# Selecting the Method to Overcome Partial and Full Multicollinearity in Binary Logistic Model

**N. Herawati, K. Nisa***, Nusyirwan**

Department of Mathematics, Faculty of Mathematics and Natural Sciences, University of Lampung, Bandar Lampung, Indonesia

**Abstract**  The aim of our study is to select the best method for overcoming partial and full multicollinearity in binary logistic model for different sample sizes. Logistic ridge regression (LRR), least absolute shrinkage and selection operator (LASSO) and principal component logistic regression (PCLR) compared to maximum likelihood estimator (MLE) using simulation data with different level of multicollinearity and different sample sizes ($n$=20, 50, 100, 200). The best method is chosen based on mean square error (MSE) values and the best model is characterized by AIC value. The results show that LRR, LASSO and PCLR surpass MLE in overcoming partial and full multicollinearity in binary logistic model. PCLR exceeds LRR and LASSO when full multicollinearity occurs in binary logistic model but LASSO and LRR are better used when partial multicollinearity exists in the model.

**Keywords**  Binary logistic model, Multicollinearity, LRR, LASSO, PCLR

## 1. Introduction

Consider that the model has the form $y_i = x_i'\boldsymbol{\beta} + \varepsilon_i$ where $x_i' = [1\ x_{i1}\ x_{2i} \dots x_{ik}]$ , $\boldsymbol{\beta}' = [\beta_0\ \beta_1\ \dots\ \beta_k]$ and dependent variables $y_i$ has value either 0 or 1. Estimating parameters in this model where the response variable is binary or multinomial is not appropriate when using the linear regression model estimation method. The linear regression model is based on a ratio scale measurement [1,2,3]. In this case logistic regression model is more suitable.

Logistic regression model is based on a logistic function to model binary dependent variables. It is a classification of individuals in different groups. Unlike multiple regression, logistic regression is much more flexible in terms of basic assumptions to be met. Logistic regression model as one of nonlinear regression model does not require liner relationship between independent and dependent variables, assumption of normal distribution and homoscedasticity in the error terms. Despite all the flexibility, the logistic regression model still requires no correlation between independent variables [4,5]. When there is a correlation between the independent variable, logistic model becomes unstable. This can cause errors in the interpretation of the relationship between the dependent and each independent variable in terms of odds ratios [6,7].

There are several methods for overcoming the problem of multicollinearity in the logistic models and have been examined by several researchers [8,9,10,11,12]. In this research, a selection of LRR, LASSO and PCLR methods was conducted in logistic model with binary responses and a set of continuous predictor variables. Each method was compared using simulation data that contains partial and full multicollinearity with different sample sizes. The best method was examined based on the minimum value of MSE and the best model is characterized by AIC value.

## 2. Logistic Regression Model

Suppose the response variable of regression application of interest has two possible outcomes or $Y_i$ is a Bernoulli random variable with the probability distribution $P(y_i = 0) = 1 - \pi_i$ and $P(y_i = 1) = \pi_i$ The probability function for each observation is $f(y_i) = \pi_i^{y_i}(1 - \pi_i)^{1-y_i}$ , i=1,2,…,$n$ [4,5,6,7]. The multiple logistic regression model of the response variable $Y = \pi(X) + \varepsilon$, with $\pi$ (**X**) is an $n$ x 1 vector and

$$\pi_i(x) = E[Y|X = x_i] = P[Y_i = 1] = \frac{\exp(\beta_0 + x_i\beta)}{1 + \exp(\beta_0 + x_i\beta)} \quad (1)$$

where $\boldsymbol{\beta}$ is a $k$ x 1 vector of estimated parameters. The logit function of $\pi_i(x)$ is $logit[\pi_i(x)] = \ln\left(\frac{\pi_i(x)}{1-\pi_i(x)}\right)$ or in linear form can be written as [3,13]:

$$L(\boldsymbol{X}) = \boldsymbol{X}\boldsymbol{\beta} \quad (2)$$

The parameters were estimated by maximizing likelihood function $L(\boldsymbol{\beta}) = \prod_{i=1}^{N} \pi_i(x_i)^{y_i} (1 - \pi_i)^{n_i - y_i}$ . When the log-likelihood is differentiated with respect to $\boldsymbol{\beta}$ equal to zero, we get

$$\widehat{\boldsymbol{\beta}}_{ML} = \left(\mathbf{X}'\widehat{\mathbf{W}}\mathbf{X}\right)^{-1}\mathbf{X}'\widehat{\mathbf{W}}\mathbf{Z} \qquad (3)$$

where $\mathbf{Z}$ is a nx1 column vector with elements $z_i = \text{logit}(\hat{\pi}_i) + \frac{y_i\hat{\pi}_i}{\hat{\pi}_i(1-\hat{\pi}_i)}$ and $\widehat{\mathbf{W}} = \text{diag}[\hat{\pi}_i(1 - \hat{\pi}_i)]$ [7].

## 2.1. Logistic Ridge Regression (LRR)

When multicollinearity exist between independent variables in the logistic model, the matrix $\mathbf{X'WX}$ is (near) singular. Using maximum likelihood method to estimate the parameters in the model is not suitable because we cannot get the inversion of the matrix. As a result, the estimation of the parameters in the logistic model using maximum likelihood method is being unstable and cannot be uniquely estimated. In this situation, the ridge regression method can be applied by using a penalty to the diagonal matrix of $\mathbf{X'WX}$ to stabilize the coefficients estimates [14,15,16]. Although this method will produce a bias in the coefficient estimates of the model, it provides a lower variance of the coefficient estimates than the unpenalized model. Ridge likelihood estimator of the logistic model is done by maximize the ridge penalized loglikelihood [17,18,19,20,21]:

$$L^{LRR}(\boldsymbol{\beta}, \lambda) = L(\boldsymbol{\beta}) - \lambda\boldsymbol{\beta}'\boldsymbol{\beta}$$
$$= \sum_{i=1}^{n}[y_i \log(\pi_i) + (1 - y_i)\log(1 + \pi_i)] - \lambda\boldsymbol{\beta}'\boldsymbol{\beta} \quad (4)$$

where the ridge penalty is the second summand (the sum of the square of the elements of $\boldsymbol{\beta}$) with $\lambda$ as penalty parameter. Because the value of the $\boldsymbol{\beta}$ equation is not linear, Newton-Raphson method is used to solve it. The solution uses and follows the iterative weighted least square algorithm to obtained the $\boldsymbol{\beta}$ estimates. The logistic ridge regression (LRR) model following [17] is:

$$\widehat{\boldsymbol{\beta}}_{LRR} = (\mathbf{X}'\widehat{\mathbf{W}}\mathbf{X} + k\mathbf{I})^{-1}\mathbf{X}'\widehat{\mathbf{W}}\mathbf{X}\,\widehat{\boldsymbol{\beta}}_{LMLE} \qquad (5)$$

with $k = \frac{1}{\widehat{\beta}^2}$ 5 and $\widehat{W}$ as in equation (3) [17].

## 2.2. Least Absolute Shrinkage and Selection Operator (LASSO)

LASSO method can be used to overcome problems in multicollinearity [22]. LASSO shrinks the coefficient parameter β which correlates to exactly zero or close to zero [23]. Lagrangian constraint (L1-norm) can be combined in a log-likelihood parameter estimation in logistic regression [24,25]. The estimation of parameters in LASSO in combining log-likelihood and Lagrangian constraints produces:

$$l(\beta) = -\sum_{i=1}^{n}[(1 - y_i)\beta'x_i' + \ln(1 + exp(-\beta'x_i'))]$$
$$- \lambda\sum_{k=1}^{p}|\beta_j| \qquad (6)$$

So we get a logistic regression parameter estimates with LASSO:

$$\widehat{\boldsymbol{\beta}}_{\lambda}^{LASSO} = argmax\left\{l(\beta) - \lambda\sum_{j=1}^{p}|\beta_j|\right\} \qquad (7)$$

$\lambda > 0$ is the tuning parameter that control the strength of penalty in the LASSO method and can be obtained using generalized cross validation [22].

## 2.3. Principal Component Logistic Regression (PCLR)

In linear regression analysis, principal component regression (PCR) is one of the methods that has been confirmed to be able to overcome the problem of multicollinearity [1,11,26,27]. PCR aims to simplify the observed variables by reducing the dimensions, where the chosen principal components must maintain as much diversity as possible. This is done by eliminating the correlation between the independent variables through the transformation of the original independent variable into a new variable that does not correlate at all. In terms of the principal component (PC) of the predictor variables, the logit transformation (2) can be written as principal component regression form as:

$$L(\boldsymbol{X}) = \boldsymbol{X}\boldsymbol{\beta} = \boldsymbol{Z}\boldsymbol{V}'\boldsymbol{\beta} = \boldsymbol{Z}\boldsymbol{\gamma} \qquad (8)$$

where $\boldsymbol{Z} = \boldsymbol{XV}$ as an $n$ x $k$ matrix whose columns are the PCs of $\boldsymbol{X}$ with $\boldsymbol{V}$ is a $k$ x $k$ matrix whose columns are the eigenvectors of the of the matrix $\boldsymbol{X'X}$ denoted by $\boldsymbol{v_j}$ with $j = 1, \dots, p$. It is obvious that $\boldsymbol{\gamma} = \boldsymbol{V^T\beta}$ can be estimated by

$$\widehat{\boldsymbol{\gamma}} = \boldsymbol{V}'\widehat{\boldsymbol{\beta}} \qquad (9)$$

The prediction equation of MLE is $\hat{Y}(x) = \hat{\pi}(x)$ with

$$\hat{\pi}(x) =$$
$$\exp(\widehat{\beta_0} + \sum_{j=1}^{k} z_j(x)\widehat{\gamma_j})/(1 + \exp(\widehat{\beta_0} + \sum_{j=1}^{k} z_j(x)\widehat{\gamma_j}))$$

where $z_j(\boldsymbol{x}) = \boldsymbol{x}'\boldsymbol{v_j}$ is the j-th PC value for a point $\boldsymbol{x}$. The logit model (8) can be expressed as

$$L(\boldsymbol{X}) = \beta_0 1 + \boldsymbol{Z_s}\boldsymbol{\gamma^s} + \boldsymbol{Z_r}\boldsymbol{\gamma^r} \qquad (10)$$

The principal component logistic regression (PCLR) model in terms of the first PC is $Y = \pi^s(X) + \varepsilon^s$ and the logit transformation, $L^s(X)$ has components $L_i{}^s = \ln([\pi_i{}^s/(1 - \pi_i{}^s)]$ is defined as $L^s(X) = \boldsymbol{W_s}\boldsymbol{\gamma^s}$. The parameter estimate of the PCLR [9] is:

$$\widehat{\boldsymbol{\beta}}^s = \boldsymbol{V_s}\widehat{\boldsymbol{\gamma}}^s \qquad (11)$$

where the subscript (s) indicates number of PCs were used in the PCLR model.

This method was introduced by Aguilera et al. [10] for solving the problem of high-dimensional multicollinear data in logistic regression of binary response variable and a set of continuous predictor variables. They showed that the PCLR model provides better estimation of model parameters compared to partial least square (PLS) logistic regression.

# 3. Methods

Illustration of the performance of LRR, LASSO PCLR methods used in this study was carried out using a simulation study to show how these methods can improve the estimation of parameters of the binary logistic model contains partial and full multicollinearity using R. Six independent variables ($p$=6) were generated using the formula $X_p = \sqrt{(1 - \rho^2)z_{ij}} + \rho z_{i(p+1)}$; $i = 1,2,\dots,n$ and $j = 1,2,\dots,6$ with $z_{ij} \sim N(0,1)$ and $\rho = 0,99$. The dependent variable Y is

generated by the binary logistic regression probability
$P(y_i = 1) = \pi_i(x) = \frac{exp(\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \cdots + \beta_p x_p)}{1 + exp(\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \cdots + \beta_p x_p)}$ with $\beta_0 = 0$ and $\beta_1 = \beta_2 = \cdots = \beta_p = 1$ respectively. Partial and full multicollinearity between independent variables were applied in the model with different sample sizes (n=20, 40, 60, 100, 200) and replicated 1000 times. Multicollinearity of the independent variables is measured by $VIF = \frac{1}{(1-R_j^2)}$ with $R_j^2$ is the coefficient of determination. The best method in estimating the parameters is evaluated using MSE with formula:

$$MSE(\widehat{\beta}) = \frac{1}{n} \sum_{l=1}^{m} \|\widehat{\beta}^{(l)} - \beta\|^2,$$

and the best model is characterized by $AIC_c = 2k - 2\ln(\hat{L})$

where $\hat{L} = p(\hat{\theta}, M)$, $\hat{\theta}$ is the value that maximize the likelihood function, n is the number of recorded measurements and k is the number of parameters estimated [28,29].

# 4. Results and Discussion

The partial and full multicollinearity of the independent variables applied in this study is shown in Table 1. First, partial multicollinearity in which correlation only applies between $X_1$ and $X_2$; second, partial multicollinearity where correlation occurs only between $X_1$, $X_2$, and $X_3$; third, full multicollinearity in all independent variables. This condition is applied to all sample sizes that are being studied.

**Table 1.** Multicollinearity in independent variables for all sample sizes studied

| Independent Variables | VIF | | |
|---|---|---|---|
| | Partial multicollinearity in $X_1, X_2$ | Partial multicollinearity in $X_1, X_2, X_3, X_4$ | Full multi-collinearity |
| $X_1$ | 20.35 | 20.35 | 20.35 |
| $X_2$ | 20.68 | 20.68 | 20.68 |
| $X_3$ | 1.62 | 23.17 | 23.17 |
| $X_4$ | 1.65 | 30.33 | 30.33 |
| $X_5$ | 1.48 | 1.28 | 31.28 |
| $X_6$ | 1.58 | 1.14 | 41.14 |

From Table 1, we can see that the VIF values are greater than 10 for all given cases in this study. It means that the independent variables seem to correlate to each other and indicate there is a multicollinearity between these variables. To select the method that is considered the best in overcoming the multicollinearity problems in this study, MSE value is used. The best method is determined from an MSE value that close to zero. The MSE values of MLE, LRR, LASSO, PCLR for partial and full multicollinearity in the model at different sample sizes are shown in Table 2.

From Table 2 where partial multicollinearity in $X_1$ and $X_2$ occurs in the model, MLE gives MSE =35608.77 for $n$=20, MSE= 1112.951, for $n$=50, MSE= 0.0820, for $n$=100, respectively. These values are far above the MSE of LRR, LASSO and PCLR which give MSE = 0.0857, 0.0170, 0.0207 for $n$=20, MSE =0.0506, 0.0085, 0.0175 for $n$=50, and MSE =0.0385, 0.0026, 0.0126 for $n$=100, respectively. Similar results are obtained when partial multicollinearity exists in $X_1$, $X_1$, $X_3$, $X_4$ and when the model contains full multicollinearity. It is obvious that MLE is unable to overcome partial and full multicollinearity between independent variables very well in logistic regression with binary responses when sample sizes are small enough. In a larger sample size ($n$=200) the MSE of MLE seems to decrease significantly, but its value still above the MSE of LRR, LASSO and PCLR. This suggests that MLE should not be used in estimating the parameters of logistic models with binary responses that have partial and full multicollinearity

on small and large sample sizes.

**Table 2.** MSE of MLE, LRR, LASSO, PCLR

| Multicolli-nearity in | | MSE | | | |
|---|---|---|---|---|---|
| | | MLE | LRR | LASSO | PCLR |
| $X_1, X_2$ | n=20 | 35608.77 | 0.0857 | 0.0170 | 0.0207 |
| | n=50 | 1112.951 | 0.0506 | 0.0085 | 0.0175 |
| | n=100 | 0.0820 | 0.0385 | 0.0026 | 0.0126 |
| | n=200 | 0.0064 | 0.0054 | 0.0017 | 0.0057 |
| | | | | | |
| $X_1, X_2, X_3, X_4$ | n=20 | 1.03E+27 | 0.0572 | 0.0495 | 0.0216 |
| | n=50 | 2.53E+25 | 0.0536 | 0.0291 | 0.0194 |
| | n=100 | 47.6825 | 0.0435 | 0.0078 | 0.0161 |
| | n=200 | 0.0162 | 0.0051 | 0.0052 | 0.0084 |
| | | | | | |
| $X_1, X_2, X_3, X_4, X_5, X_6$ | n=20 | 3.77E+26 | 0.0668 | 0.1452 | 0.0182 |
| | n=50 | 2.66E+25 | 0.0528 | 0.0664 | 0.0089 |
| | n=100 | 16867.23 | 0.0114 | 0.0149 | 0.0049 |
| | n=200 | 0.0336 | 0.0062 | 0.0092 | 0.0006 |

To provide clearer results from the LRR, LASSO and PCLR methods in overcoming partial and full multicollinearity for all sample sizes ($n$=20, 40, 60, 100, 200), we compared the MSE of the three methods separately from MLE as shown in Figure 1-3.
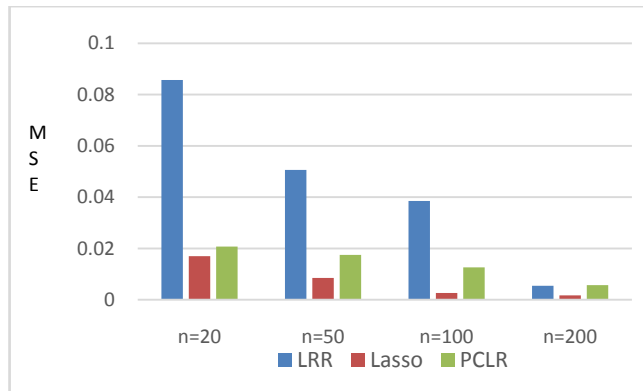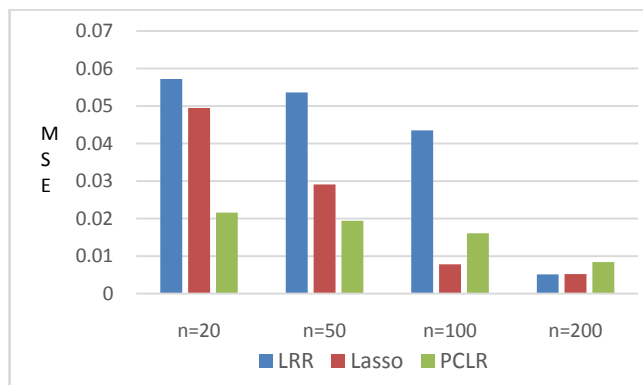
**Figure 1.** MSE of partial multicollinearity in $X_1, X_2$



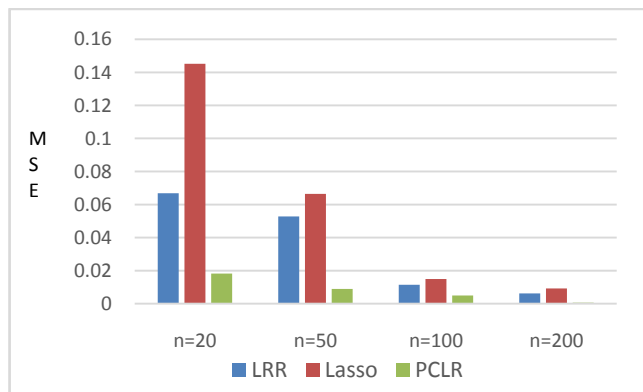**Figure 2.** MSE of Partial multicollinearity in $X_1, X_2, X_3, X_4$



**Figure 3.** MSE of full multicollinearity

Figures 1-3 shows MSE of LRR, LASSO and PCLR in conditions where the binary logistic model contains partial and full multicollinearity at different sample sizes (*n*=20, 50, 100, 200). It can be seen that MSE values of LRR, LASSO and PCLR vary depending on the number of correlated variables and sample sizes. If partial multicollinearity occurs between $X_1$ and $X_2$ in the model, LASSO gives MSE= 0.0170, 0.0085, 0.0025, and 0.0017 for *n*=20, 50, 100, and 200, respectively. These values are much lower than MSE value of LRR and PCLR. However, when partial multicollinearity occurs among $X_1, X_2, X_3$ and $X_4$ in the model, the results vary. For *n* = 20 and 50, LASSO and LRR gives lower MSE than PCLR. Conversely, for *n*= 100 and *n*=200, PCLR has the lowest MSE compared to LASSO and LRR. This suggests that when partial multicollinearity exists

in the binary logistic model, LASSO, LRR, and PCLR can be used depending on the amount of multicollinearity and sample sizes.

In situation where there is full multicollinearity between the independent variables in the model, we can see from Figure 3 that PCLR has the lowest MSE value than LRR and LASSO with MSE= 0.0182, 0.0089, 0.0049, and 0.0006 for *n*=20, 50, 100, 200, respectively. Obviously, LRR and LASSO appear unable to overcome full multicollinearity in logistic regression with binary responses. In this case PCLR exceeds LRR and LASSO for each sample size studied (*n*=20, 50. 100, 200). This indicates that PCLR is the best method for overcoming full multicollinearity in binary logistic model for all sample sizes being studied.

To determine the best model, we examine the AIC values of LRR, LASSO and PCLR as shown in Table 3.

**Table 3.** AIC of LRR, LASSO and PCLR

| Multicollinearity in | | Method | | |
|---|---|---|---|---|
| | | LRR | LASSO | PCLR |
| $X_1, X_2$ | n=20 | 23.7817 | 22.2384 | 23.3529 |
| | n=50 | 23.9976 | 23.4506 | 23.9840 |
| | n=100 | 24.8375 | 24.1487 | 24.6437 |
| | n=200 | 25.6043 | 25.2272 | 25.3984 |
| $X_1, X_2, X_3, X_4$ | n=20 | 21.9677 | 21.2499 | 19.6000 |
| | n=50 | 21.9843 | 21.7072 | 19.8603 |
| | n=100 | 22.8726 | 20.1575 | 22.5838 |
| | n=200 | 21.6062 | 23.1186 | 23.5904 |
| $X_1, X_2, X_3, X_4, X_5, X_6$ | n=20 | 20.4921 | 19.6963 | 19.5883 |
| | n=50 | 20.7754 | 20.2561 | 19.8593 |
| | n=100 | 21.8765 | 21.3905 | 20.1537 |
| | n=200 | 22.3716 | 21.9450 | 21.2761 |

Based on the AIC values from Table 3 it was found that the best model depends on the number of correlated variables in the binary logistics model and sample size. This supports the results obtained based on the MSE value.

# 5. Conclusions

We conclude from the results of this study that LRR, LASSO and PCLR surpass MLE in overcoming partial multicollinearity and full multicollinearity occur in binary logistic model. PCLR exceeds LRR and LASSO when full multicollinearity occurs in binary logistic model but LASSO and LRR are better used when partial multicollinearity exists in the model.

# ACKNOWLEDGEMENTS

assistance.

# REFERENCES

[1] D.C. Montgomery, E.A. Peck, and G.G. Vinning, Introduction to Linear Regression Analysis, New York: A Wiley Intersection Publication, 2012.

[2] N.R. Draper and H. Smith, Applied Regression Analysis, 3rd ed., New York: Wiley, 1998.

[3] R.H. Myers, Classical and Modern Regression With Application, Boston: PWSKENT publishing Company, 1990.

[4] M.H. Kutner, C.J. Nachtsheim, J. Neter, and W. Li. Applied Linear Statistical Models, Boston: McGraw-Hill, 2005.

[5] Midi, H., Sarkar, S.H., and Rana, S., 2010, Collinearity diagnostics of binary logistic regression model, Journal of Interdisciplinary Mathematics, 13(3): 253-267.

[6] D.W. Hosmer and S. Lemeshow, Applied Logistic Regression. New York: John Wiley & Sons, 2000.

[7] T.P. Ryan, Modern Regression Methods. New York: Wiley, 1997.

[8] Schaeffer, R.L., 1986, Alternative estimators in logistic regression when the data is collinear. Journal of Statistical Computation and Simulation, 25: 75-91.

[9] Aguilera, A.M., and Escabias, M., 2000, Principal component logistic regression. In: Bethlehem J.G., van der Heijden P.G.M. (eds) COMPSTAT. Physica, Heidelberg, 175-180.

[10] Aguilera, A.M., Escabias, M., and Valderrama, M.J., 2006, Using principal components for estimating logistic regression with-high-dimensional multicollinear data, Computational Statistics & Data Analysis, 50:1905-1924.

[11] T. Hastie, R. Thibsirani, and J. Friedman, The Elements of .Statistical Learning: Data Mining, Inference, and Prediction, 2nd ed., New York: Springer-Verlag, 2001.

[12] Toka, O., 2016, A Comparative Study on Regression Methods in the presence of Multicollinearity, Journal of Statisticians: Statistics and Actuarial Sciences, 2: 47-53.

[13] P. McCullagh, and J.A. Nelder, Generalized Linear Models, 2nd Ed., London: Chapman and Hall, 1989.

[14] Hoerl, A.E., 1962, Application of ridge analysis to regression problems, Chem. Eng. Prog., 58: 54-59.

[15] Hoerl, A.E. and Kennard, R.W., 1970, Ridge Regression: Biased Estimation for nonorthogonal problems, Technometrics, 12(1): 55-67.

[16] Dorugade, A.V. and Kashid, D.N., 2010, Alternative method for choosing ridge parameter for regression, Applied Mathematical Sciences, 4(9): 447-456.

[17] Schaffer, R.L., Roi, L.D., and Wolfe, R.A., 1984, A ridge logistic estimator, Communications in Statistics: Theory and Methods, 13: 99-113.

[18] Le Cessie, S. and van Houwelingen, J.C., 1992, Ridge estimators in logistic regression, Applied Statistics, 41(1), 191-201.

[19] Kibria, B.M.G., Shukur, G., Mansson, K., 2012, Performance of some logistic ridge regression estimators, Computational Economics, 40: 401-4014.

[20] Wu, J. and Asar, Y., 2016, On almost unbiased ridge logistic estimator for the logistic regression model, Hacettepe Journal of Mathematics and Statistic, 43(3): 989-998.

[21] Duffy, D.E. and Santner, T.J., 1989, On the small properties of norm − restricted maximum likelihood estimators for logistic regression models, Communications in Statistics - Theory and Methods, 18: 959-980

[22] Tibshirani, R., 1996, Regression Shrinkage and Selection via LASSO, Journal of the Royal Statistical Society, 58(1): 267-288.

[23] G. James, D. Witten, T. Hastie, and R. Tibshirani, An Introduction to Statistical Learning with Applications in R. New York: Springer-Verlag, 2013.

[24] Hastie, T., Tibshirani, R., and Wainwright, M., Statistical Learning with Sparsity, The LASSO and Generalizations. New Jersey: CRC Press, 2015.

[25] Fonti, V. and Belitser, E, 2017, Feature Selection using LASSO, p. 1-25, Research Paper In Business Analytics, VU Amsterdam.

[26] I.T. Jolliffe, Principal Component Analysis, 2nd ed., New York: Springer-Verlag.

[27] Herawati, N., Nisa, K., Setiawan, E., Nusyirwan and Tiryono, 2018, Regularized Multiple Regression Methods to Deal with Severe Multicollinearity, International Journal of Statistics and Applications, 8(4): 167-172.

[28] Akaike, H. 1973. Information theory and an extension of the maximum likelihood principle. In B.N. Petrow and F. Csaki (eds), Second International symposium on information theory (pp.267-281). Budapest: Academiai Kiado.

[29] Akaike, H. 1974. A new look at the statistical model identification. IEEE Transactions on Automatic Control, 19, 716-723.