996.
- [8] M. C. Casdagli and A. S. Weigend, “Exploring the continuum between deterministic and stochastic modeling,” in *Time Series Prediction: Forecasting the Future and Understanding the Past*, A. S. Weigend and N. A. Gershenfeld, Eds. Reading, MA: Addison-Wesley, 1994, pp. 347–367.
- [9] D. Kaplan, “Exceptional events as evidence for determinism,” *Physica D*, vol. 73, no. 1, pp. 38–48, 1994.

- [10] L. Cao, "Practical method for determining the minimum embedding dimension of a scalar time series," *Physica D*, vol. 110, pp. 43–50, 1997.
- [11] A. Savran, "Multifeedback-layer neural network," *IEEE Trans. Neural Netw.*, vol. 18, no. 2, pp. 373–384, Mar. 2007.
- [12] Z. Hou, M. Gupta, P. Nikiforuk, M. Tan, and L. Cheng, "A recurrent neural network for hierarchical control of interconnected dynamic systems," *IEEE Trans. Neural Netw.*, vol. 18, no. 2, pp. 466–481, Mar. 2007.
- [13] M. Chen, T. Gautama, M. M. Van Hulle, and D. P. Mandic, "On non-linear modular neural filters," in *Proc. IEEE Int. Conf. Acoust. Speech Signal Process. (ICASSP)*, 2005, vol. 5, pp. 317–320.
- [14] D. P. Mandic, J. Baltersee, and J. A. Chambers, "Non-linear adaptive prediction of speech with a pipelined recurrent neural network and advanced learning algorithms," in *Signal Analysis and Prediction*, A. Prochazka, J. Uhler, P. W. Rayner, and N. G. Kingsbury, Eds. Boston, MA: Birkhauser, 1998, vol. 5.

Recursive Support Vector Machines for Dimensionality Reduction

Qing Tao, Dejun Chu, and Jue Wang

Abstract—The usual dimensionality reduction technique in supervised learning is mainly based on linear discriminant analysis (LDA), but it suffers from singularity or undersampled problems. On the other hand, a regular support vector machine (SVM) separates the data only in terms of one single direction of maximum margin, and the classification accuracy may be not good enough. In this letter, a recursive SVM (RSVM) is presented, in which several orthogonal directions that best separate the data with the maximum margin are obtained. Theoretical analysis shows that a completely orthogonal basis can be derived in feature subspace spanned by the training samples and the margin is decreasing along the recursive components in linearly separable cases. As a result, a new dimensionality reduction technique based on multilevel maximum margin components and then a classifier with high accuracy are achieved. Experiments in synthetic and several real data sets show that RSVM using multilevel maximum margin features can do efficient dimensionality reduction and outperform regular SVM in binary classification problems.

Index Terms—Classification, dimensionality reduction, feature extraction, projection, recursive support vector machines (RSVMs), support vector machines (SVMs).

I. INTRODUCTION

Dimensionality reduction is an important preprocessing step in many applications of data mining, machine learning, and pattern recognition, due to the so-called curse of dimensionality [1], [2]. Now, principal component analysis (PCA, [3]) and linear discriminant analysis (LDA,

[4]) are regarded as the most fundamental and powerful tools of dimensionality reduction for extracting effective features from high-dimensional vectors of input data. From the point of view of mathematics, PCA is an orthogonal transformation of the coordinate system in which we describe our data. The new coordinate values by which we represent the data are called principal components. Usually, a small number of principal components is sufficient to account for most of the structure in the data. From the viewpoint of pattern recognition, LDA aims to find the optimal discriminant vectors (and then, an orthogonal transformation) by maximizing the ratio of the between-class distance to the within-class distance, thus achieving the maximum class discrimination. LDA is the benchmark for the linear discrimination between two classes in multidimensional space. One of the most obvious differences between PCA and LDA is that the former does not employ the labels of all samples while the latter does.

Around 1997, several comparative studies between LDA and PCA on the face recognition problems were reported independently by numerous authors [5], [6], in which LDA outperformed PCA significantly. So far, LDA has proven to be a more efficient approach for extracting features for many pattern classification problems as compared to PCA. However, there exists a serious limitation for using LDA to solve high-dimensional recognition with finite samples. Usually, LDA requires the so-called total scatter matrix to be nonsingular. In many applications, especially in face recognition, all scatter matrices in question can be singular since the data points are from a very high-dimensional space and, in general, the sample size does not exceed this dimensionality. This is known as the *singularity* or *undersampled problem* [8] and inevitably gives rise to a problem of unstable numerical computation. In recent years, many approaches have been proposed to deal with such high-dimensional undersampled problems, including null space LDA and orthogonal LDA, and their detailed computational and theoretical analysis can be seen in [9]. Recently, a recursive LDA for calculating the discriminant features was suggested in [10]. This new algorithm incorporates the same fundamental idea behind LDA of seeking the projection that best separates the data corresponding to different classes, while in contrast to regular LDA, the features are obtained recursively and the number of features that may be derived is independent of the number of the classes to be recognized. Extensive experiments of comparing the recursive LDA algorithm with the traditional approaches have been carried out on face recognition problems, in which the resulting improvement of the performances by the new feature extraction scheme is significant. Obviously, how to employ the recursive idea to get a dimensionality reduction approach without undersampled problems is very interesting.

In the last few years, there have been very significant developments in the understanding of support vector machines (SVMs) and statistical learning theory [11]–[13]. In appearance, the geometric interpretation of a linear SVM, known as the maximum margin algorithm, is very clear. In theory, increasing margin has been shown to improve the generalization performance. In [14], an SVM-like framework was established for LDA and it was proved that the general framework of LDA is based on the simplest and most intuitive LDA with zero within-class variance. Further, it can be found that LDA and SVM are closely related. Commonly, they all try to seek the projection that best separates the data in terms of a specific objective function. Along the former direction, the within-class variance is minimized while between-class variance is maximized. Along the latter, the between-class distance is maximized with the large margin while within-class distance is not considered. Since the usual dimensionality reduction technique in supervised learning is mainly based on using a small number of orthogonal

Manuscript received November 11, 2006; revised February 7, 2007; accepted July 2, 2007. This work was supported by the National Basic Research Program 2004CB318103 and the National Science Foundation of China under Grant 60575001.

Q. Tao is with the Key Laboratory of Complex Systems and Intelligence Science, Institute of Automation, Chinese Academy of Sciences, Beijing 100080, P.R. China and also with the New Star Research Institute of Applied Technology, Hefei 230031, P.R. China (e-mail: qing.tao@mail.ia.ac.cn; taoqing@gmail.com).

J. Wang is with the Key Laboratory of Complex Systems and Intelligence Science, Institute of Automation, Chinese Academy of Sciences, Beijing 100080, P.R. China (e-mail: jue.wang@mail.ia.ac.cn).

D. Chu is with the New Star Research Institute of Applied Technology, Hefei, 230031 P.R. China.

Digital Object Identifier 10.1109/TNN.2007.908267