# SIMPLE REGRESSION MODELS: A COMPARISON USING CRITERIA MEASURES

## Osuagwu Chidimma Udo[1] and Okenwe Idochi[2]

[1]Department of Statistics, Federal University of Technology, Owerri, Imo State, Nigeria.

[2]Department of Statistics, School of Applied Sciences, Ken Saro Wiwa Polytechnic, Bori, Rivers State, Nigeria.

**ABSTRACT:** *The study is on simple regression models: a comparison using criteria measures. The source of the dataset used for this study was extracted from records of the Federal Medical Centre, Owerri, Imo State, on weight of babies and hemoglobin level of mothers. The response variable is weight of babies while the explanatory variable is hemoglobin level of mothers. Eleven simple regression models—Linear, Growth, Quadratic, Polynomial, Logarithmic, Hyperbolic, Power, Exponential Growth, Square Root, Sinusoidal and Arctangent— were stated and employed for the study. For ease of data analysis, E-views package was implemented. Three model selection criteria measures for comparison, known as Akaike Information Criterion (AIC), Schwarz Information Criterion (SIC) and Hannan-Quinn Information Criterion (HQIC), were employed. The result of the study showed that, when it comes to analyzing the association between baby weight and mothers' hemoglobin levels, the exponential growth regression model performs better than the other ten models that were examined. Therefore, researchers should investigate other models that were not included in this analysis and compare the findings using goodness of fit metrics other than the criteria measures used in this work.*

**KEYWORDS:** Simple Nonlinear Regression, Simple Linear Regression, AIC, SIC, HQIC, Model Comparison.

## INTRODUCTION

While fitting a simple linear model to data is rare because most data follow nonlinear models, simple regression model fitting is typically used in many scientific domains, including pharmaceutical and biochemical test quantification (Duong & Lim, 2023). There are nonlinear models and choosing the best model for the data requires a combination of expertise, understanding of the underlying mechanism, and statistical analysis of the fitting result (Esemokumo et al., 2020). Quantifying the validity of a fit using a metric that distinguishes between "good" and "bad" fits is crucial. When performing calibration experiments for samples to be measured, many researchers typically use a common measure known as the coefficient of determination ($R^2$) employed in linear regression (Montgomery et al., 2006).

Because values between 0 and 1 make it simple to grasp how much of the variation in the data is explained by the fit, this measure is therefore particularly intuitive from a linear perspective (Chicco et al., 2021). Many scientists and academics continue to utilize $R^2$ in studies pertaining to nonlinear data processing, despite the fact that it has been proven for some time to be an inappropriate metric for nonlinear regression (Berk, 2020). This problem had been highlighted by a number of earlier descriptions of $R^2$ being useless in nonlinear fitting, but they have presumably now been forgotten (Bartlett et al., 2020). This observation may be the result of the disparities in mathematical training between researchers and trained statisticians, who frequently use statistical techniques but lack in-depth statistical understanding (Spiess & Neumeyer, 2010).

$R^2$ is not the best option in a nonlinear regime because, unlike in linear regression, the total sum-of-squares (TSS) is not equal to the regression sum-of-squares (REGSS) plus the residual sum-of-squares (RSS), and as a result, it lacks the appropriate interpretation. It has been stated that researchers arbitrarily use $R^2$ to evaluate the validity of a specific model when dealing with nonlinear data fit. One possible explanation for the prevalence of relying just on $R^2$ values to assess the validity of nonlinear models is that researchers may not be aware of this common misunderstanding.

This study only employed three criteria models known as the Akaike Information Criterion, Schwarz Information Criterion, and Hannan-Quinn Information Criterion for model selection, correct interpretation, and conclusion because using $R^2$ alone to assess the performance of nonlinear data analysis has been discouraged.

In terms of medicine, it has been demonstrated that a patient's weight and pulse rate have a linear relationship. But many researchers, particularly those in other fields where they most likely lack enough statistical skills, typically used the linear regression technique to find a relationship between these two variables without considering the nonlinear models. Because of this, the goal of this study is to compare several non-linear models with linear models in order to determine which model best fits the patient's weight and pulse rate based on the data collected for this investigation.

## METHODOLOGY

### Regression Models

Eleven Regression models were considered in this study, which are Linear, Growth, Quadratic, Polynomial, Logarithmic, Hyperbolic, Power, Exponential Growth, Square Root, Sinusoidal and Arctangent Regression models as written in Equations (1), (2), (3), (4), (5), (6), (7), (8), (9), (10) and (11) respectively:

$$Y = \lambda_0 + \lambda_1 Z + \varepsilon \tag{1}$$

$$Y = \frac{\lambda_0 Z}{\lambda_1 + Z} + \varepsilon \tag{2}$$

$$Y = \lambda_0 + \lambda_1 Z + \lambda_2 Z^2 + \varepsilon \tag{3}$$

$$Y = \lambda_0 + \lambda_1 Z + \lambda_2 Z^2 + \lambda_3 Z^3 + \varepsilon \tag{4}$$

$$Y = \lambda_0 + \lambda_1 \ln(Z) + \varepsilon \tag{5}$$

$$Y = \lambda_0 + \lambda_1 (1/Z) + \varepsilon \tag{6}$$

$$Y = \lambda_0 Z^{\lambda_1} + \varepsilon \tag{7}$$

$$Y = \lambda_0 + \exp(\lambda_1 Z) + \varepsilon \tag{8}$$

$$Y = \lambda_0 + \lambda_1 \sqrt{z} + \varepsilon \tag{9}$$

$$Y = \lambda_0 + \lambda_1 Sin(Z) + \varepsilon \tag{10}$$

$$Y = \lambda_0 + \lambda_1 \arctan(\lambda_2 Z + \lambda_3) + \varepsilon \tag{11}$$

### Simple Linear Regression

This is a regression line involving only two variables as it is applicable in this study. A widely used procedure for obtaining the regression line of Y and Z is the least square method.

The linear regression of Y on Z is stated in Equation (1)

If there are n pairs of sample observations $(Z_1, Y_1), (Z_2, Y_2), \cdots, (Z_n, Y_n),$ then we get

$$Y_i = \lambda_0 + \lambda_1 Z_i + \varepsilon_i, \ \ i = 1, 2, \cdots, n \tag{12}$$

Then seeking for the estimators $\hat{\lambda}_0$ and $\hat{\lambda}_1$ of $\lambda_0$ and $\lambda_1$ respectively in such a way that P is minimized.

Let
$$P = \sum_{i=1}^{n} \varepsilon_i^2 = \sum_{i=1}^{n} (Y_i - \lambda_0 - \lambda_1 Z_i)^2 \tag{13}$$

Differentiate (13) partially w.r.t. $\lambda_0$ and $\lambda_1$, to get Equations (14) and (15) respectively

$$\sum_{i=1}^{n} Y_i - n\lambda_0 - \lambda_1 \sum_{i=1}^{n} Z_i = 0 \qquad \dots \quad (14)$$

$$\sum Z_i Y_i - \lambda_0 \sum Z_i - \lambda_1 \sum Z_i^2 = 0 \qquad \dots \quad (15)$$

Solving Equations (14) and (15) simultaneously, we get

$$\hat{\lambda}_1 = \frac{n\Sigma Z_i Y_i - \Sigma Z_i Y_i}{n\Sigma Z_i^2 - (\Sigma Z_i)^2} \qquad \dots \quad (16)$$

$$\hat{\lambda}_0 = \frac{\Sigma Z_i^2 \Sigma Z_i - \Sigma Z_i Y_i}{n\Sigma Z_i^2 - (\Sigma Z_i)^2} \qquad \dots \quad (17)$$

The calculation is usually set out in ANOVA form as shown (see Table 1).

**Table 1: Regression ANOVA Table**

| Variance | Degree of freedom | Sum of square | Mean square |
|---|---|---|---|
| Regression | 1 | $RSS = \lambda_1 \sum zy$ | $RMS = \dfrac{RSS}{1}$ |
| Error | $n-2$ | $ESS = TSS - RSS$ | $EMS = \dfrac{ESS}{n-2}$ |
| Total | $n-1$ | $TSS = \sum y^2$ | |

In the same procedure, the parameters of other nonlinear models can be obtained.

**Akaike Information Criterion (AIC)**

The degree of goodness of fit for an assessed measurable equation is known as AIC (Maguilla et al., 2021) and it can be employed for model choice. It is scientifically characterized as:

$$AIC = \exp^{\frac{2p}{n}} \frac{\sum \hat{e}_i^2}{n} = \exp^{\frac{2p}{n}} \frac{SS_R}{n} \qquad (18)$$

where $p$ is the number of parameters with the inclusion of the intercept. Equation (18) is stated mathematically for convenience sake as:

$$\ln(AIC) = \left(\frac{2p}{n}\right) + \ln\left(\frac{SS_R}{n}\right) \qquad (19)$$

**Schwarz Information Criterion (SIC)**

The degree of goodness of fit for an evaluated measurable equation is known as SIC (Obaji & Nwagor, 2021) and it can be employed for model choice. It is mathematically characterized as:

$$SIC = n^{\frac{p}{n}} \frac{\sum \hat{e}_i^2}{n} = n^{\frac{p}{n}} \frac{SS_R}{n} \qquad (20)$$

The log of (20) gives (21):

$$\log_e(SIC) = \frac{p}{n}\log_e(n) + \log_e\left(\frac{SS_R}{n}\right)$$

(21)

## Hannan-Quinn Information Criterion (HQIC)

The degree of goodness of fit for an evaluated measurable equation is known as HQIC (Obaji & Nwagor, 2021) and it can be utilized for model choice. It is mathematically characterized as:

$$HQIC = n\ln\frac{SS_E}{n} + 2p\ln(\ln n)$$

(22)

The equation with least AIC, SIC or HQIC value is chosen as the best model.

### Analysis of Data

The dataset used for this study was extracted from the records of Federal Medical Centre, Owerri, Imo State, Nigeria and presented in Table 2.

**Table 2: Weight of Babies and Hemoglobin Level of Mothers**

| S/N | Weight of babies (Y) | Hemoglobin Level of Mothers (Z) | S/N | Weight of babies (Y) | Hamoglobin Level of Mothers (Z) |
|---|---|---|---|---|---|
| 1 | 3.6 | 14.7 | 41 | 2.8 | 7.7 |
| 2 | 3.1 | 13.6 | 42 | 3.3 | 7.9 |
| 3 | 3.7 | 12.2 | 43 | 3.1 | 8.9 |
| 4 | 3.8 | 14.8 | 44 | 3.2 | 9.4 |
| 5 | 3.0 | 11.7 | 45 | 3.2 | 5.7 |
| 6 | 3.2 | 12.1 | 46 | 3.4 | 14.7 |
| 7 | 2.9 | 7.5 | 47 | 3.0 | 10.1 |
| 8 | 3.1 | 12.5 | 48 | 2.5 | 8.9 |
| 9 | 2.5 | 11.2 | 49 | 3.6 | 9.7 |
| 10 | 2.6 | 12.7 | 50 | 2.9 | 7.4 |
| 11 | 3.7 | 12.9 | 51 | 3.2 | 9.4 |
| 12 | 2.4 | 10.8 | 52 | 2.6 | 8.4 |
| 13 | 2.6 | 11.1 | 53 | 2.3 | 5.7 |
| 14 | 2.7 | 11.6 | 54 | 2.3 | 14.7 |
| 15 | 3.7 | 12.1 | 55 | 3.0 | 13.0 |
| 16 | 3.1 | 5.5 | 56 | 2.9 | 10.1 |
| 17 | 2.8 | 10.5 | 57 | 2.9 | 7.3 |
| 18 | 3.2 | 10.9 | 58 | 4.0 | 6.3 |
| 19 | 3.0 | 10.1 | 59 | 3.4 | 9.5 |
| 20 | 2.5 | 8.9 | 60 | 3.3 | 12.3 |
| 21 | 3.6 | 9.7 | 61 | 3.3 | 10.9 |
| 22 | 2.8 | 7.4 | 62 | 2.8 | 9.9 |
| 23 | 3.2 | 9.4 | 63 | 3.3 | 10.8 |

| 24 | 2.6 | 8.4 | 64 | 3.4 | 11.5 |
|----|-----|------|----|-----|------|
| 25 | 2.3 | 5.7 | 65 | 3.2 | 10.3 |
| 26 | 2.8 | 11.7 | 66 | 2.7 | 8.9 |
| 27 | 3.2 | 13.4 | 67 | 2.9 | 9.9 |
| 28 | 2.9 | 10.1 | 68 | 3.0 | 10.7 |
| 29 | 2.7 | 7.3 | 69 | 2.8 | 7.7 |
| 30 | 4.2 | 12.3 | 70 | 3.3 | 10.9 |
| 31 | 3.4 | 9.5 | 71 | 3.1 | 8.9 |
| 32 | 3.3 | 8.3 | 72 | 3.0 | 8.3 |
| 33 | 2.9 | 10.9 | 73 | 2.5 | 8.1 |
| 34 | 2.5 | 9.9 | 74 | 3.6 | 9.7 |
| 35 | 3.3 | 10.8 | 75 | 2.9 | 7.4 |
| 36 | 3.4 | 13.5 | 76 | 3.2 | 9.5 |
| 37 | 3.2 | 13.3 | 77 | 2.6 | 8.4 |
| 38 | 2.7 | 7.9 | 78 | 2.3 | 5.7 |
| 39 | 2.9 | 9.9 | 79 | 3.8 | 14.7 |
| 40 | 3.0 | 10.7 | 80 | 3.1 | 13.0 |

**Table 3: E-views Software Output for Linear Regression Model**

| Variable | Coefficient | Std. Error | t-Statistic | Prob. |
|----------|-------------|------------|-------------|-------|
| C | 2.375630 | 0.190979 | 12.43920 | 0.0000 |
| Z | 0.066374 | 0.018385 | 3.610266 | 0.0005 |

| | | | |
|----|----|----|----|
| R-squared | 0.143177 | Mean dependent var | 3.047500 |
| Adjusted R-squared | 0.132193 | S.D. dependent var | 0.411842 |
| S.E. of regression | 0.383656 | Akaike info criterion | 0.946543 |
| Sum squared resid | 11.48099 | Schwarz criterion | 1.006094 |
| Log likelihood | -35.86174 | Hannan-Quinn criter. | 0.970419 |
| F-statistic | 13.03402 | Durbin-Watson stat | 1.712895 |
| Prob(F-statistic) | 0.000539 | | |

**Table 4: E-views Software Output for Growth Regression Model**

|  | Coefficient | Std. Error | t-Statistic | Prob. |
|---|---|---|---|---|
| C(1) | 3.802399 | 0.281716 | 13.49729 | 0.0000 |
| C(2) | 2.395918 | 0.894716 | 2.677855 | 0.0090 |

| | | | | |
|---|---|---|---|---|
| R-squared | 0.130143 | Mean dependent var | 3.047500 | |
| Adjusted R-squared | 0.118991 | S.D. dependent var | 0.411842 | |
| S.E. of regression | 0.386564 | Akaike info criterion | 0.961641 | |
| Sum squared resid | 11.65565 | Schwarz criterion | 1.021192 | |
| Log likelihood | -36.46565 | Hannan-Quinn criter. | 0.985517 | |
| Durbin-Watson stat | 1.743530 | | | |

**Table 5: E-views Computer Software for Quadratic Regression Model**

| Variable | Coefficient | Std. Error | t-Statistic | Prob. |
|---|---|---|---|---|
| C | 2.566495 | 0.674993 | 3.802253 | 0.0003 |
| Z | 0.026923 | 0.135020 | 0.199400 | 0.8425 |
| Z^2 | 0.001932 | 0.006550 | 0.294964 | 0.7688 |

| | | | | |
|---|---|---|---|---|
| R-squared | 0.144144 | Mean dependent var | 3.047500 | |
| Adjusted R-squared | 0.121914 | S.D. dependent var | 0.411842 | |
| S.E. of regression | 0.385922 | Akaike info criterion | 0.970414 | |
| Sum squared resid | 11.46804 | Schwarz criterion | 1.059740 | |
| Log likelihood | -35.81656 | Hannan-Quinn criter. | 1.006227 | |
| F-statistic | 6.484229 | Durbin-Watson stat | 1.710524 | |
| Prob(F-statistic) | 0.002497 | | | |

**Table 6: E-views Software Output for Polynomial Regression Model**

| Variable | Coefficient | Std. Error | t-Statistic | Prob. |
|---|---|---|---|---|
| C | 2.583073 | 2.453272 | 1.052909 | 0.2957 |
| Z | 0.021548 | 0.776344 | 0.027755 | 0.9779 |
| Z^2 | 0.002485 | 0.078881 | 0.031500 | 0.9750 |
| Z^3 | -1.81E-05 | 0.002578 | -0.007033 | 0.9944 |

| | | | | |
|---|---|---|---|---|
| R-squared | 0.144145 | Mean dependent var | 3.047500 | |
| Adjusted R-squared | 0.110361 | S.D. dependent var | 0.411842 | |
| S.E. of regression | 0.388452 | Akaike info criterion | 0.995413 | |
| Sum squared resid | 11.46803 | Schwarz criterion | 1.114515 | |
| Log likelihood | -35.81654 | Hannan-Quinn criter. | 1.043165 | |
| F-statistic | 4.266698 | Durbin-Watson stat | 1.709922 | |
| Prob(F-statistic) | 0.007712 | | | |

**Table 7: E-views Software Output for Logarithmic Regression Model**

| Variable | Coefficient | Std. Error | t-Statistic | Prob. |
|---|---|---|---|---|
| C | 1.618701 | 0.408810 | 3.959541 | 0.0002 |
| LOG(Z) | 0.624862 | 0.177792 | 3.514564 | 0.0007 |

| | | | | |
|---|---|---|---|---|
| R-squared | 0.136711 | Mean dependent var | 3.047500 | |
| Adjusted R-squared | 0.125643 | S.D. dependent var | 0.411842 | |
| S.E. of regression | 0.385101 | Akaike info criterion | 0.954062 | |
| Sum squared resid | 11.56764 | Schwarz criterion | 1.013612 | |
| Log likelihood | -36.16247 | Hannan-Quinn criter. | 0.977937 | |
| F-statistic | 12.35216 | Durbin-Watson stat | 1.731528 | |
| Prob(F-statistic) | 0.000737 | | | |

**Table 8: E-views Software Output for Hyperbolic Regression Model**

| Variable | Coefficient | Std. Error | t-Statistic | Prob. |
|---|---|---|---|---|
| C | 3.605930 | 0.173181 | 20.82175 | 0.0000 |
| 1/Z | -5.331010 | 1.600603 | -3.330625 | 0.0013 |

| | | | | |
|---|---|---|---|---|
| R-squared | 0.124511 | Mean dependent var | 3.047500 | |
| Adjusted R-squared | 0.113287 | S.D. dependent var | 0.411842 | |
| S.E. of regression | 0.387813 | Akaike info criterion | 0.968095 | |
| Sum squared resid | 11.73112 | Schwarz criterion | 1.027646 | |
| Log likelihood | -36.72381 | Hannan-Quinn criter. | 0.991971 | |
| F-statistic | 11.09307 | Durbin-Watson stat | 1.761924 | |
| Prob(F-statistic) | 0.001326 | | | |

**Table 9: E-views Software Output for Power Regression Model**

| | Coefficient | Std. Error | t-Statistic | Prob. |
|---|---|---|---|---|
| C(1) | 1.881995 | 0.260943 | 7.212279 | 0.0000 |
| C(2) | 0.210221 | 0.059695 | 3.521559 | 0.0007 |

| | | | | |
|---|---|---|---|---|
| R-squared | 0.138580 | Mean dependent var | 3.047500 | |
| Adjusted R-squared | 0.127536 | S.D. dependent var | 0.411842 | |
| S.E. of regression | 0.384684 | Akaike info criterion | 0.951895 | |
| Sum squared resid | 11.54260 | Schwarz criterion | 1.011446 | |
| Log likelihood | -36.07581 | Hannan-Quinn criter. | 0.975771 | |
| Durbin-Watson stat | 1.725146 | | | |

**Table 10: E-views Software Output for Exponential Growth Regression Model**

|  | Coefficient | Std. Error | t-Statistic | Prob. |
|---|---|---|---|---|
| C(1) | 1.497585 | 0.137485 | 10.89269 | 0.0000 |
| C(2) | 0.042796 | 0.008138 | 5.258567 | 0.0000 |

| | | | | |
|---|---|---|---|---|
| R-squared | 0.144073 | Mean dependent var | 3.047500 | |
| Adjusted R-squared | 0.133100 | S.D. dependent var | 0.411842 | |
| S.E. of regression | 0.383456 | Akaike info criterion | 0.945497 | |
| Sum squared resid | 11.46899 | Schwarz criterion | 1.005048 | |
| Log likelihood | -35.81990 | Hannan-Quinn criter. | 0.969373 | |
| F-statistic | 13.12928 | Durbin-Watson stat | 1.711333 | |
| Prob(F-statistic) | 0.000516 | | | |

**Table 11: E-views Software Output for Square Root Regression Model**

|  | Coefficient | Std. Error | t-Statistic | Prob. |
|---|---|---|---|---|
| C(1) | 1.745097 | 0.366927 | 4.755986 | 0.0000 |
| C(2) | 0.412191 | 0.115328 | 3.574068 | 0.0006 |

| | | | | |
|---|---|---|---|---|
| R-squared | 0.140723 | Mean dependent var | 3.047500 | |
| Adjusted R-squared | 0.129706 | S.D. dependent var | 0.411842 | |
| S.E. of regression | 0.384206 | Akaike info criterion | 0.949404 | |
| Sum squared resid | 11.51389 | Schwarz criterion | 1.008955 | |
| Log likelihood | -35.97617 | Hannan-Quinn criter. | 0.973280 | |
| F-statistic | 12.77396 | Durbin-Watson stat | 1.720299 | |
| Prob(F-statistic) | 0.000607 | | | |

**Table 12: E-views Software Output for Sinusoidal Regression Model**

|  | Coefficient | Std. Error | t-Statistic | Prob. |
|---|---|---|---|---|
| C(1) | 3.047498 | 0.046335 | 65.77086 | 0.0000 |
| C(2) | -0.007996 | 0.064980 | -0.123052 | 0.9024 |

| | | | | |
|---|---|---|---|---|
| R-squared | 0.000194 | Mean dependent var | 3.047500 |
| Adjusted R-squared | -0.012624 | S.D. dependent var | 0.411842 |
| S.E. of regression | 0.414433 | Akaike info criterion | 1.100874 |
| Sum squared resid | 13.39690 | Schwarz criterion | 1.160424 |
| Log likelihood | -42.03495 | Hannan-Quinn criter. | 1.124749 |
| F-statistic | 0.015142 | Durbin-Watson stat | 1.886517 |
| Prob(F-statistic) | 0.902383 | | |

**Table 13: E-views Software Output for Arctangent Regression Model**

|  | Coefficient | Std. Error | t-Statistic | Prob. |
|---|---|---|---|---|
| C(1) | 3.830565 | 28573.27 | 0.000134 | 0.9999 |
| C(2) | 37.15154 | 678694.2 | 5.47E-05 | 1.0000 |
| C(3) | 0.001787 | 32.61915 | 5.48E-05 | 1.0000 |
| C(4) | -0.039174 | 76.22436 | -0.000514 | 0.9996 |

| | | | | |
|---|---|---|---|---|
| R-squared | 0.143180 | Mean dependent var | 3.047500 |
| Adjusted R-squared | 0.109358 | S.D. dependent var | 0.411842 |
| S.E. of regression | 0.388671 | Akaike info criterion | 0.996540 |
| Sum squared resid | 11.48096 | Schwarz criterion | 1.115642 |
| Log likelihood | -35.86162 | Hannan-Quinn criter. | 1.044292 |
| F-statistic | 4.233357 | Durbin-Watson stat | 1.712888 |
| Prob(F-statistic) | 0.008026 | | |

**Table 14: Summary Result of Different Regression Models**

| Model | AIC | SIC | HQIC |
|---|---|---|---|
| Linear Regression | 0.9465 | 1.0061 | 0.9704 |
| Growth Regression | 0.9616 | 1.0212 | 0.9855 |
| Quadratic Regression | 0.9704 | 1.0597 | 1.0062 |
| Polynomial Regression | 0.9954 | 1.1145 | 1.0431 |
| Logarithmic Regression | 0.9541 | 1.0136 | 0.9779 |
| Hyperbolic Regression | 0.9681 | 1.0276 | 0.9920 |
| Power Regression | 0.9519 | 1.0114 | 0.9758 |
| Exponential Growth Regression | 0.9455 | 1.0050 | 0.9694 |
| Square Root Regression | 0.9494 | 1.0090 | 0.9733 |
| Sinusoidal Regression | 1.1009 | 1.1604 | 1.1247 |
| Arctangent Regression | 0.9965 | 1.1156 | 1.0443 |

**Source:** *E-views Software*

Table 14 shows that the polynomial regression model had the lowest HQIC (0.9694), SIC (1.0050), and AIC (0.9455) criteria measures. This suggests that the exponential growth regression model is the most effective model using the dataset employed in this study. The linear regression model—whose criteria scores for AIC is 0.9465, BIC is 1.0061, and HQIC is 0.9704—is the second-best model. Once more, the least performed equation is the sinusoidal regression model, which has the highest HQIC (1.1247), SIC (1.1604), and AIC (1.1009).

## CONCLUSION AND RECOMMENDATION

The result of the study showed that, when it comes to analyzing the association between baby weight and mothers' hemoglobin levels, the exponential growth regression model performs better than the other ten models that were examined. Therefore, researchers should investigate other models that were not included in this analysis and compare the findings using goodness of fit metrics other than the criteria measures used in this work.

## REFERENCES

Bartlett, P. L., Long, P. M., Lugosi, G. & Tsigler, A. (2020). Benign overfitting in linear regression. Proceedings of the National Academy of Sciences of the USA 117(48):30063–30070.

Berk, R. A. (2020). Statistical learning as a regression problem. In: statistical learning from a regression perspective. Berlin: Springer International Publishing, 1–72.

Chicco, D., Warrens, M. J. & Jurman, G. (2021). The coefficient of determination R-squared is more informative than SMAPE, MAE, MAPE, MSE and RMSE in regression analysis evaluation. *PeerJ Computer Science*, 7(2021), 1-24.

Duong, C. M. & Lim, T.T. (2023). Use of regression models for development of a simple and effective biogas decision-support tool. *Scientific Report*, 13(2023), 1-11.

Esemokumo, A. P., Bekesuoyeibo, M. & Nwobi, A. C. (2020). Model selection in bivariate regression models. *International Journal of Applied Science*, 3(4), 1-8.

Esemokumo, P. E. (2023). Asymmetric distributions and nonlinear functions in a canonical correlation analysis using simulated and real-life medical data. An unpublished PhD Thesis submitted to the department of Mathematics and Statistics, Ignatius Ajuru University of Education Rivers State.

Maguilla, E., Escudero, M., Jiménez-Lobato, V., Díaz-Lifante, Z., Andrés-Camacho, C. & Arroyo, J. (2021). Polyploidy expands the range of centaurium (Gentianaceae). *Frontiers in Plant Science*, 12(2021), 1-12.

Montgomery, D. C., Peck, E. A. & Vining, G. G. (2006). Introduction to Linear Regression Analysis. Wiley & Sons, Hoboken.

Obaji, I. & Nwagor, P. (2021). Multiple regression model selection via birth weight, mother age and gestation variables. International Journal of Statistics and Applied Mathematics, 6(6), 83-90.

Spiess, A. & Neumeyer, N. (2010). An evaluation of $R^2$ as an inadequate measure for nonlinear models in pharmacological and biochemical research: a Monte Carlo approach. *BMC Pharmacol*, 10(2010), 34-45.