

1 **Framework for Integrating an Artificial Neural Network and a Genetic**
2 **Algorithm to Develop a Predictive Model for Construction Labor Productivity**

3 Sara EBRAHIMI¹, Mohammad RAOUFI², and Aminah Robinson FAYEK³

4 ¹ MSc Student and Graduate Research Assistant, Hole School of Construction
5 Engineering, Department of Civil and Environmental Engineering, University of
6 Alberta, 7-203 Donadeo Innovation Centre for Engineering, 9211 116 Street NW,
7 Edmonton, Alberta, T6G 1H9, Canada, email: eb4@ualberta.ca

8 ² Postdoctoral Fellow, Hole School of Construction Engineering, Department of Civil
9 and Environmental Engineering, University of Alberta, 7-381 Donadeo Innovation
10 Centre for Engineering, 9211 116 Street NW, Edmonton, Alberta, T6G 1H9, Canada,
11 email: mraoufi@ualberta.ca

12 ³ Professor, Director of the Construction Innovation Centre, Tier 1 Canada Research
13 Chair in Fuzzy Hybrid Decision Support Systems for Construction, NSERC Industrial
14 Research Chair in Strategic Construction Modeling and Delivery, Ledcor Professor in
15 Construction Engineering, Hole School of Construction Engineering, Department of
16 Civil and Environmental Engineering, University of Alberta, 7-232 Donadeo
17 Innovation Centre for Engineering, 9211 116 Street NW, Edmonton, Alberta, T6G
18 1H9, Canada, PH: (780) 492-1205, email: aminah.robinson@ualberta.ca

19 **ABSTRACT**

20 Construction labor productivity (CLP) is one of the most important factors in the
21 construction industry, as it has a direct effect on a company's efficiency and
22 profitability. The accurate prediction of CLP is essential for effective decision-making
23 prior to project execution, and continuous tracking and improvement of productivity
24 over a project life cycle is necessary for its success. The objective of this paper is to
25 develop a framework to help construction organizations predict and measure
26 construction productivity, leading to improved project performance in terms of cost,
27 time, and quality. CLP is affected by numerous factors, including the high-dimensional
28 factors that result from a large number of model input variables and which often impose
29 a high computational cost and the risk of overfitting of data. Therefore, it is necessary
30 to use feature selection methods to reduce the dimensionality of CLP data. This paper
31 proposes a framework that integrates an artificial neural network (ANN) and a genetic
32 algorithm (GA) for feature selection. The proposed framework is used to develop a
33 predictive model for CLP using features selected because they provide the best
34 prediction of CLP. The ability of GAs to generate an optimal feature subset in
35 combination with the superior accuracy of ANNs is a unique advancement that this
36 framework offers for improving the prediction of labor productivity. The developed
37 model can predict productivity and specify which factors are most predictive of CLP.
38 The contributions of this paper are (1) the development of a framework that uses an
39 integrated ANN and GA as a wrapper method for selecting the features with the most
40 influence on CLP and (2) the development of an improved predictive model that can
41 be used to both predict and measure CLP.

42 INTRODUCTION

43 Since many activities in the construction industry are labor dependent, improving
44 labor productivity is key for improving project performance. Construction labor
45 productivity (CLP) significantly impacts a company's profitability and project cost,
46 and construction organizations therefore need a predictive model of activity-level CLP
47 that helps them understand which factors affect labor productivity (Moselhi and Khan
48 2012). Numerous factors, both subjective (e.g., foreman skill and task complexity) and
49 objective (e.g., crew size), have been identified that affect CLP, causing complex
50 variability (Tsehayae and Fayek 2014; Raoufi and Fayek 2018; Hamza et al. 2019).
51 Therefore, providing a predictive model for CLP requires complex mapping of the
52 affecting factors (Heravi and Eslamdoost 2015). A large number of inputs and high-
53 dimensional data may present different problems, such as reduced accuracy and
54 increased complexity (Piao and Ryu 2017). To overcome these problems and find the
55 factors with the most influence on CLP, feature selection methods are used. In data
56 mining, feature selection is a necessary preprocessing approach for identifying a
57 relevant subset for classification. The aim of feature selection is to quickly develop
58 prediction models with better performance (Piao and Ryu 2017).

59 Past research that used both filter methods and wrapper methods indicates that the
60 wrapper method produces better results in feature selection (Alolfe et al. 2009).
61 However, current CLP studies are limited in their use of feature selection methods, as
62 they use only filter methods.

63 In this paper, a framework for feature selection that uses an ANN and a GA as a
64 wrapper method is developed and used to produce a predictive model for CLP. The
65 framework integrates the ANN and the GA for feature selection in order to find the
66 optimal feature subset by minimizing the fitness error of the ANN. Then, by employing
67 the selected features in the ANN as inputs with CLP as the output, a model for
68 predicting CLP is developed.

69 This approach predicts and measures CLP using the most influential factors, which
70 are selected by employing the neural-genetic algorithm as a combination of the ANN
71 and GA for selecting the best predictive CLP factors. Using both the abovementioned
72 feature selection algorithm and a database for CLP factors enables the identification of
73 the factors with the most influence on CLP so they can be taken into account on various
74 construction projects.

75 This paper is organized as follows: First, a review of past research on feature
76 selection methods, past methods used to select CLP factors, and the integration of an
77 ANN and a GA is presented. Then, a neural-genetic algorithm for selecting the best
78 predictive factors of CLP is described. Using the selected features, a predictive model
79 for labor productivity is developed. Next, a case study and the results of implementing
80 the algorithm are presented along with the predictive model. Finally, conclusions and
81 future research are presented.

82 LITERATURE REVIEW

83 Feature selection is the process of identifying and removing irrelevant and
84 redundant data (Hall 1999). It focuses on choosing a subset of the input features that
85 efficiently represents the input data while decreasing irrelevant features or noise effects
86 and providing a relatively accurate prediction of results. The main benefits of feature
87 selection are that (1) it decreases the amount of data needed to achieve learning, (2) it

88 enhances the predictive accuracy of models, (3) it reduces model execution time
89 because there are fewer inputs, and (4) it allows learned knowledge to be easily
90 understood because it is more compact (Hall 1999). Filter and wrapper methods are the
91 main approaches for feature selection (Yao et al. 2015). Filter methods are independent
92 of learning algorithms and choose best features based on some of the statistical
93 properties of data, such as their correlation coefficients. Because a small number of
94 features are used for classification in filter methods, the computation cost is low, which
95 is the main advantage of these methods. However, a small number of features, even if
96 they are the “best” ones, does not guarantee high classification accuracy (Cover 1974).
97 Furthermore, most filter methods are only suitable for developing mathematical
98 equations by the statistical regression method (Guyon et al. 2008; Gerami Seresht and
99 Fayek 2018; Raoufi and Fayek 2018a). Wrapper methods, on the other hand, use the
100 accuracy of a learning algorithm as a criterion for selecting useful features (Yao et al.
101 2015), and they explore the feature space to score feature subsets according to their
102 predictive power. Wrapper methods are therefore a more effective means of
103 constructing a predictive model than filter methods because they are tuned to the
104 specific interaction between a learning algorithm and its training data (Ahmad et al.
105 2015; Aličković and Subasi 2017). However, their application is limited because of the
106 high computational complexity that occurs when numerous feature sets are considered
107 (Piao and Ryu. 2017).

108 In the construction discipline, there are some studies that use filter feature selection
109 for reducing the number of factors influencing CLP. Tsehayae and Fayek (2014)
110 identified a total of 169 parameters influencing CLP in building and industrial projects.
111 As the 169 input parameters and seven process variables result in a high-dimension
112 feature space, Tsehayae and Fayek (2016) used a correlation-based feature selection
113 (CFS) algorithm, which is a filter method, and proposed a predictive model that uses a
114 fuzzy inference system. Although the CFS algorithm is suitable in that it has the ability
115 to deal with a high dimension of input space and a small number of data instances,
116 using a wrapper method is more appropriate for predictive modeling using artificial
117 intelligence (AI) techniques, such as fuzzy inference systems and ANNs, because of its
118 superior performance (Piao and Ryu 2017). Therefore, this research proposes the
119 integration of an ANN and a GA as a wrapper method for CLP feature selection.

120 The integration of the ANN with the GA has not been investigated in CLP studies
121 in the construction domain. However, in other domains, such as medical research, the
122 integration of ANNs with GAs for feature selection has been studied. For example, the
123 technique of combining ANN parameters that are simultaneously optimized by the GA
124 was proposed by Verma and Zhang (2007) and Ahmad et al. (2015) to implement
125 feature selection for diagnosing breast cancer. The classification rates achieved in both
126 studies were promising and showed better results than most previous studies that used
127 filter methods.

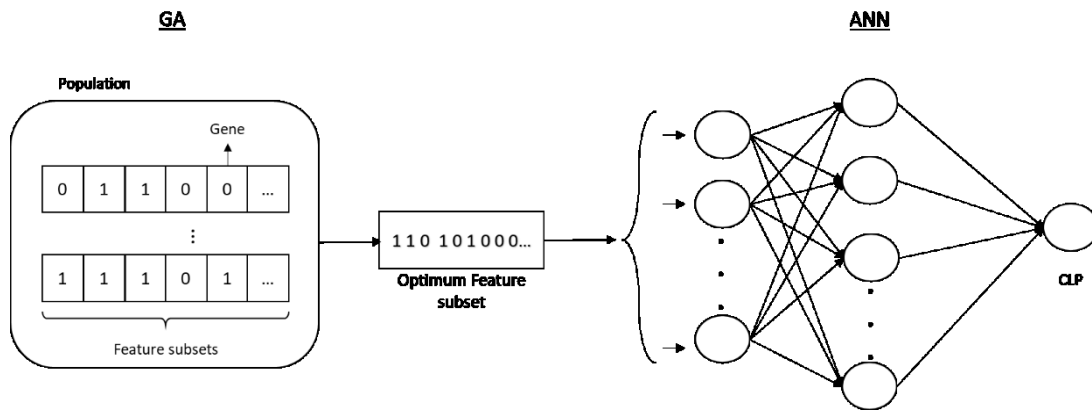
128 GAs, inspired by the natural selection process, are strong evolutionary optimization
129 algorithms that search for the best subset of system parameters for the development of
130 an accurate predictive model. GAs have been applied successfully for feature selection
131 by various researchers (Gerami Seresht and Fayek 2018). The integration of a GA and
132 an ANN (e.g., in Verma and Zhang 2007) has recently attracted wide attention. The

133 objective of this paper is to present a neural-genetic algorithm for finding the most
 134 influential features and developing a predictive model for CLP.

135 **METHODOLOGY TO DEVELOP A NEURAL-GENETIC ALGORITHM**
 136 **FOR FEATURE SELECTION**

137 This paper presents a neural-genetic algorithm for finding the most influential
 138 factors from a number of existing features that affect CLP in order to develop a model
 139 for predicting CLP.

140 Figure 1 shows the conceptual integration of a GA and an ANN, where each
 141 individual in the population indicates a candidate solution for selecting the feature
 142 subset. If there are n features affecting CLP, there are 2^n possible feature subsets. Each
 143 feature subset is called a “chromosome” and contains n genes, which can have one of
 144 two values. A value of 1 indicates that the corresponding feature has been chosen for
 145 predicting CLP, and a value of 0 means that the feature has not been selected.



146
 147

Figure 1. Conceptual integration of GA and ANN.

148 An overview of the proposed framework is shown in Figure 2, which presents the
 149 process of integrating the GA and the ANN for feature selection. There five steps to
 150 performing this integration, described in detail in the following paragraphs.

151 As shown in Figure 2, in the first step, the algorithm generates the random initial
 152 population of chromosomes. Each individual in the population represents an available
 153 solution to the feature subset selection problem.

154 In the second step, the selected features are the inputs of the ANN. In the neural
 155 network, the number of hidden layer nodes is calculated using Equation (1), as the
 156 appropriate hidden layer size is calculated based on the number of inputs and outputs
 157 (Heaton 2008).

158
$$\text{Number of nodes} = 2 * \sqrt{\text{inputs} + \text{outputs}} \quad (1)$$

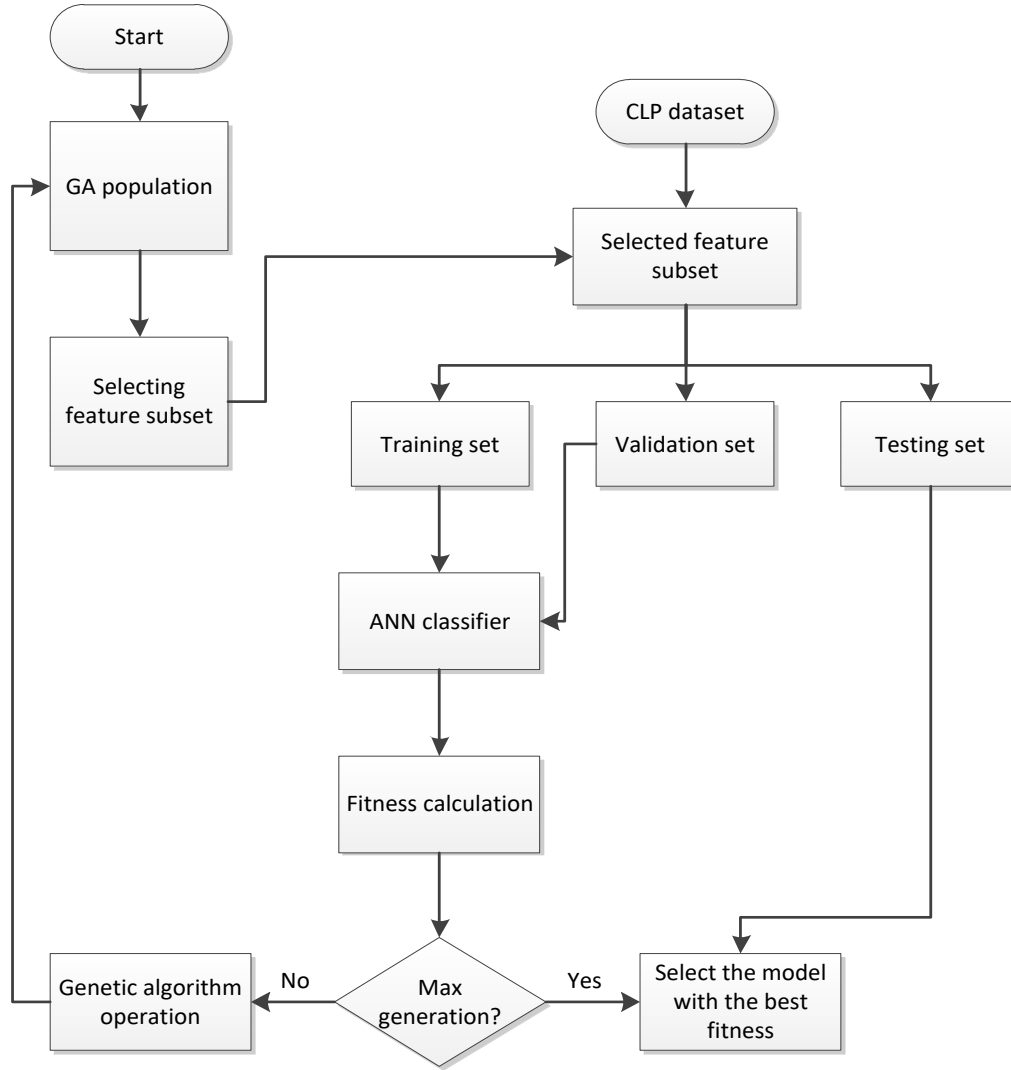
159 In the third step, the training set, validation set, and testing set are built from the
 160 CLP dataset collected by Tsehayae and Fayek (2014). Then, using the training set, the
 161 process of training the network is started. To avoid overtraining in ANN, the error of
 162 the validation set is considered during this process. If the error grows for five iterations
 163 consecutively, then the training stops.

164 The fourth step is the fitness calculation process, wherein the validation set is used
 165 to simulate the network and calculate the error by using the root mean square error
 166 (RMSE) based on Equation (2).

167

$$Error = \sqrt{\frac{1}{n} \sum_{i=1}^p (A_i - T_i)^2} \quad (2)$$

168 where p , A_i and T_i are the number of output nodes, the actual output value of the i th
 169 output node, and the target output value of the i th output node, respectively. In this
 170 paper, there is one output node, which is CLP. A better fitness of the ANN requires a
 171 smaller error.



172
 173

Figure 2. Overview of the CLP feature selection algorithm.

174

The fifth step is a GA operation, which consists of the following process:

175

- (a) Selection: A roulette wheel selection strategy is used to choose the individual probabilistically to form a parent whose number is equal to the population size minus the elitism number. If f_i is the fitness of individual i in the population, the probability of being selected is given by Equation (3), where N is the population size.

176

177

178

179

180

$$P_i = i \frac{f_i}{\sum_{j=1}^N f_j} \quad (3)$$

- 181 (b) Crossover and mutation: Crossover is the process used to specify which two
182 chromosomes will create a new offspring chromosome (Bean 1994). The
183 mutation operation changes one or more genes in a chromosome from its
184 initial state. All the chromosomes after the crossover operation will go
185 through a mutation operation and a new offspring is produced. In this
186 research, single point binary crossover and binary mutation are performed.
- 187 (c) Elitism: Another process in the GA algorithm is elitism. Elitism involves
188 copying a small proportion of the fittest candidates, unchanged, into the next
189 generation. In this paper, the three best chromosomes are selected to be part
190 of the population in the next generation.
- 191 (d) Fitness function: The GA optimization method minimizes the value of a
192 fitness function, which is shown as Z and calculated for each chromosome.
193 The Z function is defined by Equation (4), where b is a coefficient of the
194 number of selected features (nf). In this study, we consider “ b ” to be equal
195 to 0.008 in order to have fewer selected features.

196
$$Z = Error * (1 + b * nf) \quad (4)$$

197 After these processes, once the final generation with the best fitness value is
198 reached, the iteration stops, and the feature subset that is chosen as a final solution is
199 the one that is the best predictor of CLP among all feature subsets. Accordingly, a
200 predictive model for CLP can be developed, which has the minimum fitness error.

201 CASE STUDY

202 To illustrate the proposed method of feature selection and construct the predictive
203 model for CLP, a case study was conducted. Tsehayae and Fayek (2014) identified a
204 total of 96 activity-level sub-parameters. In this case study, all 20 activity-level factors,
205 identified by Tsehayae and Fayek (2014), that showed non-zero variance in data were
206 considered inputs to the feature selection algorithm (Table 1), and CLP was the output
207 of the predictive model. A total of 92 data instances were used. The aim of feature
208 selection in this case study is to identify the most influential features among the 20
209 activity-level factors to be able to quickly develop a predictive model of CLP.

210 The proposed algorithm was developed in the MATLAB 9.6 environment. The
211 backpropagation technique, which is one type of ANN learning algorithm, was used
212 due to its fast execution and simple implementation in MATLAB. The output and the
213 hidden nodes' activation functions were pure linear and hyperbolic tangent,
214 respectively. The output layer consisted of one output node, which is CLP. Table 2
215 shows the GA parameter settings, which were based on Zhuo et al. (2008). In order to
216 make data consistent across all tables, feature normalization was required. This
217 approach uses normalized data, which are real numbers in the range 0–1. In this paper,
218 70% of the CLP dataset was used for training, the next 15% was used for validation,
219 and the last 15% was used for testing. The validation set was used to calculate the
220 overall fitness of the network and choose the best network, and the testing set was used
221 to achieve the desired test accuracy of the selected neural network. The early stopping
222 method, which defines the maximum number of iterations before overfitting begins,
223 was used to avoid overfitting of the network (Ebrahimi Kahou et al. 2015; Gal and
224 Ghahramani 2016).

Table 1. CLP factors.

No	Factor
f1	Crew size
f2	Craftsperson education
f3	Craftsperson technical training
f4	Crew composition
f5	Crew experience (seniority)
f6	Number of languages spoken
f7	Cooperation among craftspersons
f8	Treatment of craftspersons by foreman
f9	Craftsperson motivation
f10	Craftsperson fatigue
f11	Team spirit of crew
f12	Fairness of work assignments
f13	Crew participation in foreman decision-making process
f14	Crew flexibility
f15	Availability of task materials
f16	Quality of task materials
f17	Sharing of tools
f18	Working condition (noise)
f19	Location of work scope (distance)
f20	Fairness in performance review of crew by foreman

Table 2. GA parameter settings.

Parameter	
Population size	50
Crossover probability	0.8
Elitism size	3
Mutation probability	0.2
Maximum generation	40

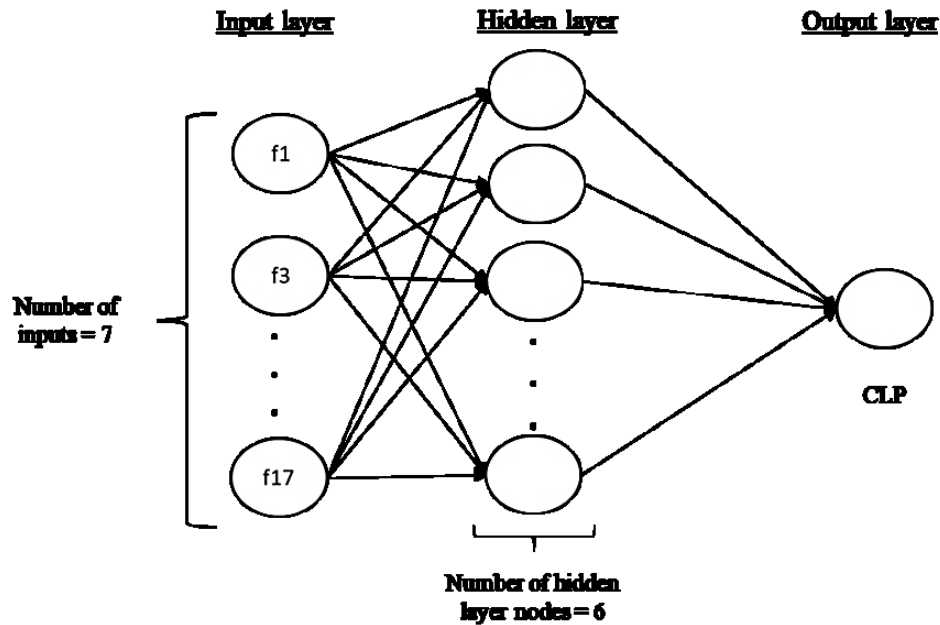
227 **RESULTS**

228 As a result of implementing the neural-genetic algorithm, a total of seven features
 229 were selected as the factors with the most influence on CLP, and they were used to
 230 develop a predictive model. These seven features were crew size (f1), craftsperson
 231 technical training (f3), cooperation among craftspersons (f7), craftsperson motivation
 232 (f9), craftsperson fatigue (f10), fairness of work assignments (f12) and sharing of tools
 233 (f17). Table 3 shows the results of implementing the neural-genetic algorithm. As
 234 shown in Table 3, the final error of fitness calculation of the ANN was 0.0107, which
 235 was calculated using Equation (2). Based on Equation (4), Z represents the minimum
 236 amount of fitness function. As the number of selected features (i.e., inputs) was seven
 237 and the number of outputs was one, the number of hidden layer nodes in the selected
 238 model was calculated to be six using Equation (1).

Table 3. Results of conducting neural-genetic algorithm.

Results	
Selected features	f1, f3, f7, f9, f10, f12, f17
Error	0.0107
Z	0.0111
Hidden layer size	6

240 These selected features were the inputs of the achieved predictive model with the
 241 best fitness, which is shown in Figure 3. The limitation of the wrapper method in terms
 242 of computational complexity in dealing with numerous feature sets does not occur in
 243 this case because the utilized CLP dataset consists of 20 features.



244
 245

Figure 3. The predictive model for CLP.

246 One of the selected features was cooperation among craftspersons (f7), which is
 247 supported by Tsehayae and Fayek (2014), whose work identified “good cooperation
 248 between craftsmen in a crew” as a top parameter having a positive effect on CLP. In
 249 addition, Jergeas (2009) found “labor relations” to be a target for CLP improvement.
 250 Tsehayae (2015) identified 27 identical input variables for CLP in a total four contexts
 251 (i.e., industrial, warehouse, high-rise, and institutional building) by using a CFS
 252 algorithm including crew size, craftsperson on-job training, craftsperson motivation,
 253 craftsperson fatigue, and fairness of work assignments, which are found in our results
 254 (f1, f3, f9, f10 and f12). Therefore, the selected features in this paper are in agreement
 255 with the results of past studies on productivity feature selection. However, the
 256 comparison of the results of this study with past literature indicates that the proposed
 257 framework can better identify predictive features of CLP. Tsehayae and Fayek (2016)
 258 obtained 2.515% as the RMSE value, while in this paper RMSE value is 1.070%.
 259 Future research will focus on collecting a larger data set from different organizations
 260 and project contexts to expand the scope of applicability of the developed algorithm.

261 CONCLUSIONS AND FUTURE RESEARCH

262 The main aim of this paper was to develop a predictive model for measuring and
263 predicting CLP using a GA algorithm and an ANN. After performing a literature review
264 to investigate past research on feature selection methods and the importance of CLP, a
265 methodology for constructing a predictive model for CLP was developed. This paper
266 illustrated the neural-genetic framework for both feature selection and presenting a
267 predictive model. By implementing the developed framework on a real case, the
268 features that were the best predictive factors of CLP were identified, and the results
269 were compared to past research and found to be consistent with previous results. The
270 achieved error in this paper indicates an improvement of the predictive model in
271 comparison to past studies. The contributions of this paper are (1) the development of
272 a framework that uses an ANN and a GA as a wrapper method for feature selection to
273 select the parameters with the most influence on CLP and (2) the development of an
274 improved model for predicting and measuring CLP. The results of this work will
275 improve the prediction of CLP. Better identification of the factors with the most
276 influence on CLP can lead to more effective management of CLP and project
277 performance. The findings of this paper also provide a basis for future research work,
278 including modeling CLP based on feature selection on all identified influencing factors
279 with actual data. Modeling multifactor construction productivity, which includes labor,
280 material, and equipment, can be done in future works by using the proposed framework.
281 Future research can also focus on implementing other AI techniques, such as the
282 combination of a GA and a neuro-fuzzy system, to focus on subjective factors affecting
283 CLP and comparing the accuracy of predictive models with different AI techniques.

284 REFERENCES

- 285 Ahmad, F., Isa, N. A. M., Hussain, Z., Osman, M. K., and Sulaiman, S. N. (2015). "A
286 GA-based feature selection and parameter optimization of an ANN in
287 diagnosing breast cancer." *Pattern Analysis and Applications*, 18(4), 861-870.
- 288 Aličković, E., and Subasi, A. (2017). "Breast cancer diagnosis using GA feature
289 selection and Rotation Forest." *Neural Computing and Applications*, 28(4),
290 753-763.
- 291 Alolfe, M. A., Mohamed, W. A., Youssef, A. M., Kadah, Y. M., and Mohamed, A. S.
292 (2009). "Feature selection in computer aided diagnostic system for
293 microcalcification detection in digital mammograms." *Proc., 2009 National
294 Radio Science Conference*, IEEE, 1-9.
- 295 Bean, J. C. (1994). "Genetic algorithms and random keys for sequencing and
296 optimization." *ORSA Journal on Computing*, 6(2), 154-160.
- 297 Cover, T. M. (1974). "The best two independent measurements are not the two best."
298 *IEEE Trans. Syst. Man Cybern.*, SMC-4(1), 116-117.
- 299 Ebrahimi Kahou, S., Michalski, V., Konda, K., Memisevic, R., and Pal, C. (2015).
300 "Recurrent neural networks for emotion recognition in video." In *Proc., the
301 2015 ACM on International Conference on Multimodal Interaction*, ACM,
302 Seattle, Washington, USA, 467-474.
- 303 Gal, Y., and Ghahramani, Z. (2016). "A theoretically grounded application of dropout
304 in recurrent neural networks." In *Proc., 30th Annual Conference on Neural
305 Information Processing Systems 2016 (NIPS 2016)*, Barcelona, Spain, 1019-
306 1027.

307 Gerami Seresht, N., and Fayek, A. R. (2018). "Dynamic modeling of multifactor
308 construction productivity for equipment-intensive activities." *J. Constr. Eng.*
309 *Manage.*, 144(9), 04018091.

310 Guyon, I., Gunn, S., Nikravesh, M., and Zadeh, L. A. (2008). "Feature extraction:
311 Foundations and applications." *Springer*.

312 Hall, M. A. (1999). "Correlation-based feature selection for machine learning."

313 Hamza, M., Shahid, S., Bin Hainin, M. R., and Nashwan, M. S. (2019). "Construction
314 labour productivity: Review of factors identified." *Int. J. Constr. Manage.*, 1-
315 13.

316 Heaton, J. (2008). "*Introduction to neural networks with Java*." Heaton Research, Inc.,
317 St. Louis, USA.

318 Heravi, G., and Eslamdoost, E. (2015). "Applying artificial neural networks for
319 measuring and predicting construction-labor productivity." *J. Constr. Eng.*
320 *Manage.*, 141(10), 04015032.

321 Jergeas, G. (2009). "Improving construction productivity on Alberta oil and gas capital
322 projects." A Report Submitted to: Alberta Finance and Enterprise.

323 Moselhi, O., and Khan, Z. (2012). "Significance ranking of parameters impacting
324 construction labour productivity." *Construction Innovation*, 12(3), 272-296.

325 Piao, Y., and Ryu, K. H. (2017). "A hybrid feature selection method based on
326 symmetrical uncertainty and support vector machine for high-dimensional data
327 classification." *Proc., Asian Conference on Intelligent Information and*
328 *Database Systems*, Springer, 721-727.

329 Raoufi, M., and Fayek, A. R. (2018a). "Key moderators of the relationship between
330 construction crew motivation and performance." *J. Constr. Eng. Manage.*,
331 144(6), 04018047.

332 Raoufi, M., and Fayek, A. R. (2018b). "Framework for identification of factors
333 affecting construction crew motivation and performance." *J. Constr. Eng.*
334 *Manage.*, 144(9), 04018080.

335 Tsehayae, A. A. (2015). "Developing and optimizing context-specific and universal
336 construction labour productivity models."

337 Tsehayae, A. A., and Fayek, A. R. (2014). "Identification and comparative analysis of
338 key parameters influencing construction labour productivity in building and
339 industrial projects." *Can. J. Civ. Eng.*, 41(10), 878-891.

340 Tsehayae, A. A., and Fayek, A. R. (2016). "Developing and optimizing context-specific
341 fuzzy inference system-based construction labor productivity models." *J.*
342 *Constr. Eng. Manage.*, 142(7), 04016017.

343 Verma, B., and Zhang, P. (2007). "A novel neural-genetic algorithm to find the most
344 significant combination of features in digital mammograms." *Applied Soft*
345 *Computing*, 7(2), 612-625.

346 Yao, J., Mao, Q., Goodison, S., Mai, V., and Sun, Y. (2015). "Feature selection for
347 unsupervised learning through local learning." *Pattern Recog. Lett.*, 53 100-
348 107.

349 Zhuo, L., Zheng, J., Li, X., Wang, F., Ai, B., and Qian, J. (2008). "A genetic algorithm
350 based wrapper feature selection method for classification of hyperspectral
351 images using support vector machine." *Proc., Geoinformatics 2008 and Joint*
352 *Conference on GIS and Built Environment*, 71471J.