A COMPARATIVE STUDY OF SOME VARIABLE SELECTION TECHNIQUES IN LOGISTIC REGRESSION

Ijomah Maxwell . A. & Nwali Obinna Department of Mathematics/Statistics University of Port Harcourt, **NIGERIA**

ABSTRACT

Several research works have studied on the performance of variable selection techniques in logistic regression but were limited to models without interaction. In this research, we considered a comparative study of some variable selection techniques in logistic regression for models with and without interaction. Newton Raphson iteration method was applied to obtain coefficients of the variables in the full model (model without interaction). The performance of each technique was judged by their Akaike Information Criterion (AIC) value and the value of the Area under Reciever's Operating Characteristic (AROC) curve. Our findings show that for models without interaction, the forward stepwise, backward stepwise and best subset methods gave same result. Also, for model with interaction, Best subset method outperformed the other two methods. The AROC also revealed that the model fitted using these three methods have an excellent discrimination ability.

1. INTRODUCTION

Logistic Regression is an approach to studying relationship among variables when the dependent variable is categorical (dichotomous, polytomous or ordinal). Binomial logistic regression or binary logistic regression is an aspect of logistic regression that deals with a dependent variable with dichotomous outcome (pass or fail, success or failure, dead or alive, etc). Statistical modeling is aimed at fitting a model with a minimized number of variables which gives a better description of the data and also results in numerical stability. Some commonly used methods for selecting variables in logistic regression include forward selection, backward elimination, stepwise selection, best subset selection, purposeful selection, tabu search, and Bayesian model averaging.

Wang et al. (2004) compared the performance of Bayesian Model Averaging method and stepwise selection method. Their work result in a conclusion that the Bayesian Model Averaging is better that the stepwise selection method. Saker at el. (2009) conducted a study which was aimed at selecting variables for fitting a model for the explanatory variable. They used both stepwise selection method and best subset selection method for variable selection. Their findings revealed that both methods gave same results, but they did not consider comparing the performance of these methods when there are interaction factors present in the model. Hosmer and Lemeshow (2000) also showed that both stepwise selection method and Best subset selection method selected same variables using the UIS data. They also did not consider variable selection when interaction factors are present.

In the absence of comparison of the performance of these three variable selection methods when

interaction factors are present in the model, we are studying the performance of these selection methods to identify the one which is more reliable in the presence of interaction using the model's AkaikeInformation Criterion (AIC) value and the Area Under the Reciever's Operating Characteristic (ROC) curveas a measure for this comparison.

This work is limited to the use of these variable selection techniques in Binary Logistic Regression and did not extend beyond models with two factor interactions.

2. METHODS

A model (without interaction) containing all variables is fitted and these three selection techniques are used to select variables that are considered important. Following this selection is a selection of variables using these three selection techniques when two factor interactions are present. The preferred technique resulting in the suitable fitted model is judged using the Akaike Information Criterion (AIC) and the Area under Receiver's Operating Characteristic (AROC) curve. We performed Newton Raphson Iteration to obtain the coefficient of the variables in the full model (model with all main effect factors). Best subset selection method was done with XLMINER and SPSS 20 was used for forward stepwise, backward stepwise selection and Area under Receiver's Operating Characteristic (AROC) curve. To compare the performance of Best Subset method to other methods, the set of variables resulting in a model with minimum AIC among the selected subsets of models is compared with the AIC resulting from the set of variables selected by other methods.

2.1 Data

The data is a primary data collected (using a questionnaire) from commercial motorcycle operator who carry out their commercial motorcycle operation within some (Aluu, Omagwa, Isiokpo and Elele) areas of Ikwerre Local Government Area. Data was collected from loading points in Aluu, Omagwa, Isiokpo and Elele using a questionnaire. Out of a total of 705 motorcycle operators who are members of the commercial motorcycle operator's union in Ikwerre Local Government Area, a total of 303 motorcycle operators took part in the survey and 274 filled their questionnaire correctly, while 29 of them had some issues with their response. The 274 motorcycle operators whose questionnaires were correctly filled were used as the cases for the study. There are seven independent variables with Crash Involvement as the dependent variable. Dummy variable coding method was used for coding the variables. The variables and variable coding are shown in Table 1.

Variable	Measurement/category H		Parameter coding	
		1	2	
Crash Involvement	Not-involved (0)	0		
	Involved (1)	1		
Motorcycle Ownership	Rented (0)	0		
	Owner-operator (1)	1		
Age	Less than 30 years (1)	0	0	

Table 1. Categorical variable coding

	30-40 years (2)	1	0
	Above 40 years (3)	0	1
Possession of Valid Driver's License	Have no valid license (0)	0	
Driver's License	Have valid license (1)	1	
Knowledge of Road Signs	Have no knowledge of road signs (0)	0	
	Have knowledge of road signs (1)	1	
Alcohol Intake	No (0)	0	
	Yes (1)	1	
Marital Status	Single (0)	0	
	Married (1)	1	
Educational Status	Have no formal education (1)	0	0
	Primary education (2)	1	0
	Secondary education and above (3)	0	1

2.2 Logistic Regression Model

Let N be the number of subjects/population in the dataset,

 $X = (X_1, X_2, X_3, \dots, X_N)^T$, where X_i is a collection of the outcomes of the *i*th subject associated with the k independent random variable and a constant ($x_{i0} = 1 \forall i$), i.e, $X_i = (x_{i0}, x_{i1}, x_{i2}, \dots, x_{ik})$

Let
$$Q = (Q_1, Q_2, Q_3, ..., Q_N)^T$$
, where

 Q_i , i = 1, 2, 3, ..., N are binomial random variables (dependent variables) with values one for success and zero for failure.

Let $X_{ij} = (X_{i1}, X_{i2}, \dots, X_{ik})^T$ be a column vector with the k independent variables as its elements. Let $n = (n_1, n_2, n_3, ..., n_N)^T$, where n_i denotes the number of observation for the i^{th} subject. The probability of success occurring in the *i*th population is

$$\pi(X_i) = E(Q_i = 1|X_i) = \frac{e^{\beta_0 x_{i0} + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_k x_{ik}}}{1 + e^{\beta_0 x_{i0} + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_k x_{ik}}}$$

(1.1)

The Multiple Logistic Regression Model is defined as $\pi(X) = \frac{e^{\beta_0 + \beta_1 X_{i_1} + \beta_2 X_{i_2} + \dots + \beta_k X_{ik}}}{1 + e^{\beta_0 + \beta_1 X_{i_1} + \beta_2 X_{i_2} + \dots + \beta_k X_{ik}}} (1.2)$

 β_0 is the constant term and β_i , j = 1, 2, ..., k is the Logistic Regression coefficient for the kth variable.

The Multiple Logistic Regression Model with two factor interaction is defined as $\pi(X) =$

 $\frac{e^{\beta_0 + \sum_{j=1}^k \beta_j X_{ij} + \sum_{j=1}^{k-1} \beta_{1j+1} X_{i1} X_{ij+1} + \dots + \beta_{k-1k} X_{ik-1} X_{ik}}{1 + e^{\beta_0 + \sum_{j=1}^k \beta_j X_{ij} + \sum_{j=1}^{k-1} \beta_{1j+1} X_{i1} X_{ij+1} + \dots + \beta_{k-1k} X_{ik-1} X_{ik}}}$

(1.3)

where β_{jj+1} , j = 1,2,3,...,k-1, is the Logistic Regression coefficient of the interaction between the (k-1)th variable and the *kth* variable.

The logit of (3.3) is defined as

$$G(X) = \beta_0 + \sum_{j=1}^k \beta_j X_{ij} + \sum_{j=1}^{k-1} \beta_{1j+1} X_{i1} X_{ij+1} + \sum_{j=2}^{k-1} \beta_{2j+1} X_{i2} X_{ij+1} + \beta_{k-1k} X_{ik-1} X_{ik}$$
(1.4)

In (1.3) and (1.4), $X_{ik-1}X_{ik} = 0$, if the (k-1)th variable and the kth variable are design variables of same factor.

2.3 Parameter estimation

The method of estimation adopted for the estimation of the Logistic Regression model parameters is the method of maximum likelihood.

Let the likelihood function be denoted by $L(\beta)$.

$$L(\beta) = \prod_{i=1}^{N} {\binom{n_i}{Q_i}} \pi(X_i)^{Q_i} \left(1 - \pi(X_i)\right)^{n_i - Q_i}$$
(1.5)

Although $\binom{n_i}{Q_i}$ is a constant term in the likelihood function and maximizing

$$L(\beta) = \prod_{i=1}^{N} \pi(X_i)^{Q_i} \left(1 - \pi(X_i)\right)^{n_i - Q_i} \text{gives same result as maximizing}$$
$$L(\beta) = \prod_{i=1}^{N} {n_i \choose Q_i} \pi(X_i)^{Q_i} \left(1 - \pi(X_i)\right)^{n_i - Q_i}$$

since
$$\binom{n_i}{Q_i}$$
 does not contain $\pi(X_i)$, for

$$n_i = 1 \text{ and } Q_i = 0 \text{ or } 1$$
, $\binom{n_i}{Q_i} = 1$.

Therefore,

$$L(\beta) = \ln \left[\prod_{i=1}^{N} \pi(X_i)^{Q_i} \left(1 - \pi(X_i) \right)^{1-Q_i} \right] (1.6)$$

$$\frac{\partial \ln L(\beta)}{\partial \beta_j} = \sum_{i=1}^{N} x_{ij} (Q_i - \pi(X_i)) \qquad (1.7)$$

$$\frac{\partial^2 \ln L(\beta)}{\partial \beta_j \partial \beta_{j'}} = -\sum_{i=1}^{N} x_{ij} x_{ij'} \pi(X_i) (1 - \pi(X_i)) (1.8)$$

The Newton-Raphson equation for iteration is defined as $\beta^{(i+1)} = \beta^{(i)} + (X^T T X)^{-1} X^T (Q - \mu_Q) (1.9)$

where
$$(\beta^{(i)})^{T} = (\beta_{0} \quad \beta_{1} \quad \beta_{2} \quad \cdots \quad \beta_{K})$$

 $Q = \begin{pmatrix} Q_{1} \\ Q_{2} \\ Q_{3} \\ \vdots \\ Q_{N} \end{pmatrix}, \mu_{Q} = \begin{pmatrix} \pi(X_{1}) \\ \pi(X_{2}) \\ \pi(X_{3}) \\ \vdots \\ \pi(X_{n}) \end{pmatrix},$
 $T = \begin{pmatrix} t_{1} \quad 0 \quad 0 \quad \cdots \quad 0 \\ 0 \quad t_{2} \quad 0 \quad \cdots \quad 0 \\ 0 \quad 0 \quad t_{3} \quad \cdots \quad 0 \\ \vdots \quad \vdots \quad \vdots \quad \cdots \quad \vdots \\ 0 \quad 0 \quad 0 \quad \cdots \quad t_{N} \end{pmatrix}$ (1.10)
where $t_{i} = \pi(X_{i})(1 - \pi(X_{i}))$

and
$$X = \begin{pmatrix} x_{10} & x_{11} & x_{12} & & x_{1k} \\ x_{20} & x_{21} & x_{22} & \cdots & x_{2k} \\ x_{30} & x_{31} & x_{32} & & x_{3k} \\ & \vdots & & \ddots & \vdots \\ x_{N0} & x_{N1} & x_{N2} & \cdots & x_{Nk} \end{pmatrix}$$

For the first step of the iteration (i = 0), we guess the elements of $\beta^{(i)}$ $(i.e.\beta^{(0)})$ and use our guess to find $\beta^{(i+1)}(i.e.\beta^{(1)})$. In the next iteration (i = 1), the result for β^{i+1} which is $\beta^{(1)}$ is substituted for $\beta^{(i)}$ in equation (1.9) and then we solve for $\beta^{(i+1)}(i.e.\beta^{(2)})$. The iteration process continues until $\beta^{(i)} \cong \beta^{i+1}$. Note that at any step of the iteration, the elements of $\beta^{(i)}$ are used to solve for $\pi(X_i)$ which is further used for finding $(X^TTX)^{-1}$ and $X^T(Q - \mu_Q)$ at that step.

2.3.1 Variance and covariance

The variance covariance matrix denoted by \mathbb{V} is defined as $\mathbb{V} = (X^T T X)^{-1}$ (1.11)

The diagonal elements in \mathbb{V} are the variance of the Logistic Regression coefficients while the off diagonal elements are the covariance. To compute t_i in T, we make use of our maximum likelihood estimates-

2.3.2 Standard Error and Confidence Interval

The standard error of the jth Logistic Regression coefficient denoted as $\widehat{SE}(\hat{\beta}_i)$ is defined as

 $\widehat{SE}(\hat{\beta}_j) = \left[Var(\hat{\beta}_j) \right]^{\frac{1}{2}} (1.12)$

The $100(1-\alpha)$ % confidence interval estimate of the *kth* Logistic Regression coefficient, denoted as *CI* is

$$\widehat{CI} = \widehat{\beta}_j \pm Z_{1-\frac{\alpha}{2}} \sqrt{Var(\widehat{\beta}_j)}$$
(1.13)

Where $Z_{1-\alpha}$ is the normal critical value of a two-tail test of size α .

2.4 Likelihood ratio test

The likelihood ratio test is used for comparing the likelihood ratio of one Logistic Regression model to another. Let G denote the statistic used for this comparison.

$$\begin{split} G &= \left[-2(loglikelihood of the reduced model - loglikelihood of the full model)\right] \\ (1.14) \\ &\Rightarrow G = -2\ln\left[\frac{likelihood of the reduced model}{likelihood of the full model}\right] \quad (1.15) \\ &\text{Supposing that} \quad X_{i1}, X_{i2}, X_{i3} \text{ are independent variables of a Logistic Regression model,} \end{split}$$

To test $H_0: \beta_1 = 0$, we fit $G(X) = \beta + \beta X + \beta X + \beta X$ and

 $G(X) = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \beta_3 X_{i3}$ and $G(X') = \beta_1 + \beta_2 X_1 + \beta_3 X_2$ and then commute

 $G(X') = \hat{\beta}_0 + \hat{\beta}_2 X_{i2} + \hat{\beta}_3 X_{i3}$ and then compute the log-likelihood of each model.

Likelihood of the full model is $\prod_{i=1}^{N} \pi(X_i)^{Q_i} (1 - \pi(X_i))^{1-Q_i}$, where

$$\pi(X_i) = \frac{e^{\beta_0 x_{i0} + \beta_1 x_{i1} + \beta_2 x_{i2} + \beta_3 x_{i3}}}{1 + e^{\beta_0 x_{i0} + \beta_1 x_{i1} + \beta_2 x_{i2} + \beta_3 x_{i3}}}$$

Also, the likelihood of the reduced model is $\prod_{i=1}^{N} \pi(X_i')^{Q_i} (1 - \pi(X_i'))^{1-Q_i}$, where

$$\pi(X_{i}') = \frac{e^{\beta_{0}x_{i0} + \beta_{2}x_{i2} + \beta_{3}x_{i3}}}{1 + e^{\beta_{0}x_{i0} + \beta_{2}x_{i2} + \beta_{3}x_{i3}}}$$
$$G = \left[-2\ln\left[\frac{\prod_{i=1}^{N} \pi(X_{i})^{Q_{i}} \left(1 - \pi(X_{i})\right)^{1 - Q_{i}}}{\prod_{i=1}^{N} \pi(X_{i}')^{Q_{i}} \left(1 - \pi(X_{i}')\right)^{1 - Q_{i}}}\right]\right]$$

If $P(\chi^2_{\nu} > G) < \alpha$, we reject the null hypothesis and conclude at α level of significance that the variable X_{i1} contributes significantly to the prediction of the dependent variable if it is included in the model.

2.5 Area under the Receiver Operating Characteristic (ROC) curve.

Classification is done using a cut-off point between zero and one which is chosen by the researcher. This cut-off point if not well selected may lead to an inaccurate classification. The area under the ROC curve gives a very good description of classification accuracy. To produce an ROC curve, we plot the probability of detecting a true signal (sensitivity) and false signal (1-specificity) for an entire range of cut-off points. The area under the ROC curve ranges from zero to one and it gives a measure of the model's ability to discriminate between those subjects who experience the outcome of interest versus those who do not.

According to Hosmer and Lemeshow (2000), if area under *ROC curve* ≈ 0.5 , this suggests no discrimination, if $0.7 \leq area under ROC curve < 0.8$, this is considered an acceptable discrimination if $0.8 \leq area under ROC curve \leq 0.9$, this is considered an excellent discrimination, if *area under ROC curve* > 0.9, this is considered an outstanding discrimination.

2.5.1 Estimate of area under ROC curve

Let N_c be the set of subjects with $Q_i = 1$ and N'_c the set of subjects with $Q_i = 0$.

Let $\mathbb{N}_c = N_c \times N'_c = \{n_1, n_1, n_3, ..., n_n\}$, where n_i in \mathbb{N}_c is a pair of an element in N_c and an element in N'_c , and n is the number of elements in N_c multiplied by the number of element in N'_c or the total number of pairs of subjects with $Q_i = 1$ and subjects with $Q_i = 0$.

For each pair, we compare the estimated odd of the subject with $Q_i = 1$ and the estimated odd for the subject with $Q_i = 0$.

Let us define $R = \{r_1, r_1, r_3, ..., r_n\}$, where r_i is a value assigned to n_i based on comparison of the estimated odd of subjects in n_i .

For each n_i : $r_i = 1$ if the estimated odd of the subject with $Q_i = 1$ is greater than the estimated odd of the subject with $Q_i = 0$; $r_i = 0.5$ if the estimated odd of the subject with $Q_i = 1$ is equal to the estimated odd of the subject with $Q_i = 0$; $r_i = 0$ if the estimated odd of the subject with $Q_i = 1$ is less than the estimated odd of the subject with $Q_i = 0$.

The estimate of the area under the ROC curve is $\frac{\sum_{i=1}^{n} r_i}{n}$.

2.6 Stepwise Logistic Regression

In regression analysis, a collection of variables (independent) are studied to know the association of such variables with a particular variable (dependent). This collection of independent variables contains important and unimportant ones. Stepwise Regression is employed to carry out a stepwise selection procedure aimed at screening this collection of variables, and fitting several Logistic Regression equations simultaneously. The selection procedure is based on a statistical algorithm which carries out a check for important variables, and either add or delete them in accordance with a decision rule. The importance of an independent variable depends on a measure of the significance of the coefficient of that variable. In Logistic Regression, significance of a variable is determined using the likelihood ratio chi-square test. The most important variable is usually the variable with the largest change in log-likelihood relative to the model not containing the variable. The steps for Forward Stepwise selection and Backward Stepwise selection are as follow:

2.6.1 Forward Stepwise selection STEP 0

Let L_0 denote the log-likelihood of the intercept only model,

 $L_{ii}^{(0)}$ denote the log-likelihood of the model containing the jth independent variable,

 $G_{ij}^{0} = -2(L_0 - L_{ij}^{(0)}), P_{ij}^{(0)} = Pr(\chi^2(v) > G_{ij}^{(0)}), \text{ where } v = 1 \text{ if } X_{ij} \text{ is continuous and } k - 1 \text{ if } X_{ij} \text{ is polytomous,}$

 $P_{ie_1}^{(0)} = min(P_{ij}^{(0)}),$

 X_{ie} , denote the most important variable at step zero,

 P_E denotes the α – *level* to judge the importance of a variable

Compute $P_{ie_1}^{(0)}$ and move to step 1 if $P_{ie_1}^{(0)} < P_E$, if otherwise, terminate the process.

STEP 1

 $L_{ie}^{(1)}$ denote the log-likelihood of the model containing X_{ie_1} ,

 $L_{ie_1 ij}^{(1)}$ denote the log-likelihood of the model containing X_{ie_1} and X_{ij} , $j = 1, 2, ..., k, j \neq e_1$

 $G_{ij}^{(1)} = -2\left(L_{ie_1}^{(1)} - L_{ie_1ij}^{(1)}\right)$, be the likelihood statistic of the model containing X_{ie_1} and X_{ij} versus the model containing only X_{ie_1} ,

$$P_{ij}^{(1)} = Pr(\chi^2(v) > G_{ij}^1),$$

 X_{ie_2} denote the variable with minimum p-value when added to the model containing X_{ie_1} , $P_{ie_2}^{(1)} = min(P_{ij}^{(1)})$. Compute $P_{ie_2}^{(1)}$ and proceed to step 2 if $P_{ie_2}^{(1)} < P_E$, stop if otherwise.

STEP 2

Let $L^2_{-ie_i}$ denote the log-likelihood of the model without X_{ie_j}

 $G_{-ie_j}^{(2)} = -2\left(L_{-ie_j}^{(2)} - L_{ie_1ie_2}^{(2)}\right) , \text{ where } L_{ie_1ie_2}^{(2)} \text{ is the likelihood of the model}$ containing X_{ie_1} and X_{ie_2} , $P_{ie_1}^{(2)} = P_{ie_1ie_2}^{(2)} = P_{ie_1ie_2$

$$P_{-ie_{j}}^{(2)} = Pr\left(\chi^{2}(v) > G_{-ie_{j}}^{2}\right)$$

Let X_{ir_2} denote the variable that yields the maximum p-value when removed from the model containing X_{ie_1} and X_{ie_2} ,

$$P_{r_2}^{(2)} = max \left(P_{-ie_1}^{(2)}, P_{-ie_2}^{(2)} \right).$$

Remove X_{ir_2} if $P_{r_2}^{(2)} > P_R$, where P_R denote the $\alpha - level$ for removal of a variable. If $P_{r_2}^{(2)} < P_R$, X_{ir_2} remains in the model and then we fit a model containing

 $X_{ie_1}, X_{ie_2}, and \; X_{ij}, j = 1, 2, \dots, k, j \neq e_1, e_2.$ Let $L_{ie_1ie_2ij}^{(2)}$ denote the log-likelihood containing X_{ie_1}, X_{ie_2} , and $X_{ij}, j = 1, 2, ..., k, j \neq e_1, e_2$, Let of the model $G_{ij}^{(2)} = -2 \left(L_{ie_1 ie_2}^{(2)} - L_{ie_1 ie_2 ij}^{(2)} \right),$ $P_{ie_3}^{(2)} = min \left(P_{ij}^{(2)} \right).$ Compute $P_{ie_{\pi}}^{(2)}$ and move to step 3 if $P_{ie_{\pi}}^{(2)} < P_{E}$, if otherwise, terminate the process.

STEP3

The procedure in STEP 3 is similar to the procedure in STEP 2. The process continues in this manner until the last step, STEP S.

STEP S

The final step occurs when all k variables have entered the model or when all the variables that are in the model have p-values to remove less than P_R , and the variables not in the model have p-value to enter greater than P_E .

2.6.2 Backward Stepwise selection **STEP 0**

Let L_f denote the log-likelihood of the full model,

 $L_{-i}^{(0)}$ denote the log-likelihood of the model not containing the jth independent variable,

 $G_{-ij}^{0} = -2(L_{-ij}^{(0)} - L_{f}),$ $P_{-ij}^{(0)} = Pr(\chi^2(v) > G_{-ij}^{(0)})$, where v = 1 if X_{ij} is continuous and k - 1 if X_{ij} is polytomous, $P_{ir_{i}}^{(0)} = max(P_{ii}^{(0)}),$

 X_{ir} , denote the most unimportant variable at step zero,

 P_R denote the α - *level* to judge the unimportance of a variable Compute $P_{ir_1}^{(0)}$ and move to step 1 if $P_{ir_1}^{(0)} > P_R$, if otherwise, terminate the process.

STEP1

 $L_{-ir_1}^{(1)}$ denote the log-likelihood of the model not containing X_{ir_1} , $L_{-ir_1ij}^{(1)}$ denote the log-likelihood of the model not containing X_{ir_1} and X_{ij} , $j = 1, 2, ..., k, j \neq r_1$ $C_{-ir_1ij}^{(1)} = -2(I_1^1, ..., I_{ij}^1)$ here the likelihood statistic of the model not containing X_{ir_1} and X_{ij} , $j = 1, 2, ..., k, j \neq r_1$ $G_{-ij}^{(1)} = -2(L_{-ir_1ij}^1 - L_{-ir_1}^1)$, be the likelihood statistic of the model not containing X_{ir_1} and X_{ij} versus the model not containing only X_{ir_i} ,

$$P_{-ij}^{(1)} = Pr(\chi^2(v) > G_{-ij}^1),$$

 $X_{ir_{s}}$ denote the variable with maximum p-value when removed from the model not containing $X_{ir_{s}}$, $P_{-ir_2}^{(\tilde{1})} = max(P_{-ij}^{(1)})$. Compute $P_{-ir_2}^{(1)}$ and proceed to step 2 if $P_{-ir_2}^{(1)} > P_R$, stop if otherwise.

STEP 2

Let $L_{in}^{(2)}$ denote the log-likelihood of the model containing X_{ink} (k = 1,2) and other variables not

removed in step one.

 $G_{ir_{k}}^{(2)} = -2\left(L_{-ir_{1}ir_{2}}^{(2)} - L_{ir_{k}}^{(2)}\right), \text{ where } L_{-ir_{1}ir_{2}}^{(2)} \text{ is the likelihood of the model not containing } X_{ir_{1}} \text{ and } X_{ir_{2}},$ $P_{ir_{k}}^{(2)} = Pr\left(\chi^{2}(v) > G_{ir_{k}}^{2}\right).$

Let X_{ie_2} denote the variable that yields the minimum p-value when entered into the model not containing X_{ir_1} and X_{ir_2} ,

 $P_{e_2}^{(2)} = \min\left(P_{ir_1}^{(2)}, P_{ir_2}^{(2)}\right).$ Enter X_{ie_2} if $P_{e_2}^{(2)} < P_E$, where P_E denote the α - *level* for entering of a variable. If $P_{e_2}^{(2)} > P_E$, X_{ie_2} is removed from the model and then we fit a model not containing $X_{ir_1}, X_{ir_2}, and X_{ij}, j = 1, 2, ..., k, j \neq e_1, e_2.$ Let $L_{-ir_1ir_2 ij}^{(2)}$ denote the log-likelihood of the model not

 $\begin{array}{l} & \text{Let} & T_{-ir_{1}ir_{2}\ ij} & \text{denote} & \text{the} & \text{tog intermodel} & \text{of} & \text{the} & \text{finder} & \text{for} \\ & \text{containing} X_{ir_{1}}, X_{ir_{2}}, \text{and} & X_{ij}, j = 1, 2, \dots, k, j \neq r_{1}, r_{2}, \\ & G_{-ij}^{(2)} = -2 \left(L_{ir_{1}ir_{2}\ ij}^{(2)} - L_{ir_{1}ir_{2}}^{(2)} \right), & P_{-ij}^{(2)} = Pr \left(\chi^{2}(v) > G_{-ij}^{(2)} \right) \\ & P_{-ir_{8}}^{(2)} = max \left(P_{-ij}^{(2)} \right). \\ & \text{Compute} & P_{-ir_{8}}^{(2)} & \text{and move to step 3 if} P_{-ir_{8}}^{(2)} > P_{R} \text{, if otherwise, terminate the process.} \end{array}$

STEP 3

The procedure in STEP 3 is similar to the procedure in STEP 2. The process continues in this manner until the last step, STEP S.

STEPS

The final step occurs when all k variables have been removed from the model or when all the variables that are not in the model have p-values to enter greater than P_E , and the variables in the model have p-value to be removed less than P_R .

2.7 Best subset Logistic Regression

Best subset logistic regression is another method used for variable selection. This method of variable selection searches for the best model among models with equal number of variables, based on some criterion (Mallow's C_q , AR^2 , etc.).

Best subset selection method is used to select the model with the minimum Cq or the maximum AR^2 from the set of models with one variable, two variables, three variables, etc.

Hosmer et al. (1989) showed how to conduct best subset selection in logistic regression using any software capable of performing best subset linear regression analysis when weights are involved. To conduct best subset logistic regression using a linear regression program, it is required that we already know the coefficients of the logistic regression variables.

$$Z = \begin{pmatrix} z_1 \\ z_2 \\ z_3 \\ \vdots \\ z_N \end{pmatrix}$$
$$\hat{\beta} = (X^T T X)^{-1} X^T T (X \hat{\beta} + T^{-1} (Q - \mu_Q))$$
$$= (X^T T X)^{-1} X^T T Z, \text{ where}$$

$$Z = (X\hat{\beta} + T^{-1}(Q - \mu_Q))$$

$$Z_i = (1 \quad x_{i1} \quad x_{i2} \quad \dots \quad x_{ik}) \begin{pmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \\ \hat{\beta}_2 \\ \vdots \\ \hat{\beta}_k \end{pmatrix} + \frac{Q_i - \pi(X_i)}{\pi(X_i)(1 - \pi(X_i))}$$

Let $X = (\mathbb{X}_1 \ \mathbb{X}_2)$, where \mathbb{X}_1 is an $N \times q + 1$ design matrix for the subset of q variables in the model and the constant, and \mathbb{X}_2 is an $N \times k - q$ matrix containing the remaining k - q. Let us also partition the vector of coefficients for \mathbb{X}_1 and \mathbb{X}_1 as $\beta^T = (\beta_1^T \ \beta_2^T)$. The resulting information matrix (I) is $I = \begin{pmatrix} I_{11} & I_{12} \\ I_{21} & I_{22} \end{pmatrix}$, where $I_{11} = (\mathbb{X}_1^T T \mathbb{X}_1)$, $I_{12} = (\mathbb{X}_1^T T \mathbb{X}_2) = I_{21}^T$, and $I_{22} = (\mathbb{X}_2^T T \mathbb{X}_2)$. The estimator of the coefficient vector β_1 obtained from the linear regression of Z on \mathbb{X}_1 using weighted matrix V is $\widetilde{\beta_1} = (\mathbb{X}_1^T T \mathbb{X}_1)^{-1} \mathbb{X}_1^T T Z (1.16)$ The vector of fitted values for \mathbb{X}_1 obtained using linear regression is $\widetilde{Z}(q) = \mathbb{X}_1 \widetilde{\beta_1}$. (1.17) The residual sum of squares for the fitted model containing variables in \mathbb{X}_1 is

$$SSE(q) = [Z - \tilde{Z}(q)]^T T[Z - \tilde{Z}(q)]$$

= $Z^T T Z - \tilde{\beta_1}^T \mathbb{X}_1^T T \mathbb{X}_1 \tilde{\beta_1} (3.38)$
$$SSE(k) = [Z - \hat{Z}(k)]^T T[Z - \hat{Z}(k)]$$

= $Z^T T Z - \beta^T X^T T X \hat{\beta} (1.18)$

The Mallow's Cq statistic for a particular subset of q variables, using linear regression is

$$C_q = \frac{SSE(q)}{(SSE(k)/N-k-1)} + 2(q+1) - N(1.19)$$

2.8 Akaike Information Criterion (AIC)

Akaike Information Criterion (AIC) is a measure which enables the comparison of a set of statistical models to identify the model which minimizes information lost. From a set of statistical models, AIC gives an estimate of the quality of each of the model relative to each of the other model.

AIC is defined as

AIC = 2k - 2loglikelihood, where k is the number of estimated parameters in the model Assuming that there are five models for comparison, if we define the AIC's of these models as AIC_1 , AIC_2 , AIC_3 , AIC_4 , and AIC_5 , then the model having min(AIC_1 , AIC_2 , AIC_3 , AIC_4 , AIC_5) is the best model among these five models.

3. RESULTS

3.1 The Logistic Regression Model with only the constant term

The Logistic Regression model with only the constant term is defined as

 $\pi(X) = \frac{e^{\beta_0}}{1 + e^{\beta_0}}$

3.1.1 Parameter Estimation

$$\begin{aligned} \widehat{\pi}(X_i) &= \frac{e^{\beta_0}}{1 + e^{\beta_0}} = \frac{\sum_{i=1}^n Q_i}{\sum_{i=1}^n n_i} = \frac{204}{274} = 0.7445 \text{ for all } i = 1, 2, \dots, 274 \\ \mathbf{1} - \widehat{\pi}(X_i) &= 0.2555 \text{ for } i = 1, 2, \dots, 274 \text{ , } \widehat{\beta_0} = \\ \ln \frac{0.7445}{0.2555} &= 1.0695 \\ \mathbf{v}\widehat{\alpha}r(\widehat{\beta_0}) &= \frac{1}{\sum_{i=1}^n \widehat{\pi}(X_i)(1 - \widehat{\pi}(X_i))} = 0.01919, \quad \widehat{SE}(\widehat{\beta_0}) = \sqrt{var(\beta_0)} \approx 0.139 \\ \widehat{CI} &= 1.07 \pm 1.96 \times 0.139 = 1.07 \pm 0.27244 = (0.79756, 1.34244) \end{aligned}$$

3.1.2 Model Likelihood

$$\begin{aligned} \mathbf{L}(\boldsymbol{\beta}_{0}) &= \prod_{i=1}^{n} \pi(X_{i})^{Q_{i}} \left(1 - \pi(X_{i})\right)^{1 - Q_{i}} \\ &= 2.365 \times 10^{-68} \\ &\ln L(\boldsymbol{\beta}_{0}) = \ln(2.365 \times 10^{-68}) = -155.7060 , -2 \ln L(\boldsymbol{\beta}_{0}) = 311.4119 \end{aligned}$$

3.2 Logistic Regression model without interaction

The Logistic regression model without interaction, described by Equation 1.2 is defined as follow

- X_{i1} Motorcycle Ownership (MO)
- X_{i2} –Possesion of valid driver's license (POVDL)
- X_{i3} Alcohol intake (AI)
- X_{i4} Knowledge of Road Signs (KORS)
- X_{i5} Marital Status (MS)
- $X_{i6} (30 40 \text{ years})$ (AGE1)
- X_{i7} (above 40 years) (AGE2)
- X_{i8} Primary education (EDU1),

 X_{i9} –Secondary education and above (EDU2)

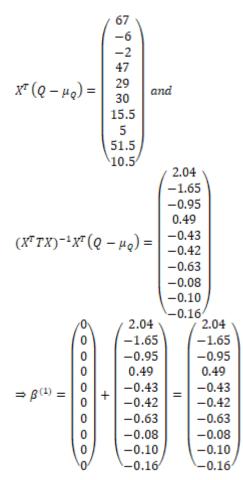
$$\pi(X) = \frac{e^{\beta_0 + \beta_1 X_{i_1} + \beta_2 X_{i_2} + \beta_3 X_{i_3} + \beta_4 X_{i_4} + \beta_5 X_{i_5} + \beta_6 X_{i_6} + \beta_7 X_{i_7} + \beta_8 X_{i_8} + \beta_9 X_{i_9}}{1 + e^{\beta_0 + \beta_1 X_{i_1} + \beta_2 X_{i_2} + \beta_3 X_{i_3} + \beta_4 X_{i_4} + \beta_5 X_{i_5} + \beta_6 X_{i_6} + \beta_7 X_{i_7} + \beta_8 X_{i_8} + \beta_9 X_{i_9}}}$$

3.2.1 Parameter estimation

3.2.1.1 Logistic Regression coefficient estimation using Newton-Raphson iteration method. $\beta^{(i)} = \beta^{(i-1)} + (X^T T X)^{-1} X^T (Q - \mu_Q)$. Let $\beta^{(0)}$ be a zero vector, $\beta^{(1)} = \beta^{(0)} + (X^T T X)^{-1} X^T (Q - \mu_Q)$.

$$\begin{aligned} \beta^{(3)} &= \beta^{(3)} + (X^{T}X)^{-1}X^{T}(Q - \mu_{Q}) \\ \pi(X_{i}) &= \frac{e^{(0)x_{i0} + (0)x_{i1} + (1)x_{i2} + (0)x_{i3} + (0)x_{i4} + (0)x_{i5} + (0)x_{i6} + (0)x_{i7} + (0)x_{i8}(0)x_{i9}}{1 + e^{(0)x_{i0} + (0)x_{i1} + (1)x_{i2} + (0)x_{i3} + (0)x_{i4} + (0)x_{i5} + (0)x_{i6} + (0)x_{i7} + (0)x_{i8}(0)x_{i9}}, t_{i} \\ &= \pi(X_{i})(1 - \pi(X_{i})) \end{aligned}$$

(X	$^{T}TX) =$			$ \begin{array}{c} \sum_{i=1}^{274} x_{i2} t_i \\ \sum_{i=1}^{274} x_{i1} x_{i2} t_i \\ \sum_{i=1}^{274} x_{i2}^2 t_i \end{array} \\ \end{array} \\ \begin{array}{c} \sum_{i=1}^{274} x_{i2} x_{i9} t_i \end{array} $	$ \begin{array}{c} \sum_{i=1}^{274} x_{ii} \\ \cdots & \sum_{i=1}^{274} x_{i1} \\ x_{i2} \\ \sum_{i=1}^{274} x_{i2} \\ \vdots \\ \cdots & \sum_{i=1}^{274} x_{i9} \end{array} $	$t_{i9}t_i$ $t_{i9}t_i$
	/68.50	17.509.00			6.50049.75	13.75
	17.50	17.505.00	0 9.25010.25	8.7504.250	1.50011.50	3.500
	9.000	5.0009.00	0 4.5005.500	4.0001.000	1.0005.950	2.750
	38.50	9.2504.50	0 38.5021.00	22.0013.00	3.75026.70	8.250
_	37.50	10.255.50	0 21.0037.50	19.0011.75	2.50026.75	8.000
_	38.00	8.7504.00	0 22.0019.00	38.0017.75	6.50026.50	9.250
	22.75	4.2501.00	0 13.0011.75	17.7522.75	0.00017.25	4.500
	6.500	1.5001.00	0 3.7506.500	6.5000.000	6.5003.500	1.750
	49.75	11.505.75	0 26.7526.50	26.5017.25	3.50049.75	0.000
	\13.75	3.5002.75	0 8.0009.250	9.2504.500	1.7500.000	13.75/



	/ 3.139 \		/ 3.513 \
	-2.233		-2.455
	-1.308		-1.436
	0.757		0.859
$\beta^{(2)} =$	-0.729	β ⁽³⁾ =	-0.849
p =	-0.544	р —	-0.593
	-1.141		-1.352
	-0.784		-0.985
	-0.322		0.359
	-0.460		-0.513
	/ 3.558 \		/ 3.559 \
	-2.482		-2.482
	-1.455		-1.455
	0.870		0.870
$\beta^{(4)} =$	-0.864	$\beta^{(5)} =$	-0.865
p	-0.605	P** -	-0.606
	-1.376		-1.377
	-1.004		-1.005
	-0.359		0.359
	-0.511		-0.512

Since $\beta^{(4)} \cong \beta^{(5)}$ the iteration converges at the fifth iteration and $\beta^{(5)}$ is the vector of our maximum likelihood estimates of the Logistic Regression Model coefficients.

	/β ₀ \		/ 3.559 \
	β_1		-2.482
	β ₂		-1.455
_	β3	=	0.870
	β_4		-0.865
~	β ₅		-0.606
	β ₆		-1.377
	β7		-1.005
	β ₈		0.359
	\ β ₀∕		-0.512

3.2.1.2 Variance and Standard error of estimate

At fifth iteration. the 35.060 13.4366.0680 16.98021.575 22.68614.000 4.500024.717 7.4340 13.436 13.4363.2580 7.42707.5750 6.61502.5310 1.32708.7910 2.5090 6.0680 3.25806.0680 3.0680 3.6790 2.89600.3530 0.92903.9070 1.7530 16.980 7.42703.0680 16.98010.696 11.3286.5860 2.410011.407 4.0620 $(X^T T X) =$ 21.575 7.57503.6790 10.69621.575 13.0848.3670 1.897014.988 4.8050 22.686 6.61502.8960 11.32813.084 22.68611.793 4.499015.816 5.7080 14.000 2.53100.3530 6.58608.3670 11.79314.000 0.000010.674 2.9930 4.5000 1.32700.9290 2.41001.8970 4.49900.0000 4.50002.4360 1.2350 24.717 8.79103.9070 11.40714.988 15.81610.674 2.436024.717 0.0000 -0.365 70.6040 ^ 0.13900.0440 ^ 0.0390 - 0.106 ^ 0.044 - 0.086 ^ - 0.134 - 0.386 0.1390 0.15400.0070 0.0170 0.0210 0.0110 0.0520 0.0460 0.0470 0.0480 0.0040 0.00700.2280 0.0010 0.0060 0.0120 0.0550 0.0110 -0.025 -0.043 $\Rightarrow (X^T T X)^{-1} = \begin{bmatrix} 0.0070 & 0.00700.2280 & 0.0010 & 0.0080 & 0.0120 & 0.0550 & 0.0110 & -0.025 \\ 0.0390 & 0.01700.0010 & 0.1170 - 0.006 & -0.007 - 0.003 & -0.008 & 0.0020 \\ 0.1060 & 0.02100.0060 & 0.00600.1280 & 0.0080 & 0.0170 & 0.0370 & 0.0080 \\ 0.0440 & 0.01100.0120 & 0.00700.0080 & 0.1810 - 0.071 & -0.106 - 0.037 \\ 0.0860 & 0.05200 & 0.550 & 0.00200 & 0.170 & 0.02710 & 0.0120 & 0.0120 & 0.0120 & 0.0120 & 0.0010 & 0$ -0.0070.0000 -0.056

Ť	$\begin{pmatrix} VAR(\beta_0) \\ VAR(\beta_1) \\ VAR(\beta_2) \\ VAR(\beta_3) \\ VAR(\beta_4) \\ VAR(\beta_5) \\ VAR(\beta_5) \\ VAR(\beta_6) \\ VAR(\beta_7) \\ VAR(\beta_8) \\ VAR(\beta_8$))))) = =	=	0.604 0.154 0.228 0.117 0.128 0.181 0.201 0.378 0.426 0.530	
î	$\begin{array}{c} VAR(\beta_9)\\ SE(\beta_0)\\ SE(\beta_1)\\ SE(\beta_2)\\ SE(\beta_3)\\ SE(\beta_4)\\ SE(\beta_4)\\ SE(\beta_5)\\ SE(\beta_6)\\ SE(\beta_7)\\ SE(\beta_8)\\ SE(\beta_9)\\ SE(\beta_9) \end{array}$	=		0.777 0.392 0.478 0.342 0.358 0.425 0.448 0.615 0.653 0.728	

3.2.2 Variable selection for model without interaction

For Forward and Backward stepwise logistic regression, $P_E = 0.15$ and $P_R = 0.2$.

Table 2. Summary of Mallow's Cq, AIC and AROC for selected models (model without interaction)

Selection Method	Variables in the model	Cq	AIC	AROC
Forward stepwise	$\begin{array}{c} X_{i1}, X_{i2}, X_{i3}, \\ X_{i4}, X_{i5}, X_{i6}, X_{i7} \end{array}$	6.500	238.493	0.845
Backward stepwise	$\begin{array}{c} X_{i1}, X_{i2}, X_{i3}, \\ X_{i4}, X_{i5}, X_{i6}, X_{i7} \end{array}$	6.500	238.493	0.845
Best subset	$\begin{array}{c} X_{i1}, X_{i2}, X_{i3}, \\ X_{i4}, X_{i5}, X_{i6}, X_{i7} \end{array}$	6.500	238.493	0.845

3.3 The Logistic regression model with two factor interaction

The Logistic regression model with interaction, described by Equation 1.3 is defined as follow

 X_{i1} – Motorcycle Ownership(MO)

 X_{i2} –Possesion of Valid Driver's License (POVDL)

 X_{i3} – Alcohol intake (AI)

 X_{i4} – Knowledge of Road Signs (KORS)

 X_{i5} – Marital Status (MS)

 $X_{i6} - (30 - 40 \text{ years}) \text{ (AGE1)}$

 X_{i7} – (above 40 years) (AGE2)

$$\pi(X) = \frac{e^{\beta_0 + \sum_{j=1}^7 \beta_j X_{ij} + \sum_{j=1}^6 \beta_{1j+1} X_{i1} X_{ij+1} + \sum_{j=2}^6 \beta_{2j+1} X_{i2} X_{ij+1} + \dots + \beta_{67} X_{i6} X_{i7}}{1 + e^{\beta_0 + \sum_{j=1}^7 \beta_j X_{ij} + \sum_{j=1}^6 \beta_{1j+1} X_{i1} X_{ij+1} + \sum_{j=2}^6 \beta_{2j+1} X_{i2} X_{ij+1} + \dots + \beta_{67} X_{i6} X_{i7}}, X_{i6} X_{i7} = 0$$

3.3.1 Variable selection for models with two factor interaction.

For Forward and Backward stepwise logistic regression, $P_E = 0.15$ and $P_R = 0.2$.

Selection Method	Variables in the model	AIC	AROC
Forward stepwise	$\begin{array}{c} X_{i1}, \ X_{i3}, \ X_{i5}, \ X_{i1} X_{i4}, \\ X_{i2} X_{i3}, \ X_{i3} X_{i6}, \ X_{i4} X_{i7} \\ X_{i5} X_{i6} \end{array}$	226.512	0.862
Backward stepwise	$\begin{array}{c} X_{i1}, \ X_{i2}, \ X_{i4}, \ X_{i6}, \ X_{i7}, \\ X_{i2} X_{i3}, \ X_{i3} X_{i5}, \ X_{i4} X_{i6} \end{array}$	220.568	0.868
Best subset	$\begin{array}{c} X_{i1}, X_{i3}, X_{i4}, X_{i5}, X_{i6}, X_{i7}, \\ X_{i2}X_{i3}, X_{i4}X_{i6} \end{array}$	220.185	0.871

4. DISCUSSION

The data analyzed contains seven dependent variables with 274 observations. The independent variables are Motorcycle Ownership (X_{i1}) , Possession of Valid Driver's License (X_{i2}) , Alchohol Intake (X_{i3}) , Knowledge of Road Signs (X_{i4}) Marital Status (X_{i5}) , AGE1 [30-40 (X_{i6})], AGE2 [above 40 (X_{i7})], EDU1 [Primary education (X_{i8})], EDU2 [Secondary education and above (X_{i9})]. We used the Newton Ralphson iteration method to obtain coefficients of the variables in the full model (model without interaction factors). We also performed variable selection to fit a model not containing interaction factors and a model containing two factor interactions using Forward Stepwise, Backward Stepwise, and Best Subset methods of variable selection.

4.1 Model without interaction

The variables considered for selection are $X_{i1}, X_{i2}, X_{i3}, X_{i4}, X_{i5}, X_{i6}, X_{i7}, X_{i8}$, and X_{i9} . The three selection techniques employed excluded X_{i8} and X_{i9} . Table 2 gave a summary of the Mallow's Cq, AIC and AROC for selected models. The information on Table 2 show that Best subset method, Forward stepwise method and backward stepwisemethod selected same set of variables $(X_{i1}, X_{i2}, X_{i3}, X_{i4}, X_{i5}, X_{i6}, X_{i7})$ [Mallow's Cq (6.500) AIC (238.493), and AROC (0.845)]. This shows that forward stepwise method, backward stepwise method and best subset method has equal performance for selecting variables when interaction factors are not present.

The resulting AROC values for the models in Table 2 are all between 0.8 and 0.9, which indicates that the fitted models have excellent discrimination ability.

4.2 Model with interaction

For model with interaction, there are differences in the variables selected by the three variable selection methods adopted in this research work.

Best Subset method selected 25 models from which we chose the model [variable in the model $(X_{i1}, X_{i3}, X_{i4}, X_{i5}, X_{i6}, X_{i7}, X_{i2}X_{i3}, X_{i4}X_{i6})$] with the smallest AIC value (220.185) (see Table 5).

Backward Stepwise method selected a model [variable in the model (X_{i1}, X_{i3}, X_{i4}, X_{i6}, X_{i7},

 $X_{i2}X_{i3}$, $X_{i3}X_{i5}$, $X_{i4}X_{i6}$] with larger AIC value (220.568) compared to the AIC value (220.185) of the selected model with minimum AIC value among models selected using Best subset method (see Table 3).

Forward Stepwise method selected a set of variables resulting in a model [variable in the model $(X_{i1}, X_{i3}, X_{i5}, X_{i1}X_{i4}, X_{i2}X_{i3}, X_{i3}X_{i6})$

$$X_{i4}X_{i7}, X_{i5}X_{i6})$$
]

with the largest AIC value (226.512) and the smallest AROC value (0.862) (see Table 3).

The model with the smallest AIC value among models selected using Best Subset method have AIC value smaller than the model fitted with the variables selected using Forward Stepwise selection method and Backward Stepwise selection method. Table 3 shows the AIC values and the AROC values for the models selected by each selection technique.

From Table3, the method with the best performance for fitting a model containing two factor interactions is the Best subset method which fitted a model [variables in the model $(X_{i1}, X_{i3}, X_{i4}, X_{i5}, X_{i6}, X_{i7}, X_{i2}X_{i3}, X_{i4}X_{i6})$] with the minimum AIC value (220.185) and maximum AROC value (0.871) when compared to that of Backward Stepwise selection method and Forward Stepwise method.

5. CONCLUSION

In conclusion,

- The three methods of variable selection considered in this study have same performance in selecting variables for fitting a model without interaction.
- Best subset method outperformed Backward Stepwise method and Forward stepwise method in selecting variables for fitting a model with two factor interaction.
- All the models selected have excellent discrimination ability.

REFERENCE

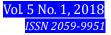
- 1. Agresti, A. (2002). Categorical data analysis, Second Edition. John Wiley & Sons, Inc., Hoboken, New Jersey.
- 2. Berry, Michael J. A. (1997). Data Mining Techniques for Marketing, Sales and Customer Support. Wiley. p. 10.
- 3. Biondo S., Ramos E., Deiros M. et al. (2000). Prognostic factors for mortality in left colonic peritonitis: a new scoring system, J. Am. Coll. Surg. Vol. 191, № 6. p. 635-642.
- Boyd, C. R., Tolson, M. A., and Copes, W. S. (1987). Evaluating trauma care: The TRISS method. Trauma Score and the Injury Severity Score. The Journal of trauma. 27 (4): 370–378. PMID 3106646. doi:10.1097/00005373-198704000-00005.
- Bursac, Z., Gauss, C. H., Williams, D. K., and Hosmer, D. W. (2008). Purposeful selection of variables in logistic regression. Source Code for Biology and Medicine 2008, 3:17 doi:10.1186/1751-0473-3-17
- 6. Cox, D. R. & Snell, E. J. (1989). The analysis of binary data, Second Edition, London: Chapman and Hall.
- 7. Dabbour, E. (2012). Using Logistic Regression to identify risk factors causing rollover

collisions. International Journal for Traffic and Transport Engineering, 2(4): 372 – 379, DOI: p:/http/dx.doi.org/10.7708/ijtte.2012.2(4).07

- Dobson, A. J. (2002). An Introduction to Generalized Linear Models. 2nded. Boca Raton, FL:Chapman& Hall/CRCSchafer, J. L. (2001). Lecture Notes for Statistics 544: Categorical Data Analysis I, Fall 2001. Penn State Univ. http://www.stat.psu.edu/_jls/
- 9. Draper, N. R. and Smith, H. (1981). Applied Regression Analysis. 2nded. New York: John Wiley.
- 10. Hosmer, D. W., Jovanovic, B and Lemeshow, S. (1989). Best Subsets Logistic Regression. Biometrics, Vol. 45, No. 4 (Dec., 1989), pp. 1265-1270.
- 11. Hosmer, D. W. &Lemeshow, S. (2000). Applied Logistic Regression, Second edition. John Wiley & Sons, Inc. Hoboken, New Jersey.
- 12. Kleinbaum, D. G. & Klein, M. (2010). Logistic regression a self learning text (3rd ed.). Springer New York Dordrecht Heidelberg London.
- 13. Kologlu M., Elker D., Altun H., Sayek I. (2001). Validation of MPI and OIA II in two different groups of patients with secondary peritonitis, Hepato-Gastroenterology. Vol. 48, № 37. pp. 147–151.
- 14. Lee K. I. and Koval J. J. (1997). Determination of the best significance level in forward stepwise logistic regression, Communications in Statistics Simulation and Computation, 26:2, 559-575.
- Lihui, Z., Chen, Y. &Schaffner, D. (2001), Comparison of logistic regression and linear regression in modeling percentage data. Appl. Environ. Microbiol, 67(5):2129. DOI: 10.1128/AEM.67.5.2129-2135,2001.
- Marshall J. C., Cook D. J., Christou N. V. et al. (1995). Multiple Organ Dysfunction Score: A reliable descriptor of a complex clinical outcome, Crit. Care Med. Vol. 23. – pp. 1638–1652.
- 17. Nagelkerke, N. J. (1991). A note on general definition of coefficient of determination. Biometrika, 78: 691-692.
- 18. Nelder, J. A. and Wedderburn, R. W. M. (1972). Generalized Linear Models, Journal of the Royal Statistical Society, Series A, vol. 135, pp.370-384.
- Oluwadiya, K. S., Kolawole, I. K., Adegbehingbe, O. O., Olasinde, A. A., Olaide A., &Uwaezuoke S. C. (2008). Motorcycle crash characteristics in Nigeria: Implication for control. Accident Analysis and Prevention 41 (2009) 294-298.
- Pacheco, J., Casado, S., Núñez, L. (2009). A variable selection method based on Tabu search for logistic regression models. European Journal of Operational Research 199 (2009) 506–511
- 21. Palei, S. K. Das, S. K. (2009). Logistic regression model for prediction of roof fall risks in bord and pillar workings in coal mines: An approach. Safety Science. 47: 88–96. doi:10.1016/j.ssci.2008.01.002.Quin, G. &Keough. M. (2001). Generalized Linear Models and logistic regression. Experimental Design and Data Analysis for Biologists, Cambridge University Press.
- 22. Pohar, M. Blas, M. Turk, S. (2004). Comparison of Logistic Regression and Linear Discriminant Analysis: A Simulation Study. (PDF). Metodološkizvezki. 1 (1).
- 23. Sarkar, S.K. and Habshah, M. (2009). Optimization Techniques for Variable Selection in Binary Logistic Regression Model Applied to Desire for Children Data. Journal of Mathematics and Statistics 5 (4): 387-394, 2009. ISSN 1549-3644.
- 24. Shuhua, H. (2007). Akaike Information Criterion. Center for research in scientific

computation, North Carolina State University, Raleigh, NC.

- 25. Strano, M., Colosimo, B.M. (2006). Logistic regression analysis for experimental determination of forming limit diagrams. International Journal of Machine Tools and Manufacture. 46 (6): 673–682. doi:10.1016/j.ijmachtools.2005.07.005.
- Wang, D., Zhang, W. and Bakhai, A. (2004). Comparison of Bayesian model averaging and stepwise methods for model selection in logistic regression. Statist. Med. 2004; 23:3451–3467, Wiley InterScience (www.interscience.wiley.com). DOI: 10.1002/sim.1930



APPENDIX

Table 4. Summary of all selected models (model without interaction) using best subset method Table 5. Summary of all selected models (model with interaction) using best subset method

k-1	Variable in the model	RSS		-2LL	AIC
1	X _{i1}	287.434	L_q 12.7613	261.788	265.788
2	X _{i1} , X _{i5}	280.269	7.7136	253.913	259.913
3	$X_{i1}, X_{i5}, X_{i2}X_{i4}$	276.164	5.6748	240.251	248.251
4	$X_{i1}, X_{i3}, X_{i2}X_{i3}, X_{i4}X_{i7}$	271.112	2.7049	229.406	239.406
5	$X_{i1}, X_{i3}, X_{i2}X_{i3}, X_{i4}X_{i7}, X_{i5}X_{i6}$	268.323	1.9618	219.652	231.652
6	$\begin{array}{c} X_{i1}, X_{i3}, X_{i5}, X_{i1}X_{i4}, \\ X_{i2}X_{i3}, X_{i3}X_{i6} \end{array}$	264.976	0.6691	213.569	227.569
7	$X_{i1}, X_{i3}, X_{i4}, X_{i6}, X_{i7},$	262.185	-0.0770	204.704	220.704
8	$\frac{X_{i2}X_{i3}, X_{i4}X_{i6}}{X_{i1}, X_{i3}, X_{i4}, X_{i5}, X_{i6}, X_{i7},}$ $\frac{X_{i2}X_{i3}, X_{i4}X_{i6}}{X_{i2}X_{i3}, X_{i4}X_{i6}}$	260.019	-0.2074	202.185	220.185
9	$X_{i1}, X_{i3}, X_{i4}, X_{i5}, X_{i6}, X_{i7}, X_{i2}X_{i3}, X_{i4}X_{i6}, X_{i3}X_{i6}$	258.287	0.0887	200.393	220.393
10	$\begin{array}{c} X_{i1}, X_{i3}, X_{i4}, X_{i5}, X_{i6}, X_{i7}, \\ X_{i2}X_{i3}, X_{i4}X_{i6}, X_{i3}X_{i6}, X_{i1}X_{i7} \end{array}$	257.148	0.9684	199.213	221.213
11	$\begin{array}{c} X_{i1}, \ X_{i3}, \ X_{i4}, \ X_{i5}, \ X_{i6}, \ X_{i7}, \\ X_{i2}X_{i3}, \ X_{i4}X_{i6}, \ X_{i3}X_{i6}, \ X_{i1}X_{i6}, \ X_{i1}X_{i7} \end{array}$	255.605	1.4503	197.522	221.522
12	$\begin{array}{c} X_{i1}, X_{i3}, X_{i4}, X_{i5}, X_{i6}, X_{i7}, \\ X_{i2}X_{i3}, X_{i4}X_{i6}, X_{i3}X_{i6}, \ X_{i1}X_{i6}, \ X_{i1}X_{i7}, \ X_{i4}X_{i7} \end{array}$	255.063	2.9166	196.951	222.951
13	$\begin{array}{c} X_{i1}, X_{i3}, X_{i4}, X_{i6}, X_{i7}, \\ X_{i1}X_{i6}, X_{i1}X_{i7}, X_{i2}X_{i3}, X_{i2}X_{i6}, X_{i2}X_{i7}, X_{i3}X_{i5}, \\ X_{i3}X_{i6}, X_{i4}X_{i6} \end{array}$	254.576	4.4381	196.453	224.453
14	$\begin{array}{c} X_{i1}, X_{i3}, X_{i4}, X_{i5}, X_{i6}, X_{i7}, \\ X_{i1}X_{i6}, X_{i1}X_{i7}, X_{i2}X_{i3}, X_{i2}X_{i6}, \ X_{i2}X_{i7}, \ X_{i3}X_{i6}, \\ X_{i4}X_{i6}, \ X_{i4}X_{i7} \end{array}$	254.015	5.8857	195.241	225.241
15	$\begin{array}{c} X_{i1}, X_{i3}, X_{i4}, X_{i6}, X_{i7}, \\ X_{i1}X_{i6}, X_{i1}X_{i7}, X_{i2}X_{i3}, X_{i2}X_{i6}, \ X_{i2}X_{i7}, \ X_{i3}X_{i5}, X_{i3}X_{i6}, \\ X_{i4}X_{i5}, X_{i4}X_{i6}, \ X_{i4}X_{i7} \end{array}$	253.519	7.3980	194.796	226.796
16	$\begin{array}{c} X_{i1}, X_{i3}, X_{i4}, X_{i6}, X_{i7}, \\ X_{i1}X_{i6}, X_{i1}X_{i7}, X_{i2}X_{i3}, X_{i2}X_{i6}, \\ X_{i2}X_{i7}, X_{i3}X_{i4}, X_{i3}X_{i5}, X_{i3}X_{i6}, X_{i4}X_{i5}, X_{i4}X_{i6}, X_{i4}X_{i7} \end{array}$	253.264	9.1474	194.486	228.486
17	$\begin{array}{c} \begin{array}{c} X_{i1}, X_{i2}, X_{i3}, X_{i4}, X_{i6}, X_{i7}, \\ X_{i1}X_{i2}, X_{i1}X_{i6}, X_{i1}X_{i7}, X_{i2}X_{i3}, X_{i2}X_{i6}, X_{i2}X_{i7}, \\ X_{i3}X_{i5}, X_{i3}X_{i6}, X_{i4}X_{i5}, X_{i4}X_{i6}, X_{i4}X_{i7} \end{array}$	252.854	10.7437	194.266	230.266
18	$\begin{array}{c} \begin{array}{c} X_{i1}, X_{i2}, X_{i3}, X_{i4}, X_{i6}, X_{i7}, \\ X_{i1}X_{i2}, X_{i1}X_{i6}, X_{i1}X_{i7}, X_{i2}X_{i3}, X_{i2}X_{i6}, \\ X_{i2}X_{i7}, X_{i3}X_{i4}, X_{i3}X_{i5}, X_{i3}X_{i6}, X_{i4}X_{i5}, X_{i4}X_{i6}, X_{i4}X_{i7} \end{array}$	252.599	12.4923	193.936	231.936

European Journal of Mathematics and Computer Science

Vol. 5 No. 1, 2018 ISSN 2059-9951

19	$X_{i1}, X_{i2}, X_{i3}, X_{i4}, X_{i6}, X_{i7},$	252.455	14.3511	193.743	233.743
	$X_{i1}X_{i2}, X_{i1}X_{i6}, X_{i1}X_{i7}, X_{i2}X_{i3}, X_{i2}X_{i6}, X_{i2}X_{i7},$				
	$X_{i3}X_{i4}, X_{i3}X_{i7},$				
	$X_{i3}X_{i5}, X_{i3}X_{i6}, X_{i4}X_{i5}, X_{i4}X_{i6}, X_{i4}X_{i7}$				
20	$X_{i1}, X_{i2}, X_{i3}, X_{i4}, X_{i6}, X_{i7},$	252.323	16.2217	193.692	235.692
	$X_{i1}X_{i2}, X_{i1}X_{i3}, X_{i1}X_{i6}, X_{i3}X_{i4}, X_{i1}X_{i7}, X_{i3}X_{i7}, X_{i2}X_{i3}, X_{i7}$	i			
	, $X_{i2}X_{i7}$, $X_{i3}X_{i5}$, $X_{i3}X_{i6}$, $X_{i4}X_{i5}$, $X_{i4}X_{i6}$, $X_{i4}X_{i7}$				
21	$X_{i1}, X_{i2}, X_{i3}, X_{i4}, X_{i5}, X_{i6}, X_{i7},$	252.198	18.0990	193.659	237.659
	$X_{i1}X_{i2}, X_{i1}X_{i3}, X_{i1}X_{i6}, X_{i3}X_{i4}, X_{i1}X_{i7}, X_{i3}X_{i7}, X_{i2}X_{i3}, X_{i7}$	i			
	, $X_{i2}X_{i7}$, $X_{i3}X_{i5}$, $X_{i3}X_{i6}$, $X_{i4}X_{i5}$, $X_{i4}X_{i6}$, $X_{i4}X_{i7}$				
22	$X_{i1}, X_{i2}, X_{i3}, X_{i4}, X_{i5}, X_{i6}, X_{i7},$	252.114	20.0161	193.612	239.612
	$X_{i1}X_{i2}, X_{i1}X_{i3}, X_{i1}X_{i6}, X_{i3}X_{i4}, X_{i1}X_{i7}, X_{i3}X_{i7}, X_{i2}X_{i3}, X_{i7}$	i			
	, $X_{i2}X_{i7}$, $X_{i3}X_{i5}$, $X_{i3}X_{i6}$, $X_{i4}X_{i5}$, $X_{i4}X_{i6}$, $X_{i4}X_{i7}$				
	X _{i5} X _{i6}				
23	$X_{i1}, X_{i2}, X_{i3}, X_{i4}, X_{i5}, X_{i6}, X_{i7},$	252.100	22.002	193.582	241.582
	$X_{i1}X_{i2}, X_{i1}X_{i3}, X_{i1}X_{i6}, X_{i3}X_{i4}, X_{i1}X_{i7}, X_{i3}X_{i7}, X_{i2}X_{i3}, X_{i7}$				
	, $X_{i2}X_{i7}$, $X_{i3}X_{i5}$, $X_{i3}X_{i6}$, $X_{i4}X_{i5}$, $X_{i4}X_{i6}$,				
24	X _{i4} X _{i7} , X _{i5} X _{i6}	252.009	24.0002	102 5 (0	242.56
24	$X_{i1}, X_{i2}, X_{i3}, X_{i4}, X_{i5}, X_{i6}, X_{i7}, X_{i1}X_{i2}$	252.098	24.0003	193.560	243.56
	$X_{i1}X_{i3}, X_{i1}X_{i4}, X_{i1}X_{i6}, X_{i3}X_{i4}, X_{i1}X_{i7}, X_{i3}X_{i7}, X_{i2}X_{i3}, X_{i7}, X_{i2}X_{i3}, X_{i7}, X_{i7}X_{i7}, X_{i7}X_{i7}$				
	$, X_{i2}X_{i7}, X_{i3}X_{i5}, X_{i3}X_{i6}, X_{i2}X_{i2}, X_{i3}X_{i6}, X_{i3}X_{i6}, X_{i3}X_{i6}, X_{i3}X_{i6}, X_{i3}X_{i6}, X_{i3}X_{i6}, X_{i3}X_{i6}, X_{i3}X_{i6}, X_{i3}X_{i6}, X_{i6}X_{i6}, X_{i6}X_{i6}X_{i6}, X_{i6}X_{i6}X_{i6}, X_{i6}X_{i6}X_{i6}, X_{i6}X_{i6}X_{i6}X_{i6}, X_{i6}$				
25	$\frac{X_{i4}X_{i6}, X_{i4}X_{i5}, X_{i4}X_{i7}, X_{i5}X_{i6}}{X_{i1}, X_{i2}X_{i3}, X_{i4}, X_{i5}, X_{i6}, X_{i7}, X_{i1}X_{i2}, X_{i1}X_{i3},}$	252.098	26	193.549	245.549
25	$X_{i1}X_{i2}X_{i3}X_{i4}X_{i5}X_{i6}X_{i6}X_{i7}X_{i1}X_{i2}X_{i1}X_{i3}, X_{i1}X_{i5}X_{i1}X_{i6}, X_{i3}X_{i4}, X_{i1}X_{i7}, X_{i3}X_{i7}, X_{i2}X_{i3}, X_{i7}, X_{i7}X_{i7}, X_{i$		20	175.5 17	213.319
	$X_{i2}X_{i7}, X_{i3}X_{i5}, X_{i3}X_{i6},$				
	$X_{i4}X_{i6}, X_{i4}X_{i5}, X_{i4}X_{i7}, X_{i5}X_{i6}$				
			21		

k-1	Variable in the model	SSE	Cq	-2LL	AIC
1	<i>X</i> _{<i>i</i>1}	293.9354	25.749	261.788	265.788
2	X_{i1}, X_{i6}	288.0576	21.834	251.658	257.658
3	X_{i1}, X_{i2}, X_{i6}	279.6204	15.345	241.144	249.144
4	$X_{i1}, X_{i2}, X_{i3}, X_{i6}$	274.6962	12.391	234.867	244.867
5	$X_{i1}, X_{i2}, X_{i3}, X_{i5}, X_{i6}$	270.1564	9.823	230.078	242.078
6	$X_{i1}, X_{i2}, X_{i3}, X_{i4}, X_{i5}, X_{i6}$	265.3246	6.961	224.898	238.898
7	$X_{i1}, X_{i2}, X_{i3}, X_{i4}, X_{i5}, X_{i6}, X_{i7}$	262.8789	6.500	222.493	238.493
8	$X_{i1}, X_{i2}, X_{i3}, X_{i4}, X_{i5}, X_{i6}, X_{i7}, X_{i9}$	262.6844	8.305	222.296	240.296
9	$X_{i1}, X_{i2}, X_{i3}, X_{i4}, X_{i5}, X_{i6}, X_{i7}, X_{i8}, X_{i9}$	262.3816	10	221.985	241.985