

A Comprehensive Evaluation of Random Vector Functional Link Networks

Le Zhang, P.N. Suganthan*

School of Electric and Electronic Engineering,
Nanyang Technological University,
Singapore, 639798; Email: epnsugan@ntu.edu.sg

Abstract

With randomly generated weights between input and hidden layers, a random vector functional link network is a universal approximator for continuous functions on compact sets with fast learning property. Though it was proposed two decades ago, the classification ability of this family of networks has not been fully investigated yet. Through a very comprehensive evaluation by using 121 UCI datasets, the effect of bias in the output layer, direct links from the input layer to the output layer and type of activation functions in the hidden layer, scaling of parameter randomization as well as the solution procedure for the output weights are investigated in this work. Surprisingly, we found that the direct link plays an important performance enhancing role in RVFL, while the bias term in the output neuron had no significant effect. The ridge regression based closed-form solution was better than those with Moore-Penrose pseudoinverse. Instead of using a uniform randomization in $[-1, +1]$ for all datasets, tuning the scaling of the uniform randomization range for each dataset enhances the overall performance. Six commonly used activation functions were investigated in this work and we found that *hardlim* and *sign* activation functions degenerate the overall performance. These basic conclusions can serve as general guidelines for designing RVFL networks based classifiers.

Keywords: Random vector functional link networks, ridge regression, Moore-Penrose pseudoinverse, data classification.

1. Introduction

Single layer feedforward neural networks (SLFN) have been widely applied to solve problems such as classification and regression because of their universal approximation capability [14, 20, 31, 17]. Conventional methods for training SLFN are back-propagation based learning algorithms [10, 7]. These iterative methods suffer from slow convergence, getting trapped in a local minimum and being sensitivity to learning rate setting. Random Vector Functional Link Networks (RVFL), shown in Fig. 1, which is a randomized version of the functional link neural network [25, 8], shows that actual values of the weights from the input layer to hidden layer can be randomly generated in a suitable domain and kept fixed in the learning stage. Independently developed method in [35] also belongs to the family of randomized methods for training artificial neural networks with randomized input layer weights. This method [35] does not have direct links between the inputs and the outputs whereas RVFL has highly beneficial direct links.

RVFL was proposed in [28]. Learning and generalization characteristics of RVFL were discussed in [26]. In [17], Igelnik and Pao proved that the RVFL network is a universal approximator for a continuous function on a bounded finite dimensional set with a closed-form solution. From then on, RVFL has been employed to solve problems in diverse domains. A dynamic step-wise updating algorithm was proposed to update the output weights of the RVFL on-the-fly in [5] for both a new added pattern and a new added enhancement node. The RVFL network was investigated in [37] in the context of modeling and control. They [37] suggested to combine unsupervised placement of network nodes to the input data density with subsequent supervised or reinforcement learning of the linear parameters of the approximator. Modelling conditional probabilities with RVFL was reported in [15].

RVFL can also be combined with other learning methods. In [6], RVFL was combined with statistical hypothesis testing and self-organization of a number of enhancement nodes to generate a new learning system called a statistical self-organizing learning system (SSOLS) for remote sensing applications. In [16], expectation maximization was

combined with RVFL to improve its performance. RVFL has also been investigated in ensemble learning framework. In [1], decorrelated RVFL ensemble was introduced based on the negative correlation learning. RVFL based multi-source data ensemble for clinker free lime content estimation in rotary kiln sintering processes can be found [21]. RVFL has also been widely applied to solve real-life problems. In [30], the authors reported the performance of a holistic-styled word-based approach to off-line recognition of English language script. Radial basis function neural net and RVFL were combined. Their approach, named as density-based random-vector functional-link net (DBRVFLN), was helpful in improving the performance of the word recognition. In [29], RVFL was used in MPEG-4 coder. In [38] RVFL was applied for pedestrian detection based on combination of multi-feature. In [39], RVFL was combined with Adaboost in the pedestrian detection system. In [23], the authors investigated the performance of hardware implementation methods for RVFL. In [34], distributed learning of RVFL was proposed where training data is distributed under a decentralized information structure.

Consider an RVFL as demonstrated in Fig. 1. As mentioned before, the weights a_{ij} from the input to the enhance-

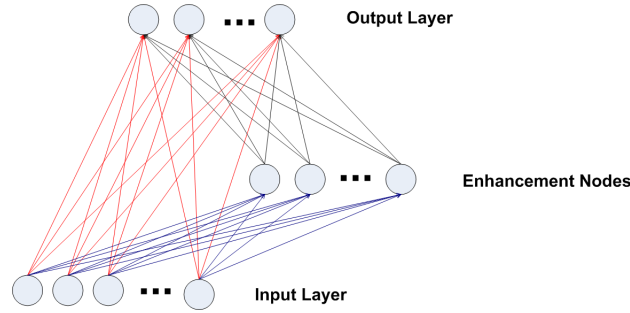


Figure 1: The structure of RVFL. The input features are firstly transformed into the enhanced features by the enhancement nodes. Input weights and biases of the enhancement nodes are randomly generated. At the output layer, all the enhanced and original features are concatenated and fed into output neurons.

ment nodes are randomly generated such that the activation functions $g(a_j^T x + b_j)$ are not all saturated. Following the approach in [1], all the weights are generated with the a uniform distribution within $[-S, +S]$ in this work, where S is a scale factor to be determined during the parameter tuning stage for each dataset. For RVFL, only the output weights β need to be determined by solving the following problem:

$$t_i = d_i^T \beta, \quad i = 1, 2, \dots, P \quad (1)$$

where P is the number of data samples, t is the target and d is the vector version of the concatenation of the original features as well as the random features.¹ Directly solving the problem in Eq. (1) may lead to over-fitting. In practice, a regularization on the solution such as regularized least square or preference of the solution with smaller norm [3] can be adopted to obtain the solution. RVFL can be roughly divided into 2 classes based on the algorithm to obtain the output weights. One is iterative RVFL, which obtains the output weights in an iterative manner based on the gradient of the error function. The other one is closed-form based RVFL, which obtains the output weights in a single-step. The present work focuses on the closed-form based RVFL because of its efficiency. A straightforward solution within a single learning step can be achieved by the pseudo-inverse [17, 27], among which Moore-Penrose pseudoinverse, $\beta = D^+ T$, where D and T are the matrix versions of the features and targets by stacking the features and targets of all data samples, is most commonly used. Another alternative is the $L2$ norm regularized least square (or ridge regression), which solves the following problem:

$$\sum_i (t_i - d_i^T \beta)^2 + \lambda \|\beta\|^2; i = 1, 2, \dots, P \quad (2)$$

¹For notational simplicity, we use the same formulation for all cases no matter whether there are biases in the output neurons since representing the features for the output neurons with $d = [d, 1]$ is equivalent to having a bias term in the output neurons.

The solution is given by $\beta = D(D^T D + \lambda I)^{-1} T$, where λ is the regularization parameter to be tuned.

Though there are many RVFL variants in the literature, some core features of RVFL remain unchanged. In this work, We choose the closed-form based RVFL and the following issues are investigated by using 121 UCI datasets as done in [11]

1. Effect of direct links from the input layer to the output layer.
2. Effect of the bias in the output neuron.
3. Performance of 6 commonly used activation functions as summarized in Table 1.
4. Performance of Moore-Penrose pseudoinverse and ridge regression (or regularized least square solutions) for the computation of the output weights.
5. Effect of range for randomly generated parameters in hidden neurons.

Issues 1 – 4 in the above list are discussed in Subsection 2.3 while issue 5 is discussed in Subsection 2.5.

Table 1: Activation functions used in this work. s and y are the inputs and outputs, respectively.

activation function	formulation
sigmoid	$y = \frac{1}{1+e^{-s}}$
sine	$y = \text{sine}(s)$
hardlim	$y = (\text{sign}(s) + 1)/2$
tribas	$y = \max(1 - s , 0)$
radbas	$y = \exp(-s^2)$
sign	$y = \text{sign}(s)$

2. Evaluation Protocol

2.1. Datasets

All 121 datasets are from the UCI repository [22]. The details of the datasets are summarized in Table 2.

We follow the same procedure as in [11]. Randomized stratified sampling is employed to make sure one training and one test set are generated (each with 50% of the available patterns), where each class has the same number of training and test patterns. Parameter tuning is performed on this couple of sets to identify parameters with the best performance on the test set. There are two parameters in the present work. One is the number of hidden neurons, which is tuned over 3: 203 with a step-size of 20 [11]. The other one is λ in ridge regression in Eq. (2), which is set to be 2^C and C is $-5 : 1 : 14$ [11]. Then, with the selected values for the tunable parameters, a 4-fold cross validation is developed using the whole data. However, for some datasets where the training-testing partition is already available (such as annealing and audiology-std, among others), the classifier is trained on the predefined training set and evaluated on the test set. In this case, the test result is calculated on the test set [11]. Each input feature is normalized by removing the mean value and dividing by its l_2 norm.

2.2. Different RVFL Configurations

We evaluate 48 different closed-form based RVFL configurations listed as follows:

1. RVFL with and without bias in the output neuron.
2. RVFL with and without direct link from input layer to output layer.
3. The performance of 6 commonly used activation functions as summarized in Table 1.
4. RVFL with Moore-Penrose pseudoinverse and RVFL with ridge regression.

Table 2: Datasets used in this work

Datasets	Patterns	Features	Classes	Datasets	Patterns	Features	Classes
abalone	4177	8	3	monks-1	124	6	2
ac-inflam	120	6	2	monks-2	169	6	2
acute-nephritis	120	6	2	monks-3	3190	6	2
adult	48842	14	2	mushroom	8124	21	2
annealing	798	38	6	musk-1	476	166	2
arrhythmia	452	262	13	musk-2	6598	166	2
audiology-std	226	59	18	nursery	12960	8	5
balance-scale	625	4	3	oocMerl2F	1022	25	3
balloons	16	4	2	oocMerl4D	1022	41	2
bank	45211	17	2	oocTris2F	912	25	2
blood	748	4	2	oocTris5B	912	32	3
breast-cancer	286	9	2	optical	3823	62	10
bc-wisc	699	9	2	ozone	2536	72	2
bc-wisc-diag	569	30	2	page-blocks	5473	10	5
bc-wisc-prog	198	33	2	parkinsons	195	22	2
breast-tissue	106	9	6	pendigits	7494	16	10
car	1728	6	4	pima	768	8	2
ctg-10classes	2126	21	10	pb-MATERIAL	106	4	3
ctg-3classes	2126	21	3	pb-REL-L	103	4	3
chess-krvk	28056	6	18	pb-SPAN	92	4	3
chess-krvkp	3196	36	2	pb-T-OR-D	102	4	2
congress-voting	435	16	2	pb-TYPE	105	4	6
conn-bench-sonar	208	60	2	planning	182	12	2
conn-bench-vowel	528	11	11	plant-margin	1600	64	100
connect-4	67557	42	2	plant-shape	1600	64	100
contrac	1473	9	3	plant-texture	1600	64	100
credit-approval	690	15	2	post-operative	90	8	3
cylinder-bands	512	35	2	primary-tumor	330	17	15
dermatology	366	34	6	ringnorm	7400	20	2
echocardiogram	131	10	2	seeds	210	7	3
ecoli	336	7	8	semeion	1593	256	10
energy-y1	768	8	3	soybean	307	35	18
energy-y2	768	8	3	spambase	4601	57	2
fertility	100	9	2	spect	80	22	2
flags	194	28	8	spectf	80	44	2
glass	214	9	6	st-aus-credit	690	14	2
haberman-survival	306	3	2	st-german-credit	1000	24	2
hayes-roth	132	3	3	st-heart	270	13	2
heart-cleveland	303	13	5	st-image	2310	18	7
heart-hungarian	294	12	2	st-landsat	4435	36	6
heart-switzerland	123	12	2	st-shuttle	43500	9	7
heart-va	200	12	5	st-vehicle	846	18	4
hepatitis	155	19	2	steel-plates	1941	27	7
hill-valley	606	100	2	synthetic-control	600	60	6
horse-colic	300	25	2	teaching	151	5	3
ilpd-indian-liver	583	9	2	thyroid	3772	21	3
image-segmentation	210	19	7	tic-tac-toe	958	9	2
ionosphere	351	33	2	titanic	2201	3	2
iris	150	4	3	rains	10	28	2
led-display	1000	7	10	twonorm	7400	20	2
lenses	24	4	3	vc-2classes	310	6	2
letter	20000	16	26	vc-3classes	310	6	3
libras	360	90	15	wall-following	5456	24	4
low-res-spect	531	100	9	waveform	5000	21	3
lung-cancer	32	56	3	waveform-noise	5000	40	3
lymphography	148	18	4	wine	179	13	3
magic	19020	10	2	w-qu-a-red	1599	11	6
mammographic	961	5	2	w-qu-a-white	4898	11	7
miniboone	130064	50	2	yeast	1484	8	10
molec-biol-promoter	106	57	2	zoo	101	16	7
molec-biol-splice	3190	60	3				

Details of the datasets. Some keys are: ac=inam=acute-inamnation, bc=breastcancer, congress-vot= congressional-voting, ctg=cardiotocography, conn-benchsonar/ vowel= connectionist-benchmark-sonar-mines-rocks/vowel-deterding, pb=pittsburg-bridges, st=statlog, aus=australian, vc=vertebral-column, w-qu-a=wine-quality.

2.3. Results and Discussion

Due to page limits, the detailed accuracy for each method is omitted in this paper and it can be downloaded from the authors' homepage². In the following, we summarize the overall performance of the variants based on their rank on each dataset. The most straightforward way to compare classifiers is to compute the average accuracies over all datasets. However, their averages are meaningless if the results on different datasets are not comparable. Moreover, averaging of accuracies is also susceptible to outliers. They allow a classifier's excellent performance on one dataset to compensate for poor performance. Further, a total failure on one problem can submerge fair results on most others [9]. Hence, in this study, we follow the method in [9], and use the rank of each classifier to reflect its performance [11]. This approach ranks the algorithms for each data set separately, the best performing algorithm getting the rank of 1, the second best rank 2..., average ranks are assigned in case of ties. Based on the 121 rankings, we summarize the overall ranking of each method in Table 3.

In order to give a detailed analysis of the results, we follow the method in [9] to test the significance of their differences. The statistical test is based on Friedman test. The Friedman test [12, 13] is a non-parametric equivalent of the repeated-measures ANOVA. It ranks the algorithms for each dataset separately, (the best performing algorithm getting the rank of 1, the second best rank 2 and so on), as shown in Table 3. In case of ties, average ranks are assigned. Let r_i^j be the rank of the j^{th} of k algorithms on the i^{th} of N datasets. The Friedman test compares the average ranks of algorithms, $R_j = \sum_i r_i^j$. The null-hypothesis states that all the algorithms are equivalent and so their ranks R_j should be equal. Let N and k denotes the number of algorithms and datasets respectively, when N and k are large enough, the Friedman statistic.

$$\chi_F^2 = \frac{12N}{k(k+1)} \left[\sum_j R_j^2 - \frac{k(k+1)^2}{4} \right], \quad (3)$$

is distributed according to χ_F^2 with $k-1$ degrees of freedom under the null-hypothesis. In that case, Friedman's χ_F^2 is undesirably conservative. A better statistic is:

$$F_F = \frac{(N-1)\chi_F^2}{N(k-1) - \chi_F^2}, \quad (4)$$

which is distributed according to the F -distribution with $k-1$ and $(k-1)(N-1)$ degrees of freedom.

If the null-hypothesis is rejected, which means the differences among the algorithms are statistically significant, the Nemenyi post-hoc test [24] can be used to check whether the performance of two among k classifiers is significantly different. If the corresponding average ranks of two different algorithms differ by at least the critical difference

$$CD = q_\alpha \sqrt{\frac{k(k+1)}{6N}}, \quad (5)$$

the performances of them are considered as significantly different. In Eq. (5), critical values q_α are based on the Studentized range statistic divided by $\sqrt{2}$. α is the significance level and is set to be 0.05 in this work.

2.3.1. Effect of Direct Link and Bias

We report in this section the effect of the direct link and bias in various RVFL configurations. In order to have reasonable comparisons, we keep other issues (the activation functions and the solution for the output weights) fixed. In Table 3, each column stands for an RVFL variant based on direct link and bias. Interestingly, we find that the superiority of the bias cannot be observed. However, the direct link is much more important than the bias term. It can be easily demonstrated that the $5^{th} - 6^{th}$ columns have lower ranks than those in $3^{rd} - 4^{th}$ columns. The direct link from the input layer to the output layer can serve as a regularization for the randomization thereby making the RVFL to achieve overall better performance than the RVFL variants without the direct link.

In our case ($k = 4, N = 121$), the critical values for the F distribution with 3 and 360 degrees of freedom is 2.6 and CD for the Nemenyi test is 0.4264. So, for those rows with F value larger than 2.6, their difference is statistically significant. Further in those rows, the pair of columns whose ranks differ by more than 0.4264 are statistically significantly different. In Table 4, we summarize the statistical significance of the rank differences for those rows where the F value is larger than 2.6.

²<http://www.ntu.edu.sg/home/EPNSugan/>

Table 3: Average rank values ¹ based on classification accuracies of RVFL variants.

solution & activation function		-bias, -link ²	+bias, -link ³	-bias, +link ⁴	+bias, +link ⁵	F value ⁶
Ridge Regression	sigmoid	2.7975	2.7025	2.3306	2.1694	6.7829
	sine	2.5868	2.6033	2.4174	2.3926	0.8843
	hardlim	3.0785	3.0702	1.9091	1.9421	43.0537
	tribas	2.6901	2.6736	2.2810	2.3554	3.3333
	radbas	2.6612	2.6405	2.3182	2.3802	2.2775
	sign	3.0455	3.0372	1.9545	1.9628	36.7475
Moore-Penrose pseudoinverse	sigmoid	2.5868	2.5620	2.3926	2.4587	0.5639
	sine	2.5702	2.6612	2.2769	2.4917	1.9702
	hardlim	3.0785	3.0785	1.9732	1.9298	43.8826
	tribas	2.5620	2.6116	2.4008	2.4256	0.7646
	radbas	2.5000	2.7066	2.4008	2.3926	1.5576
	sign	3.0950	3.0289	1.9404	1.9256	40.6768

¹ Each row represents a distinct comparison. Lower value indicates higher accuracy.

² RVFL without bias and without direct link.

³ RVFL with bias and without direct link.

⁴ RVFL without bias and with direct link.

⁵ RVFL with bias and with direct link.

⁶ Statistic value derived from Eq. (4).

2.3.2. Activation Function

In this section, we give a detailed analysis of different activation functions. Following the similar procedure introduced in section 2.3 to test the significance of difference of classifier’s performance, it is easy to get the critical value for F distribution with 5 and 600 degrees of freedom as 2.21 and the CD for the following Nemenyi test as 0.6196.

From Table 5, we can see that the *radbas* activation function always achieves better performance in all cases. *hardlim* and *sign* activation functions lead to the penultimate worst and the worst performances, respectively. For Ridge Regression based methods, there is a clear pattern

$$radbas > sine > tribas > sig > hardlim > sign \quad (6)$$

where “>” means that the method on the left performs better than the method on the right. However, for the Moore-Penrose pseudoinverse based methods, *tribas* always achieves the 4th rank. We also find that there is no clear superiority between *sigmoid* and *sine* activation functions.

2.3.3. Closed-form Solution

Both the Ridge Regression and Moore-Penrose pseudoinverse lead to closed-form solution for RVFL. In order to investigate the effect of these two methods on different RVFL variants, we keep the activation functions and the network structure (bias, direct link) to be the same and compare the performances. Hence, it leads to 24 pairs of comparisons in total. Sign test [36] is employed to test the statistical significance since each comparison only involves two classifiers.

If the two algorithms compared are equivalent as assumed under the null-hypothesis, each should win on approximately $N/2$ out of N datasets. The number of wins is distributed according to the binomial distribution. For a greater number of datasets, the number of wins is under the null-hypothesis distributed according to $N(N/2, \sqrt{N}/2)$,

Table 4: Result of statistical test¹. Statistically significant differences are marked with $--^2$ or $++^3$.

(a) Ridge regression with “sigmoid”

	-b,-d	+b,-d	-b,+d	+b,+d
-b,-d			--	--
+b,-d				--
-b,+d	++			
+b,+d	++	++		

(b) Ridge regression with “hardlim”

	-b,-d	+b,-d	-b,+d	+b,+d
-b,-d			--	--
+b,-d			--	--
-b,+d	++	++		
+b,+d	++	++		

(c) Ridge regression with “sign”

	-b,-d	+b,-d	-b,+d	+b,+d
-b,-d			--	--
+b,-d				--
-b,+d	++			
+b,+d	++	++		

(d) M-P⁴ pseudoinverse with “hardlim”

	-b,-d	+b,-d	-b,+d	+b,+d
-b,-d			--	--
+b,-d				--
-b,+d	++			
+b,+d	++	++		

(e) M-P pseudoinverse with “sign”

	-b,-d	+b,-d	-b,+d	+b,+d
-b,-d			--	--
+b,-d				--
-b,+d	++			
+b,+d	++	++		

¹ “+”, “-”, “b” and “d” have the same meaning as in Table 3.

² “--” means the method in the row is statistically significantly worse than the method in the column.

³ “++” means the method in the row is statistically significantly better than the method in the column.

⁴ “M-P” stands for Moore-Penrose.

which allows for the use of z -test: if the number of wins is at least $N/2 + 1.96\sqrt{N}/2$ (or, for a quick rule of a thumb, $N/2 + \sqrt{N}$), the algorithm is significantly better with $p < 0.05$. In this case, if one algorithm wins more than 71.28 times on 121 datasets, then it is considered as statistically significantly better than the other one. These cases are highlighted in Table 6. Table 6 summarizes the results for each pair of comparisons. The entry in each column represents the number of times ridge regression (pseudoinverse) is better than pseudoinverse (ridge regression) for the same activation function. For example, the first column means ridge regression is better than Moore-Penrose pseudoinverse in 59 of 121 datasets and worse in 57 of 121 datasets. Generally, ridge regression leads to a better performance for almost all cases.

2.4. Overall Comparison

In this section, we present an overall comparison of the RVFL variants. Since we have already found that *hardlim* and *sign* activation functions consistently performed poorly, RVFL with these two activation functions will be excluded in this comparison. Hence, an overall comparison with 36 RVFL variants is presented in Table 7. In the same way, pairs of methods whose ranks differ by more than 4.5589 are statistically significantly different.

2.5. Range of the random parameters

In [17], the authors indicate that the performance of the RVFL may depend on the ranges of uniformly distributed random weights. However, this issue has been untouched in the literature to the best of the authors’ knowledge. In this work, we investigate this issue by introducing one scaling factor S to control the ranges of the randomization. This process can also shed light on the effect of saturation of hidden neurons in RVFL. Based on the performance in previous section, we choose ridge regression based RVFL with *radbas* activation function because it achieves the best performance among all variants.

Table 5: Statistical significance test for different activation functions. The values in bracket stands for their average rank. Lower values stands for better performance. The mark \checkmark indicates that these two methods are statistically significantly different.

(a) -bias, -link, Ridge regression							(b) +bias, -link, Ridge regression						
	sigmoid (3.18)	sine (2.67)	hardlim (4.83)	tribas (2.95)	radbas (2.48)	sign (4.89)		sigmoid (3.16)	sine (2.69)	hardlim (4.80)	tribas (2.98)	radbas (2.49)	sign (4.88)
sigmoid (3.18)			\checkmark		\checkmark	\checkmark	sigmoid (3.16)			\checkmark		\checkmark	\checkmark
sine (2.67)			\checkmark			\checkmark	sine (2.69)			\checkmark			\checkmark
hardlim (4.83)	\checkmark	\checkmark		\checkmark	\checkmark		hardlim (4.80)	\checkmark	\checkmark		\checkmark	\checkmark	
tribas (2.95)			\checkmark			\checkmark	tribas (2.98)			\checkmark			\checkmark
radbas (2.48)	\checkmark		\checkmark			\checkmark	radbas (2.49)	\checkmark		\checkmark			\checkmark
sign (4.89)	\checkmark	\checkmark		\checkmark	\checkmark		sign (4.88)	\checkmark	\checkmark		\checkmark	\checkmark	

(c) -bias, +link, Ridge regression							(d) +bias, +link, Ridge regression						
	sigmoid (3.45)	sine (2.91)	hardlim (4.36)	tribas (2.98)	radbas (2.76)	sign (4.55)		sigmoid (3.37)	sine (2.98)	hardlim (4.37)	tribas (3.02)	radbas (2.71)	sign (4.54)
sigmoid (3.45)			\checkmark		\checkmark	\checkmark	sigmoid (3.37)			\checkmark		\checkmark	\checkmark
sine (2.91)			\checkmark			\checkmark	sine (2.98)			\checkmark			\checkmark
hardlim (4.36)	\checkmark	\checkmark		\checkmark	\checkmark		hardlim (4.37)	\checkmark	\checkmark		\checkmark	\checkmark	
tribas (2.98)			\checkmark			\checkmark	tribas (3.02)			\checkmark			\checkmark
radbas (2.76)	\checkmark		\checkmark			\checkmark	radbas (2.71)	\checkmark		\checkmark			\checkmark
sign (4.55)	\checkmark	\checkmark		\checkmark	\checkmark		sign (4.54)	\checkmark	\checkmark		\checkmark	\checkmark	

(e) -bias, -link, Moore-Penrose pseudoinverse							(f) +bias, -link, Moore-Penrose pseudoinverse						
	sigmoid (2.79)	sine (2.82)	hardlim (4.68)	tribas (3.36)	radbas (2.60)	sign (4.75)		sigmoid (2.73)	sine (2.96)	hardlim (4.68)	tribas (3.34)	radbas (2.59)	sign (4.69)
sigmoid (2.79)			\checkmark			\checkmark	sigmoid (2.73)			\checkmark			\checkmark
sine (2.82)			\checkmark			\checkmark	sine (2.96)			\checkmark			\checkmark
hardlim (4.68)	\checkmark	\checkmark		\checkmark	\checkmark		hardlim (4.68)	\checkmark	\checkmark		\checkmark	\checkmark	
tribas (3.36)			\checkmark		\checkmark	\checkmark	tribas (3.34)			\checkmark		\checkmark	\checkmark
radbas (2.60)			\checkmark	\checkmark		\checkmark	radbas (2.59)	\checkmark		\checkmark			\checkmark
sign (4.75)	\checkmark	\checkmark		\checkmark	\checkmark		sign (4.69)	\checkmark	\checkmark		\checkmark	\checkmark	

(g) -bias, +link, Moore-Penrose pseudoinverse							(h) +bias, +link, Moore-Penrose pseudoinverse						
	sigmoid (3.10)	sine (2.95)	hardlim (4.19)	tribas (3.56)	radbas (2.93)	sign (4.27)		sigmoid (3.17)	sine (3.05)	hardlim (4.08)	tribas (3.58)	radbas (2.99)	sign (4.12)
sigmoid (3.10)			\checkmark			\checkmark	sigmoid (3.17)			\checkmark			\checkmark
sine (2.95)			\checkmark			\checkmark	sine (3.05)			\checkmark			\checkmark
hardlim (4.19)	\checkmark	\checkmark		\checkmark	\checkmark		hardlim (4.08)	\checkmark	\checkmark		\checkmark	\checkmark	
tribas (3.56)			\checkmark			\checkmark	tribas (3.58)			\checkmark			\checkmark
radbas (2.93)			\checkmark			\checkmark	radbas (2.99)			\checkmark			\checkmark
sign (4.27)	\checkmark	\checkmark		\checkmark	\checkmark		sign (4.12)	\checkmark	\checkmark		\checkmark	\checkmark	

In previously section, the random weights are generated with uniform distribution in $[-1, 1]$, as done exactly in [35], while the biases are in $[0, 1]$. In this section, the random weights and biases are generated with uniform distribution in $[-S, S]$ and $[0, S]$ respectively, where S is a positive scaling factor. In this work, we set $S = 2^t$, t is set to be $-5:0.5:5$. The performances 21 RVFL configurations with direct links and bias are summarized in Fig. 2. It is obvious that all RVFL variants perform poorly when the range of the random parameters becomes either too large or too small. Another interesting conclusion is the commonly adopted approach that $S = 1$ for the randomization may not lead to the optimal performance.

For RVFL without direct link, setting $t > 0$ for scaling factor S to increase the discrimination power of the features in the hidden neurons may make more neurons to saturate. This can be compensated by either having more hidden neurons or the direct link from the input layer to the output layer. On the other hand, setting $t < 0$ for scaling factor S to reduce the possibility of neuronal saturation may reduce the discrimination power of the features in the hidden

Table 6: Comparisons between ridge regression and Moore-Penrose pseudoinverse. The entry in each column represents the number of times ridge regression (pseudoinverse) is better than pseudoinverse (ridge regression) for the same activation function. Statistically significant columns are highlighted.

(a) -bias, -link,

	sigmoid	sine	hardlim	tribas	radbas	sign
Ridge Regression	59	61	80	78	63	80
Moore-Penrose pseudoinverse	57	56	36	38	54	37

(b) +bias, -link,

	sigmoid	sine	hardlim	tribas	radbas	sign
Ridge Regression	57	70	81	77	63	81
Moore-Penrose pseudoinverse	58	48	36	40	54	37

(c) -bias, +link,

	sigmoid	sine	hardlim	tribas	radbas	sign
Ridge Regression	64	57	68	76	58	63
Moore-Penrose pseudoinverse	52	56	46	42	59	51

(d) +bias, +link,

	sigmoid	sine	hardlim	tribas	radbas	sign
Ridge Regression	68	61	67	78	63	63
Moore-Penrose pseudoinverse	49	56	47	39	53	51

Table 7: Overall Comparison based on overall ranks of 32 RVFL variants ¹. Lower values in rank reflects better performance.

Method	RR, -b,-d,sigmoid,	RR, +b,-d,sigmoid,	RR, -b,+d,sigmoid,	RR, +b,+d,sigmoid
Rank	18.3182	18.0579	17.0000	16.6860
method	RR, -b,-d,sine,	RR,+b,-d,sine,	RR,-b,+d,sine,	RR,+b,+d,sine
Rank	15.1157	15.2603	14.5455	14.9421
method	RR, -b,-d,tribas,	RR,+b,-d,tribas,	RR,-b,+d,tribas,	RR,+b,+d,tribas
Rank	16.7810	16.7479	15.1446	15.5702
method	RR, -b,-d,radbas,	RR,+b,-d,radbas,	RR,-b,+d,radbas,	RR,+b,+d,radbas
Rank	13.9545	14.0083	13.2686	13.4256
Method	MP, -b,-d,sigmoid,	MP, +b,-d,sigmoid,	MP, -b,+d,sigmoid,	MP, +b,+d,sigmoid
Rank	17.0124	16.7149	16.5124	16.9298
method	MP, -b,-d,sine,	MP,+b,-d,sine,	MP,-b,+d,sine,	MP,+b,+d,sine
Rank	17.3223	18.2066	15.7107	16.3678
method	MP, -b,-d,tribas,	MP,+b,-d,tribas,	MP,-b,+d,tribas,	MP,+b,+d,tribas
Rank	20.8636	20.9876	19.7934	20.1612
method	MP, -b,-d,radbas,	MP,+b,-d,radbas,	MP,-b,+d,radbas,	MP,+b,+d,radbas
Rank	15.4504	16.5785	15.0496	15.5124

¹ RR and MP stand for ridge regression and Moore-Penrose pseudoinverse, respectively. “+”, “-”, “b” and “d” are the same meaning as in Table 3.

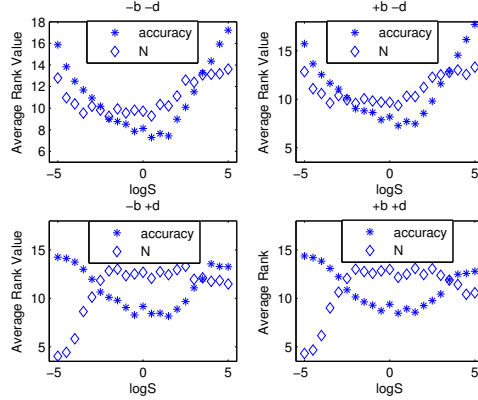


Figure 2: Performance of RVFL based for different ranges of randomization. Smaller rank indicates better accuracy and less number of hidden neurons. N stands for the number of hidden neurons corresponding to the testing accuracy used in the ranking.

neurons. Again, this can be compensated to some degree by having more hidden neurons or the direct link from the input layer to output layer. In Table 8 and Table 9, we also present rank values based on accuracies and number of hidden neurons of different RVFL variants for different scaling factor S with $t = -1.5 : 0.5 : 1.5$ since their performance is consistently better than others as indicated in Fig. 2. The average rank values based on accuracies of RVFL when all parameters (N, λ, S) are tuned are presented in Table 10.

Table 8: Average rank values of RVFL based on accuracies for different scale factor values over 121 datasets. Each column stands for a distinct comparison. Lower rank means higher accuracy. All RVFLs use ridge regression solution and *radbas* activation function.

Method	$S = 2^{-1.5}$	$S = 2^{-1}$	$S = 2^{-0.5}$	$S = 1$	$S = 2^{0.5}$	$S = 2^1$	$S = 2^{1.5}$
-b,-d	2.6446	2.5992	2.5992	2.6612	2.7231	2.6322	2.6653
+b,-d	2.6529	2.5785	2.6446	2.6405	2.6736	2.6860	2.7355
-b,+d	2.4298	2.4752	2.3140	2.3182	2.3430	2.4215	2.3182
+b,+d	2.2727	2.3471	2.4421	2.3802	2.2603	2.2603	2.2810

Table 9: Average rank values of RVFL based on the number of hidden neurons for different scale factor values over 121 datasets. Each column stands for a distinct comparison. Lower rank means less number of hidden neurons. All RVFLs use ridge regression solution and *radbas* activation function.

Method	$S = 2^{-1.5}$	$S = 2^{-1}$	$S = 2^{-0.5}$	$S = 1$	$S = 2^{0.5}$	$S = 2^1$	$S = 2^{1.5}$
-b,-d	2.6818	2.6529	2.6322	2.5868	2.6653	2.6736	2.6942
+b,-d	2.6488	2.7025	2.5785	2.5785	2.6488	2.6570	2.6281
-b,+d	2.3760	2.3140	2.3760	2.4091	2.3636	2.3926	2.3099
+b,+d	2.2934	2.3306	2.4132	2.4256	2.3223	2.2769	2.3678

Results in Table 8 and Table 10 are consistent with the previous subsections which clearly indicate the advantage of the direct link. Moreover, RVFL with direct link achieves better accuracy with less number of hidden neurons. The

Table 10: Average rank values of RVFL based on accuracies when all parameters (N , λ , S) are tuned. Lower rank means higher accuracy. All RVFLs use ridge regression solution and *radbas* activation function.

	-b,-d	+b,-d	-b,+d	+b,+d
Rank	2.8140	2.6653	2.3430	2.1777

direct link from input layer to output layer serves as a standing out regularization and make a high chance for RVFL to achieve a better performance with high possibility than those without direct link. That is, with smaller number of random hidden neurons, direct link in RVFL leads to a thinner and simpler model than those without. For a given set of observations or data, there is always an infinite number of possible hypotheses fit the same data. According to the Occams Razor principle, one should choose from a set of otherwise equivalent models of a given phenomenon the simplest one. For example, it is possible for us to further increase the performance of RVFL if we enlarge the number of random hidden neurons. “Super flat” RVFL models (i.e. with a large number of hidden neurons) are more likely to overfit the available data. This is also in line with the PAC learning theory [18] that advocates for learning with lower complexity models. Hence, the tuning range for the number of hidden neurons can be relatively narrower for RVFL with direct links.

3. Concluding Remarks

In this work we presented extensive and comprehensive evaluation of variants of RVFL with closed-form solution by using 121 UCI datasets [11]. The conclusion of our investigations are as follows:

1. The effect of the direct links from the input layer to the output layer. It turns out that the direct links lead to better performance than those without in all cases as seen in Table 4.
2. The effect of the bias in the output layer. It turns out that the bias term in the output neurons only has mixed effects on the performance, as it may or may not improve performance. Hence, bias can be a tunable network configuration depending on the specific problem.
3. Effect of scaling the randomization range of input weights and biases. We show scaling down the randomization range of input weights and biases to avoid saturating the neurons may risk at degenerating the discrimination power of the random features. However, this can be compensated by having more hidden neurons or direct link. Scaling the randomization range of input weights and biases up to enhance the discrimination power of the random features may risk saturating the neurons. Again, this can be compensated by having more hidden neurons or combining with the direct link from the input to the output layer. However, for reasons explained in Subsection 2.5, we prefer lower model complexity.
4. The performance of 6 commonly used activation functions summarized in Table 1. It turns out that *radbas* function always leads to a better performance. *hardlim* and *sign* activation functions lead to penultimate worst and worst performances, respectively.
5. The performance of Moore-Penrose pseudoinverse and ridge regression (or regularized Least Square) solutions for the output weights. It turns out that with one more parameter (λ in Eq. (2)) to tune, ridge regression based RVFL shows better performance than the Moore-Penrose pseudoinverse based RVFL.

This work sets a basis for future research for random vector functional link network. Future studies and developments of RVFL may include:

1. Performance of RVFL ensemble. Neural network has low bias and high variance [4]. Hence, the performance of RVFL can be significantly improved by ensemble methods.
2. Performance of kernel methods with RVFL. Recent research [32, 33, 2] shows the success of random features for large-scale kernel machines. Hence, it is worthy to investigate kernel machines with random features extracted from RVFL.

3. Performance of deep RVFL structure. Recent work in computer vision and machine learning community demonstrates the success of deep neural networks [19]. Hence, how to design a good deep RVFL structure for a specific problem is an open problem now.

Acknowledgement

The authors would like to thank the Guest Editors and the reviewers for their valuable comments. In particular, authors thank the managing Guest Editor Associate Professor Dianhui Wang for suggesting us to investigate the scaling of randomization. Results presented in Section 2.5 show overall performance enhancement due to tuning the scaling of randomization.

References

- [1] Alhamdoosh, M. and Wang, D. (2014). Fast decorrelated neural network ensembles with random weights. *Information Sciences*, 264:104–117.
- [2] An, S., Liu, W., and Venkatesh, S. (2007). Face recognition using kernel ridge regression. In *2007 IEEE Conference on Computer Vision and Pattern Recognition*, pages 1–7. IEEE.
- [3] Bartlett, P. L. (1998). The sample complexity of pattern classification with neural networks: the size of the weights is more important than the size of the network. *IEEE Transactions on Information Theory*, 44(2):525–536.
- [4] Breiman, L. (1998). Arcing classifier (with discussion and a rejoinder by the author). *The Annals of Statistics*, 26(3):801–849.
- [5] Chen, C. P. and Wan, J. Z. (1999). A rapid learning and dynamic stepwise updating algorithm for flat neural networks and the application to time-series prediction. *IEEE Transactions on Systems, Man, and Cybernetics, Part B: Cybernetics*, 29(1):62–72.
- [6] Chi, H.-M. and Ersoy, O. K. (2005). A statistical self-organizing learning system for remote sensing classification. *IEEE Transactions on Geoscience and Remote Sensing*, 43(8):1890–1900.
- [7] Cun, L., Boser, B., Denker, J. S., Henderson, D., Howard, R. E., Hubbard, W., and Jackel, L. D. (1990). Handwritten digit recognition with a back-propagation network. In *Advances in Neural Information Processing Systems*, pages 396–404.
- [8] Dehuri, S. and Cho, S.-B. (2010). A comprehensive survey on functional link neural networks and an adaptive pso–bp learning for cflnn. *Neural Computing and Applications*, 19(2):187–205.
- [9] Demšar, J. (2006). Statistical comparisons of classifiers over multiple data sets. *Journal of Machine Learning Research*, 7:1–30.
- [10] Denker, J. S., Gardner, W., Graf, H. P., Henderson, D., Howard, R., Hubbard, W., Jackel, L. D., Baird, H. S., and Guyon, I. (1989). Neural network recognizer for hand-written zip code digits. In *Advances in Neural Information Processing Systems*, pages 323–331.
- [11] Fernández-Delgado, M., Cernadas, E., Barro, S., and Amorim, D. (2014). Do we need hundreds of classifiers to solve real world classification problems? *Journal of Machine Learning Research*, 15(1):3133–3181.
- [12] Friedman, M. (1937). The use of ranks to avoid the assumption of normality implicit in the analysis of variance. *Journal of the American Statistical Association*, 32(200):675–701.
- [13] Friedman, M. (1940). A comparison of alternative tests of significance for the problem of m rankings. *Annals of Mathematical Statistics*, 11(1):86–92.
- [14] Hornik, K., Stinchcombe, M., and White, H. (1989). Multilayer feedforward networks are universal approximators. *Neural Networks*, 2(5):359–366.
- [15] Husmeier, D. and Taylor, J. G. (1997). Modelling conditional probabilities with committees of RVFL networks. In *Artificial Neural Networks—ICANN’97*, pages 1053–1058. Springer.
- [16] Husmeier, D. and Taylor, J. G. (1998). Neural networks for predicting conditional probability densities: Improved training scheme combining EM and RVFL. *Neural Networks*, 11(1):89–116.
- [17] Igel'nik, B. and Pao, Y.-H. (1995). Stochastic choice of basis functions in adaptive function approximation and the functional-link net. *IEEE Transactions on Neural Networks*, 6(6):1320–1329.
- [18] Kearns, M. J. and Vazirani, U. V. (1994). *An introduction to computational learning theory*. MIT press.
- [19] Krizhevsky, A., Sutskever, I., and Hinton, G. E. (2012). Imagenet classification with deep convolutional neural networks. In *Advances in Neural Information Processing Systems*, pages 1097–1105.
- [20] Leshno, M., Lin, V. Y., Pinkus, A., and Schocken, S. (1993). Multilayer feedforward networks with a nonpolynomial activation function can approximate any function. *Neural Networks*, 6(6):861–867.
- [21] Li, W., Wang, D., and Chai, T. (2015). Multisource data ensemble modeling for clinker free lime content estimate in rotary kiln sintering processes. *IEEE Transactions on Systems, Man, and Cybernetics: Systems*, 45(2):303–314.
- [22] Lichman, M. (2013). UCI machine learning repository. <http://archive.ics.uci.edu/ml>.
- [23] Martínez-Villena, J. M., Rosado-Muñoz, A., and Soria-Olivas, E. (2014). Hardware implementation methods in random vector functional-link networks. *Applied intelligence*, 41(1):184–195.
- [24] Nemenyi, P. (1963). *Distribution-free Multiple Comparisons*. Princeton University.
- [25] Pao, Y. (1989). Adaptive pattern recognition and neural networks.
- [26] Pao, Y.-H., Park, G.-H., and Sobajic, D. J. (1994). Learning and generalization characteristics of the random vector functional-link net. *Neurocomputing*, 6(2):163–180.
- [27] Pao, Y.-H. and Phillips, S. M. (1995). The functional link net and learning optimal control. *Neurocomputing*, 9(2):149–164.
- [28] Pao, Y.-H., Phillips, S. M., and Sobajic, D. J. (1992). Neural-net computing and the intelligent control of systems. *International Journal of Control*, 56(2):263–289.

- [29] Park, G. H., Lee, Y. J., and LeClair, S. R. (2000). Intelligent rate control for MPEG-4 coders. *Engineering Applications of Artificial Intelligence*, 13(5):565–575.
- [30] Park, G. H. and Pao, Y. H. (2000). Unconstrained word-based approach for off-line script recognition using density-based random-vector functional-link net. *Neurocomputing*, 31(1):45–65.
- [31] Park, J. and Sandberg, I. W. (1991). Universal approximation using radial-basis-function networks. *Neural Computation*, 3(2):246–257.
- [32] Rahimi, A. and Recht, B. (2007). Random features for large-scale kernel machines. In *Advances in Neural Information Processing Systems*, pages 1177–1184.
- [33] Saunders, C., Gammerman, A., and Vovk, V. (1998). Ridge regression learning algorithm in dual variables. In *Proceedings of the 15th International Conference on Machine Learning*, pages 515–521. Morgan Kaufmann.
- [34] Scardapane, S., Wang, D., Panella, M., and Uncini, A. (2015). Distributed learning for random vector functional-link networks. *Information Sciences*, 301:271–284.
- [35] Schmidt, W. F., Kraaijveld, M., and Duin, R. P. (1992). Feedforward neural networks with random weights. In *11th IAPR International Conference on Pattern Recognition*, pages 1–4. IEEE.
- [36] Sheskin, D. J. (2003). *Handbook of parametric and nonparametric statistical procedures*. CRC Press.
- [37] Tyukin, I. and Prokhorov, D. (2009). Feasibility of random basis function approximators for modeling and control. In *IEEE International Symposium on Intelligent Control*, pages 1391–1396.
- [38] Wang, Z., Yoon, S., Xie, S. J., Lu, Y., and Park, D. S. (2013). Random vector functional-link net based pedestrian detection using multi-feature combination. In *2013 6th International Congress on Image and Signal Processing (CISP)*, volume 2, pages 773–777. IEEE.
- [39] Wang, Z., Yoon, S., Xie, S. J., Lu, Y., and Park, D. S. (2014). A high accuracy pedestrian detection system combining a cascade adaboost detector and random vector functional-link net. *The Scientific World Journal*, 2014.