

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31

**Multicollinearity Effect in Regression Analysis: A Feed Forward Artificial Neural Network**

**Approach**

**ABSTRACT**

In this study we compared the performance of Ordinary Least Squares Regression (OLSR) and the Artificial Neural Network (ANN) in the presence of multicollinearity using two datasets – a real life insurance data and a simulated data – to know which of the methods, models a highly correlated dataset better using the Root Mean Square Error (RMSE) as the performance measure. The ANN performed better than the OLSR model for all the different ANN models except the models with nine and ten nodes in the hidden layer for the real life data. The network with four hidden nodes was the best model. For the simulated data, the ANN model with two hidden nodes gave us the least RMSE when compared to the OLSR model and the other ANN models in the testing set. The network with two hidden nodes modelled the data very well. In the presence of multicollinearity, ANN model achieves a better fit and forecast than the OLSR.

**Keywords:** Multicollinearity; Ordinary Least Squares; Artificial Neural Network; Root Mean Square Error.

**1.0 INTRODUCTION**

In modelling a linear relationship between a dependent variable and one or more independent variables, OLSR is being used to estimate the parameters of the model by minimizing the Residual Sum of Squares. The OLSR gives an unbiased estimate of the regression coefficients, it is very easy to compute and interpret. Though OLSR is preferred, it can only yield the best results when some assumptions are satisfied; There must be a linear relationship between each of the independent variables and the dependent variable, the independent variables must not be highly correlated, the variance of the error must be constant, the errors must not be correlated and there must not be an outlier. Most real life data does not always satisfy all the assumptions of the OLSR and if we insist on using the OLSR method to estimate the parameters, we will not be able to achieve a better fit for the data and a good prediction with the model. Wonsuk *et al.* (2014), Onoghojobi *et al.* (2016) and Olewuezi *et al.* (2016) studied the nature of

32 multicollinearity, their consequences, how to detect them and some remedial measures that can  
33 be taken to get a good estimate of the regression coefficients.

34 In this study, we considered a solution to the OLSR method when the multicollinearity  
35 assumption is not satisfied. In the presence of multicollinearity, it is impossible to estimate the  
36 unique effects of individual variables in the regression equation, the variance and covariance of  
37 the Least Squares (LS) estimates become too large though still the Best Linear Unbiased  
38 Estimator (BLUE), most of the regression coefficients are not significant and there is a high  $R^2$   
39 value even though the t values for most of the regression coefficients are small. Multicollinearity  
40 becomes one of the serious problems in linear regression analysis, Yazid and Mowafaq (2009).  
41 Many attempts have been made to improve the OLSR estimation procedure, some of which are  
42 Ridge Regression Onoghojobi *et al.* (2016), Latent Root Regression, Partial Least Squares  
43 Olewuezi *et al.* (2016), Principal Component Regression Onoghojobi *et al.* (2016) etc. and more  
44 recently, machine learning method which have smaller Mean Square Error (MSE) than the  
45 OLSR method Zou *et al.* (2008).

46 Artificial Neural Network (ANN) is an example of a machine learning method that evolved from  
47 the idea of simulating the human brain Zou *et al.* (2008). They are networks of simple processing  
48 elements called neurons or nodes. The ANN models complex nonlinear relationships between  
49 the predictor variables and the response with great flexibility by defining input neurons – nodes –  
50 which are the predictor variables, a hidden layer with a number of nodes connected to each of the  
51 input nodes and lastly, an output layer with one or more nodes. The theoretical advantage of  
52 ANNs is that the relationship between the variables need not to be specified in advance since the  
53 method establishes the relationship through a learning process. The model learns the relationship  
54 from the data used to train it. The ANNs do not also require any assumptions about the

55 underlying population distribution. Hsiao-Tien (2008) and Kumar (2005) compared the  
56 performance of ANN and OLSR. Hsiao-Tien (2008) compared OLSR and ANN models with  
57 seven explanatory variables of corporation's feature and three external macro-economic control  
58 variables to analyse the important determinants of capital structures of the high-tech and  
59 traditional industries respectively in Taiwan. He used the RMSE as the criterion to know the best  
60 model. The ANN model achieved a better fit and forecast than the OLSR model as it had the  
61 least RMSE. Ramirez *et al.* (2000) also compared ANN and OLSR model. They found out that  
62 ANN method performed better than the OLSR, although both methods showed good  
63 performance for daily rainfall. This paper compares the performance of ANN and OLSR method  
64 to model a highly correlated real life Nigerian Insurance Company's data and a simulated data.

## 65 **2.0 METHODOLOGY**

66 The OLSR and ANN were used to model the two datasets to know which of the methods, models  
67 a highly correlated dataset better using the RMSE as the performance measure. The model with  
68 the least RMSE is chosen as the best model. Correlation coefficient is used to test for  
69 multicollinearity in the two datasets. There is high multicollinearity in the data if the correlation  
70 coefficient is high (i.e. greater than 0.8 or less than -0.8).

### 71 **2.1 ORDINARY LEAST SQUARES METHOD.**

72 Let us consider the standard model for Multiple Regression Analysis

$$73 \quad Y = X\beta + \varepsilon \quad (1)$$

74 where

75  $Y$  is  $(n \times 1)$  vector of the dependent variable.

76  $X$  is  $(n \times p)$  matrix of independent variables.

77  $\beta$  is  $(p \times 1)$  vector of regression parameters.

78  $\varepsilon$  is  $(n \times 1)$  vector of errors.

79 From equation (1), we have

$$80 \quad \varepsilon^T \varepsilon = (Y - X\beta)^T (Y - X\beta) \quad (2)$$

81 This term is differentiated with respect to  $\beta$  and set equal to 0 to obtain an estimate of  $\beta$   
82 provided the inverse of  $X^T X$  exists and is unique. We therefore have:

$$83 \quad \hat{\beta}_{OLS} = (X^T X)^{-1} X^T Y \quad (3)$$

84 where  $\hat{\beta}_{OLS}$  is  $p \times 1$  vector of OLSR estimated parameters.

## 85 2.2 ARTIFICIAL NEURAL NETWORK (ANN)

86 The ANN models complex nonlinear relationships between the predictor variables and the  
87 response with great flexibility by defining input neurons – nodes – which are the predictor  
88 variables, a hidden layer with a number of nodes connected to each of the input nodes and lastly,  
89 an output layer with one or more nodes. An activation function is applied to both the hidden and  
90 the output layers. The connections between the nodes (input nodes and the hidden layer nodes)  
91 are assigned weights. The weights are the parameters the Neural Network estimates, and they are  
92 chosen so as to minimize a pre-defined loss function. Neural Network tries to minimize the  
93 difference between the observed responses and the output. Figure 1 is an example of an Artificial  
94 Neural Network. Three layers of nodes are defined – an input layer that comprises of three input  
95 nodes and a bias node, a single hidden layer and an output layer.

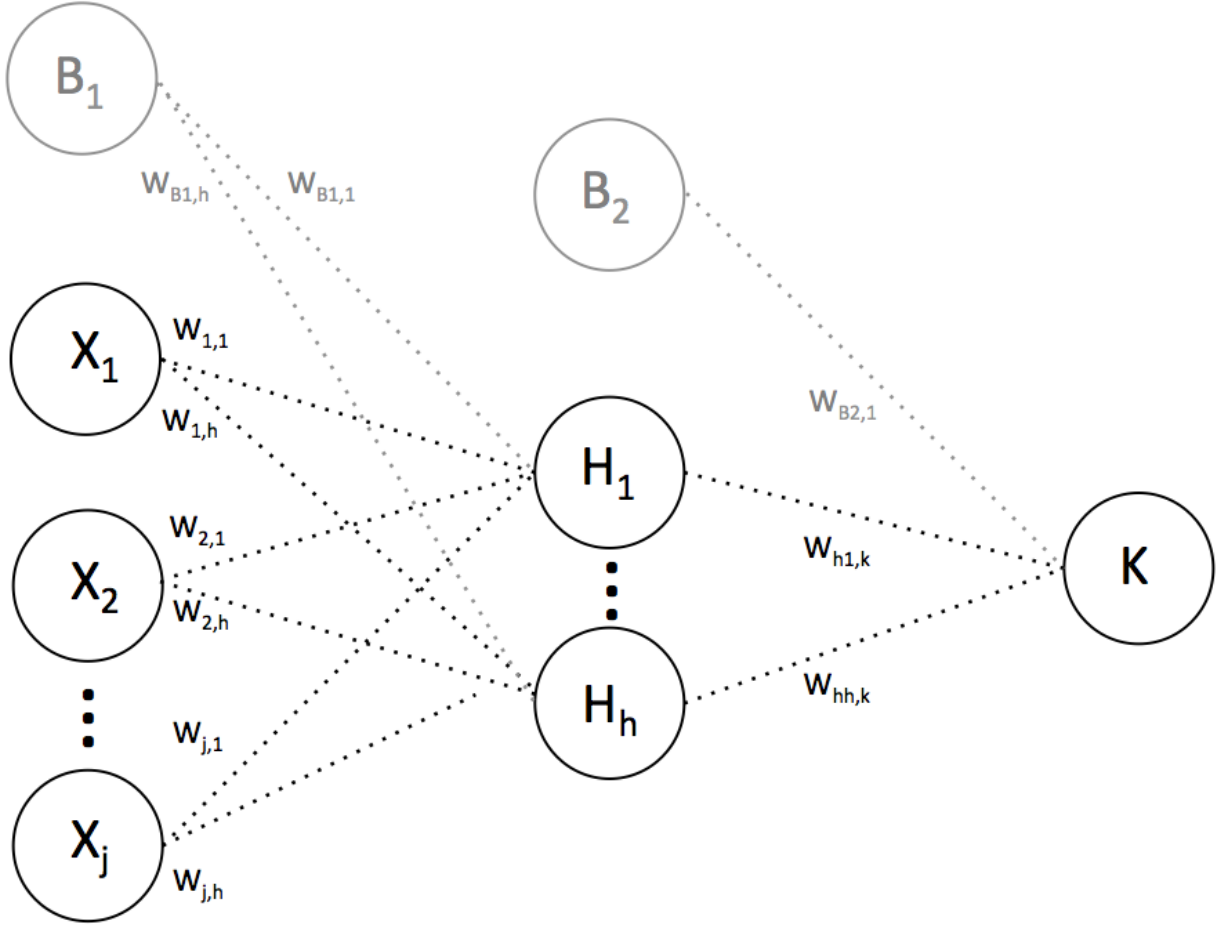
96 Let X represent the inputs,

97 H, K and B represent the hidden, output and bias nodes respectively, and

98 W represent the weights.

99 The weights in the bias nodes can be interpreted similarly to an intercept in a linear regression.

100 **Figure 1.** Single hidden layer feed forward neural network.



101

102 A Feed Forward Neural Network (FFNN) is a uni-directional connection from the input to the  
 103 hidden layer and from the hidden to the output layer. A mathematical representation of a single  
 104 layer FFNN is given in equation (4).

105 
$$\hat{y}_k(x_i, w) = \Phi_0\left(\alpha_k + \sum_{h=1}^H w_{hk} \Phi_h\left(\alpha_h + \sum_{j=1}^J w_{jh} x_{ij}\right)\right) \quad (4)$$

106 That is, the sum of the product of the weights  $w_{jh}$  and the inputs  $x_{ij}$  plus a bias term  $\alpha_h$  gives us  
 107 a node in the hidden layer. An activation function is applied to each node. An activation function  
 108 also called a threshold or transfer function is a non-linear transformation applied to the input.  
 109 The sum is taken over the hidden neurons H of the product of the transformed input and the  
 110 weights  $w_{hk}$  plus a bias term  $\alpha_k$ . A final transformation  $\Phi_0$  is applied to the output. We have  
 111 different activation function for both the transmission from the input units to the hidden units and  
 112 from the hidden units to the output units, namely: linear activation function, unit step activation

113 function, rectified linear unit activation function, hyperbolic tangent activation function, sigmoid  
114 activation function, logistic activation function, etc.

115 The error function is minimized to get an estimate of the weight  $w$  for both the input and the  
116 hidden layer. The commonly used error function has been the quadratic error function while  
117 cross-entropy error function is more suitable for binary classification. The Quadratic error  
118 function  $E_Q$  and cross entropy error function  $E_C$  are given below

$$119 \quad E_Q = \sum_{k=1}^K \sum_{i=1}^n (\hat{y}_k(x_i, w) - y_{ik})^2 \quad (5)$$

$$120 \quad E_C = - \sum_{k=1}^K \sum_{i=1}^n y_{ik} \log \hat{y}_k(x_i, w) + (1 - y_{ik}) \log[1 - \hat{y}_k(x_i, w)] \quad (6)$$

121 The more the nodes in the hidden layer, the more complicated non-linear relationship can be  
122 modelled. Increasing the nodes in the hidden layer also increases the likelihood of training an  
123 over fitted model. A model is over fitted when it does not generalize well to new observations  
124 though it will still perform very well on the training set.

125 The dataset is divided into two - the training set and the testing set. It should be noted that 70%  
126 of the data were used as the training set and the other 30% as the testing set. The training set is  
127 used to train the network and the optimally performing hyper parameters are identified. The final  
128 model performance is then tested using the testing set.

### 129 **3.0 EMPIRICAL ILLUSTRATION**

#### 130 **3.1 Illustration 1**

131 The real life data used in this study is a secondary data on Nigerian Insurance Expenditure from  
132 1996 to 2011. Claims, fire, accident, motor, employers, marine, miscellaneous were used as the  
133 independent variables with total expenditure as the dependent variable. That is,

$X_1 = \text{Claims}$   
 $X_2 = \text{Fire}$   
 $X_3 = \text{Accident}$   
 $X_4 = \text{Motor}$   
 $X_5 = \text{Employers}$   
 $X_6 = \text{Marine}$   
 $X_7 = \text{Miscellaneous}$   
 $Y = \text{Total Expenditure}$

134

135 **3.1.1 Test for Multicollinearity**

136 The correlation matrix was used to test for multicollinearity. Table 1 below is the correlation  
 137 matrix for the data.

138 **Table 1:** Correlation Matrix for the Real Life Data

1	0.800667	0.972993	0.9843	0.931225	0.954349	0.818016099
0.800667	1	0.664078	0.802241	0.559947	0.623471	0.330163173
0.972993	0.664078	1	0.956963	0.96467	0.970487	0.888117153
0.9843	0.802241	0.956963	1	0.900635	0.917368	0.763616208
0.931225	0.559947	0.96467	0.900635	1	0.991298	0.948715023
0.954349	0.623471	0.970487	0.917368	0.991298	1	0.930325184
0.818016	0.330163	0.888117	0.763616	0.948715	0.930325	1

139

140 From Table 1, there is high multicollinearity in the data since most of the independent variables  
 141 are highly correlated.

142 **3.1.2 Ordinary Least Squares Regression for the Real Life Data**

143 **Table 2:** Overall Fit for the Life Data

Multiple R	0.996305
R Square	0.992623
Adjusted R Square	0.986169
Standard Error	2917.984
Observations	16

**Table 3:** ANOVA Table for the Real Life Data

	<i>df</i>	<i>SS</i>	<i>MS</i>	<i>F</i>	<i>p-value</i>	<i>Sig</i>
Regression	7	9.17E+09	1.31E+09	153.789	6.84E-08	Yes
Residual	8	68117033	8514629			
Total	15	9.23E+09				

**Table 4:** OLSR Parameter Estimates for the Real Life Data

	<i>Coeff</i>	<i>std err</i>	<i>t stat</i>	<i>p-value</i>	<i>lower</i>	<i>Upper</i>
Intercept	1346.411	2425.234	0.555167	0.593959	-4246.19	6939.011915
Claims	2.867821	2.354802	1.217861	0.257975	-2.56236	8.298002766
Fire	-3.78929	3.977479	-0.95269	0.368643	-12.9614	5.382793653
Accident	-2.21072	2.9512	-0.74909	0.475249	-9.0162	4.594762006
Motor	0.778697	2.750741	0.283086	0.784298	-5.56452	7.121915705
Employers	-17.4773	57.36528	-0.30467	0.768395	-149.762	114.8072724
Marine	-1.36654	3.986157	-0.34282	0.740566	-10.5586	7.825552293
Miscellaneous	-3.19616	3.414654	-0.93601	0.376657	-11.0704	4.67804833

144

145 Although the R-square value (0.9926) for the model is high, all of the regression coefficients are  
 146 not significant since their p values > 0.05 at 0.05 level of significance and their confidence  
 147 intervals are large. This contradiction is as a result of the assumption of multicollinearity not  
 148 being satisfied.

### 149 3.1.3 Artificial Neural Network for the Real Life Data

150 Ten ANN models with different number of nodes in the hidden layer were trained. We used 1 to  
 151 10 nodes in the hidden layers to know which one of them will yield the best estimate of the  
 152 parameters of the network using the RMSE as the performance measure. The logistic activation  
 153 function was used for the transmission from input units to hidden units and the linear activation  
 154 function was used for the transmission from hidden units to output units. The quadratic error  
 155 function was used to determine the weights of the network. Table 5 below gives us the summary  
 156 of the result.

**Table 5:** RMSE statistics for the Real Life Data

	OLSR	ANN1	ANN2	ANN3	ANN4	ANN5	ANN6	ANN7	ANN8	ANN9	ANN10
<b>RMSE (testing)</b>	3188.1	1875.1	2522.9	2105.9	1350.7	2030.2	1526.8	1608.8	2633.7	4512.5	3387.2



<b>RMSE (training)</b>	2786.6	1982.0	2347.9	1781.8	1493.6	2376.6	1663.1	953.1	1916.6	1259.6	1470.3
----------------------------	--------	--------	--------	--------	--------	--------	--------	-------	--------	--------	--------

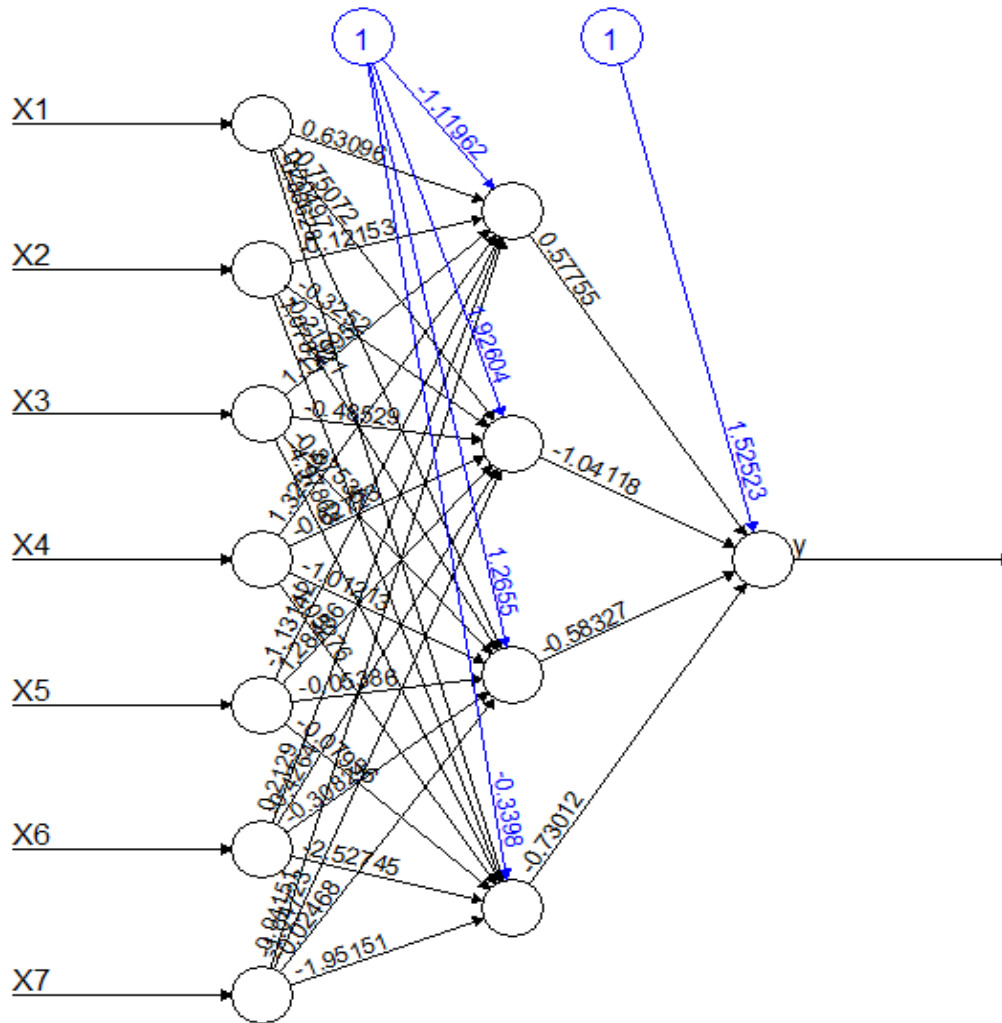
158

159 From Table 5, the ANN models had a lesser RMSE than the OLSR model for all the different  
160 models except the models with nine and ten hidden nodes, the ANN models with nine and ten  
161 hidden nodes over fitted the training set. It modelled well the training set but could not predict  
162 well the testing set. It was observed that ANN performed better than the OLSR model for all the  
163 different ANN models except the models with nine and ten hidden nodes. The network with four  
164 hidden nodes modelled the data very well than the other ANN models. It did not over fit the  
165 training set and also predicts well the testing set. It has the least RMSE when used on the testing  
166 set. Below is the estimate of the parameters of the ANN model with four hidden nodes and the  
167 graphical representation of the model.

168 Intercept.to.1layhid1 -1.119624897  
169 x1.to.1layhid1 0.630956846  
170 x2.to.1layhid1 -0.121533513  
171 x3.to.1layhid1 1.126946335  
172 x4.to.1layhid1 1.325447728  
173 x5.to.1layhid1 -1.131424271  
174 x6.to.1layhid1 0.212902953  
175 x7.to.1layhid1 -0.041508678  
176 Intercept.to.1layhid2 1.926037650  
177 x1.to.1layhid2 -0.750716229  
178 x2.to.1layhid2 -0.325200164  
179 x3.to.1layhid2 -0.485294045  
180 x4.to.1layhid2 -0.627234654  
181 x5.to.1layhid2 1.284860662  
182 x6.to.1layhid2 -0.426399785  
183 x7.to.1layhid2 -1.047225798  
184 Intercept.to.1layhid3 1.265501183  
185 x1.to.1layhid3 0.204966188  
186 x2.to.1layhid3 -0.219213203  
187 x3.to.1layhid3 -0.875358244  
188 x4.to.1layhid3 -1.012126040  
189 x5.to.1layhid3 -0.053855979  
190 x6.to.1layhid3 -0.308241782  
191 x7.to.1layhid3 -0.024684527  
192 Intercept.to.1layhid4 -0.339802202  
193 x1.to.1layhid4 -0.886280005  
194 x2.to.1layhid4 1.078712791  
195 x3.to.1layhid4 -4.978679098  
196 x4.to.1layhid4 -1.052764500  
197 x5.to.1layhid4 -0.079364780  
198 x6.to.1layhid4 -2.527453252  
199 x7.to.1layhid4 -1.951514060  
200 Intercept.to.y 1.525232613  
201 1layhid1.to.y 0.577550927  
202 1layhid2.to.y -1.041176296  
203 1layhid3.to.y -0.583267342

205 1layhid4.to.y -0.730115784

206 **Figure 2.** Single Hidden Layer Feed Forward Neural Network for the Real Life Data.



Error: 0.009966 Steps: 90

207

### 208 3.2 Illustration 2

209 We simulated a correlated data with six independent variables and one dependent variable  
210 replicated 150 times.

#### 211 3.2.1 Test for Multicollinearity

212 The correlation matrix was used to test for multicollinearity. Below is the correlation matrix for  
213 the simulated data.

214

215 **Table 6:** Correlation Matrix for the Simulated Data

1	0.98	0.93	0.89	0.87	0.83
0.98	1	0.95	0.90	0.91	0.77
0.93	0.95	1	0.96	0.84	0.71
0.89	0.90	0.96	1	0.80	0.69
0.87	0.91	0.84	0.80	1	0.67
0.83	0.77	0.71	0.69	0.67	1

216

217 From the correlation matrix, high multicollinearity was observed in the data since most of the  
 218 independent variables were highly correlated.

219 **3.2.2 Ordinary Least Squares Regression for the Simulated Data**

220 **Table 7:** Overall Fit for the Simulated Data

Multiple R	0.693028
R Square	0.480288
Adjusted R Square	0.458482
Standard Error	9.589799
Observations	150

**Table 8:** ANOVA Table for the Simulated Data

				Alpha	0.05	
	<i>Df</i>	<i>SS</i>	<i>MS</i>	<i>F</i>	<i>p-value</i>	<i>sig</i>
Regression	6	12153.3	2025.549	22.0254	3.01E-18	yes
Residual	143	13150.89	91.96425			
Total	149	25304.18				

**Table 9:** OLSR Parameter Estimates for the Simulated Data

	<i>Coeff</i>	<i>std err</i>	<i>t stat</i>	<i>p-value</i>	<i>lower</i>	<i>upper</i>
Intercept	-92.7866	159.8667	-0.5804	0.562559	-408.794	223.2207
X <sub>1</sub>	-1.4703	5.114539	-0.28748	0.774165	-11.5802	8.639566
X <sub>2</sub>	5.355521	6.111732	0.876269	0.382353	-6.72549	17.43653
X <sub>3</sub>	-1.18115	4.103448	-0.28784	0.773883	-9.29241	6.930103
X <sub>4</sub>	2.382363	2.891732	0.823853	0.411395	-3.3337	8.098428
X <sub>5</sub>	4.736788	2.019845	2.345124	0.020396	0.744176	8.7294

221  $X_6$                       -0.656   1.542828   -0.4252   0.671333   -3.7057   2.393691

222 From Table 9, all of the regression coefficients except  $X_5$  are not significant since their p values  
 223  $> 0.05$  at 0.05 level of significance and the confidence intervals are large. This again, is as a  
 224 result of the assumption of multicollinearity not being satisfied.

225 **3.2.3 Artificial Neural Network for the Simulated Data**

226 Table 10 below gives us the summary of the result.

227 **Table 10:** RMSE statistics for the Simulated Data

	OLSR	ANN1	ANN2	ANN3	ANN4	ANN5	ANN6	ANN7	ANN8	ANN9	ANN10
<b>RMSE (testing)</b>	13.5	13.9	13.1	31.5	16.5	18.9	20.5	45.0	19.1	20.3	31.8
<b>RMSE (training)</b>	13.1	12.9	11.4	9.5	9.9	8.6	8.2	7.8	6.3	4.1	3.6

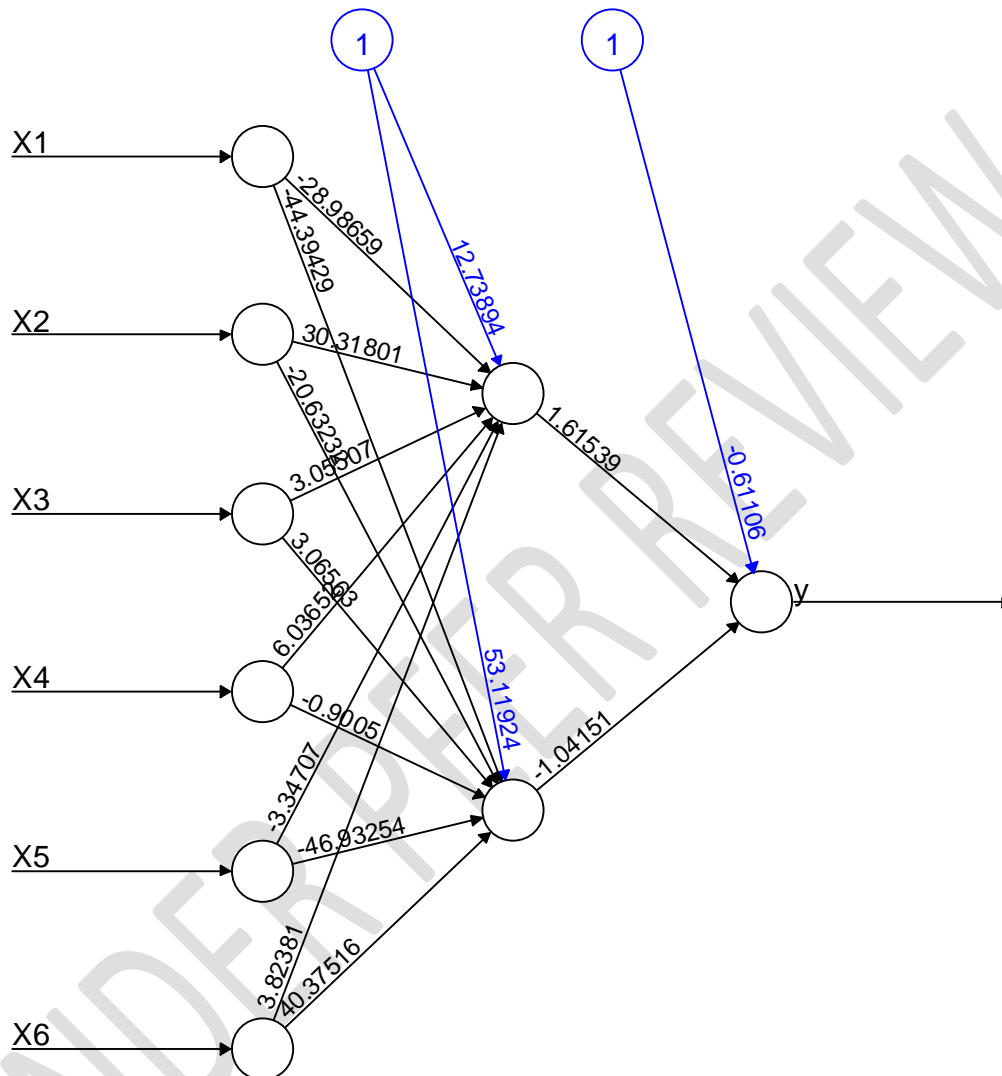
229 The ANN model with two hidden nodes gave us the least RMSE when compared to the OLSR  
 230 model and the other ANN modelled with one, three, four, five, six, seven, eight, nine and ten  
 231 hidden nodes in the testing set. The network with two hidden nodes modelled the data very well.  
 232 It did not over fit the training set and also predicted well the testing set. Below is the estimate of  
 233 the parameters of the ANN model with two hidden nodes and the graphical representation of the  
 234 model.

235 Intercept.to.1layhid1   1.273894e+01  
 236 x1.to.1layhid1        -2.898659e+01  
 237 x2.to.1layhid1        3.031801e+01  
 238 x3.to.1layhid1        3.055071e+00  
 239 x4.to.1layhid1        6.036519e+00  
 240 x5.to.1layhid1        -3.347070e+00  
 241 x6.to.1layhid1        3.823806e+00  
 242 Intercept.to.1layhid2   5.311924e+01  
 243 x1.to.1layhid2        -4.439429e+01  
 244 x2.to.1layhid2        -2.063232e+01  
 245 x3.to.1layhid2        3.065635e+00  
 246 x4.to.1layhid2        -9.005025e-01  
 247 x5.to.1layhid2        -4.693254e+01  
 248 x6.to.1layhid2        4.037516e+01  
 249 Intercept.to.y        -6.110575e-01  
 250 1layhid1.to.y        1.615394e+00  
 251 1layhid2.to.y        -1.041508e+00  
 252

253

254

255 **Figure 3.** Single Hidden Layer Feed Forward Neural Network for the Simulated Data.



Error: 20.051264 Steps: 13557

256

#### 257 **4.0 CONCLUSION**

258 Correlation coefficient was used to test for multicollinearity in the two data set and both the real  
259 life and simulated data failed to satisfy the multicollinearity assumption. The ANN models had a  
260 lesser RMSE than the OLSR model for all the different models except the models with nine and

261 ten nodes in the hidden layer for the real life data, the ANN models with nine and ten hidden  
262 nodes over fitted the training set. The network with four hidden nodes had the least RMSE when  
263 used on the testing set. It did not over fit the training set and also predicted well the testing set.

264 For the simulated data, the ANN model with two hidden nodes gave us the least RMSE when  
265 compared to the OLSR model and the other ANN models with one, three, four, five, six, seven,  
266 eight, nine and ten hidden nodes in the testing set. The network with two hidden nodes modelled  
267 the data very well. It did not over fit the training set and also predicted well the testing set.

268 When there is multicollinearity, it is advisable to use the ANN to model the data since unlike the  
269 OLSR method, it has no assumption that must be satisfied and it achieves a better fit and forecast  
270 than the OLSR in the presence of multicollinearity as seen from this study using a real life and a  
271 simulated data.

## 272 REFERENCES

273 Hsiao-Tien Pao (2008). A Comparison of Neural Network and Multiple Regression Analysis in  
274 Modelling Capital Structure. *Expert Systems with Applications*, 35, 720-727.

275 Kumar, U. A. (2005). Comparison of Neural Networks and Regression Analysis: A New Insight.  
276 *Expert Systems with Applications*, 29(2), 424–430.

277 Olewuezi, N.P., Onoghojobi, B., & Bartholomew, D.C. (2016). “Estimation of Nonorthogonal  
278 Problem Using Time Series Dataset”. *Journal of the Nigerian Association of*  
279 *Mathematical Physics*, Vol. 34, (March, 2016), pp: 141 – 150.

280 Onoghojobi, B., Olewuezi, N. P. & Obite, C. P. (2016). “Comparison of Linear Prediction  
281 Methods in Time Series”. *Journal of the Nigerian Association of Mathematical Physics*,  
282 vol. (34), pp 151-156.

283 Ramírez, M. C. V., de Campos Velho, H. F., and Ferreira, N. J. (2005). Artificial neural network  
284 technique for rainfall forecasting applied to the São Paulo region, *J. Hydrol.*, 301, 146–  
285 162.

286 Wonsuk Y., Robert M., Sejong B., Karan S., Qinghua P. & James W. L (2014). “A Study of  
287 Effects of Multicollinearity in the Multivariable Analysis. *International Journal of*  
288 *Applied Science and Technology*, 4(5), 9-19.

- 289 Yazid, M. A. and Mowafaq, M. A. (2009). “A Monte Carlo comparison between ridge and  
290 principal components regression methods”, *Applied Mathematical Sciences*, Vol. 3, Vol.  
291 42, pp. 2085 – 2098.
- 292 Zou J., Han Y., and So S.S (2008). Overview of Artificial Neural Networks, *Methods Mol Biol*,  
293 458, pp. 15-23.

UNDER PEER REVIEW