1
2
# Cascade neural network algorithm with analytical connection weights determination for modelling operations and energy applications

3    Zhengxu Wang[1#], Waqar Ahmed Khan[2*#], Hoi-Lam Ma[3] & Xin Wen[2]

4    The performance and learning speed of the Cascade Correlation neural network
5    (CasCor) may not be optimal because of redundant hidden units' in the cascade
6    architecture and the tuning of connection weights. This study explores the
7    limitations of CasCor and its variants and proposes a novel constructive neural
8    network (CNN). The basic idea is to compute the input connection weights by
9    generating linearly independent hidden units from the orthogonal linear
10   transformation, and the output connection weights by connecting hidden units in a
11   linear relationship to the output units. The work is unique in that few attempts have
12   been made to analytically determine the connection weights on both sides of the
13   network. Experimental work on real energy application problems such as
14   predicting powerplant electrical energy, predicting seismic hazards to prevent fatal
15   accidents and reducing energy consumption by predicting building occupancy
16   detection shows that analytically calculating the connection weights and
17   generating non-redundant hidden units improves the convergence of the network.
18   The proposed CNN is compared with that of the state-of-the-art machine learning
19   algorithms. The work demonstrates that proposed CNN predicts a wide range of
20   applications better than other methods.

21   Keywords: **energy management; forecasting; machine learning; neural networks;**
22   **sustainability**

## 23  1. Introduction

24   The Cascade Correlation learning algorithm (CasCor) has been extensively applied in

[1] School of Business Administration, Institute of Supply Chain Analytics, Dongbei University of Finance and Economics, Dalian, People's Republic of China.
[2] Department of Industrial and Systems Engineering, The Hong Kong Polytechnic University, Hung Hom, Hong Kong.
[3] Department of Supply Chain and Information Management, The Hang Seng University of Hong Kong, Shatin, Hong Kong.
[*] Correponding author.
[#] Zhengxu Wang and Waqar Ahmed Khan contributed equally to this work.

25    many application areas (Heidari et al. 2018; Chung, Ma, and Chan 2017) due to its self-

26    organizing neural network property, and in many cases, it is considered to be more

27    powerful than the standard multilayer perceptron (Qiao et al. 2016; Hunter et al. 2012).

28    The selection of a neural network (NN) depends upon the application area (Wang et al.

29    2018; Deng et al. 2019) to achieve better and faster convergence. Learning NNs by

30    gradient algorithms along with too many hyperparameters may make the network more

31    complex, causing the generalization performance to converge at a suboptimal solution

32    (Liew, Khalil-Hani, and Bakhteri 2016; Kapanova, Dimov, and Sellier 2018).

33    Backpropagation (BP) gradient descent is a well-known learning algorithm for NNs

34    (Rumelhart, Hinton, and Williams 1986), but it faces the problem of local minima if the

35    global minima is far away, and the learning speed is highly influenced by gradient

36    iteration and the learning rate hyperparameter (Hecht-Nielsen 1989). To address the

37    backpropagation neural network (BPNN) slowness and topology problem, the self-

38    organizing quick prop (QP) CasCor was formulated (Fahlman and Lebiere 1990).

39       The QP can reduce the error of CasCor to a small value, but it does not guarantee

40    that the network performance will be satisfactory (Hwarng 2005; Hunter et al. 2012) due

41    to its chaotic behaviour and numerical instability. QP during weight updating takes much

42    larger steps based on previous and current gradients to moves faster towards the minimum

43    of the function (Fahlman 1988). The current gradient may be larger or smaller and in the

44    same or opposite direction to the previous gradient. The larger and opposite gradient will

45    cause the algorithm to cross the minimum of the function and needs to be brought back.

46    This may cause the QP to behave chaotically across the minimum valley of the function.

47    Banerjee et al. (2011) explained that QP becomes numerically unstable if the current

48    gradient is very close or equal to the previous gradient. If the difference between current

49  and previous gradient becomes zero, the weight difference will also become zero and the

50  QP formula will remain zero permanently, even if the gradient changes.

51  Due to its widespread popularity and the recent increase in interest for self-

52  organizing neural networks (Khan et al. 2019a, 2019b), researchers are extensively

53  focused on improving the existing CasCor. Huang, Song, and Wu (2012) proposed an

54  orthogonal least squares algorithm for training cascade neural networks (OLSCN) by

55  explaining that a larger network size causes lowering the generalization performance of

56  CasCor. Besides, the covariance objective function efforts to adjust the input connection

57  weights cannot assure maximum error reduction on the addition of a new hidden unit.

58  The repeatedly tuning of connection weights, before and after hidden unit generation,

59  causes the network to be more time-consuming. However, Qiao et al. (2016) explained

60  that the new objective function formulated along with the modified Newton method by

61  OLSCN may make mistakes during linear dependencies among variables and results in

62  local minimum with slow convergence. A Faster Cascade Neural Network (FCNN) was

63  proposed to address the CasCor and OLSCN generalization and convergence issues.

64  FCNN selects linearly independent input units one by one by the Gram-Schmidt

65  Orthogonalization method and candidate units by the modified index (MI) formulated

66  objective function. It assures that the selected candidate unit (hidden unit) may have the

67  largest contribution in the existing candidate pools but cannot guarantee that the next

68  expected candidate unit (hidden unit) error reduction will be maximized. For the sake of

69  simplicity, in this paper CasCor, OLSCN and FCNN are referred to as CasCor and its

70  variants because of similar network structure, unless specified.

71  This paper proposes a novel Cascade Principal Component Least Squares Neural

72  Network Learning Algorithm (CPCLS) to address the convergence limitations of CasCor

73  and its variants. The main contributions are listed below:

74    • The linear dependence among input units and/or hidden units can be avoided by

75       transforming a set of correlated units orthogonally into linearly independent units.

76       The cascade architecture can be made better by connecting hidden units (or layers)

77       to the output units that may have no linear dependence with each other. Similarly,

78       the input unit's direct linear connection to the output units can be avoided to get

79       rid of the input unit's dependency.

80    • The best least-squares solution can be achieved by connecting only newly added

81       linearly independent (no multicollinearity) hidden layer to the output units and

82       eliminating previous output connections (hidden units).

83    • Multiple hidden units can be generated in the hidden layer to make the

84       convergence faster.

85    The advancement in information technology has enabled industries to create a

86    model of products and processes from high dimensional data to benefit production

87    research (Kusiak 2020). Traditional models based on mathematical formulations and

88    physical approaches advantageous to provide a physical understanding of the system.

89    However, in real practices, mathematical models may be inaccurate and difficult to adopt

90    because of ignoring nonlinearities (Wang et al. 2019), unable to understand symbolic

91    data, need of prior expert knowledge, and maybe not well suited to represent relationships

92    among variables (Kuo and Kusiak 2019).

93    Nowadays, the availability of high dimensional data has made it possible to

94    extract useful information, rather than physical measurement or manual work that may

95    cause subjective judgment or fatigues (Kim et al. 2019), to facilitate in making real-time

96    decisions, time and cost-saving (Q. Liu et al. 2019). It is considered that the application

97    of machine learning compared to mathematical modelling is likely to be more beneficial

98  in improving production research (Kusiak 2020; Kuo and Kusiak 2019; Lv et al. 2020;

99  Y. Liu et al. 2019). The machine learning that has gained significant interest in the

100  literature include NNs and its variants (Kumar, Singh, and Singh 2020; Ertuğrul 2018;

101  Bansal et al. 2019; Zou et al. 2018; Lorencin et al. 2019; Grasso, Luchetta, and Manetti

102  2018; Nayyeri et al. 2018), support vector machine (SVM) (Bansal et al. 2019), decision

103  tree (DT) (Mantas et al. 2019; Bansal et al. 2019; Candanedo and Feldheim 2016), naïve

104  Bayes (NB) (Bansal et al. 2019), metaheuristics search algorithms and its variants (Bansal

105  et al. 2019; Aljarah, Faris, and Mirjalili 2018), random forest (RF) (Mantas et al. 2019;

106  Candanedo and Feldheim 2016), ensembles (Mantas et al. 2019), gradient boosting

107  machine (Candanedo and Feldheim 2016), regression and its variants (Lorencin et al.

108  2019), and linear discriminant analysis (LDA) (Candanedo and Feldheim 2016).

109  Compared to other machine learning algorithms, NNs is widely adopted because

110  of its superior performance and universal approximation ability (Wang et al. 2019).

111  Usually, the application of NNs in production research involves learning of the

112  connection weights by either BP or random generation with a lot of hyperparameter

113  tuning (Chien, Lin, and Lin 2020) which makes learning complicated and challenging

114  (Kusiak 2020; Solimanpur, Vrat, and Shankar 2004). According to the best of our

115  knowledge, insufficient attempts have been made to improve the NNs performance and

116  speed by analytically calculating connection weights on both sides of the network with a

117  small number of hyperparameters initialization. The novelty of the proposed algorithm

118  exists in its improved cascade architecture by connecting linearly independent hidden

119  layer to the output units and analytically calculating connection weights. This may

120  facilitate to predict a wide range of applications with less human intervention.

121  This work applies the proposed CPCLS algorithm and made a state-of-the-art

122  comparison with other machine learning algorithms to predict health sciences,

123 engineering, marine, food products, forestry, and energy application problems. Better

124 generalization performance and faster learning speed of CPCLS give insight that NNs

125 based model prediction capability can be made better by analytically calculating

126 connection weights rather than BP or random generation. Moreover, in current practice,

127 the majority of the production research is focused on solving problems belonging to a

128 single application. This limits the proposed method, in real practice, to a single industry

129 or single business function. The better performance of CPCLS on a wide range of

130 applications give managerial insight that it can be practiced in general and able to handle

131 industrial and business function problems on an integrated platform. Furthermore, the

132 cascade architecture of CPCLS helps to eliminate the problem of "what-if" of fixed

133 topology BPNNs for determining hidden units and layers that involves human

134 interventions and simultaneously affect decision making. The CPCLS can facilitate in

135 optimizing the operations by providing predictive advice and may derive the decision-

136 making process by building greater confidence in prediction from historical data rather

137 than mathematical formulation or manual work.

138   This paper is a revised and extended version of that of Khan, Chung, and Chan

139 (2018). In this extended version, the property of maximum error reduction of the CPCLS

140 is explained by supporting statements, lemmas, theoretical analysis, and remarks and

141 further demonstrated by experimental work. The rest of the paper is structured as follows.

142 In Section 2, CasCor and its variants with convergence limitations, Orthogonal linear

143 transformation (OLT) and Ordinary Least Squares (OLS) are briefly explained. Section

144 3 presents the novel CPCLS. Section 4 describes the state-of-the-art comparison. Section

145 5 concludes the paper.

## 2. Existing learning methodologies

### 2.1. CasCor and its variants with convergence drawbacks

CasCor initializes by linearly connecting the input units to the output units and tuning randomly generated output connection weights by the QP learning algorithm. When training converges, hidden units are added one by one to discover nonlinear patterns in the problem. The candidate units are added to select the hidden unit, having the property of maximum error reduction. The candidate units receive the input connections from input units and any pre-existing hidden units. The aim is to maximize the covariance $S$ between network error and the candidate units by the gradient ascent. When $S$ stops improving, the candidate unit with the maximum value of the $S$ is chosen as the hidden unit and is linked to the output units by the output connection weights, while incoming connections are kept frozen. Again, the output connection weights are trained by the QP and this procedure continues till the error converges. Figure 1 illustrates the architecture of CasCor.

Huang, Song, and Wu (2012) explained that the $S$ objective function to maximize the correlation between the hidden unit and network error cannot assure a maximum error reduction with the addition of new hidden unit to the network. Secondly, the output training is repeatedly performed after every hidden unit generation which increases the computational burden. OLSCN was proposed to overcome the above disadvantages which lead CasCor to slow convergence and poor generalization performance. The OLSCN reformulated new objective function based on the OLS for input training which was further optimized by the second order modified Newton method. Qiao et al. (2016) supported the work of Huang, Song, and Wu (2012) and concluded that the CasCor objective function cannot guarantee a maximum error reduction and repeatedly output training can be more time-consuming. However, Qiao et al. (2016) argued that the

171  OLSCN may result in a local minimum, slow convergence, and a computational burden

172  by updating the weights of the hidden units by the modified Newton method. In addition,

173  linear independence of the input units and the hidden units are necessary for QR

174  factorization and the newly formulated objective function, respectively. FCNN was

175  proposed to address the generalization performance and learning speed of CasCor and

176  OLSCN.

177      In Theorem 3.1 (Qiao et al. 2016) of FCNN, it is explained that one or more

178  candidate units in the pool may be linearly independent to the input and any pre-existing

179  hidden units. However, the column matrix of hidden units may not necessarily full rank

180  due to the random generation of input weights. Therefore, MI was proposed to evaluate

181  the candidate unit among the pool of candidates. The candidate unit with the maximum

182  contribution to the sum of squares error (SSE) is added to the network which specifies

183  linearly independence of the candidate unit, however, network optimal error

184  minimization capability cannot be guaranteed. For instance, if among a pool of candidate

185  units, fewer candidate units are linearly independent than the chances of getting the

186  largest contributed MI also decreases. Secondly, the selected candidate unit (hidden unit)

187  may have the property of maximum error reduction capability among the existing

188  candidate pool which cannot assure that the next expecting candidate unit (hidden unit)

189  error reduction will be maximized. This may cause the network to generate some hidden

190  units with less error minimization capability. Eventually, more hidden units need to be

191  added by randomly generating input weight and bias which may make the network more

192  complex. For better understanding, Figure 9 (Qiao et al. 2016) in experimental work

193  illustrates the same problem of not achieving maximal error reduction by FCNN at each

194  hidden unit. It can be seen that error reduction by adding a new hidden unit is not smooth

195 and the objective of maximum error reduction by next newly added hidden units is not

196 achieved. This may result in redundant hidden units with minor effect on the convergence.

197 *2.2. OLT and OLS*

198 This section describes the existing methodologies that assist proposed CPCLS to

199 analytically calculate the connection weights for achieving maximum error reduction on

200 each hidden layer generation. Consider a training data sample with $(X, Y)$, where $X$ is

201 the input unit matrix of $m \times n$ and $Y$ is the output unit matrix of $m \times q$ with hidden units

202 matrix $H$ of $m \times p$. The input connection weights matrix of $n \times p$ is exemplified by $W$,

203 whereas, the output connection weights matrix of $p \times q$ are exemplified by $\boldsymbol{\beta}$.

204 OLT generates new $p$-features space of linearly independent $H$ by orthogonally

205 transforming $n$-features $X$ (Jolliffe 2002). It helps to reduce the dimensionality of the

206 correlated $X$ by determining the unknown components $W$, with each component

207 explaining the amount of variance in the data. OLT initializes by determining the

208 covariance matrix $S$ of equal dimension $n \times n$ matrix, with diagonal numbers indicating

209 covariance for the same feature and each number indicating the covariance between $n$-

210 features of $X$, to compute the eigenvalue $\lambda$ and its corresponding eigenvector:

$$S = \frac{1}{m-1}(X - \overline{x})^T(X - \overline{x}) \tag{1}$$

211 where $\overline{x} = \sum_{i=1}^{m} x_i$, with each quantity indicating the mean of $n$ features.

212 The eigenvector, explaining the coordinate system for the new $p$-features by

213 decreasing dimensions equal to or less than $n$-features, selection is based on the $\lambda$ value.

214 The $\lambda$ is computed from the $S$ matrix:

$$|S - \lambda I| = 0 \tag{2}$$

215  The corresponding eigenvector based on highest $\lambda$ can be determined by

216 computing the component $W$:

$$(S - \lambda I)W = 0 \tag{3}$$

217  The matrix $W$ linearly transforms $n$-features $X$ into new linearly independent $p$-

218 features $H$ :

$$H = XW \tag{4}$$

219  OLS reduces the estimation error between the predicted $\widehat{Y}$ and the observed $Y$

220 variables by determining the unknown parameter $\beta$ (Goldberger 1964):

$$Y = H\beta + e \tag{5}$$

221  OLS theory is used for determining $\beta$ by:

$$\beta = (H^T H)^{-1} H^T Y \tag{6}$$

222 where $(H^T H)^{-1} H^T$ is the Moore Penrose pseudo-inverse of matrix $H$ . For better

223 convergence, there should be no linear dependence among $H$.

224  In the last step, the $\widehat{Y}$ is determined:

$$\widehat{Y} = H\beta \tag{7}$$

225  Better network convergence can be achieved by optimally calculating the

226 connection weights in the forward step. Equations (3) and (6) play a key role in

227 determining the connection weights for the novel CPCLS.

228 **3. Proposed CPCLS learning algorithm**

229 Like CasCor and its variants, which have a similar network structure, CPCLS also works

230 on two concepts of cascade architecture and learning. Figure 2 illustrates CPCLS

231 architecture, which is an improved form of CasCor and its variants. *Firstly*, CPCLS

232 connects input units to the output units by the linearly independent hidden units to avoid

233     the linear dependency of the input units. *Secondly*, more than single hidden units can be

234     generated in the hidden layer to achieve faster convergence. *Thirdly*, the newly generated

235     hidden layer is only linked to the output units, and earlier connections are removed to

236     avoid the linear dependence of the hidden units among the hidden layers. In learning,

237     CasCor repeatedly tunes the connection weights in forward and backward steps by the

238     gradient method, while its variants either perform the gradient method or randomly

239     generate the input weights, which can take more time, and it is equally problematic to

240     control convergence. CPCLS eliminates the need for random generation and gradient

241     methods by analytically computing the connection weights in the forward step.

242     ### *3.1 Supporting statement and lemma*

243     Statement 1: (Jolliffe 2002) OLT: The $X$ values of $n$-features are orthogonally

244     transformed into a linearly independent $H$ of $p$-features by determining the eigenvalue $\lambda$

245     and its eigenvector $W$ from the input covariance $S$.

246        Remark 1: Statement 1 implies that the hidden units generated are linearly

247     independent (uncorrelated) because of the OLT of the input features.

248        Lemma 1: (Huang, Zhu, and Siew 2006) Given a standard Single hidden Layer

249     Feedforward Network (SLFN) with $N$ hidden nodes and activation function $g: R \rightarrow R$,

250     which is infinitely differentiable in any interval, for $N$ arbitrary distinct samples $(x_i, y_i)$,

251     where $x_i \in \mathbf{R^n}$ and $y_i \in \mathbf{R^m}$, for any $w_i$ and $b_i$ randomly chosen from any intervals of $\mathbf{R^n}$

252     and $\mathbf{R}$, respectively, according to any continuous probability distribution, then with

253     probability one, the hidden layer output matrix $H$ of the SLFN is invertible and

254     $||H\beta - Y|| = 0$.

255        Remark 2: Lemma 1 implies that the hidden units need to be linearly independent

256     with a probability of one to obtain the best least-squares solution of $Y = H\beta$.

257    Remark 3: (Goldberger 1964) According to ordinary least squares theory, the

258    smallest error $\left\|\widehat{Y} - Y\right\| = 0$ can be achieved by calculating $\boldsymbol{\beta} = (\boldsymbol{H}^T\boldsymbol{H})^{-1}\boldsymbol{H}^T\boldsymbol{Y}$ such

259    that there exists no multicollinearity (linearly dependence) among the hidden units.

260    *3.2 Input connection weights $\boldsymbol{W}$ determination*

261    Based on the above supporting statement, lemma, and remarks, the CPCLS can achieve

262    a best least-squares unique solution by the orthogonal transformation of the input and pre-

263    existing hidden units. CPCLS initializes by defining number $N$ of $\boldsymbol{H}$ in the first hidden

264    layer such that $p \leq n$. Initially, $\boldsymbol{X}$ is indirectly connected to $\boldsymbol{Y}$ through $\boldsymbol{H}$ to avoid input

265    feature linear dependence. For $\boldsymbol{W}$ determination, the eigendecomposition of $\boldsymbol{S}$ (1)

266    generates $\lambda$ (2) and the highest $\lambda$ values explaining maximum variance in data are used to

267    determine the eigenvectors (3). The determined eigenvectors are referred to as $\boldsymbol{W}$.

268    Knowing $\boldsymbol{X}$ and $\boldsymbol{W}$, the value of $\boldsymbol{H}$ is computed as:

$$\boldsymbol{H} = \emptyset(\boldsymbol{XW}) \tag{8}$$

269    where $\emptyset(z)$ can be any differentiable or nondifferentiable continuous activation function.

270    *3.3 Output connection weights $\boldsymbol{\beta}$ determination*

271    The second step is to compute the $\boldsymbol{\beta}$ by considering the linear relationship of $\boldsymbol{H}$ to $\boldsymbol{Y}$. The

272    Moore Penrose pseudo-inverse of $\boldsymbol{H}$ is calculated and its product with $\boldsymbol{Y}$ is used to

273    calculate $\boldsymbol{\beta}$ (6). The linear conversion of $\boldsymbol{H}$ through $\boldsymbol{\beta}$ generates $\widehat{Y}$ (7). The algorithm

274    aims to efficiently converge the network by minimizing the error function $E$ faster:

$$E = \frac{1}{m}\sum_{i=1}^{m}\left(\widehat{Y}_i - Y_i\right)^2 \tag{9}$$

275        If $E$ is a smaller amount than the described target error $e$, the CPCLS loop will

276    terminate, else a new $H$ will be generated until the required convergence is reached.

### 277 *3.4 Newly added hidden layer connection to the output layer*

278    In the proceeding hidden layers, the newly added $H_k$ *(k=1, 2, 3...)* receives all incoming

279    connections from $X$ and any preexisting hidden layers $H_{k-1}$, $H_{k-2}$ and so on, whereas,

280    the output layer receives connections from only the newly added $H_k$ and diminishes its

281    previous connections i.e. $H_{k-1}$, $H_{k-2}$ and so on. Connecting the previously added hidden

282    layers to the output units plays no significant role in the network. It may only add burden

283    to the network by connecting linearly dependent and redundant hidden units which can

284    reduce the generalization performance, as well as learning speed. Each newly added $H_k$

285    adds its non-linearity based on the variance in $X$ and previously added $H_{k-1}$, $H_{k-2}$ and

286    so on. This can be expressed in term of error minimization as:

$$E^{LHL} = (\beta H_k) - Y \tag{10}$$

287    where $E^{LHL}$ is network error by connecting only the newly added hidden layer to the

288    output layer. The newly added $H_k$ is of a higher level which has learned from the

289    orthogonal linear transformation of both $X$ and previously added $H_{k-1}$, $H_{k-2}$ and so on,

290    and represents the maximum variance of the network in that it guarantees the convergence

291    of the CPCLS.

292        Suppose symmetric matrix $S$ has two different eigenvalues $\lambda_1$ and $\lambda_2$

293    corresponding to eigenvectors $w_1$ and $w_2$ in matrix $W$ respectively. Two vectors can be

294    considered orthogonal if their inner product is zero, such as: $w_1 . w_2 = 0$ or $w_1^T w_2 = 0$.

295    where $w_1^T$ is the transpose of $w_1$.

296    We have:

$$S w_1 = \lambda_1 w_1 \tag{11}$$

297    and

$$Sw_2 = \lambda_2 w_2 \tag{12}$$

298    To prove that $w_1$ and $w_2$ are orthogonal:

299    $$\lambda_1(w_1.w_2) = (\lambda_1 w_1).w_2 = (Sw_1).w_2 = (Sw_1)^T w_2 = w_1^T S^T w_2$$

300    $$= w_1^T Sw_2 = w_1^T \lambda_2 w_2 = \lambda_2(w_1^T w_2) = \lambda_2(w_1.w_2)$$

301    $S = S^T$ because $S$ is a symmetric matrix. From mathematical work, we have:

$$\lambda_1(w_1.w_2) = \lambda_2(w_1.w_2) \tag{13}$$

$$(\lambda_1 - \lambda_2)(w_1.w_2) = 0 \tag{14}$$

302    Since $\lambda_1 - \lambda_2 \neq 0$, because both are different. So, we have:

$$w_1.w_2 = 0 \tag{15}$$

303    which means that eigenvectors $w_1$ and $w_2$ are orthogonal to each other in matrix $W$, i.e.,

304    $w_1 \perp w_2$. This orthogonal property of $W$ causes $X$ and preexisting $H_{k-1}$ to orthogonally

305    linearly transform into linearly independent $H_k$. Suppose if two hidden unit vectors are

306    generated in $H_k$ such that the $h_{k_1}$ is generated from $w_1$ and $h_{k_2}$ is generated from $w_2$,

307    then they can also be considered orthogonal, i.e., $h_{k_1} \perp h_{k_2}$. The proof supports Lemma

308    1 and guarantees the convergence of CPCLS because of the $H_k$ generated are invertible

309    and hence $\left\| (\beta H_k) - Y \right\| = 0$.

310         However, if all (every previous and newly) hidden layers are connected to the

311    output layer, we have:

$$E^{AHL} = \left( \beta(H_k + H_{k-1} + H_{k-2} + \cdots + H_1) \right) - Y \tag{16}$$

312    where $E^{AHL}$ is the network error by connecting all the hidden layers to the output layer.

313    According to Remarks 1 and Lemma 1, the hidden units in multiple hidden layers may

314    create linear dependency and redundancy in that it will avoid the best least square solution

315    assumption. Suppose if two hidden unit vectors are generated in $H_{k-1}$ such that $h_{k-1_1}$ is

316      generated from $w_{k-1_1}$ and $h_{k-1_2}$ is generated from $w_{k-1_2}$ and two hidden unit vectors

317      are generated in $\boldsymbol{H_k}$ such that $h_{k_1}$ is generated from $w_{k_1}$ and $h_{k_2}$ is generated from $w_{k_2}$

318      than there is a chance that it may or may not be orthogonal, i.e., $\boldsymbol{H_{k-1}} \perp \boldsymbol{H_k}$ or

319      $\boldsymbol{H_{k-1}} \not\perp \boldsymbol{H_k}$ . In the latter case, it may void the assumption that the $\boldsymbol{H}$ generated are

320      invertible and hence $\left\|\left(\boldsymbol{\beta}(\boldsymbol{H_k} + \boldsymbol{H_{k-1}} + \boldsymbol{H_{k-2}} + \cdots + \boldsymbol{H_1})\right) - \boldsymbol{Y}\right\| \neq \boldsymbol{0}$.

321          Hidden units are generated from the eigenvalue and corresponding eigenvector;

322      therefore, the new hidden units feature generation will always be less than or equal to the

323      input units and the previously hidden unit features $\boldsymbol{X} = (\boldsymbol{X}, \boldsymbol{H})$, such that $p \leq n$. Jolliffe

324      and Cadima (2016) stated that the eigenvalues having cumulative percentage variance

325      (CPV) of 70% are commonly used to extract eigenvectors. However, Jolliffe and Cadima

326      (2016) further added that there may circumstances in which the last few eigenvalues may

327      be also of interest in explaining more variance in the data. Researchers (Jolliffe and

328      Cadima 2016; Tortorella et al. 2016) in their work recommended selecting eigenvalues

329      giving a CPV greater than 70% to a maximum of 99%. The experimental work has been

330      performed to study the effect of hidden unit selection on generalization performance and

331      learning speed.

332      *3.5 CPCLS hyperparameters*

333      CPCLS initializes with a small number of hyperparameters i.e. $\boldsymbol{H}$ and e, in comparison

334      with other fixed and constructive topology algorithms i.e. learning rate, hidden nodes,

335      candidate units, etc. This makes learning simple.

336      **Algorithm CPCLS**

337      Given a training set $(\boldsymbol{X}, \boldsymbol{Y})$ with input unit matrix $\boldsymbol{X}$ be $m \times n$, output unit matrix $\boldsymbol{Y}$ be

338      $m \times q$, hidden unit matrix $\boldsymbol{H}$ be $m \times p$, and target error $e$:

339    Step 1) **Initialization:** Define the initial number $N$ of $\mathbf{H}$ in a first hidden layer such that

340    $p \leq n$

341    Step 2) **Learning Step:**

342    While $E > e$

343    a) Determine the $\mathbf{W}$ matrix of $n \times p$:

344        1. Calculate the $\mathbf{S}$ matrix of $n \times n$ from $n$ features $\mathbf{X}$:

345
$$\mathbf{S} = \frac{1}{m-1}(\mathbf{X} - \bar{\mathbf{x}})^T(\mathbf{X} - \bar{\mathbf{x}})$$

346
$$\bar{x} = \frac{1}{m}\sum_{i=1}^{m} x_i$$

347        2. Select $\lambda$ with the highest values to calculate the eigenvectors. The calculated $N$

348           eigenvectors are considered as $\mathbf{W}$ for $\mathbf{H}$:

349
$$|\mathbf{S} - \lambda\mathbf{I}| = 0$$

350
$$(\mathbf{S} - \lambda\mathbf{I})\mathbf{W} = 0$$

351    b) Take $\emptyset$ of $\mathbf{X}$ and $\mathbf{W}$ to generate $\mathbf{H}$:

352
$$\mathbf{H} = \emptyset(\mathbf{XW})$$

353    c) Determine the $\boldsymbol{\beta}$ matrix of $p \times q$:

354
$$\boldsymbol{\beta} = (\mathbf{H}^T\mathbf{H})^{-1}\mathbf{H}^T\mathbf{Y}$$

355    d) Calculate $\widehat{\mathbf{Y}}$:

356
$$\widehat{\mathbf{Y}} = \mathbf{H}\boldsymbol{\beta}$$

357    e) Calculate $E$:

358
$$E = \frac{1}{m} \sum_{i=1}^{m} \left( \hat{Y}_i - Y_i \right)^2$$

359    f) Combine the columns of $\boldsymbol{H}$ with $\boldsymbol{X}$:

360
$$\boldsymbol{X} = (\boldsymbol{X}, \boldsymbol{H})$$

361    g) increase the number of $\boldsymbol{H}$ by $N'$ in the proceeding hidden layers such that $p \leq n$:

362
$$N = N + N'$$

363    end

## 4. Experimental study

365    The comparative study of the proposed algorithm CPCLS with state-of-the-art machine

366    learning algorithms was conducted to demonstrate its effectiveness. The experimental

367    work was performed in Netmaker v0.9.5.2 and Anaconda Spyder Python v3.2.6. The

368    experimental work of CPCLS, BPNN, and self-adaptive extreme learning machine

369    (SaELM) (Wang et al. 2016) were performed in Python, whereas, the CasCor work was

370    performed in the built-in powerful Netmaker C-programming code. Generally speaking,

371    experimental work in the two programming codes will not affect the comparative study

372    because C programming is considered much faster than Python. The dataset was

373    normalized in the range [0,1] for both input and output and sigmoid activation function

374    $\emptyset(z) = 1/(1 + e^{-z})$ was used in the hidden units of the algorithms.

375        The experimental work was divided into three parts: real-world applications

376    prediction, energy applications prediction and studying the CPCLS hidden units and

377    layers characteristics followed by further discussion. Table 1 shows the most popular and

378    widely used dataset in machine learning extracted from UCI (Dua and Graff 2019). The

379    number of hidden units in hidden layers of the CPCLS was set to (2,2), (4,3), (2,1), (2,2),

380    (5,5) for real-world applications such as abalone, airfoil self-noise, forest fires, breast

381    cancer, wine respectively, and (2,2), (2,7), (1,1) for energy applications such as combined

382    cycle power plant, occupancy detection, seismic bumps respectively. The number of

383    CasCor candidate units was set to 3 Nos. The number of hidden units for stochastic

384    gradient descent BPNN was decided by a trial and error approach in the range 5-25 and

385    the hidden units with optimal results are reported. The minimum, maximum and interval

386    hidden units for SaELM was set to 5, 500 and 10 respectively with width factor $Q=2$ and

387    scale factor $L=4$.

388          Tables 2, 3 and 4 show the average best results of 25 trials obtained by the machine

389    learning algorithms. The testing RMSE/accuracy represents the generalization

390    performance, and the learning time represents the learning speed of the algorithms, while

391    the mean and stdev in the table refer to the average and standard deviation results of 25

392    trials. The performance criteria for regression problems and classification problems are

393    RMSE and percentage accuracy respectively.

394    *4.1 Real-world applications prediction*

395    Table 2 shows the prediction results of real-world applications. The proposed CPCLS

396    algorithm was able to achieve a better generalization performance and learning speed in

397    all cases as compared to CasCor, BPNN, and SaELM. The best results in terms of

398    generalization and learning speed are highlighted in bold and underlined in Table 2. For

399    an in-depth understanding of the convergence rate during each hidden layer, Figure 3

400    illustrates the CPCLS convergence rate of 25 trials for the Abalone dataset. It can be

401    observed that the convergence rate of CPCLS during each hidden layer addition is smooth

402    and stable.

403          CPCLS performance comparison has also been made with CasCor variants to

404    demonstrate its effectiveness. Due to the limitation caused by the unavailability of the

405  original programming code of OLSCN and FCNN, the simulation results of selected real-

406  world problems representing both algorithms are taken from their original source papers.

407  To make the comparison more valuable and to get better insights, the CPCLS simulation

408  is carried out by considering all test conditions mentioned in the original paper of OLSCN

409  and FCNN. Table 3 shows the dataset description, algorithms comparison in terms of

410  generalization performance and learning speed. It can be observed that CPCLS

411  generalization performance and learning speed averaged over 25 trials are better with

412  more improved results compared to FCNN and OLSCN.

413  *4.2 Energy applications prediction*

414  To further validate the performance, a comparative study was performed on energy-based

415  problems. The most demanding energy applications are:

416  (1) *Combined cycle power plant:* A combined cycle power plant is used to generate

417  electricity from gas turbines and consequently uses the waste energy in a steam

418  turbine to improve the efficiency of the electrical output. The attributes that

419  considerably affect the performance of gas turbine are atmospheric pressure

420  (millibar), temperature (°C) and relative humidity (%), whereas, the attributes that

421  affect the performance of the steam turbine are exhaust steam pressure (cm Hg).

422  The dataset contains an hourly average of attributes (atmospheric pressure,

423  temperature and relative humidity, exhaust steam pressure) to predict the net

424  hourly electrical energy (MW) of the powerplant.

425  (2) *Seismic bumps:* Seismic hazard prediction is a challenging application area in

426  coal mining. The purpose is to detect the possibilities of the occurrence of rock

427  bursts from seismic activity. The task is to classify high energy seismic bumps as

428  "hazardous" and "non-hazardous" from attributes such as possible seismic hazard,

429        seismic energy, pulses, energy deviation, number of seismic bumps with different

430        energy levels, total and maximum energy recorded for seismic bumps.

431    (3) *Occupancy detection:* Predicting occupancy detection in an office building is

432        attracting significant interest in reducing energy consumption. Various

433        measurements of light energy (Lux), temperature (°C), relative humidity (%),

434        humidity ratio (kgwater-vapor/kg/air), and $CO_2$ (ppm) along with the time are

435        used to classify whether the room is occupied or not.

436        Table 4 shows the performance of various machine learning algorithms for energy

437  application prediction. For the combined cycle power plant, CPCLS was able to achieve

438  a better performance of 0.0545 RMSE in a learning time of 2.96s compared to CasCor of

439  0.0573 in 29.69s, BPNN of 0.0577 in 59.54s, and SaELM of 0.0547 in 7.09s. For seismic

440  bumps, the generalization accuracy of CPCLS and BPNN is the same with the advantage

441  of CPCLS in that it took 0.01s compared to BPNN of 1.06s. The CPCLS demonstrated

442  its effectiveness by achieving a performance accuracy of 93.83% in a learning time of

443  0.01s compared to CasCor of 92.98% in 29.86s, BPNN of 93.83% in 1.06s, and SaELM

444  of 93.44% in 1.87s respectively. Similar to the combined cycle power plant and seismic

445  bumps, CPCLS also efficiently predicted occupancy detection. CPCLS achieved a better

446  performance accuracy of 99.05% in learning time of 3.95s compared to CasCor of 98.97%

447  in 31.54s, BPNN of 98.98% in 75.33s, and SaELM of 99.03% in 17.64s respectively. The

448  standard deviation of the generalization performance and learning time are also lower

449  which demonstrates the stable results of CPCLS.

*4.3 Connecting hidden layers to the output layer and varying hidden unit size in*

*the hidden layer of CPCLS*

*4.3.1 Varying hidden unit sizes in the hidden layers*

For CPCLS, the selection of hidden units in the first hidden layer and proceeding hidden layers is only a single hyperparameter that needs to be defined based on the eigenvalue and corresponding eigenvector. For illustration, experimental work has been performed by taking the example of the abalone dataset. The abalone dataset consists of 9 input attributes with bias. This implies that a lower and higher combination can be (1,1) and (9,9) respectively with a total of 81 combinations.

Figures 4, 5 and 6 show the generalization performance, learning speed and number of hidden layers for different combinations. The horizontal axis concerns the addition of hidden units in the first layer and the right legend concerns the addition of hidden units in the proceeding layers. Figure 4 illustrates that the generalization performance is stable for a maximum number of combinations. The minimum 0.0748RMSE and maximum 0.0774RMSE were achieved by (5,2) and (4,2) combinations respectively. Furthermore, a lower combination (1,1) achieved 0.0765RMSE and higher combination (9,9) achieved 0.0755RMSE. The (5,2), (4,2), (1,1) and (9,9) hidden units are generated from the eigenvalue CPV of (99.65%,96.94%), (99.19%,96.94%), (71.26%,72.92%) and (100%,100%) respectively. The minimum and maximum RMSE combination, and lower and higher hidden unit combinations give insight that hidden units generated based on eigenvalue explaining CPV $\lambda > 70\%$ are helpful in achieving better generalization performance. However, as shown in Figure 5, the learning time was 2.03s with (1,1) as compared to 0.03s for (9,9). The increase in learning time happens because of the higher computational burden by hidden layers. Figure 6 illustrates that hidden layers reach to 45 Nos. for lower combination (1,1)

475     compared to 4 Nos. for higher combination (9,9). The findings support the existing work

476     and recommend generating hidden units in the layers having eigenvalue explaining CPV

477     $\lambda > 70\%$. Based on our experimental work, it is recommended that the CPV should not

478     be greater than 99% because many of the last few eigenvalues may have approximately

479     zero variability. The zero variability eigenvalues may create a problem of overfitting

480     which needs to be avoided.


481     *4.3.2 The effect of hidden layers connection to the output layer*

482     Experimental work has been performed to study the effect of hidden layers connection to

483     output layer by considering both cases for CPCLS:

484       (1) Connecting the last hidden layer to the output layer (LHL)

485       (2) Connect all hidden layers to the output layer (AHL)

486       The work was performed on artificial nonlinear SinC function regression task,

487     generating 4,000 observations in the range [-20,20], by changing the data random state

488     from 0 to 100 with an interval of 5 and data test size from 30% to 70% with an interval

489     of 5%. This makes a total of 21 trials with different random states and 9 trials with

490     different test sizes. The 21 trials with different random states were performed by keeping

491     the constraint of test size equal to 50%. The best result by the random state was selected

492     to perform 9 trials by varying the test sizes.

493       Table 5 shows the generalization performance and learning speed of both cases.

494     Figures 7 and 8 illustrate the generalization performance and learning speed of both cases

495     for each random state and for each test size respectively. Both figures show that the

496     generalization performance becomes worse in most cases for AHL. Compared to AHL,

497     the generalization results of LHL are more stable with minimal deviation. Similarly, the

498     learning time increases for AHL compared to LHL. To avoid an increase in further

499     learning time, the algorithm for AHL needs to stop early when there is no further decrease

500     in error, and the training time is about five times more than LHL.

501        The difference in Figure 9 illustrates that AHL is unable to correctly predict the

502     SinC function, whereas LHL, (the original CPCLS), has predicted accurately all data

503     points of the function.

504     *4.4 Further discussion*

505     The better generalization performance and faster learning speed of CPCLS on real-world

506     and energy problems compared to CasCor, BPNN, SaELM, OLSCN, and FCNN

507     demonstrate its effectiveness. However, comparison with state-of-the-art machine

508     learning algorithms is important to build greater confidence in the application of CPCLS.

509     Table 6 shows the comparison of CPCLS with popular machine learning algorithms. The

510     comparative study gives an important insight that CPCLS generalization performance in

511     solving various real-world and energy problems is better compared to other machine

512     learning results, that are published recently in the literature. This finding supports that

513     CPCLS is a promising machine learning tool that can be practiced in general to improve

514     various operations of production research.

515        In real practice, the work is beneficial in numerous manners. Taking the example

516     of breast cancer, the CPCLS correctly classified its class as malignant or benign. It is

517     important to avoid misclassification of malignant cancer as benign because it can cause

518     human death. In engineering, the aviation sector works on zero-defect philosophy. Better

519     prediction of airfoils noise by CPCLS can facilitate in improving aircraft efficiency and

520     reduce environmental pollution. CPCLS efficiently prediction of marine species ages

521     rather than a microscope measurement can facilitate in avoiding subjective judgment and

522     fatigue. Besides, the application of CPCLS in predicting possible future hazards can help

523     to protect food products and the wastage of natural resources.

524        The better prediction results of CPCLS for energy applications such as predicting

525     electrical energy of powerplant and reducing energy consumption by accurately

526     predicting building occupancy detection can help in designing better energy management

527     systems. Moreover, predicting seismic hazards by CPCLS as hazardous and non-

528     hazardous can prevent fatal accidents.

529     **5. Conclusions**

530     In this paper, a novel learning algorithm called CPCLS is proposed. Unlike other cascade

531     algorithms, in this approach, hidden units are linearly generated by orthogonal linear

532     transformation and only the last hidden layer is connected to the output layer. It was

533     theoretically and experimentally verified that the hidden units generated in the respective

534     hidden layer are inevitable (i.e. linearly independent) which guarantees CPCLS

535     convergence. Connecting only the last hidden layer to the output layer eventually

536     improves the performance and increase the learning speed because all the hidden units

537     are orthogonal.

538        Compared to the state-of-the-art machine learning algorithms, the proposed

539     CPCLS achieved better generalization performance and learning speed in various

540     prediction tasks. Experimental work also demonstrated that connecting only the last

541     hidden layer rather than all the hidden layers to the output layer creates less burden on

542     the network and significantly improves convergence.

543        The major contributions and findings are: i) The CPCLS provides new insight into

544     existing algorithms by analytically calculating connection weights on both sides of the

545     network rather than gradient iteration or random generation, ii) In CPCLS, the generated

546     hidden units are inevitable ensuring that convergence will be optimal, iii) CPCLS

547     initialize with small number of hyperparameters, such as only defining number of hidden

548     units in the layer, iv) Compared to the existing works, this study provides insight that

549    avoiding direct linear connection of the input layer to the output layer and connecting

550    only newly added hidden layer to the output layer reduces network burden and improves

551    convergence, and v) In current practice, majority of research or models are proposed for

552    specific applications. The better performance of CPCLS, on various applications, in

553    comparison with state-of-the-art machine learning algorithms demonstrate that CPCLS

554    can be practiced in general for prediction of regression and classification tasks to make

555    better-informed decisions.

556         The implications are: i) In the proposed CPCLS, the experimental work was

557    performed on the OLT of the covariance matrix. Other than the covariance matrix, single

558    value decomposition and the correlation matrix can also be applied for OLT. Future work

559    may include studying the application of single value decomposition and correlation

560    matrix and their performance on the CPCLS, ii) Besides, the experimental work is limited

561    to the application of commonly used sigmoid activation function. Other than sigmoid

562    function, the effect of various other activation functions on the performance of CPCLS

563    needs to be explored in future work.

564    **References**

565    Aljarah, I., H. Faris, and S. Mirjalili. 2018. "Optimizing connection weights in neural

566        networks using the whale optimization algorithm." *Soft Computing* 22 (1): 1-15.

567        doi: 10.1007/s00500-016-2442-1.

568    Banerjee, P., V. S. Singh, K. Chatttopadhyay, P. C. Chandra, and B. Singh. 2011.

569        "Artificial neural network model as a potential alternative for groundwater

570        salinity forecasting." *Journal of Hydrology* 398 (3-4):212-20.

571    Bansal, P., S. Gupta, S. Kumar, S. Sharma, and S. Sharma. 2019. "MLP-LOA: A

572        metaheuristic approach to design an optimal multilayer perceptron." *Soft*

573        *Computing* 23 (23): 12331-45.

574    Candanedo, L. M., and V. Feldheim. 2016. "Accurate occupancy detection of an office
575        room from light, temperature, humidity and CO2 measurements using statistical
576        learning models." *Energy and Buildings* 112: 28-39.

577    Chien, C. F., Y. S. Lin, and S. K. Lin. 2020. "Deep Reinforcement Learning for
578        Selecting Demand Forecast Models to Empower Industry 3.5 and an Empirical
579        Study for a Semiconductor Component Distributor." *International Journal of*
580        *Production Research* (in press) doi: 10.1080/00207543.2020.1733125.

581    Chung, S. H., H. L. Ma, and H. K. Chan. 2017. "Cascading delay risk of airline
582        workforce deployments with crew pairing and schedule optimization." *Risk*
583        *Analysis* 37 (8):1443-58. doi: 10.1111/risa.12746.

584    Deng, C., J. Miao, Y. Ma, B. Wei, and Y. Feng. 2019. "Reliability analysis of chatter
585        stability for milling process system with uncertainties based on neural network
586        and fourth moment method." *International Journal of Production Research* (in
587        press). doi: 10.1080/00207543.2019.1636327.

588    Dua, D., and C. Graff. 2019. "UCI Machine Learning Repository." Accessed 2019-03-
589        28. http://archive.ics.uci.edu/ml.

590    Ertuğrul, Ö. F. 2018. "A novel type of activation function in artificial neural networks:
591        Trained activation function." *Neural Networks* 99: 148-57.

592    Fahlman, S. E., and C. Lebiere. 1990. The cascade-correlation learning architecture.
593        Paper presented at the Advances in neural information processing systems.

594    Fahlman, S. E. 1988. "An empirical study of learning speed in back-propagation
595        networks." In. Pittsburgh PA 15213: School of Computer Science, Carnegie
596        Mellon University.

597    Goldberger, A. S. 1964. "Classical linear regression." In *Econometric theory*, 156-212.
598        New York: John Wiley & Sons.

599    Grasso, F., A. Luchetta, and S. Manetti. 2018. "A multi-valued neuron based complex

600        ELM neural network." *Neural Processing Letters* 48 (1): 389-401.

601    Hecht-Nielsen, R. 1989. Theory of the backpropagation neural network. Paper

602        presented at the International Joint Conference on Neural Networks.

603    Heidari, A., V. G. Agelidis, J. Pou, J. Aghaei, and A. M. Y. M. Ghias. 2018. "Reliability

604        worth analysis of distribution systems using cascade correlation neural

605        networks." *IEEE Transactions on Power Systems* 33 (1):412-20.

606    Huang, G., S. Song, and C. Wu. 2012. "Orthogonal least squares algorithm for training

607        cascade neural networks." *IEEE Transactions on Circuits and Systems I:*

608        *Regular Papers* 59 (11):2629-37.

609    Huang, G. B., Q. Y. Zhu, and C. K. Siew. 2006. "Extreme learning machine: theory and

610        applications." *Neurocomputing* 70 (1-3):489-501.

611    Hunter, D., H. Yu, M. S. Pukish III, J. Kolbusz, and B. M. Wilamowski. 2012.

612        "Selection of proper neural network sizes and architectures—A comparative

613        study." *IEEE Transactions on Industrial Informatics* 8 (2):228-40.

614    Hwarng, H. B. 2005. "Simultaneous identification of mean shift and correlation change

615        in AR(1) processes." *International Journal of Production Research* 43

616        (9):1761-83. doi: 10.1080/00207540512331311822.

617    Jolliffe, I. T., and J. Cadima. 2016. "Principal component analysis: a review and recent

618        developments." *Philosophical Transactions of the Royal Society A:*

619        *Mathematical, Physical and Engineering Sciences* 374 (2065):20150202.

620    Jolliffe, I. T. 2002. *Principal component analysis*. New York: Springer-Verlag.

621    Kapanova, K. G., I. Dimov, and J. M. Sellier. 2018. "A genetic approach to automatic

622        neural network architecture optimization." *Neural Computing and Applications*

623        29 (5):1481-92.

624　Khan, W. A., S. H. Chung, M. U. Awan, and X. Wen. 2019a. "Machine learning

625　　　facilitated business intelligence (Part I): Neural networks learning algorithms

626　　　and applications." *Industrial Management & Data Systems* 120 (1):164-95.

627　Khan, W. A., S. H. Chung, M. U. Awan, and X. Wen. 2019b. "Machine learning

628　　　facilitated business intelligence (Part II): Neural networks optimization

629　　　techniques and applications." *Industrial Management & Data Systems* 120

630　　　(1):128-63.

631　Khan, W. A., S. H. Chung, and C. Y. Chan. 2018. Cascade Principal Component Least

632　　　Squares Neural Network Learning Algorithm. Paper presented at the 2018 24th

633　　　International Conference on Automation and Computing (ICAC).

634　Kim, B., Y. S. Jeong, S. H. Tong, and M. K. Jeong. 2019. "A generalised uncertain

635　　　decision tree for defect classification of multiple wafer maps." *International*

636　　　*Journal of Production Research* (in press). doi: 10.1080/00207543.2019.1637035.

637　Kumar, S., J. Singh, and O. Singh. 2020. "Ensemble-based extreme learning machine

638　　　model for occupancy detection with ambient attributes." *International Journal of*

639　　　*System Assurance Engineering and Management* (in press).

640　Kuo, Y. H., and A. Kusiak. 2019. "From data to big data in production research: The

641　　　past and future trends." *International Journal of Production Research* 57 (15-16):

642　　　4828-53. doi: 10.1080/00207543.2018.1443230.

643　Kusiak, A. 2020. "Convolutional and generative adversarial neural networks in

644　　　manufacturing." *International Journal of Production Research* 58 (5): 1594-1604.

645　　　doi: 10.1080/00207543.2019.1662133.

646　Kwok, T. Y., and D. Y. Yeung. 1997. "Constructive algorithms for structure learning in

647　　　feedforward neural networks for regression problems." *IEEE Transactions on*

648　　　*Neural Networks* 8 (3):630-45.

649    Liew, S. S., M. Khalil-Hani, and R. Bakhteri. 2016. "An optimized second order

650        stochastic learning algorithm for neural network training." *Neurocomputing*

651        186:74-89.

652    Liu, Q., H. Zhang, J. Leng, and X. Chen. 2019. "Digital Twin-Driven Rapid

653        Individualised Designing of Automated Flow-Shop Manufacturing System."

654        *International Journal of Production Research* 57 (12): 3903–19. doi:

655        10.1080/00207543.2018.1471243.

656    Liu, Y., L. Wang, X. V. Wang, X. Xu, and L. Zhang. 2019. "Scheduling in Cloud

657        Manufacturing: State-of-the-Art and Research Challenges." *International*

658        *Journal of Production Research* 57 (15–16): 4854–79. doi:

659        10.1080/00207543.2018.1449978.

660    Lorencin, I., N. Anđelić, V. Mrzljak, and Z. Car. 2019. "Genetic algorithm approach to

661        design of multi-layer perceptron for combined cycle power plant electrical

662        power output estimation." *Energies* 12 (22): 435201-26.

663    Lv, J., T. Peng, Y. Zhang, and Y. Wang. 2020. "A novel method to forecast energy

664        consumption of selective laser melting processes." *International Journal of*

665        *Production Research* (in press). doi: 10.1080/00207543.2020.1733126.

666    Mantas, C. J., J. G. Castellano, S. M. García, and J. Abellán. 2019. "A comparison of

667        random forest based algorithms: Random credal random forest versus oblique

668        random forest." *Soft Computing* 23 (21): 10739-54.

669    Nayyeri, M., H. S. Yazdi, A. Maskooki, and M. Rouhani. 2018. "Universal

670        approximation by using the correntropy objective function." *IEEE Transactions*

671        *on Neural Networks and Learning Systems* 29 (9): 4515-21.

672   Qiao, J., F. Li, H. Han, and W. Li. 2016. "Constructive algorithm for fully connected

673        cascade feedforward neural networks." *Neurocomputing* 182:154-64. doi:

674        10.1016/j.neucom.2015.12.003.

675   Rumelhart, D. E., G. E. Hinton, and R. J. Williams. 1986. "Learning representations by

676        back-propagating errors." *Nature* 323 (6088):533-6. doi: 10.1038/323533a0.

677   Solimanpur, M., P. Vrat, and R. Shankar. 2004. "Feasibility and Robustness of

678        Transiently Chaotic Neural Networks Applied to the Cell Formation Problem."

679        *International Journal of Production Research* 42 (6): 1065–82. doi:

680        10.1080/00207543.2004.10750072.

681   Tortorella, G. L., G. A. Marodin, D. D. C. Fettermann, and F. S. Fogliatto. 2016.

682        "Relationships between lean product development enablers and problems."

683        *International Journal of Production Research* 54 (10):2837-55.

684   Wang, G. G., M. Lu, Y. Q. Dong, and X. J. Zhao. 2016. "Self-adaptive extreme learning

685        machine." *Neural Computing and Applications* 27 (2):291-303.

686   Wang, H., S. Ding, D. Wu, Y. Zhang, and S. Yang. 2018. "Smart connected electronic

687        gastroscope system for gastric cancer screening using multi-column

688        convolutional neural networks." *International Journal of Production Research*

689        (in press). doi: 10.1080/00207543.2018.1464232.

690   Wang, Z., H. Ma, H. Chen, B. Yan, and X. Chu. 2019. "Performance degradation

691        assessment of rolling bearing based on convolutional neural network and deep

692        long-short term memory network." *International Journal of Production Research*

693        (in press). doi: 10.1080/00207543.2019.1636325.

694   Zou, W., F. Yao, B. Zhang, and Z. Guan. 2018. "Improved meta-ELM with error

695        feedback incremental ELM as hidden nodes." *Neural Computing and*

696        *Applications* 30 (11): 3363-70.