

Parsimonious Random Vector Functional Link Network for Data Streams

Mahardhika Pratama^a, Plamen P. Angelov^b, Edwin Lughofer^c, Deepak Puthal^d

^a*School of Computer Science and Engineering, Nanyang Technological University, Singapore, 639798, Singapore*

^b*School of Computing and Communication, Lancaster University, Lancaster, UK*

^c*Department of Knowledge-based Mathematical Systems, Johannes Kepler University, Linz, Austria*

^d*School of Electrical and Data Engineering, University of Technology, Sydney, Australia*

Abstract

The majority of the existing work on random vector functional link networks (RVFLNs) is not scalable for data stream analytics because they work under a batched learning scenario and lack a self-organizing property. A novel RVFLN, namely the parsimonious random vector functional link network (pRVFLN), is proposed in this paper. pRVFLN adopts a fully flexible and adaptive working principle where its network structure can be configured from scratch and can be automatically generated, pruned and recalled from data streams. pRVFLN is capable of selecting and deselecting input attributes on the fly as well as capable of extracting important training samples for model updates. In addition, pRVFLN introduces a non-parametric type of hidden node which completely reflects the real data distribution and is not constrained by a specific shape of the cluster. All learning procedures of pRVFLN follow a strictly single-pass learning mode, which is applicable for online time-critical applications. The advantage of pRVFLN is verified through numerous simulations with real-world data streams. It was benchmarked against recently published algorithms where it demonstrated comparable and even higher predictive accuracies while imposing the lowest complexities.

Email addresses: `d.pratama@ntu.edu.sg` (Mahardhika Pratama), `p.angelov@lancaster.ac.uk` (Plamen P. Angelov), `edwin.lughofer@jku.at` (Edwin Lughofer), `Deepak.Puthal@uts.edu.au` (Deepak Puthal)

Keywords: Random Vector Functional Link, Evolving Intelligent System, Online Learning, Online Identification, Randomized Neural Networks

1. Introduction

For decades, research in artificial neural networks has mainly investigated the best way to determine network-free parameters, which produces a model with low generalization error. Various approaches were proposed, but a large volume of work is based on a first or second-order derivative approach in respect to the loss function. Due to the rapid technological progress in data storage, capture, and transmission, the machine learning community has encountered an information explosion, which calls for scalable data analytics. Significant growth of the problem space has led to a scalability issue for conventional machine learning approaches, which require iterating entire batches of data over multiple epochs. This phenomenon results in a strong demand for a simple, fast machine learning algorithm to be well-suited for deployment in numerous data-rich applications [1]. This provides a strong case for research in the area of randomness in neural networks [2, 3], which was very popular in the late 80s and early 90s. This concept offers an algorithmic framework, which allows them to generate most of the network parameters randomly while still retaining reasonable performance [3]. One of the most prominent examples of randomness in neural networks is the random vector functional link network (RVFLN) which features solid universal approximation theory under strict conditions [4].

Due to its simple but sound working principle, randomness in neural networks has regained its popularity in the current literature [5, 6, 7, 8, 9]. Nonetheless, the vast majority of work in the literature suffers from the issue of complexity which makes their computational complexity and memory burden prohibitive for data stream analytics since their complexities are manually determined and rely heavily on expert domain knowledge. These works present a model with a fixed size which lacks of adaptive mechanism to encounter changing training patterns in the data streams. The random selection of network parameters often causes the network complexity to go beyond what is necessary due to the existence of superfluous hidden nodes which contribute little to the generalization performance [25]. Although the universal approximation capability of such an approach is assured only when sufficient complexity is selected, choosing a suitable complexity for a given problem entails expert-domain knowledge and is problem-dependent.

35 A novel RVFLN, namely the parsimonious random vector functional link
 36 network (pRVFLN), is proposed. pRVFLN combines the simple and fast
 37 working principles of RFVLN where all network parameters but the output
 38 weights are randomly generated with no tuning mechanism for hidden nodes.
 39 It characterises the online and adaptive nature of evolving intelligent systems.
 40 pRVFLN is capable of tracking any variations of data streams no matter how
 41 slow, rapid, gradual, sudden or temporal the drifts in data streams because it
 42 can initiate its learning structure from scratch with no initial structure and its
 43 structure is self-evolved from data streams in the one-pass learning mode by
 44 automatically adding, pruning and recalling its hidden nodes [10]. Further-
 45 more, it is compatible for online real-time deployment because data streams
 46 are handled without revisiting previously seen samples. pRVFLN is equipped
 47 with a hidden node pruning mechanism which guarantees a low structural
 48 burden and the rule recall mechanism which aims to address cyclic concept
 49 drift. pRVFLN incorporates a dynamic input selection scenario which makes
 50 possible the activation and deactivation of input attributes on the fly and
 51 an online active learning scenario which rules out inconsequential samples
 52 from the training process. pRVFLN is a plug-and-play learner where a single
 53 training process encompasses all learning scenarios in a sample-wise man-
 54 ner without pre-and/or post-processing steps. pRVFLN offers at least four
 55 novelties: 1) it introduces the interval-valued data cloud paradigm which is
 56 an extension of the data cloud in [11]. This modification aims to induce
 57 robustness in dealing with data uncertainty caused by noisy measurement,
 58 noisy data, etc. Unlike conventional hidden nodes, the interval-valued data
 59 cloud is parameter-free and requires no parametrization. It evolves naturally,
 60 similar to real data distribution; 2) an online active learning scenario based
 61 on the sequential entropy method (SEM) is proposed. The SEM is derived
 62 from the concept of neighbourhood probability [12] but here the concept of
 63 the data cloud is integrated. The data cloud concept simplifies the sample
 64 selection process because the neighbourhood probability is inferred with ease
 65 from the activation degree of the data cloud; 3) pRVFLN is capable of au-
 66 tomatically generating its hidden nodes on the fly with the help of a type-2
 67 self-constructing clustering (T2SCC) mechanism [13, 14]. This rule growing
 68 process differs from existing approaches because the hidden nodes are cre-
 69 ated from the rule growing condition, which considers the locations of the
 70 data samples in the input space; 4) pRVFLN is capable of carrying out an
 71 online feature selection process, borrowing several concepts of online feature
 72 selection (OFS) [15]. The original version [15] is generalized here since it

73 is originally devised for linear regression and calls for some modification to
74 be a perfect fit for pRVFLN. The prominent trait of this method lies in a
75 flexible online feature selection scenario, which makes it possible to select or
76 deselect input attributes on demand by assigning crisp weights (0 or 1) to
77 input features.

78 The efficacy of pRVFLN was thoroughly evaluated using numerous real-
79 world data streams and was benchmarked against recently published algo-
80 rithms in the literature, with pRVFLN demonstrating a highly scalable ap-
81 proach for data stream analytics while retaining acceptable generalization
82 performance. An analysis of the robustness of random intervals was per-
83 formed. It is concluded that random regions should be carefully selected
84 and should be chosen close to the true operating regions of a system being
85 modelled. Moreover, we also present a sensitivity analysis of the predefined
86 threshold and study the effect of learning components. A supplemental doc-
87 ument containing additional numerical studies is also provided in ¹ and the
88 MATLAB codes of pRVFLN have been made publicly available in ² to help
89 further study. Key mathematical notations are listed in Table 1.

90 The rest of this paper is structured as follows: the network architecture of
91 pRVFLN is outlined in Section 2; the algorithmic development of pRVFLN is
92 detailed in Section 3; proof of concept is outlined in Section 4; and conclusions
93 are drawn in the last section of this paper.

94 2. Related Work

95 The concept of randomness in neural networks was initiated by Broom-
96 head and Lowe in their work on radial basis function networks (RBFNs)
97 [3]. A closed pseudo-inverse solution can be formulated to obtain the output
98 weights of the RBFN and the centres of RBF units can be randomly sampled
99 from data samples. This work later was generalized in [16], where the centre
100 of the RBF neurons can be sampled from an independent distribution of the
101 training data. The randomness in neural networks was substantiated by the
102 findings of White [9], who developed a statistical test on hidden nodes. It
103 was found that some nonlinear structures in the mapping function can be
104 neglected without substantial loss of accuracy. In [9], the input weights of

¹https://www.dropbox.com/s/lytpt4huqyoqa6p/supplemental_document.docx?dl=0

²<https://www.dropbox.com/sh/zbm54epk8trlnq9/AACgxLHt5Nsy7MISbgXcdkTba?dl=0>

105 the hidden layers are randomly chosen. It is shown that the input weights
106 are not sensitive to the overall learning performance.

107 A prominent contribution was made by Pao et al. with the random vector
108 functional link network (RVFLN) [17]. This work presents a specific case of
109 the functional link neural network [18], which embraces the concept of ran-
110 domness in the functional link network. Note that a closed pseudo-inversion
111 solution can be also defined for the RVFLN in lieu of the conjugate gradient
112 (CG) approach. The universal approximation capability of the RVFLN is
113 proven in [4] by formalising the Monte Carlo method approximating a limit-
114 integral representation of a function. To attain the universal approximation
115 capability, the hidden node should be chosen as either absolutely integrable or
116 differentiable function. In practise, the region of random parameters should
117 also be chosen carefully and the number of hidden nodes should be sufficiently
118 large. There also exists another research direction in this area, namely reser-
119 voir computing (RC), which puts forward a recurrent network architecture in
120 order to take into account temporal dependencies between subsequent pat-
121 terns and in order to avoid dependencies on time-delayed input attributes
122 [19]. RC is constructed with a fixed number of recurrent layers and adopts
123 the concept of randomness in neural network where all the parameters are
124 randomly generated except the output weight. A comprehensive survey of
125 randomness in neural network can be found in [2, 6].

126 Since the last decade, RVFLN has transformed into one of the most vi-
127 brant fields in the neural network community as evidenced by its numerous
128 extensions and variations. The vast majority of RNNs in the literature are
129 not compatible with online real-time learning situations because it requires a
130 complete dataset to be collected after which it performs a one-shot learning
131 process based on a closed pseudo-inverse solution. This issue led to the de-
132 velopment of online learning in RVFLNs, which follows a single-pass learning
133 concept [20, 21]. Nevertheless, this work still lacks the capability to cope
134 with changing training patterns because they are built upon a fixed network
135 structure which cannot evolve in accordance with up-to-date data trends.
136 Several concepts of dynamic structure were offered in [22] and [8] by putting
137 forward the notion of a growing structure. Notwithstanding their dynamic
138 natures, concept drift remains an uncharted territory in these works because
139 all parameters are chosen at random without paying close attention to the
140 true data distribution. RC aims to address temporal system dynamics [19]
141 but still does not consider a possible dramatic change of system behaviour.
142 To the best of our knowledge, existing RC algorithms still suffers from the

Table 1: Key Mathematical Notations

Symbol	Description
$A_t \in \mathfrak{R}^n$	The input weight vector
β_t	The output of expansion layer
$X_t \in \mathfrak{R}^n$	The input attribute
$T_t \in \mathfrak{R}^n$	The target attribute
$x_e \in \mathfrak{R}^{2n+1}$	The expanded input vector
$w_i \in \mathfrak{R}^{2n+1}$	The output weight vector
B_t	The network bias
$\tilde{G}_{i,temporal}$	The interval-valued temporal firing strength
$q \in \mathfrak{R}^m$	The design factor
$\lambda \in \mathfrak{R}^R$	The recurrent weight vector
$\tilde{\mu}_i \in \mathfrak{R}^n$	The interval-valued local mean
$\tilde{\Sigma}_i \in \mathfrak{R}^n$	The interval-valued mean square length
$\delta_i \in \mathfrak{R}^n$	The uncertainty factor
$H(N X_n)$	The entropy of neighborhood probability
$I_c(\tilde{\mu}_i, X_t)$	The input coherence
$O_c(\tilde{\mu}_i, X_t)$	The output coherence
$\zeta()$	The correlation measure
$\zeta(\tilde{G}_{i,temp}, T_t)$	The mutual information between $i - th$ rule and the target concept
Ψ_i	The output covariance matrix

143 absence of self-organizing mechanism. The problem of uncertainty is another
144 open issue in the existing literature since most work utilises a crisp activation
145 function generating certain activation degrees. Such functions lack a degree
146 of tolerance against imprecision, inaccuracy and uncertainty in the training
147 data. It is worth noting that uncertainty occurs for a number of reasons:
148 noisy measurement, noisy data, false sensor reading, etc.

149 3. Basic Concepts

150 This section outlines the foundations of pRVFLN encompassing the basic
151 concept of RVFLN [17], the use of the Chebyshev polynomial as the func-

152 tional expansion block [23] and the concept of data clouds [24].

153 3.1. Random Vector Functional Link Network

154 The idea of RVFLN was proposed by Pao in [17] and is one of the forms of
155 the functional link network combined with the random vector approach [18].
156 It starts with the fact that while the network parameters are set as random
157 pairings of points, the training set can be still learned very well, although it
158 does not remove the inherent nature of random process. It features the en-
159 hancement node performing the nonlinear transformation of input attributes
160 as well as the direct connection of input attributes to the output node. The
161 activation degree of the enhancement node along with the input attributes
162 is combined with a set of output weights to generate the final network out-
163 put. The RVFLN only leaves the weight vector to be fine-tuned during the
164 training process while the other parameters are randomly sampled from a
165 carefully selected scope. Suppose that there are J enhancement nodes and
166 N input attributes, the size of the output weight vector is $W \in \mathfrak{R}^{(J+N)}$. The
167 quadratic optimization problem is then formulated as follows:

$$E = \frac{1}{2P} \sum_{p=1}^P (t^{(p)} - B^t d^{(p)})^2 \quad (1)$$

168 where $B \in \mathfrak{R}^{(N+J)}$ is the output weight vector containing the N-dimensional
169 original input vector also in addition to the weight values. $d^{(p)}$ is the output
170 of the enhancement node. The RVFLN is similar to a single hidden layer
171 feedforward network except for the fact that the hidden node functions as an
172 enhancement of the input feature and there exists direct connection from the
173 input layer to the output layer. The steepest descent approach can be used to
174 fine-tune the output weight vector. If matrix inversion using pseudo-inverse
175 is feasible, a closed-form solution can be formulated. The generalization
176 performance of RVFLN was examined in [17] where RVFL can be trained
177 rapidly with ease. The RVFLNs convergence is also guaranteed to be attained
178 within a number of iterations.

179 The RVFL can be modified by incorporating the idea of the functional
180 link network [23]. That is, the hidden node or the enhancement node is
181 replaced by the functional expansion block generating a set of linearly in-
182 dependent functions of the entire input pattern. The functional expansion
183 block can be formulated as trigonometric expansion [25], Chebyshev expan-
184 sion, legendre expansion, etc. [23] but our scope of discussion is limited to

185 the Chebyshev expansion only due to its relevance to pRVFLN. Given the N -
 186 dimensional input vector $X = [x_1, x_2, \dots, x_N] \in \mathfrak{R}^{1 \times N}$ and its corresponding
 187 m -dimensional target vector $Y = [y_1, y_2, \dots, y_m] \in \mathfrak{R}^{1 \times m}$, the output of RVFLN
 188 with the Chebyshev functional expansion block is expressed as follows:

$$y = \sum_{j=1}^{2N+1} B_j \phi_j(A_N X_N + b_N) \quad (2)$$

189 where B_j is the output weight vector and $\phi_j()$ is the Chebyshev functional
 190 expansion mapping the N -dimensional input attribute and the input weight
 191 vector to the higher $2N + 1$ expansion space. As with the original RVFLN,
 192 the output weight vector can be learned using any optimization method. The
 193 $2N + 1$ here results from the utilisation of the Chebyshev series up to the
 194 second order. The Chebyshev series is mathematically written as follows:

$$T_{n+1} = 2xT_n(x) - T_{n-1}(x) \quad (3)$$

195 If we are only interested in the Chebyshev series up to the second order,
 196 this results in $T_0(x) = 1, T_1(x) = x, T_2(x) = 2x^2 - 1$. The advantage of the
 197 Chebyshev functional link compared to other popular functional links such as
 198 trigonometric [25], legendre, power function, etc. [23] lies in its simplicity of
 199 computation. The Chebyshev function scatters fewer parameters to be stored
 200 into memory than the trigonometric function, while the Chebyshev function
 201 has a better mapping capability than the other polynomial functions of the
 202 same order. In addition, the polynomial power function is not robust against
 203 an extrapolation case.

204 3.2. Data Cloud

205 The concept of the data cloud offers an alternative to the traditional
 206 cluster concept where it is not shape-specific and evolves naturally in ac-
 207 cordance with the true data distribution. It is also easy to use because it
 208 is non-parametric and does not require any parameterization. This strat-
 209 egy is desirable because parameterization per scalar variable often calls for
 210 complex high-level approximation and/or optimization. This approach was
 211 inspired by the idea of RDE and was integrated in the context of the TSK
 212 fuzzy system [11, 24]. Unlike a conventional fuzzy system where a degree
 213 of membership is defined by a point-to-point distance, the data cloud com-
 214 putes an accumulated distance of the point of interest to all other points in

215 the data cloud without physically keeping all data samples in the memory
 216 similar to the local data density. This notion has a positive impact on the
 217 memory and space complexity because the number of network parameters
 218 significantly reduces. The data cloud concept is formally written as:

$$\gamma_k^i = \frac{1}{1 + \|x_k - \mu_k^L\|^2 + \Sigma_k^L - \|\mu_k^L\|^2} \quad (4)$$

219 where γ_k^i denotes the i -th data cloud at the k -th observation. The data cloud
 220 evolves by updating the local mean μ_k^L and square length of i -th local region
 221 Σ_k^L as follows:

$$\mu_k^L = \left(\frac{M_k^i - 1}{M_k^i}\right)\mu_{k-1}^L + \frac{x_k}{M_k^i}, \mu_1^L = x_1 \quad (5)$$

222

$$\Sigma_k^L = \frac{M_k^i - 1}{M_k^i}\Sigma_{k-1}^L + \frac{\|x_k\|^2}{M_k^i}, \Sigma_1^L = \|x_1\|^2 \quad (6)$$

223 It is worth noting that these two parameters correspond to statistics of the
 224 i -th data cloud and are computed recursively with ease using standard re-
 225 cursive formulas. They do not impose a specific optimization or a specific
 226 setting to be performed to adjust their values.

227 4. Network Architecture of pRVFLN

228 pRVFLN utilises a local recurrent connection at the hidden node which
 229 generates the spatiotemporal activation degree. This recurrent connection
 230 is realized by a self-feedback loop of the hidden node which memorizes the
 231 previous activation degree and outputs a weighted combination between pre-
 232 vious and current activation degrees spatiotemporal firing strength. In the
 233 literature, there exist at least three types of recurrent network structures re-
 234 ferring to its recurrent connections: global [26, 27], interactive [25], and local
 235 [28], but the local recurrent connection is deemed to be the most compati-
 236 ble recurrent type in our case because it does not harm the local property,
 237 which assures stability when adding, pruning and fine-tuning hidden nodes.
 238 pRVFLN utilises the notion of the functional-link neural network where the
 239 expansion block is created by the Chebyshev polynomial up to the second
 240 order. Furthermore, the hidden layer of pRVFLN is built upon an interval-
 241 valued data cloud [11] where we integrate the idea of an interval-valued local
 242 mean into the data cloud.

243 Suppose that a pair of data points (X_t, T_t) is received at t -th time instant
 244 where $X_t \in \mathfrak{R}^n$ is an input vector and $T_t \in \mathfrak{R}^m$ is a target vector, while n
 245 and m are respectively the number of input and output variables. Because
 246 pRVFLN works in a strictly online learning environment, it has no access
 247 to previously seen samples, and a data point is simply discarded after being
 248 learned. Due to the pre-requisite of an online learner, the total number of
 249 data N is assumed to be unknown. The output of pRVFLN is defined as
 250 follows:

$$y_o = \sum_{i=1}^R \beta_i \tilde{G}_{i,temporal}(A_t X_t + B_t), \tilde{G}_{temporal} = [\underline{G}, \overline{G}] \quad (7)$$

251 where R denotes the number of hidden nodes and β_i stands for the i -th
 252 output of the functional expansion layer, produced by weighting the weight
 253 vector with an extended input vector $\beta_i = x_e^T w_i$. $x_e \in \mathfrak{R}^{(2n+1) \times 1}$ is an
 254 extended input vector resulting from the functional link neural network based
 255 on the Chebyshev function up to the second order [23] as shown in (3) and
 256 $w_i \in \mathfrak{R}^{(2n+1) \times 1}$ is a connective weight of the i -th output node. The definition
 257 of β_i is rather different from its common definition in the literature because
 258 it adopts the concept of the expansion block, mapping a lower dimensional
 259 space to a higher dimensional space with the use of certain polynomials. This
 260 paradigm produces the extended input vector x_e as follows:

$$\nu_{p+1}(x) = 2x_j \nu_p(x_j) - \nu_{p-1}(x_j) \quad (8)$$

261 where $\nu_0(x_j) = 1, \nu_1(x_j) = x_j, \nu_2(x_j) = 2x_j^2 - 1$. Suppose that three input
 262 attributes are given $X = [x_1, x_2, x_3]$, the extended input vector is expressed as
 263 the Chebyshev polynomial up to the second order $x_e = [1, x_1, \nu_2(x_1), x_2, \nu_2(x_2),$
 264 $x_3, \nu(x_3)]$. Note that the term 1 here represents an intercept of the output
 265 node to avoid going through the origin, which may risk an untypical gradient.
 266 $A_t \in \mathfrak{R}^n$ is an input weight vector randomly generated from a certain range.
 267 B_t is removed for simplicity. $\tilde{G}_{i,temporal}$ is the i -th interval-valued data cloud,
 268 triggered by the upper and lower data cloud $\underline{G}_{i,temporal}, \overline{G}_{i,temporal}$. Note that
 269 recurrence is not seen in (7) because pRVFLN makes use of local recurrent
 270 layers at the hidden node. By expanding the interval-valued data cloud [29],
 271 the following is obtained:

$$y_o = \sum_{i=1}^R (1 - q_o) \beta_i \overline{G}_{i,temporal} + \sum_{i=1}^R q_o \beta_i \underline{G}_{i,temporal} \quad (9)$$

where $q \in \mathfrak{R}^m$ is a design factor to reduce an interval-valued function to a crisp one [29]. It is worth noting that the upper and lower activation functions $\underline{G}_{i,temporal}, \overline{G}_{i,temporal}$ deliver spatiotemporal characteristics as a result of a local recurrent connection at the i -th hidden node, which combines the spatial and temporal firing strength of the i -th hidden node. These temporal activation functions output the following.

$$\begin{aligned}\underline{G}_{i,temporal}^t &= \lambda_i \underline{G}_{i,spatial}^t + (1 - \lambda_i) \underline{G}_{i,temporal}^{t-1}, \\ \overline{G}_{i,temporal}^t &= \lambda_i \overline{G}_{i,spatial}^t + (1 - \lambda_i) \overline{G}_{i,temporal}^{t-1}\end{aligned}\tag{10}$$

272 where $\lambda \in \mathfrak{R}^R$ is a weight vector of the recurrent link. The local feedback
 273 connection here feeds the spatiotemporal firing strength at the previous time
 274 step $\tilde{G}_{i,temporal}^{t-1}$ back to itself and is consistent with the local learning princi-
 275 ple. This trait happens to be very useful in coping with the temporal system
 276 dynamic because it functions as an internal memory component which mem-
 277 orizes a previously generated spatiotemporal activation function at $t - 1$.
 278 Also, the recurrent network is capable of overcoming over-dependency on
 279 time-delayed input features and lessens strong temporal dependencies of sub-
 280 sequent patterns. This trait is desired in practise since it may lower the input
 281 dimension, because prediction is done based on the most recent measurement
 282 only. Conversely, the feedforward network often relies on time-lagged input
 283 attributes to arrive at a reliable predictive performance due to the absence
 284 of an internal memory component. This strategy at least entails expert
 285 knowledge for system order to determine the suitable number of delayed
 286 components.

287 The hidden node of the pRVFLN is an extension of the cloud-based hidden
 288 node, where it embeds an interval-valued concept to address the problem of
 289 uncertainty [30]. Instead of computing an activation degree of a hidden node
 290 to a sample, the cloud-based hidden node enumerates the activation degree
 291 of a sample to all intervals in a local region on-the-fly. This results in local
 292 density information, which fully reflects real data distributions. This concept
 293 was defined in AnYa [11, 24]. This concept is also the underlying component
 294 of AutoClass and TEDA-Class [31], all of which come from Angelovs sound
 295 work of RDE [24]. This paper aims to modify these prominent works to the
 296 interval-valued case. Suppose that N_i denotes the support of the i -th data
 297 cloud, an activation degree of i -th cloud-based hidden node refers to its local

298 density estimated recursively using the Cauchy function:

$$\tilde{G}_{i,spatial} = \frac{1}{1 + \sum_{k=1}^{N_i} \left(\frac{\tilde{x}_k - x_t}{N_i} \right)}, \quad \tilde{x}_k = [\underline{x}_{k,i}, \bar{x}_{k,i}], \quad \tilde{G}_{i,spatial} = [\underline{G}_{i,spatial}, \bar{G}_{i,spatial}] \quad (11)$$

where \tilde{x}_k is k -th interval in the i -th data cloud and x_t is t -th data sample. It is observed that (11) requires the presence of all data points seen so far. Its recursive form is formalised in [24] and is generalized here to the interval-valued case:

$$\begin{aligned} \bar{G}_{i,spatial} &= \frac{1}{1 + \|A_t^T x_t - \bar{\mu}_{i,N_i}\|^2 + \bar{\Sigma}_{i,N_i} - \|\bar{\mu}_{i,N_i}\|^2}, \\ \underline{G}_{i,spatial} &= \frac{1}{1 + \|A_t^T x_t - \underline{\mu}_{i,N_i}\|^2 + \underline{\Sigma}_{i,N_i} - \|\underline{\mu}_{i,N_i}\|^2} \end{aligned} \quad (12)$$

where $\underline{\mu}_i, \bar{\mu}_i$ signify the upper and lower local means of the i -th cloud:

$$\begin{aligned} \underline{\mu}_{i,N_i} &= \left(\frac{N_i - 1}{N_i} \right) \underline{\mu}_{i,N_i-1} + \frac{x_{i,N_i} - \Delta_i}{\|N_i\|}, \quad \underline{\mu}_{i,1} = x_{i,1} - \Delta_i, \\ \bar{\mu}_{i,N_i} &= \left(\frac{N_i - 1}{N_i} \right) \bar{\mu}_{i,N_i-1} + \frac{x_{i,k} + \Delta_i}{\|N_{i,k}\|}, \quad \bar{\mu}_{i,1} = x_{i,1} + \Delta_i \end{aligned} \quad (13)$$

where Δ_i is an uncertainty factor of the i -th cloud, which determines the degree of tolerance against uncertainty. The uncertainty factor creates an interval of the data cloud, which controls the degree of tolerance for uncertainty. It is worth noting that a data sample is considered as a population of the i -th cloud when resulting in the highest density. Moreover, $\bar{\Sigma}_{i,N_i}, \underline{\Sigma}_{i,N_i}$ are the upper and lower mean square lengths of the data vector in the i -th cloud as follows:

$$\begin{aligned} \underline{\Sigma}_{i,N_i} &= \left(\frac{N_i - 1}{N_i} \right) \underline{\Sigma}_{i,N_i-1} + \frac{\|x_{i,N_i}\|^2 - \Delta_i}{\|N_i\|}, \quad \underline{\Sigma}_{i,1} = \|x_{i,1}\|^2 - \Delta_i, \\ \bar{\Sigma}_{i,N_i} &= \left(\frac{N_i - 1}{N_i} \right) \bar{\Sigma}_{i,N_i-1} + \frac{\|x_{i,N_i}\|^2 + \Delta_i}{\|N_i\|}, \quad \bar{\Sigma}_{i,1} = \|x_{i,1}\|^2 + \Delta_i \end{aligned} \quad (14)$$

299 Although the concept of the cloud-based hidden node was generalized in
300 TeDaClass [32] by introducing the eccentricity and typicality criteria, the
301 interval-valued idea is uncharted in [32]. Note that the Cauchy function is

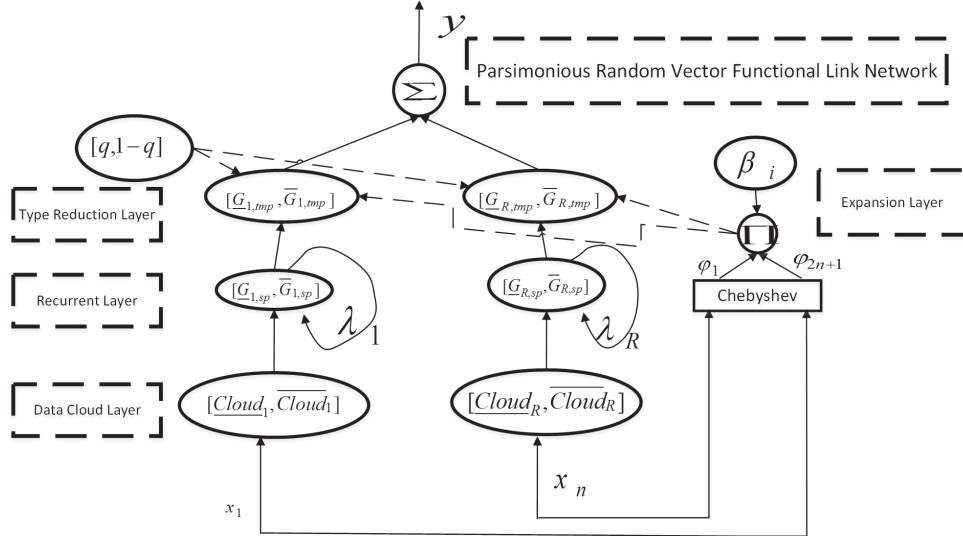


Figure 1: Network Architecture of pRVFLN

302 asymptotically a Gaussian-like function, satisfying the activation function
 303 requirement of the RVFLN to be a universal approximator.

304 Unlike conventional RVFLNs, pRVFLN puts into perspective a nonlinear
 305 mapping of the input vector through the Chebyshev polynomial up to the
 306 second order. Note that recently developed RVFLNs in the literature mostly
 307 are designed with a zero-order output node [5, 6, 7, 8]. The functional ex-
 308 pansion block expands the output node to a higher degree of freedom, which
 309 aims to improve the local mapping aptitude of the output node. pRVFLN
 310 implements the random learning concept of the RVFLN, in which all pa-
 311 rameters, namely the input weight A , design factor q , recurrent link weight
 312 λ , and uncertainty factor Δ , are randomly generated. Only the weight
 313 vector is left for parameter learning scenario w_i . Since the hidden node is
 314 parameter-free, no randomization takes place for hidden node parameters.
 315 The network structure of pRVFLN and the interval-valued data cloud are
 316 depicted in Fig. 1 and 2 respectively.

317 5. Learning Policy of pRVFLN

318 This section discusses the learning policy of pRVFLN. Section 5.1 out-
 319 lines the online active learning strategy, which deletes inconsequential sam-
 320 ples. Samples, selected in the sample selection mechanism, are fed into the

321 learning process of pRVFLN. Section 5.2 deliberates the hidden node growing
 322 strategy of pRVFLN. Section 5.3 elaborates the hidden node pruning and re-
 323 call strategy, while Section 5.4 details the online feature selection mechanism.
 324 Section 5.5 explains the parameter learning scenario of pRVFLN. Algorithm
 325 1 shows the pRVFLN learning procedure.

326 5.1. Online Active Learning Strategy

327 The active learning component of the pRVFLN is built on the extended
 328 sequential entropy (ESEM) method, which is derived from the SEM method
 329 [12]. The ESEM method makes use of the entropy of the neighborhood prob-
 330 ability to estimate the sample contribution. The underlying difference from
 331 its predecessor [12] lies in the integration of the data cloud paradigm, which
 332 greatly relieves the effort in finding the neighborhood probability because the
 333 data cloud itself is inherent with the local data density, taking into account
 334 the influence of all samples in a local region. Furthermore, it handles the
 335 regression problem which happens to be more challenging than the classifi-
 336 cation problem because the sample contribution is estimated in the absence
 337 of a decision boundary. To the best of our knowledge, only Das et al. [33]
 338 address the regression problem, but they still employ a fully supervised tech-
 339 nique because their method depends on the hinge error function to evaluate
 340 the sample contribution. The concept of neighborhood probability refers to
 341 the probability of an incoming data stream sitting in the existing data clouds:

$$P(X_i \in N_i) = \frac{\sum_{k=1}^{N_i} \frac{M(X_t, x_k)}{N_i}}{\sum_{i=1}^R \sum_{k=1}^{N_i} \frac{M(X_t, x_k)}{N_i}} \quad (15)$$

342 where X_T is a newly arriving data point and x_n is a data sample, associated
 343 with the i -th rule. $M(X_T, x_k)$ stands for a similarity measure, which can
 344 be defined as any similarity measure. The bottleneck is however caused by
 345 the requirement to revisit already seen samples. This issue can be tackled
 346 by formulating the recursive expression of (15). In the context of the data
 347 cloud, this issue becomes even simpler, because it is derived from the idea
 348 of local density and is computed based on the local mean [11]. (15) is then
 349 written as follows:

$$P(X_i \in N_i) = \frac{\Lambda_i}{\sum_{i=1}^R \Lambda_i} \quad (16)$$

350 where Λ_i is a type-reduced activation degree $\Lambda_i = (1 - q)\overline{G}_{i,spatial} + q\underline{G}_{i,spatial}$.
 351 Once the neighbourhood probability is determined, its entropy is formulated
 352 as follows:

$$H(N|X_i) = - \sum_{i=1}^R P(X_i \in N_i) \log P(X_i \in N_i) \quad (17)$$

353 Algorithm 1. Learning Architecture of pRVFLN

Algorithm 1: Parsimonious Random Vector Functional Link Network

Given a data tuple at t -th time instant $(X_t, T_t) = (x_1, \dots, x_n, t_1, \dots, t_m)$, $X_t \in \mathbb{R}^n, T_t \in \mathbb{R}R^m$; set predefined parameters α_1, α_2

*/*Step 1: Online Active Learning Strategy/**

For $i=1$ to R **do**

Calculate the neighborhood probability (8) with spatial firing strength (4)

End For

Calculate the entropy of neighborhood probability (8) and the ESEM (10)

IF (34) **Then**

*/*Step 2: Online Feature Selection/**

IF $Partial=Yes$ **Then**

Execute Algorithm 3

Else IF

Execute Algorithm 2

End IF

*/*Step 3: Data Cloud Growing Mechanism/**

For $j=1$ to n **do**

Compute $\xi(x_j, T_0)$

End For

For $i=1$ to R **do**

Calculate input coherence (12)

For $o=1$ to m **do**

Calculate $\xi(\tilde{\mu}_i, T_0)$

End For

*/*Step 4: Data Cloud Pruning Mechanism/**

For $i=1$ to R **do**

For $o=1$ to m **do**

Calculate $\xi(\tilde{G}_{i,temp}, T_0)$

End For

IF (19) **Then**

Discard i -th data cloud

End IF

End For

*/*Step 5: Adaptation of Output Weight/**

For $i=1$ to R **do**

Update output weights using FWGRLS

End For

354

355 The entropy of the neighbourhood probability measures the uncertainty
 356 induced by a training pattern. A sample with high uncertainty should be
 357 admitted for the model update, because it cannot be well-covered by an
 358 existing network structure and learning such a sample minimises uncertainty.
 359 A sample is to be accepted for model updates, provided that the following
 360 condition is met:

$$H \geq thres \quad (18)$$

361 where $thres$ is an uncertainty threshold. This parameter is not fixed dur-
 362 ing the training process, rather it is dynamically adjusted to suit the learning
 363 context. The threshold is set as $thres_{N+1} = thres_N(1 \pm inc)$, where it aug-
 364 ments $thres_{N+1} = thres_N(1 + inc)$ when a sample is admitted for the training
 365 process, whereas it decreases $thres_{N+1} = thres_N(1 - inc)$ when a sample is
 366 ruled out for the training process. inc here is a step size, set at $inc = 0.01$.
 367 This simply follows its default setting in [21].

368

369 5.2. Hidden Node Growing Strategy

370 pRVFLN relies on the T2SCC method to grow interval-valued data clouds
 371 on demand. This notion is extended from the so-called SCC method [14, 13]
 372 to adapt to the type-2 hidden node working framework. The significance of
 373 the hidden nodes in pRVFLN is evaluated by checking its input and output
 374 coherence through an analysis of its correlation to existing data clouds and
 375 the target concept. Let $\tilde{\mu}_i = [\underline{\mu}_i, \bar{\mu}_i] \in \mathfrak{R}^{1 \times n}$ be a local mean of the i -th
 376 interval-valued data cloud (5), $X_t \in \mathfrak{R}^n$ is an input vector and $T_t \in \mathfrak{R}^n$ is a
 377 target vector, the input and output coherence are written as follows:

$$I_c(\tilde{\mu}_i, X_t) = (1 - q)\zeta(\bar{\mu}_i, X_t) + q\zeta(\underline{\mu}_i, X_t) \quad (19)$$

$$378 \quad O_c(\tilde{\mu}_i, X_t) = (\zeta(X_t, T_t) - \zeta(\tilde{\mu}_i, T_t)), \quad \zeta(\tilde{\mu}_i, T_t) = (1 - q)\zeta(\bar{\mu}_i, T_t) + q\zeta(\underline{\mu}_i, T_t) \quad (20)$$

where $\zeta()$ is the correlation measure. Both linear and non-linear correlation
 measures are applicable here. However, the non-linear correlation measure
 is rather hard to deploy in the online environment, because it usually calls
 for the Discretization or Parzen Window method. The Pearson correlation
 measure is a widely used correlation measure but it is insensitive to the scaling
 and translation of variables as well as being sensitive to rotation [34]. The
 maximal information compression index (MCI) is one attempt to tackle these

problems and it is used in the T2SCC to perform the correlation measure $\zeta()$ [34]:

$$\zeta(X_1, X_2) = \frac{1}{2}(\text{var}(X_1) + \text{var}(X_2) - \sqrt{(\text{var}(X_1) + \text{var}(X_2))^2 - 4\text{var}(X_1)\text{var}(X_2)(1 - \rho(X_1, X_2)^2)}) \quad (21)$$

$$\rho(X_1, X_2) = \frac{\text{cov}(X_1, X_2)}{\sqrt{\text{var}(X_1)\text{var}(X_2)}} \quad (22)$$

379 where (X_1, X_2) are substituted with $(\bar{\mu}_i, X_t), (\underline{\mu}_t, X_t), (\bar{\mu}_i, T_t), (\underline{\mu}_t, T_t), (X_t, T_t)$
 380 to calculate the input and output correlation (19), (20). respectively stand
 381 for the variance of X, covariance of X_1 and X_2 , and Pearson correlation
 382 index of X_1 and X_2 . The local mean of the interval-valued data cloud rep-
 383 represents a data cloud because it represents a point with the highest density.
 384 In essence, the MCI method indicates the amount of information compres-
 385 sion when ignoring a newly observed sample. The MCI method features
 386 the following properties: 1) $0 \leq \zeta(X_1, X_2) \leq 0.5(\text{var}(X_1) + \text{var}(X_2))$, 2)
 387 a maximum correlation is given by $\zeta(X_1, X_2) = 0$, 3) a symmetric prop-
 388 erty $\zeta(X_1, X_2) = \zeta(X_2, X_1)$, 4) it is invariant against the translation of the
 389 dataset, and 5) it is also robust against rotation.

390 The input coherence explores the similarity between new data and ex-
 391 isting data clouds directly, while the output coherence focusses on their dis-
 392 similarity indirectly through a target vector as a reference. The input and
 393 output coherence formulates a test that determines the degree of confidence
 394 in the current hypothesis:

$$I_c(\tilde{\mu}_i, X_t) \leq \alpha_1, O_c(\tilde{\mu}_i, X_t) \geq \alpha_2 \quad (23)$$

395 where $\alpha_1 \in [0.001, 0.01], \alpha_2 \in [0.01, 0.1]$ are predefined thresholds. If a hy-
 396 pothesis meets both conditions, a new training sample is assigned to a data
 397 cloud with the highest input coherence i^* . Accordingly, the number of in-
 398 tervals N_{i^*} , local mean and square length $\tilde{\mu}_{i^*}, \tilde{\Sigma}_{i^*}$ are updated respectively
 399 with (21) and (22) as well as $N_{i^*} = N_{i^*} + 1$. A new data cloud is introduced,
 400 provided that the existing hypotheses do not pass either condition (23), that
 401 is, one of the conditions is violated. This situation reflects the fact that a new
 402 training pattern conveys significant novelty, which has to be incorporated to
 403 enrich the scope of the current hypotheses. Note that if a larger α_1 is spec-
 404 ified, fewer data clouds are generated and vice versa, whereas if a larger α_2

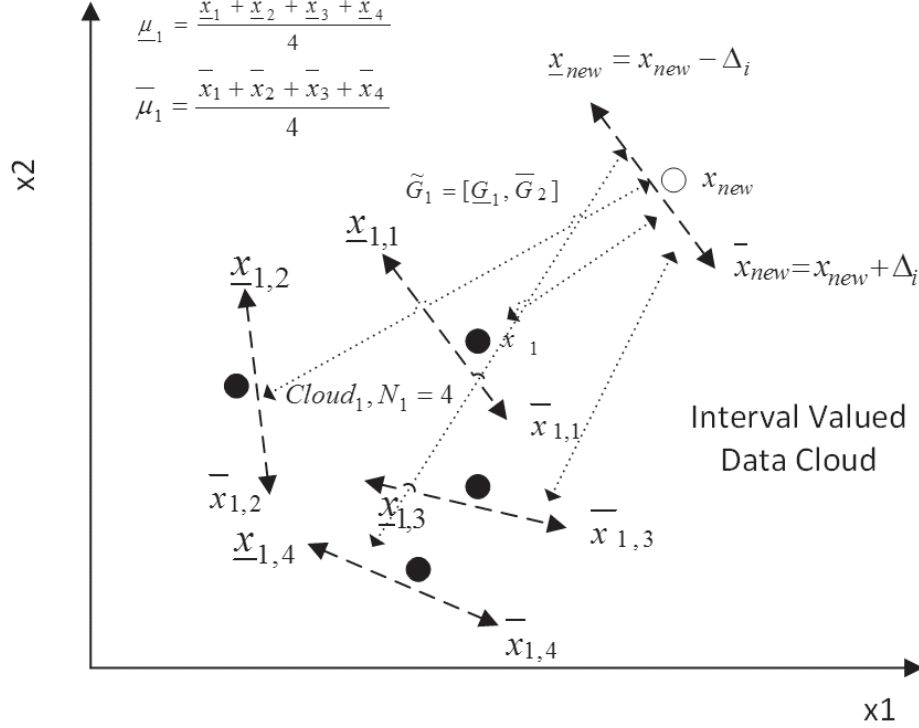


Figure 2: Interval Valued Data Cloud

405 is specified, larger data clouds are added and vice versa. The sensitivity of
 406 these two parameters is studied in the section V.E of this paper. Because a
 407 data cloud is non-parametric, no parameterization is committed when adding
 408 a new data cloud. The output node of a new data cloud is initialised:

$$W_{R+1} = W_{i^*}, \Psi_{R+1} = \bar{\omega}I \quad (24)$$

409 where $\bar{\omega} = 10^5$ is a large positive constant. The output node is set as the
 410 data cloud with the highest input coherence because this data cloud is the
 411 closest one to the new data cloud. Furthermore, the setting of covariance
 412 matrix Ψ_{R+1} leads to a good approximation of the global minimum solution
 413 of batched learning, as proven mathematically in [35].

414 5.3. Hidden Node Pruning and Recall Strategy

415 pRVFLN incorporates a data cloud pruning scenario, termed the type-
 416 2 relative mutual information (T2RMI) method. This method was firstly

417 developed in [36] for the type-1 fuzzy system. This method is convenient to
 418 apply here because it estimates mutual information between a data cloud and
 419 a target concept by analysing their correlation. Hence, the MCI method (21),
 420 (22) is valid to measure the correlation between two variables. Although this
 421 method has been well-established [36], to date, its effectiveness in handling
 422 data clouds and a recurrent structure as implemented in pRVFLN is an open
 423 question. Unlike both the RMI method that applies the classic symmetrical
 424 uncertainty method, the T2RMI method is formalised using the MCI method
 425 as follows:

$$\zeta(\tilde{G}_{i,temp}, T_t) = q\zeta(\underline{G}_{i,temp}, T_t) + (1 - q)\zeta(\overline{G}_{i,temp}, T_t) \quad (25)$$

426 where $\underline{G}_{i,temp}$, $\overline{G}_{i,temp}$ are respectively the lower and upper tempo-
 427 ral activation functions of the i -th rule. The temporal activation function
 428 is included in (25) rather than the spatial activation function in order to
 429 account for the inter-temporal dependency of subsequent training samples.
 430 The MCI method is chosen here because it possesses a significantly lower
 431 computational burden than the symmetrical uncertainty method but it is
 432 still more robust than a linear Pearson correlation index. A data cloud is
 433 deemed inconsequential, if the following is met:

$$\zeta_i < mean(\zeta_i) - 2std(\zeta_i) \quad (26)$$

434 where $mean(\zeta_i)$, $std(\zeta_i)$ are respectively the mean and standard deviation
 435 of the MCI during its lifespan. This criterion aims to capture an obsolete
 436 data cloud which does not keep up with current data distribution due to
 437 possible concept drift, because it computes the downtrend of the MCI values
 438 during its lifespan. It is worth mentioning that mutual information between
 439 hidden nodes and the target variable is a reliable indicator for changing data
 440 distributions because it monitors significance of a local region with respect
 441 to the recent data context.

442 The T2RMI method also functions as a rule recall mechanism to cope with
 443 cyclic concept drift. Cyclic concept drifts frequently happen in relation to the
 444 weather, customer preferences, electricity power consumption problems, etc.
 445 all of which are related to seasonal change. This points to a situation where
 446 a previous data distribution reappears in the current training step. Once
 447 pruned by the T2RMI, a data cloud is not forgotten permanently and is
 448 inserted into a list of pruned data clouds $R^* = R^* + 1$. In this case, its local
 449 mean, square length, population, an output node, and output covariance

450 matrix $\tilde{\mu}_{R^*}, \tilde{\Sigma}_{R^*}, N_{R^*}, \beta_{R^*}, \Psi_{R^*}$, are retained in memory. Such data clouds
 451 can be reactivated in the future, whenever their validity is confirmed by an
 452 up-to-date data trend. It is worth noting that adding a completely new data
 453 cloud when observing a previously learned concept catastrophically erases the
 454 learning history. A data cloud is recalled subject to the following condition:

$$\max_{i^*=1,\dots,R^*}(\zeta_{i^*}) > \max_{i=1,\dots,R}(\zeta_i) \quad (27)$$

455 This situation reveals that a previously pruned data cloud is more relevant
 456 than any existing ones. This condition pinpoints that a previously learned
 457 concept reappears again. A previously pruned data cloud is then regenerated
 458 as follows:

$$\tilde{\mu}_{R+1} = \tilde{\mu}_{R^*}, \tilde{\Sigma}_{R+1} = \tilde{\Sigma}_{R^*}, N_{R+1} = N_{R^*}, \beta_{R+1} = \beta_{R^*}, \Psi_{R+1} = \Psi_{R^*} \quad (28)$$

459 Although previously pruned data clouds are stored in memory, all previously
 460 pruned data clouds are excluded from any training scenarios except (18).
 461 Unlike its predecessors [10], this rule recall scenario is completely independent
 462 from the growing process (please refer to Algorithm 1).

463 5.4. Online Feature Selection Strategy

464 A prominent work, namely online feature selection (OFS), was developed
 465 in [15]. The appealing trait of OFS lies in its aptitude for flexible feature
 466 selection, as it enables the provision of different combinations of input at-
 467 tributes in each episode by activating or deactivating input features (1 or 0)
 468 in accordance to the up-to-date data trend. Furthermore, this technique is
 469 also capable of handling partial input attributes which are fruitful when the
 470 cost of feature extraction is too expensive. OFS is generalized here to fit the
 471 context of pRVFLN and to address the regression problem.

472 We start our discussion from a condition where a learner is provided with
 473 full input variables. Suppose that B input attributes are to be selected in
 474 the training process and $B < n$, the simplest approach is to discard the input
 475 features with marginal accumulated output weights $\sum_{i=1}^R \sum_{j=1}^2 \beta_{i,j}$ and maintain
 476 only B input features with the largest output weights. Note that the second
 477 term $\sum_{j=1}^2$ is required because of the extended input vector $x_e \in \mathfrak{R}^{(2n+1)}$. The
 478 rule consequent informs a tendency or orientation of a rule in the target space

479 which can be used as an alternative to gradient information [35]. Although it
480 is straightforward to use, it cannot ensure the stability of the pruning process
481 due to a lack of sensitivity analysis of the feature contribution. To correct
482 this problem, a sparsity property of the L1 norm can be analyzed to exam-
483 ine whether the values of n input features are concentrated in the L1 ball.
484 This allows the distribution of the input values to be checked to determine
485 whether they are concentrated in the largest elements and that pruning the
486 smallest elements wont harm the models accuracy. This concept is actualized
487 by first inspecting the accuracy of pRVFLN. The input pruning process is
488 carried out when the system error is large enough $T_t - y_t > \kappa$. Nevertheless,
489 the system error is not only large in the case of underfitting, but also in
490 the case of overfitting. We modify this condition by taking into account the
491 evolution of system error $|\bar{e}_t + \sigma_t| > \kappa|\bar{e}_{t-1} + \sigma_{t-1}|$ which corresponds to the
492 global error mean and standard deviation. The constant κ is a predefined
493 parameter and fixed at 1.1. The output nodes are updated using the gradient
494 descent approach and then projected to the L2 ball to guarantee a bounded
495 norm. Algorithm 2 details the algorithmic development of pRVFLN.

496

Algorithm 2. GOFS using full input attributes
Input: α learning rate, χ regularization factor, B the number of features
to be retained
Output: selected input features $X_{t,selected} \in \mathfrak{R}^{1 \times B}$
For $t=1, \dots, T$
Make a prediction y_t
IF $|\bar{e}_t + \sigma_t| > 1.1|\bar{e}_{t-1} + \sigma_{t-1}|$ // for regression $\hat{o} = \max_{o=1, \dots, m} (y_o) \neq T_t$ or //
497 for classification
 $\beta_i = \beta_i - \chi\alpha \frac{\partial E}{\partial \beta_i}$, $\beta_i = \min(1, \frac{1/\sqrt{\chi}}{\|\beta_i\|_2})\beta_i$
Prune input attributes X_t except those of B largest $\sum_{i=1}^R \sum_{j=1}^2 \beta_{i,j}$
Else
 $\beta_i = \beta_i - \chi\alpha\beta_i$
End IF
End FOR

498

499 where α, χ are respectively the learning rate and regularization factor. We
500 assign $\alpha = 0.2, \chi = 0.01$ following the same setting [15]. The optimization
501 procedure relies on the standard mean square error (MSE) as the objective

502 function and utilises the conventional gradient descent scenario:

$$\frac{\partial E}{\partial \beta_i} = (T_t - y_t) \left\{ \sum_{i=1}^R (1 - q) \overline{G}_{i,temporal} + \sum_{i=1}^R q \underline{G}_{i,temporal} \right\} \quad (29)$$

503 Furthermore, the predictive error has been theoretically proven to be bounded
504 in [17] and the upper bound is also found. One can also notice that the GOFs
505 enables different feature subsets to be elicited in each training observation t .

506 A relatively unexplored area of existing online feature selection is a situa-
507 tion where a limited number of features is accessible for the training process.
508 To actualise this scenario, we assume that at most B input variables can
509 be extracted during the training process. This strategy, however, cannot be
510 done by simply acquiring any B input features, because this scenario risks
511 having the same subset of input features during the training process. This
512 problem is addressed using the Bernaoulli distribution with confidence level
513 to sample B input attributes from n input attributes $B < n$. Algorithm 3
514 displays the feature selection procedure.

515

Algorithm 3. GOFS using partial input attributes

Input: α learning rate, χ regularization factor, B the number of features to be retained, ϵ confidence level

Output: selected input features $X_{t,selected} \in \mathfrak{R}^{1 \times B}$

For $t=1, \dots, T$

Sample γ from Bernaoulli distribution with confidence level ϵ

IF $\gamma_t = 1$

Randomly select B out of n input attributes $\tilde{X}_t \in \mathfrak{R}^{1 \times B}$

End IF

Make a prediction y_t

516 **IF** $|\bar{e}_t + \sigma_t| > 1.1|\bar{e}_{t-1} + \sigma_{t-1}|$ // for regression $\hat{o} = \max_{o=1, \dots, m} (y_o) \neq T_t$ or //
for classification

$\hat{X}_t = \tilde{X}_t / (B/n\epsilon) + (1 - \epsilon)$

$\beta_i = \beta_i - \chi\alpha \beta_i - \alpha\chi \frac{\partial E}{\partial \beta_i}, \beta_i = \min(1, \frac{1/\sqrt{\chi}}{\|\beta_i\|_2})\beta_i$

Prune input attributes X_t except those of B largest $\sum_{i=1}^R \sum_{j=1}^2 \beta_{i,j}$

Else

$\beta_{i,t} = \beta_{i,t-1}$

End IF

517 **End FOR**

518 As with Algorithm 2, the convergence of this scenario has been theoreti-
519 cally proven and the upper bound is derived in [17]. One must bear in mind
520 that the pruning process in Algorithm 1 and 2 is carried out by assigning
521 crisp weights (0 or 1), which fully reflect the importance of the input features.

522 5.5. Random Learning Strategy

523 pRVFLN adopts the random parameter learning scenario of the RVFLN,
524 leaving only the output nodes W to be analytically tuned with an online
525 learning scenario, whereas others, namely A_t, q, λ, Δ , can be randomly gener-
526 ated without any tuning process. To begin the discussion, we recall the
527 output expression of pRVFLN as follows:

$$y_o = \sum_{i=1}^R \beta_i \tilde{G}_{i,temporal}(X_t; A_t, q, \lambda, \Delta) \quad (30)$$

528 Referring to the RVFLN theory, the activation function $\tilde{G}_{i,spatial}$ should be
 529 either integrable or differentiable:

$$\int_R G^2(x)dx < \infty, \text{ or } \int_R [G'(x)]^2 dx < \infty \quad (31)$$

530 Furthermore, a large number of hidden nodes R is usually needed to ensure
 531 adequate coverage of data space because hidden node parameters are chosen
 532 at random [37]. Nevertheless, this condition can be relaxed in the pRVFLN,
 533 because the data cloud growing mechanism, namely the T2SCC method,
 534 partitions the input region in respect to real data distributions. The data
 535 cloud-based neurons are parameter-free and thus do not require any param-
 536 eterization, which often calls for a high-level approximation or complicated
 537 optimization procedure. Other parameters, namely A_t, q, λ, Δ , are randomly
 538 chosen, and their region of randomisation should be carefully selected. Re-
 539 ferring to [38], the parameters are sampled randomly from the following.

$$\begin{cases} b = -w_0 y_0 - \mu_0 \\ w_0 = \alpha w_0; w_0 \in [0; \Omega] \times [-\Omega; \Omega]^{d-1} \\ y_0 \in I^d \\ \mu_0 \in [-2d\Omega, 2d\Omega] \end{cases} \quad (32)$$

540 where u, Ω, α are probability measures. Nevertheless, this strategy is im-
 541 possible to implement in online situations because it often entails a rigorous
 542 trial-error process to determine these parameters. Most RVFLNs work sim-
 543 ply by following Schmidt et al.s strategy [9], that is, setting the region of
 544 random parameters in the range of $[-1,1]$.

545 Assuming that a complete dataset $\Xi = [X, T] \in \mathfrak{R}^{N \times (n+m)}$ is observable, a
 546 closed-form solution of (23) can be defined to determine the output weights .
 547 Although the original RVFLN adjusts the output weight with the conjugate
 548 gradient (CG) method, the closed-form solution can still be utilised with
 549 ease [4]. The obstacle for the use of pseudo-inversion in the original work
 550 was the limited computational resources in 90's. Although it is easy to use
 551 and ensures a globally optimum solution, this parameter learning scenario
 552 however imposes revisiting preceding training patterns which are intractable
 553 for online learning scenarios. pRVFLN employs the FWGRLS method [39]
 554 to adjust the output weight. As the FWGRLS approach has been detailed
 555 in [39], it is not recounted here.

Table 2: Details of Experimental Procedure

Section	Mode	Number of Runs	Benchmark Algorithm	Predefined Parameters	Number of Samples	Number of Input attributes
A (Nox Emission)	Direct Partition	10 times	PANFIS, GENEFIS, eTS, simpTS, DFNN, GDFNN, FAOS-PFNN, ANFIS, BARTFIS	$\alpha_1 = 0.002, \alpha_2 = 0.02$	826	170
	Cross Validation	5 times in each fold	DNNE, Online RVFL, Batch RVFL	$\alpha_1 = 0.002, \alpha_2 = 0.02$		
B (Tool Condition Monitoring)	Direct Partition	10 times	PANFIS, GENEFIS, eTS, simpTS, DFNN, GDFNN, FAOS-PFNN, ANFIS, BARTFIS	$\alpha_1 = 0.002, \alpha_2 = 0.02$	630	12
	Cross Validation	5 times in each fold	DNNE, Online RVFL, Batch RVFL	$\alpha_1 = 0.002, \alpha_2 = 0.02$		
C (Nox Emission, Tool Condition Monitoring)	Cross Validation	5 times in each fold	N/A	$\alpha_1 = 0.002, \alpha_2 = 0.02$	As above	As above
D (Mackey Glass)	Direct Partition	10 times	N/A	$\alpha_1 = 0.002, \alpha_2 = 0.02$	3500	4
E (BJ gas furnace)	Direct Partition	10 times	N/A	N/A	290	2

556 *5.6. Robustness of RVFLN*

557 The network parameters are usually sampled uniformly within a range of
558 $[-1,1]$ in the literature. A new finding of Li and Wang in [22] exhibits that
559 randomly generating network parameters with a fixed scope $[-\alpha, \alpha]$ does not
560 ensure a theoretically feasible solution or often the hidden node matrix is
561 not full rank. Surprisingly, the hidden node matrix was not invertible in all
562 their case studies when randomly sampling network parameters in the range
563 of $[-1,1]$ and far better numerical results were achieved by choosing the scope
564 $[-200,200]$. This trend was consistent with different numbers of hidden nodes.
565 How to properly select scopes of random parameters and its corresponding
566 distribution still require in-depth investigation [9]. In practice, a pre-training
567 process is normally required to arrive at a decent scope of random parameters.
568 Note that the range of random parameters by Igel'nik and Pao [4] is still at
569 the theoretical level and does not touch the implementation issue. We study
570 different random regions in Section 5.5 to see how pRVFLN behaves under
571 variations of the scope of random parameters.

572 **6. Numerical Examples**

573 This section presents the numerical validation of our proposed algorithm
574 using case studies and comparisons with prominent algorithms in the litera-
575 ture. Two numerical examples, namely modelling of Nox emissions from a car
576 engine and tool condition monitoring in the ball-nose end milling process, are
577 presented in this section, and the other three numerical examples are placed
578 in the supplemental document ¹ to keep the paper compact. Our numerical

¹https://www.dropbox.com/s/lytpt4huqyoqa6p/supplemental_document.docx?dl=0

Table 3: Prediction of Nox emissions Using Time-Series Mode

Model	RMSE	Node	Input	Runtime	Network	Samples
pRVFLN (P)	0.04	2	5	0.24	22	510
pRVFLN (F)	0.05	2	5	0.56	22	515
eT2Class	0.045	2	170	17.98	117304	667
PANFIS	0.052	5	170	3.37	146205	667
GENEFIS	0.048	2	2	0.41	18	667
Simp_eTS	0.14	5	170	5.5	1876	667
BARTFIS	0.11	4	4	2.55	52	667
DFNN	0.18	548	170	4332.9	280198+NS	667
GDFNN	0.48	215	170	2144.1	109865	667
eTS	0.38	27	170	1098.4	13797	667
FAOS-PFNN	0.06	6	170	14.8	2216+NS	667
ANFIS	0.15	2	170	100.41	17178	667

579 studies were carried out under two scenarios: the time-series scenario and the
580 cross-validation (CV) scenario. The time-series procedure orderly executes
581 data streams according to their arrival and partitions data streams into two
582 parts, namely training and testing. Simulations were repeated 10 times and
583 numerical results were averaged from 10 runs to arrive at conclusive findings.
584 In the time-series mode, pRVFLN was compared against 10 state-of-the-art
585 evolving algorithms: eT2Class [40], BARTFIS [41], PANFIS [42], GENEFIS
586 [43], eTS [44], simp_eTS [45], DFNN [46], GDFNN [47], FAOSPFNN [48],
587 ANFIS [49]. The CV scenarios were taken place in our experiment in order
588 to follow the commonly adopted simulation environment of other RVFLNs in
589 the literature where five runs per each fold were undertaken. pRVFLN was
590 benchmarked against the decorrelated neural network ensemble (DNNE) [5],
591 online and batch versions of RVFL [9]. The pRVFLN MATLAB codes are
592 provided in ² while the MATLAB codes of DNNE and RVFL are available
593 online ^{3,4}. Comparisons were performed against five evaluation criteria: accu-

²<https://www.dropbox.com/sh/zbm54epk8trlnq9/AACgxLHt5Nsy7MISbgXcdkTba?dl=0>

³<http://homepage.cs.latrobe.edu.au/dwang/html/DNNEweb/index.html>

⁴<http://ispac.ing.uniroma1.it/scardapane/software/lynx/>

Table 4: Prediction of Nox emissions Using CV Mode

Model	NRMSE	Node	Input	Runtime	Network	Samples
pRVFLN (P)	0.12±0.04	2	5	1.1	22	699
pRVFLN (F)	0.12±0.03	2	5	5.98	22	630.7
DNNE	0.14±0	50	170	0.81	43600+NS	744
Online RVFL	0.52±0.02	100	170	0.03	87200	744
Batch RVFL	0.59±0.05	100	170	0.003	87200+NS	744

594 racy, data clouds, input attribute, runtime, network parameters. The scope
595 of random parameters followed Schmidt’s suggestion [9] where the scope of
596 random parameters was [-1,1] but we insert the analysis of robustness in
597 part C which provides additional results with different random regions and
598 illustrates how the scope of random parameters influences the final numerical
599 results. The effect of the individual learning component to the end results and
600 the influence of user-defined predefined thresholds are analysed in Section D
601 and E. Furthermore, additional numerical results across different problems
602 are provided in the supplemental document ¹. To allow a fair comparison, all
603 consolidated algorithms were executed in the same computational resources
604 under the MATLAB environment. Details of the experimental procedure are
605 tabulated in Table 2.

606 6.1. Modeling of Nox Emissions from a Car Engine

607 This section demonstrates the efficacy of the pRVFLN in modeling Nox
608 emissions from a car engine [50]. This real-world problem is relevant to
609 validate the learning performance, not only because it features noisy and
610 uncertain characteristics similar to the nature of a car engine, it also char-
611 acterizes high dimensionality, containing 170 input attributes. That is, 17
612 physical variables were captured in 10 consecutive measurements. Further-
613 more, different engine parameters were applied to induce changing the system
614 dynamics to simulate real driving actions across different road conditions. In
615 the time-series procedure, 826 data points were streamed to consolidated al-
616 gorithms, where 667 samples were set as training samples, and the remainder
617 were fed for testing purposes. 10 runs were carried out to attain consistent

¹https://www.dropbox.com/s/lytpt4huqyoqa6p/supplemental_document.docx?dl=0

618 numerical results. In the CV procedure, the experiment was run under the
619 10-fold CV, and each fold was repeated five times similar to the scenario
620 adopted in [26]. This strategy checks the consistency of the RVFLNs learn-
621 ing performance because it adopts the random learning scenario and avoids
622 data order dependency. Table 3 and 4 exhibit the consolidated numerical
623 results of benchmarked algorithms.

Table 5: Tool Wear Prediction Using Time Series Mode

Model	RMSE	Node	Input	Runtime	Network	Samples
pRVFLN (P)	0.14	2	8	0.34	34	157
pRVFLN (F)	0.19	2	8	0.15	34	136
eT2Class	0.16	4	12	1.1	1260	320
Simp_eTS	0.22	17	12	1.29	437	320
eTS	0.15	7	12	0.56	187	320
BARTFIS	0.16	6	12	0.43	222	320
PANFIS	0.15	3	12	0.77	507	320
GENEFIS	0.13	42	12	0.88	507	320
DFNN	0.27	42	12	2.41	1092+NS	320
GDFNN	0.26	7	12	2.54	259+ NS	320
FAOS-PFNN	0.38	7	12	3.76	1022+NS	320
ANFIS	0.16	8	12	0.52	296+ NS	320

Table 6: Tool wear prediction using CV Mode

Model	NRMSE	Node	Input	Runtime	Network	Samples
pRVFLN (P)	0.16±0.01	1.3±0.2	8	0.2	22.1	245.8
pRVFLN (F)	0.17±0.06	1.92±0.1	8	0.3	32.6	432.1
DNNE	0.11±0	50	12	0.79	3310+NS	571.5
Online RVFL	0.16±0.01	100	12	0.02	1400	571.5
Batch RVFL	0.19±0.04	100	12	0.002	1400+NS	571.5

624 It is evident that pRVFLN outperformed its counterparts in all the evalua-
625 tion criteria except GENEFS for the number of input attributes and network

626 parameters. It is worth mentioning however that in the other three criteria:
627 predictive accuracy, execution time, and number of training samples, the
628 GENEFIS was inferior to ours. pRVFLN is equipped with an online active
629 learning strategy, which discards superfluous samples. This learning mod-
630 ule had a significant effect on predictive accuracy. Furthermore, pRVFLN
631 utilizes the GOFs method, which is capable of coping with the curse of di-
632 mensionality. Note that the unique feature of the GOFs method is that
633 it allows different feature subsets to be picked up in every training episode
634 which avoids the catastrophic forgetting of obsolete input attributes, tem-
635 porarily inactive due to changing data distributions. The GOFs can handle
636 partial input attributes during the training process and results in the same
637 level of accuracy as that of the full input attributes. The use of full input
638 attributes slowed down the execution time because it needed to deal with
639 170 input variables first, before reducing the input dimension. In this case
640 study, we selected five input attributes to be kept for the training process.
641 Our experiment showed that the number of selected input attributes is not
642 problem-dependent and is able to be set as any desirable number in most
643 cases. We did not observe significant performance difference when using ei-
644 ther the full input mode or partial input mode. On the other hand, consistent
645 numerical results were achieved by pRVFLN, although the pRVFLN is built
646 on the random vector functional link algorithm, as observed in the CV ex-
647 perimental scenario. In addition, pRVFLN produced the most encouraging
648 performance in almost all evaluation criteria except computational speed be-
649 cause other RVFLNs implement less comprehensive training procedures than
650 pRVFLN. Note that although DNNE attained the highest accuracy in the
651 CV mode, it imposed considerable memory and space complexities.

652 *6.2. Tool Condition Monitoring of High-Speed Machining Process*

653 this section presents a real-world problem from a complex manufacturing
654 process (Courtesy of Dr. Li Xiang, Singapore) [51]. The objective of this
655 case study is to perform predictive analytics of the tool wear in the ball-
656 nose end milling process frequently found in the metal removal process of
657 the aerospace industry. In total, 12 time-domain features were extracted
658 from the force signal and 630 samples were collected during the experiment.
659 Concept drift in this case study resulted from changing surface integrity,
660 tool wear degradation as well as varying machining configurations. For the
661 time-series experimental procedure, the consolidated algorithms were trained
662 using data from cutter A, while the testing phase exploited data from cutter

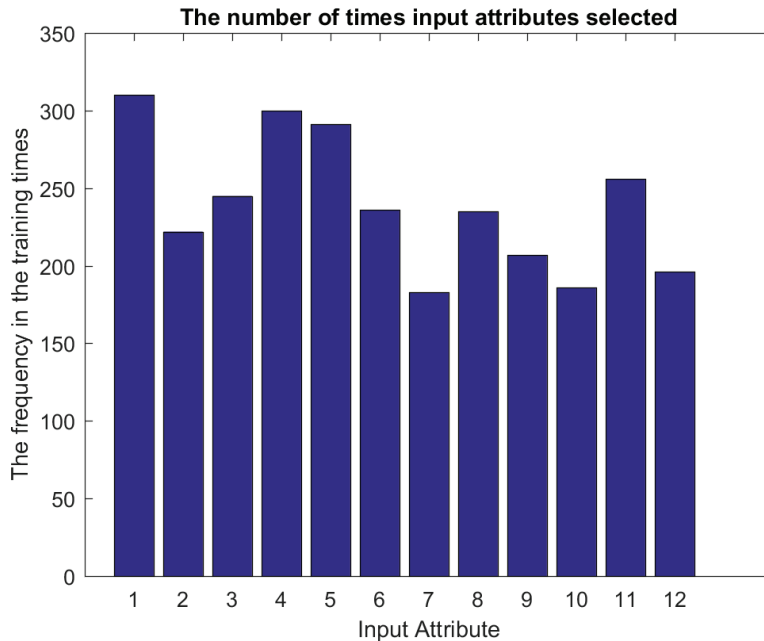


Figure 3: The frequency of input features

663 B. This process was repeated 10 times to achieve valid numerical results.
 664 For the CV experimental procedure, the 10-fold CV process was undertaken
 665 where each fold was undertaken five times to arrive at consistent findings.
 666 Tables 5 and 6 report the average numerical results across all folds. Fig. 3
 667 depicts how many times input attributes are selected during one fold of the
 668 CV process.

669 It is observed from Table 5 and 6 that pRVFLN evolved the lowest struc-
 670 tural complexities while retaining a high accuracy. It is worth noting that
 671 although the DNNE exceeded pRVFLN in accuracy, it imposed considerable
 672 complexity because it is an offline algorithm revisiting previously seen data
 673 samples and adopts an ensemble learning paradigm. The efficacy of the on-
 674 line sample selection strategy was seen, as it led to a significant reduction of
 675 training samples to be learned during the experiment. Using partial input
 676 information led to subtle differences to those with the full input information.
 677 It is seen in Fig. 3 that the GOFS selected different feature subsets in ev-
 678 ery training episode. Additional numerical examples, sensitivity analysis of

679 predefined thresholds and analysis of learning modules are given in ¹.

680 *6.3. Analysis of Robustness*

681 this section aims to numerically validate our claim in section III. E that a
682 range $[-1,1]$ does not always ensure the production of a reliable model [22, 51].
683 Additional numerical results with different intervals of random parameters
684 are presented. Four intervals, namely $[0,0.1]$, $[0,0.5]$, $[0,0.8]$, $[0,3]$, $[0,5]$, $[0,10]$
685 were tried for two case studies in section IV.A and IV.B. Our experiments
686 were undertaken in the 10-fold CV procedure as done in previous sections.
687 Table 7 displays the numerical results.

688 For the tool wear case study, the best-performing model was generated
689 by the range $[0,0.1]$. The higher the range of the model, the more inferior the
690 model. It went to the point where a model was no longer stable under the
691 range $[0,3]$. On the other side, the range $[0,0.5]$ induced the best-performing
692 model with the highest accuracy while evolving comparable network com-
693 plexity for the Nox emission case study. The higher the scope led to the
694 deterioration of numerical results. Moreover, the range $[0,0.1]$ did not deliver
695 a better accuracy than the range $[0,0.5]$ since this range did not generate
696 diverse enough random values. These numerical results are interpreted from
697 the nature of pRVFLN a clustering-based algorithm. The success of pRVFLN
698 is mainly determined from the compatibility of the zone of influence of hid-
699 den nodes to a real data distribution, and its performance worsens when the
700 scope is remote from the true data distribution. This finding is complemen-
701 tary to Li and Wang [22] where it relies on a sigmoid-based RVFL network,
702 and the scope of random parameters can be outside the applicable operating
703 intervals. Its predictive performance is set by its approximation capability
704 in the output space. It is worth-stressing that network parameters are ran-
705 domly generated in a positive range since the uncertainty threshold setting
706 the footprint of uncertainty is also chosen at random. Having negative values
707 for this parameter causes invalid interval definitions and poor performance
708 is returned as a result.

709 *6.4. Contributions of Learning Components*

710 This section demonstrates the efficacy of each learning module of the
711 pRVFLN and analyses to what extent this learning module contributes to the

¹https://www.dropbox.com/s/lytpt4huqyoqa6p/supplemental_document.docx?dl=0

712 resultant learning performance. The experiment was undertaken using the
713 Mackey-Glass time series problem, a control model of the production of white
714 blood cells. This problem features the chaotic characteristic, whose nonlinear
715 oscillations are well-accepted as a model of most psychological processes.
716 This problem is described by the following mathematical model:

$$\frac{dx(t)}{dt} = \frac{ax(t - \tau)}{1 + x^{10}(t - \tau)} - bx(t) \quad (33)$$

717 where $a = 0.2$, $b = 0.1$ and $\tau = 85$. The chaotic characteristic is attributed
718 by $\tau \geq 17$. The nonlinear dependence of this problem is built upon the
719 series-parallel identification model as follows:

$$x(t + 85) = f(x(t), x(t - 6), x(t - 12), x(t - 18)) \quad (34)$$

720 3000 training data from the range of [201,3200] and 500 testing data from
721 the range of [5001,5500] were generated using the fourth order Range-Kutta
722 method. The effect of each learning component was investigated by studying
723 the learning performance of pRVFLN under five learning configurations: A)
724 this configuration refers to pRVFLN with a feedforward network architec-
725 ture; B) the online active learning part is deactivated; C) we switch off the
726 online feature selection; D) pRVFLN is executed with the absence of hidden
727 node pruning and recall mechanisms. The numerical results are summarised
728 in Table 8. As with previous case studies, the learning performance was ex-
729 amined against five learning criteria: NDEI, hidden nodes, execution time,
730 training samples, and input attributes. Our simulation was done under the
731 time-series mode with 10 runs.

732 It is obvious from Table 6 that each learning module played a critical
733 role in the learning performance of pRVFLN. Without the recurrent connec-
734 tion, the predictive accuracy of pRVFLN slightly deteriorated and this also
735 increased the number of training samples to be seen in the training process.
736 The difference in performance was negligible and imposed pRVFLN to see all
737 the samples when the active learning strategy was shelved for the training
738 process. This fact substantiates the efficacy of the active learning scenario in
739 extracting important data points for the training process. The absence of the
740 hidden node pruning mechanism triggered the increase of hidden nodes to
741 be evolved during the training process and only minimally affected the pre-
742 dictive accuracy. Moreover, pRVFLN was capable of learning data streams
743 with partial input information as well as with full input information.

744 *6.5. Sensitivity Analysis of Predefined Thresholds*

745 This section examines the impact of two predefined thresholds, namely
746 α_1, α_2 , on the overall learning performance of pRVFLN. Intuitively, one can
747 envisage that the higher the value of α_1 , the fewer the data clouds are added
748 during the training process and vice versa, whereas the higher the value of
749 α_2 , the higher the number of data clouds are generated. To further confirm
750 this aspect, the sensitivity of these parameters is analysed using the box
751 Jenkins (BJ) gas furnace problem. The BJ gas furnace problem is a popular
752 benchmark problem in the literature, where the goal is to model the CO2
753 level in off gas based on two input attributes: the methane flow rate $u(n)$,
754 and its previous one-step output $t(n-1)$. From the literature, the best input
755 and output relationship of the regression model is known as $\hat{y}(n) = f(u(n-4), t(n-1))$.
756 290 data points were generated from the gas furnace, 200 of
757 which were assigned as the training samples, and the remainder were utilised
758 to validate the model. α_1 was varied in the range of [0.002,0.004,0.006,0.008],
759 while α_2 was assigned the values of [0.02,0.04,0.06,0.08]. Two tests were
760 carried out to test their sensitivity. That is, α_1 was fixed at 0.002, while
761 setting different values of α_2 , whereas α_2 was set at 0.02, while varying α_1 .
762 Moreover, our simulation followed the time-series mode with 10 repetitions
763 as aforementioned. The learning performance of pRVFLN was evaluated
764 against four criteria: non-dimensional error index (NDEI), number of hidden
765 nodes, execution time, number of training samples, and number of network
766 parameters. The results are reported in Table 9.

767 Referring to Table 9, it can be observed that pRVFLN can achieve satis-
768 factory learning performance while demanding very low network, computa-
769 tional, and sample complexities. Allocating different values of α_1, α_2 did not
770 cause significant performance deterioration, where the NDEI, runtime and
771 the number of samples were stable in the range of [0.27,0.38], [0.5,0.79], and
772 [10,30] respectively. Note that the slight variation in these learning perfor-
773 mances was attributed to the random learning algorithm of pRVFLN. On the
774 other hand, the number of hidden nodes and parameters remained constant
775 at 2 and 10 respectively and were not influenced by a variation of the two
776 predefined thresholds. It is worth mentioning that the data cloud-based hid-
777 den node of pRVFLN incurred modest network complexity because it did not
778 have any parameters to be memorised and adapted. In all the simulations
779 in this paper, α_1 and α_2 were fixed at 0.02 and 0.002 respectively to ensure
780 a fair comparison with its counterparts and to avoid a laborious pretraining
781 step in finding suitable values for these two parameters.

782 **7. Conclusion**

783 A novel random vector functional link network, namely the parsimonious
784 random vector functional link network (pRVFLN), is proposed. pRVFLN
785 aims to provide a concrete solution to the issue of data stream by putting
786 into perspective a synergy between adaptive and evolving characteristics and
787 fast and easy-to-use characteristics of RVFLN. pRVFLN is a fully evolving
788 algorithm where its hidden nodes can be automatically added, pruned
789 and recalled dynamically while all network parameters except the output
790 weights are randomly generated with the absence of any tuning mechanism.
791 pRVFLN is fitted by the online feature selection mechanism and the online
792 active learning scenario which further strengthens its aptitude in processing
793 data streams. Unlike conventional RVFLNs, the concept of interval-valued
794 data clouds is introduced. This concept simplifies the working principle of
795 pRVFLN because it neither requires any parameterization per scalar variables
796 nor follows a pre-specified cluster shape. It features an interval-valued
797 spatiotemporal firing strength, which provides the degree of tolerance for uncertainty.
798 Rigorous case studies were carried out to numerically validate the efficacy of pRVFLN
799 where pRVFLN delivered the most encouraging performance. The ensemble version of pRVFLN
800 will be the subject of our future investigation which aims to further improve the predictive
801 performance of pRVFLN.
802

803 **ACKNOWLEDGEMENT**

804 The third author acknowledges the support of the Austrian COMET-K2
805 program of the Linz Center of Mechatronics (LCM), funded by the Austrian
806 federal government and the federal state of Upper Austria. We thank Dr. D.
807 Wang for his suggestion pertaining to robustness issue of RVFLN Mr. MD
808 Meftahul Ferdous for his assistance for Latex typesetting of our manuscript.

809 **References**

- 810 [1] M. Pratama, S. G. Anavatti, J. Lu, Recurrent classifier based on an
811 incremental metacognitive-based scaffolding algorithm, *IEEE Transactions on Fuzzy Systems* 23 (2015) 2048–2066.
812
- 813 [2] S. Scardapane, D. Wang, Randomness in neural networks: an overview,
814 *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*
815 7 (2017).

- 816 [3] D. S. Broomhead, D. Lowe, Multi-variable functional interpolation and
817 adaptive networks, *Complex Systems* 2 (1998) 321–355.
- 818 [4] B. Igelnik, Y.-H. Pao, Stochastic choice of basis functions in adaptive
819 function approximation and the functional-link net, *IEEE Transactions*
820 *on Neural Networks* 6 (1995) 1320–1329.
- 821 [5] M. Alhamdoosh, D. Wang, Fast decorrelated neural network ensembles
822 with random weights, *Information Sciences* 264 (2014) 104–117.
- 823 [6] L. Zhang, P. N. Suganthan, A survey of randomized algorithms for
824 training neural networks, *Information Sciences* 364 (2016) 146–155.
- 825 [7] L. Zhang, P. N. Suganthan, A comprehensive evaluation of random
826 vector functional link networks, *Information Sciences* 367 (2016) 1094–
827 1105.
- 828 [8] J. Zhao, Z. Wang, F. Cao, D. Wang, A local learning algorithm for
829 random weights networks, *Knowledge-Based Systems* 74 (2015) 159–
830 166.
- 831 [9] W. F. Schmidt, M. A. Kraaijveld, R. P. Duin, Feedforward neural net-
832 works with random weights, in: *Pattern Recognition, 1992. Vol. II. Con-*
833 *ference B: Pattern Recognition Methodology and Systems, Proceedings.,*
834 *11th IAPR International Conference on, IEEE*, pp. 1–4.
- 835 [10] M. Pratama, S. G. Anavatti, M. Joo, E. D. Lughofer, pclass: an effective
836 classifier for streaming examples, *IEEE Transactions on Fuzzy Systems*
837 23 (2015) 369–386.
- 838 [11] P. Angelov, R. Yager, A new type of simplified fuzzy rule-based system,
839 *International Journal of General Systems* 41 (2012) 163–185.
- 840 [12] S. Xiong, J. Azimi, X. Z. Fern, Active learning of constraints for semi-
841 supervised clustering, *IEEE Transactions on Knowledge and Data En-*
842 *gineering* 26 (2014) 43–54.
- 843 [13] R.-F. Xu, S.-J. Lee, Dimensionality reduction by feature clustering for
844 regression problems, *Information Sciences* 299 (2015) 42–57.

- 845 [14] J.-Y. Jiang, R.-J. Liou, S.-J. Lee, A fuzzy self-constructing feature clus-
846 tering algorithm for text classification, *IEEE transactions on knowledge*
847 *and data engineering* 23 (2011) 335–349.
- 848 [15] J. Wang, P. Zhao, S. C. Hoi, R. Jin, Online feature selection and its
849 applications, *IEEE Transactions on Knowledge and Data Engineering*
850 26 (2014) 698–710.
- 851 [16] D. Lowe, Adaptive radial basis function nonlinearities, and the prob-
852 lem of generalisation, in: *Artificial Neural Networks, 1989.*, First IEE
853 *International Conference on (Conf. Publ. No. 313)*, IET, pp. 171–175.
- 854 [17] Y.-H. Pao, G.-H. Park, D. J. Sobajic, Learning and generalization char-
855 acteristics of the random vector functional-link net, *Neurocomputing* 6
856 (1994) 163–180.
- 857 [18] Y.-H. Pao, Y. Takefuji, Functional-link net computing: theory, system
858 architecture, and functionalities, *Computer* 25 (1992) 76–79.
- 859 [19] M. Lukoševičius, H. Jaeger, Reservoir computing approaches to re-
860 current neural network training, *Computer Science Review* 3 (2009)
861 127–149.
- 862 [20] L. Xie, Y. Yang, Z. Zhou, J. Zheng, M. Tao, Z. Man, Dynamic neural
863 modeling of fatigue crack growth process in ductile alloys, *Information*
864 *Sciences* 364 (2016) 167–183.
- 865 [21] I. Zliobaite, A. Bifet, B. Pfahringer, G. Holmes, Active learning with
866 drifting streaming data, *IEEE transactions on neural networks and*
867 *learning systems* 25 (2014) 27–39.
- 868 [22] M. Li, D. Wang, Insights into randomized algorithms for neural net-
869 works: Practical issues and common pitfalls, *Information Sciences* 382
870 (2017) 170–178.
- 871 [23] J. C. Patra, A. C. Kot, Nonlinear dynamic system identification using
872 chebyshev functional link artificial neural networks, *IEEE Transactions*
873 *on Systems, Man, and Cybernetics, Part B (Cybernetics)* 32 (2002) 505–
874 511.

- 875 [24] P. Angelov, *Autonomous learning systems: from data streams to knowl-*
876 *edge in real-time*, John Wiley & Sons, 2012.
- 877 [25] Y.-Y. Lin, J.-Y. Chang, C.-T. Lin, Identification and prediction of dy-
878 *dynamic systems using an interactively recurrent self-evolving fuzzy neural*
879 *network*, *IEEE Transactions on Neural Networks and Learning Systems*
880 *24* (2013) 310–321.
- 881 [26] C.-F. Juang, C.-T. Lin, A recurrent self-organizing neural fuzzy infer-
882 *ence network*, *IEEE Transactions on Neural Networks* *10* (1999) 828–
883 845.
- 884 [27] C.-F. Juang, A tsk-type recurrent fuzzy network for dynamic systems
885 *processing by neural network and genetic algorithms*, *IEEE Transactions*
886 *on Fuzzy Systems* *10* (2002) 155–170.
- 887 [28] C.-F. Juang, Y.-Y. Lin, C.-C. Tu, A recurrent self-evolving fuzzy neural
888 *network with local feedbacks and its application to dynamic system*
889 *processing*, *Fuzzy Sets and Systems* *161* (2010) 2552–2568.
- 890 [29] R. H. Abiyev, O. Kaynak, Type 2 fuzzy neural structure for identifica-
891 *tion and control of time-varying plants*, *IEEE Transactions on Industrial*
892 *Electronics* *57* (2010) 4147–4159.
- 893 [30] Q. Liang, J. M. Mendel, Interval type-2 fuzzy logic systems: theory and
894 *design*, *IEEE Transactions on Fuzzy systems* *8* (2000) 535–550.
- 895 [31] P. Angelov, D. Kangin, X. Zhou, D. Kolev, Symbol recognition with
896 *a new autonomously evolving classifier autotool*, in: *Evolving and*
897 *Adaptive Intelligent Systems (EAIS)*, 2014 IEEE Conference on, IEEE,
898 pp. 1–7.
- 899 [32] D. Kangin, P. Angelov, J. A. Iglesias, *Autonomously evolving classifier*
900 *tedaotool*, *Information Sciences* *366* (2016) 1–11.
- 901 [33] A. K. Das, K. Subramanian, S. Sundaram, An evolving interval type-2
902 *neurofuzzy inference system and its metacognitive sequential learning*
903 *algorithm*, *IEEE Transactions on Fuzzy Systems* *23* (2015) 2080–2093.
- 904 [34] P. Mitra, C. Murthy, S. K. Pal, *Unsupervised feature selection using*
905 *feature similarity*, *IEEE transactions on pattern analysis and machine*
906 *intelligence* *24* (2002) 301–312.

- 907 [35] E. Lughofer, *Evolving fuzzy systems-methodologies, advanced concepts*
908 *and applications*, volume 53, Springer, 2011.
- 909 [36] H. Han, X.-L. Wu, J.-F. Qiao, Nonlinear systems modeling based on self-
910 organizing fuzzy-neural-network with adaptive computation algorithm,
911 *IEEE transactions on cybernetics* 44 (2014) 554–564.
- 912 [37] I. Y. Tyukin, D. V. Prokhorov, Feasibility of random basis function ap-
913 proximators for modeling and control, in: *Control Applications,(CCA)*
914 *& Intelligent Control,(ISIC)*, 2009 IEEE, IEEE, pp. 1391–1396.
- 915 [38] J.-Y. Li, W. Chow, B. Igelnik, Y.-H. Pao, Comments on” stochastic
916 choice of basis functions in adaptive function approximation and the
917 functional-link net” [with reply], *IEEE Transactions on Neural Networks*
918 8 (1997) 452–454.
- 919 [39] M. Pratama, J. Lu, E. Lughofer, G. Zhang, S. Anavatti, Scaffolding
920 type-2 classifier for incremental learning under concept drifts, *Neuro-*
921 *computing* 191 (2016) 304–329.
- 922 [40] M. Pratama, J. Lu, G. Zhang, Evolving type-2 fuzzy classifier, *IEEE*
923 *Transactions on Fuzzy Systems* 24 (2016) 574–589.
- 924 [41] R. J. Oentaryo, M. J. Er, S. Linn, X. Li, Online probabilistic learning
925 for fuzzy inference system, *Expert Systems with Applications* 41 (2014)
926 5082–5096.
- 927 [42] M. Pratama, S. G. Anavatti, P. P. Angelov, E. Lughofer, Panfis: a novel
928 incremental learning machine, *IEEE Transactions on Neural Networks*
929 *and Learning Systems* 25 (2014) 55–68.
- 930 [43] M. Pratama, S. G. Anavatti, E. Lughofer, Genefis: toward an effective
931 localist network, *IEEE Transactions on Fuzzy Systems* 22 (2014) 547–
932 562.
- 933 [44] P. P. Angelov, D. P. Filev, An approach to online identification of
934 takagi-sugeno fuzzy models, *IEEE Transactions on Systems, Man, and*
935 *Cybernetics, Part B (Cybernetics)* 34 (2004) 484–498.
- 936 [45] P. Angelov, D. Filev, Simpl_ets: a simplified method for learning evol-
937 ving takagi-sugeno fuzzy models, in: *Fuzzy Systems, 2005. FUZZ’05. The*
938 *14th IEEE International Conference on*, IEEE, pp. 1068–1073.

- 939 [46] S. Wu, M. J. Er, Dynamic fuzzy neural networks-a novel approach
940 to function approximation, *IEEE Transactions on Systems, Man, and*
941 *Cybernetics, Part B (Cybernetics)* 30 (2000) 358–364.
- 942 [47] S. Wu, M. J. Er, Y. Gao, A fast approach for automatic generation
943 of fuzzy rules by generalized dynamic fuzzy neural networks, *IEEE*
944 *Transactions on Fuzzy Systems* 9 (2001) 578–594.
- 945 [48] N. Wang, M. J. Er, X. Meng, A fast and accurate online self-organizing
946 scheme for parsimonious fuzzy neural networks, *Neurocomputing* 72
947 (2009) 3818–3829.
- 948 [49] J.-S. Jang, Anfis: adaptive-network-based fuzzy inference system, *IEEE*
949 *transactions on systems, man, and cybernetics* 23 (1993) 665–685.
- 950 [50] E. Lughofer, V. Macián, C. Guardiola, E. P. Klement, Identifying static
951 and dynamic prediction models for nox emissions with evolving fuzzy
952 systems, *Applied Soft Computing* 11 (2011) 2487–2500.
- 953 [51] M. Pratama, M. J. Er, X. Li, R. J. Oentaryo, E. Lughofer, I. Arifin, Data
954 driven modeling based on dynamic parsimonious fuzzy neural network,
955 *Neurocomputing* 110 (2013) 18–28.

Table 7: Analysis of Robustness

Scope	Criteria	Tool Wear	Nox emission
[0,0.1]	RMSE	0.13±0.008	1.9±0
	Node	1.8±0.25	1
	Input	8	5
	Runtime	0.2±0.1	0.1±0.02
	Network	30.9	11
	Samples	503.1	1
[0,0.5]	RMSE	0.14±0.02	0.1±0.01
	Node	1.92±0.2	1.98±0.14
	Input	8	5
	Runtime	0.18±0.008	5.7±0.3
	Network	32.6	21.8
	Samples	571.5	743.4
[0,0.8]	RMSE	0.47±0.42	0.18±0.3
	Node	1.4±0.05	1.96±0.19
	Input	8	5
	Runtime	0.19±0.13	5.56±0.96
	Network	23.8	21.6
	Samples	385.1	711.24
[0,3]	RMSE	Unstable	Unstable
	Node		
	Input		
	Runtime		
	Network		
	Samples		
[0,5]	RMSE	Unstable	Unstable
	Node		
	Input		
	Runtime		
	Network		
	Samples		
[0,10]	RMSE	Unstable	Unstable
	Node		
	Input		
	Runtime		
	Network		
	Samples		

Table 8: Analysis of Learning Components

ALGORITHMS	HIDDEN NODES	INPUT	RUNTIME	PARAMETERS	NDEI	TRAINING SAMPLES
pRVFLN	1	2	0.7	5	0.46	2474.3
(A)	1.1	2	0.7	5.5	0.51	2480.1
(B)	1	2	0.7	5	0.46	3000
(C)	1	2	0.7	5	0.47	2474.3
(D)	1	2	0.7	5	0.51	2348.1

Table 9: Sensitivity Analysis

PARAMETERS	NDEI	HN	RUNTIME	NP
$\alpha_1 = 0.002$	0.3	19.3	0.52	96.5
$\alpha_1 = 0.004$	0.3	19.3	0.49	96.5
$\alpha_1 = 0.006$	0.3	35.9	0.67	179.5
$\alpha_1 = 0.008$	0.3	7.3	0.4	36.5
$\alpha_2 = 0.02$	0.3	17	0.44	85
$\alpha_2 = 0.04$	0.31	143	1.41	715
$\alpha_2 = 0.06$	0.32	196.3	2.01	981.5
$\alpha_2 = 0.08$	0.32	196.3	2.01	981.5