# Developing Almost and Modified Almost Unbiased Estimators to Handle Multicollinearity Problem in Logistic Regression Model

*Humam A. Abdulrazzaq*[1,*]*, Raoudha Zine*[2]*, Mustafa I. Alheety*[3]

[1] Department of Mathematics, College of Education for Pure Sciences, University of Anbar, Iraq
[2] Laboratory of Probability and Statistics,Faculty of Science of Sfax, University of Sfax, Tunisia
[3] Department of Mathematics, College of Education for Pure Sciences, University of Anbar, Iraq

**Abstract:** This paper introduces two biased estimators to avoid problems arising from multicollinearity in the logistic regression model. We investigated the theoretical excellence of the proposed estimators according to the mean square error matrix (MSE) and the scalar mean square error (MSE) criterion. We found that they have the superiority than some existing estimators. Moreover, we run the simulation study, which depended on the simulated MSE (SMSE), squared bias (SB) and generalized cross validation (GCV) as criteria to compare the estimators. The simulation results showed that the proposed estimators have the superiority than the estimators under comparison at several factors and at the same time, they work well at the high level of correlation. In addition, we investigated the behavior of the present estimators applying the real data. Under this trend, the results were consistent with the theoretical results.

**Keywords:** Maximum likelihood estimator; Multicollinearity; AL estimator; Mean squared error matrix.

**Mathematics Subject Classification** Primary: 62J12. 26A25; 26A35.

## 1 Introduction

Logistic Regression Model (LRM) is widely used in the social sciences, in economic research and in the medical fields [7]. The presence of multicollinearity in logistic regression model can pose challenges in accurately estimating the regression parameters. Multicollinearity refers to a high degree of correlation between independent variables, which can inflate the standard errors of the model parameters and lead to inaccurate estimation results. To address this issue, researchers have proposed biased estimators, such as the Logistic Ridge, Logistic Liu, and estimators with two biasing parameters. These biased estimators aim to mitigate the impact of multicollinearity and provide more stable parameter estimates in logistic regression models. These biased estimators can help overcome the sensitivity of parameter estimates to multicollinearity by introducing a controlled bias in the estimation process. By applying these biased estimators, researchers can obtain more reliable and robust estimates of the regression parameters, even in the presence of multicollinearity. The almost unbiased estimation procedure offers a solution to the multicollinearity problem by incorporating bias into the estimation process, which helps stabilize parameter estimates and improves the accuracy of the LRM. Despite this, researchers are still working on developing biased estimators to address the issue of multicollinearity, such that, the estimators take into account the correlation among independent variables and adjust the parameter estimates accordingly.

### 1.1 Literature Review

Since biased estimators aim to mitigate the impact of multicollinearity and provide more stable parameter estimates in LRM, [15] developed ridge logistic regression which was the most widely used estimator for LRM. Furthermore, a

---

* Corresponding author e-mail: ahhumam1@uoanbar.edu.iq

modified logistic ridge regression estimator was introduced by [14], to deal with the multicollinearity problem. Based on the fact that ridge regression did not completely overcome the problem of ill-conditioning, Liu-type logistic estimator was defined by [9]. [4,5] introduced new biased estimators depending on Liu – type estimator. [17] proposed the modified almost unbiased Liu logistic estimator, while [10] suggested a modified estimator depending on ridge logistic estimator. [11] proposed the modified ridge type logistic estimator, as well as [13] introduced a new alternative method based on particle swarm optimization to estimate the (k, d) pair in Liu-type logistic estimator, simultaneously. Moreover, in the study of Varathan [18] a modified almost unbiased ridge logistic estimator was proposed. Also, [7] proposed a new estimator as a general estimator which includes other biased estimators.

The motivation of this paper can be given as follows: [1] introduced a new biased estimator called (AL) estimator for linear regression model. Till now, no researchers have tried to reduce the bias of the AL estimator or to use AL estimator for other regression models like logistic, Poisson, etc. Therefore, in this paper, two estimators are constructed based on the AL estimator after transformation to LRM. The suggested estimators are called almost unbiased AL estimator and modified almost unbiased AL estimator.

The rest of this paper is organized as follows: The methodology used in this paper is given in Section 2. The conditions for superiority of the proposed estimators over the existing estimators are found with respect to the matrix of mean square error (MSEM) and scalar mean square error (SMSE) criteria and are given in Section 3. In Section 4, the simulation study has been conducted to investigate the performance of the proposed estimators in the SMSE sense. An application is given in Section 5. Finally, the conclusion is given in Section 6.

## 2  Methodology

### 2.1  The Logistic Regression Model

The form of logistic regression model (LRM) is defined as:

$$\gamma_i = \mu_i + \varepsilon_i, \ \ i = 1, 2, \ldots, n \tag{1}$$

where $\gamma_i$ distributed as Bernoulli distribution, that is $\gamma_i \sim B(\mu_i)$ where $\mu_i$ is given as:

$$\mu_i = \frac{exp(x_i'\beta)}{1 + exp(x_i'\beta)} \ ; i = 1, 2, \ldots, n \tag{2}$$

where $x_i$ is the ith row of an $n \times (p+1)$ the design matrix X with n data pints and p independent explanatory variables and $\beta$ is a $(p+1) \times 1$ vector of coefficients and $\varepsilon_i$ is independent and distributed such that $E(\varepsilon_i) = 0$ and $\text{Var}(\varepsilon_i) = \mu_i(1 - \mu_i) = \sigma_i$.

The most commonly used method of estimating $\beta$ is the maximum likelihood estimation (MLE) method to maximize the log-likelihood $l(\beta)$:

$$l(\beta) = \sum_{i=1}^{n} y_i x_i' \beta - ln\left[1 + e^{x_i'\beta}\right].$$

The MLE estimator of $\beta$ is computed by setting the first derivative of $l(\beta)$ with respect to $\beta$ to zero. Therefore, the MLE estimator is obtained by solving the following equation:

$$\frac{\partial l(\beta)}{\partial \beta} = \sum_{i=1}^{n}\left[y_i - \frac{e^{x_i'\beta}}{1 + e^{x_i'\beta}}\right] x_i = \sum_{i=1}^{n} [y_i - \mu_i] x_i = 0.$$

The iteratively weighted least squares (IWLS) is applied to obtain the solution to Equation $\frac{\partial l(\beta)}{\partial \beta} = 0$. The MLE estimator of $\beta$ is estimated by applying the IWLS algorithm as follows [13]:

$$\widehat{\beta} = \left(X^t \widehat{U} X\right)^{-1} X^t \widehat{W} Z = S^{-1} X^t \widehat{W} Z, \tag{3}$$

where $S = \left(X^t \widehat{U} X\right)$ ; $\widehat{U} = \text{diag}(\widehat{\mu}_i(1 - \widehat{\mu}_i))$ , $\widehat{W} = \text{diag}[\widehat{\mu}_i(1 - \widehat{\mu}_i)]$ and Z is a column vector such that the $i^{th}$ element is $\log it(\widehat{\mu}_i) + \frac{\gamma_i - \widehat{\mu}_i}{\widehat{\mu}_i(1 - \widehat{\mu}_i)}$. The MLE is asymptotically unbiased where $E\left(\widehat{\beta}\right) = E(\arg\max_{\beta} l(\beta))$, and by differentiate the $l(\beta)$ with respect to $\beta$ and take the expectation on both side, that is:

$$E\left(\frac{\partial l(\beta)}{\partial \beta}\right) = \frac{\partial}{\partial \beta} E(l(\beta)).$$

So, $E(l(\beta))$ should be equal to maximum value of the $l(\beta)$, which is achieved at the $\beta$. The covariance matrix of asymptotically normally distributed $\widehat{\beta}$ is defied by the inverse of the Hessian matrix, $X^t \widehat{U} X$ which is given by the following equation [13]:

$$\text{Cov}\left(\widehat{\beta}\right) = S^{-1}. \tag{4}$$

When the Hessian matrix is not invertible, this leads to problems [7], where the variance of MLE will be large and as a result for that, the confidence interval will be large also. For this reason, the LRM suffers from unstably in case there is a strong dependence among independent variables.

## 2.2 The Proposed Estimators

Alheety and Gore [1] suggested a biased estimator called AL estimator (ALE) for linear regression model by augmenting the equation $mX^t X \widehat{\beta}_{\text{OLSE}} = \beta + \varepsilon^t$ to $Y = X\beta + \varepsilon$ and then they used the least square method to get the following form:

$$\widehat{\beta}_{\text{AL}} = (1+m)(I + \left(X^t X\right)^{-1})^{-1} \widehat{\beta}_{\text{OLSE}} \tag{5}$$

Where, $\widehat{\beta}_{\text{OLSE}} = (X^t X)^{-1} X^t Y$, $0 \le m \le 1$.
Now, if we convert the estimator in 5 to LRM, the ALE will take the following form:

$$\widehat{\beta}_{\text{ALL}} = (1+m)(I + S^{-1})^{-1} \widehat{\beta} = Q_m \widehat{\beta}, \tag{6}$$

where $Q_m = (1+m)(I + S^{-1})^{-1}$ and we will refer to it as AL logistic estimator (ALLE). Many researchers including [2,3], used a method to decrease bias in biased estimators. This method aims at making a slight increase in variance to achieve biased estimators with minimal bias according to the mean square error criterion. Such biased estimators are referred to as "almost unbiased estimators ".

Due to the limited research regarding this type of estimator, we propose a novel almost unbiased AL logistic estimator. To derive this estimator, we first provide the following definition:

**Definition 1.**[20] Suppose $\beta^*$ is a biased estimator of parameter vector $\beta$, and if the bias vector of $\beta^*$ is given by $Bias(\beta^*) = E(\beta^*) - \beta = R\beta$, where $R$ is a matrix, which shows that $E(\beta^*) - R\beta = \beta$, then the estimator $\beta^{**} = \beta^* - R\beta^* = (I - R)\beta^*$ is called the almost unbiased estimator based on the biased estimator $\beta^*$.

Through addition and subtraction of the matrix $S^{-1}$ in $(I + mI)$ that given in $\widehat{\beta}_{\text{ALL}}$, the ALLE can be rewritten as:

$$\widehat{\beta}_{\text{ALL}} = (I + mI)\left(I + S^{-1}\right)^{-1} \widehat{\beta} = \left[I + \left(mI - S^{-1}\right)\left(I + S^{-1}\right)^{-1}\right] \widehat{\beta}.$$

So,

$$E\left(\widehat{\beta}_{\text{ALL}}\right) = \beta + \left(I + S^{-1}\right)^{-1}\left(mI - S^{-1}\right)\beta.$$

Therefore, the bias is:

$$B\left(\widehat{\beta}_{\text{ALL}}\right) = E\left(\widehat{\beta}_{\text{ALL}}\right) - \beta = \left(I + S^{-1}\right)^{-1}\left(mI - S^{-1}\right)\beta$$

According to Definition 1, the almost unbiased AL logistic estimator (AUALLE) is given as follows:

$$\begin{aligned}
\widehat{\beta}_{\text{AUALL}} &= \left[I - \left(mI - S^{-1}\right)\left(I + S^{-1}\right)^{-1}\right] \widehat{\beta}_{\text{ALL}} \\
&= \left[I - \left(mI - S^{-1}\right)\left(I + S^{-1}\right)^{-1}\right]\left[I + \left(mI - S^{-1}\right)\left(I + S^{-1}\right)^{-1}\right] \widehat{\beta} \\
&= \left[I - \left(I + S^{-1}\right)^{-2}\left(mI - S^{-1}\right)^2\right] \widehat{\beta} \\
&= W_m \widehat{\beta},
\end{aligned} \tag{7}$$

Where, $W_m = \left[I - \left(I + S^{-1}\right)^{-2}\left(mI - S^{-1}\right)^2\right].$

By using $\widehat{\beta}_{\text{ALL}}$ instead of $\widehat{\beta}$ in 7, a new biased estimator is proposed which is called as the modified almost unbiased AL logistic estimator (MAUALLE) and defined as:

$$\widehat{\beta}_{\text{MAUALL}} = W_m \widehat{\beta}_{\text{ALL}} = W_m Q_m \widehat{\beta} = T_m \widehat{\beta}, \tag{8}$$

Where, $T_m = W_m Q_m = (1+m)\left[I - \left(I + S^{-1}\right)^{-2}\left(mI - S^{-1}\right)^2\right]\left(I + S^{-1}\right)^{-1}$

## 3 The MSE Comparison

The asymptotic scalar mean squared error (SMSE) and the asymptotic matrix mean squared error (MSEM) of an estimator $\widehat{\widetilde{\beta}} = Z\widehat{\beta}$, where $Z$ is a matrix with proper order, and are defined in [7] as:

$$\text{MMSE}\left(\widehat{\widetilde{\beta}}\right) = E\left(\widehat{\widetilde{\beta}} - \beta\right)\left(\widehat{\widetilde{\beta}} - \beta\right)' = Z\left(\widehat{\beta} - \beta\right)\left(\widehat{\beta} - \beta\right)' Z' + (Z-I)\beta'(Z-I)'$$

$$\text{SMSE}\left(\widehat{\widetilde{\beta}}\right) = E\left(\widehat{\widetilde{\beta}} - \beta\right)'\left(\widehat{\widetilde{\beta}} - \beta\right) = \left(\widehat{\beta} - \beta\right)' Z' Z\left(\widehat{\beta} - \beta\right) + (Z-I)'\beta'\beta(Z-I)$$

Note that there is a relationship between MMSE and SMSE criteria, where $\text{SMSE}\left(\widehat{\widetilde{\beta}}\right) = \text{tr}\left(\text{MMSE}\left(\widehat{\widetilde{\beta}}\right)\right)$ and tr is a trace of a square matrix. Therefore, the MSEM of MLE is:

$$\text{MMSE}\left(\widehat{\beta}\right) = S^{-1} \tag{9}$$

Consider spectral decomposition of the matrix $S$. Let $\alpha = P'\beta$, $\Lambda = \text{diag}\ (\lambda_1, \lambda_2, \ldots, \lambda_{p+1}) = P'\left(X'\widehat{W}X\right)P$, where $\lambda_1 \geq \lambda_2 \geq \ldots \geq \lambda_{p+1} \geq 0$ are the eigenvalues of $X'\widehat{W}X$, and $P$ is the matrix whose columns are the eigenvectors of $S$. Since $\widehat{\beta}$ is asymptotically unbiased and ALLE depends on $\widehat{\beta}$, the asymptotic properties of the ALLE are derived as follows:

The expectation:     $E\left(\widehat{\beta}_{\text{ALL}}\right) = Q_m\beta \tag{10}$

The covariance:     $\text{Cov}\left(\widehat{\beta}_{\text{ALL}}\right) = Q_m S^{-1} Q_m' \tag{11}$

The bias:     $B\left(\widehat{\beta}_{\text{ALL}}\right) = (Q_m - I)\beta \tag{12}$

The MSEM:     $\text{MMSE}\left(\widehat{\beta}_{\text{ALL}}\right) = Q_m S^{-1} Q_m' + (Q_m - I)\beta\beta'(Q_m - I)' \tag{13}$

And the SMSE:     $\text{SMSE}\left(\widehat{\beta}_{\text{ALL}}\right) = \sum_{i=1}^{p+1} \frac{(1+m)^2 \lambda_i}{(\lambda_i + 1)^2} + \sum_{i=1}^{p+1} \left[\frac{m\lambda_i - 1}{\lambda_i + 1}\right]^2 \alpha_i^2. \tag{14}$

The asymptotic properties of AUALLE are obtained as in Equations 15-19, respectively:

$$E\left(\widehat{\beta}_{\text{AUALL}}\right) = W_m\beta \tag{15}$$

$$\text{Cov}\left(\widehat{\beta}_{\text{AUALL}}\right) = W_m S^{-1} W_m' \tag{16}$$

$$B\left(\widehat{\beta}_{\text{AUALL}}\right) = (W_m - I)\beta = -\left(I + S^{-1}\right)^{-2}\left(mI - S^{-1}\right)^2 \beta \tag{17}$$

$$\text{MMSE}\left(\widehat{\beta}_{\text{AUALL}}\right) = W_m S^{-1} W_m' + (W_m - I)\beta\beta'(W_m - I)', \quad And \tag{18}$$

$$\text{SMSE}\left(\widehat{\beta}_{\text{AUALL}}\right) = \text{tr}\left(\text{MMSE}\left(\widehat{\beta}_{\text{AUALL}}\right)\right) = \sum_{i=1}^{p+1} \frac{1}{\lambda_i}\left[1 - \frac{(m\lambda_i - 1)^2}{(\lambda_i + 1)^2}\right]^2 + \sum_{i=1}^{p+1}\left[\frac{m\lambda_i - 1}{\lambda_i + 1}\right]^4 \alpha_i^2 \tag{19}$$

Also, the asymptotic properties of MAUALLE are obtained as in Equations 20 - 24, respectively:

$$E\left(\widehat{\beta}_{\text{MAUALL}}\right) = T_m\beta \tag{20}$$

$$\text{Cov}\left(\widehat{\beta}_{\text{MAUALL}}\right) = T_m S^{-1} T_m' \tag{21}$$

$$B\left(\widehat{\beta}_{\text{MAUALL}}\right) = (T_m - I)\beta \tag{22}$$

$$\text{MMSE}\left(\widehat{\beta}_{\text{MAUALL}}\right) = T_m S^{-1} T_m' + (T_m - I)\beta\beta'(T_m - I)' \tag{23}$$

Consequently, the SMSE is obtained as:

$$\text{SMSE}\left(\widehat{\beta}_{\text{MAUALL}}\right) = \sum_{i=1}^{p} \frac{(m+1)^2}{\lambda_i}\left[1 - \frac{(m\lambda_i - 1)^2}{(\lambda_i + 1)^2}\right]^2 + \sum_{i=1}^{p}\left[\frac{(m+1)\lambda_i}{\lambda_i + 1}\left[1 - \frac{(m\lambda_i - 1)^2}{(\lambda_i + 1)^2}\right] - 1\right]^2 \alpha_i^2 \qquad (24)$$

The new estimators are proposed to reduce the bias of ALLE estimator as well as to reduce the matrix mean square error and scaler mean square error. Therefore, in the following section we compare the new estimators with ALLE and then AUALLE with MAUALLE respectively.

### 3.1 Bias Comparison of Estimators

In this subsection, we will use the quadratic form of bias to compare the new estimators with the ALLE estimator.

**Theorem 1.** *Let $\|.\|$ denotes the norm of a vector, then in logistic regression model, the following inequality is held:*

$$\left\|B\left(\widehat{\beta}_{AUALL}\right)\right\|^2 < \left\|B\left(\widehat{\beta}_{ALL}\right)\right\|^2 \text{ for } 0 < m < 1.$$

*Proof.*

$$\begin{aligned}
\left\|B\left(\widehat{\beta}_{\text{ALL}}\right)\right\|^2 - \left\|B\left(\widehat{\beta}_{\text{AUALL}}\right)\right\|^2 &= \beta'\left(I + S^{-1}\right)^{-2}\left(mI - S^{-1}\right)^2\beta - \beta'\left(I + S^{-1}\right)^{-4}\left(mI - S^{-1}\right)^4\beta \\
&= \alpha'\left(I + \Lambda^{-1}\right)^{-2}\left(mI - \Lambda^{-1}\right)^2\alpha - \alpha'\left(I + \Lambda^{-1}\right)^{-4}\left(mI - \Lambda^{-1}\right)^4\alpha \\
&= \alpha' H \alpha
\end{aligned}$$

Where ,

$$H = \left(I + \Lambda^{-1}\right)^{-2}\left(mI - \Lambda^{-1}\right)^2 - \left(I + \Lambda^{-1}\right)^{-4}\left(mI - \Lambda^{-1}\right)^4 = \text{diag}\left(\frac{(\lambda_i + 1)^2(m\lambda_i - 1)^2 - (m\lambda_i - 1)^4}{(\lambda_i + 1)^4}\right)$$

$$= \text{diag}\left(\frac{(m\lambda_i - 1)^2}{(\lambda_i + 1)^4}\left[(\lambda_i + 1)^2 - (m\lambda_i - 1)^4\right]\right)$$

As we observe;

$$(\lambda_i + 1)^2 - (m\lambda_i - 1)^4 = \left[(\lambda_i + 1) - (m\lambda_i - 1)\right]\left[(\lambda_i + 1) + (m\lambda_i - 1)\right] = \left[(1-m)\lambda_i + 2\right](1+m)\lambda_i.$$

Therefore, $H$ is positive definite for $0 < m < 1$ and for that, $\alpha' H \alpha$ is positive definite. The proof is completed. To facilitate comparative analysis aimed at evaluating the efficacy of the suggested estimators, it is important to consider the following lemmas:

**Lemma 1.**[8] *Let $N$ be a positive definite matrix (pd), namely $(N > 0)$ and let $c$ be a nonzero vector then $N - cc'$ is nonnegative definite; namely $(N - cc' > 0)$ if and only if $c'N^{-1}c < 1$.*

**Lemma 2.**[19] *Suppose that $Q$ is a positive definite matrix and $N$ is a nonnegative definite matrix (NND), namely $N \geq 0$. Then*

$$Q - N \geq 0 \iff \lambda_{\max}\left(NQ^{-1}\right) \leq 1,$$

*where $\lambda_{\max}\left(NQ^{-1}\right)$ is the largest eigenvalue of the matrix $NQ^{-1}$.*

**Lemma 3.**[6] *Let $\widehat{\alpha}_i \quad i = 1,2$ be two competing homogeneous linear estimators of $\alpha$. Suppose that $D = Cov\left(\widehat{\alpha}_1\right) - Cov\left(\widehat{\alpha}_2\right)$ is a positive definite, where $Cov\left(\widehat{\alpha}_i\right)$, i=1,2 is the covariance matrix of $\widehat{\alpha}_i$ and $b_i = Bias\left(\widehat{\alpha}_i\right)$, consequently. Then $\Delta = MSEM\left(\widehat{\alpha}_1\right) - MSEM\left(\widehat{\alpha}_2\right) = D + b_1'b_1 - b_2'b_2 \geq 0$ if and only if $b_2'\left(D + b_1'b_1\right)b_2 < 1$, where $MSEM\left(\widehat{\alpha}_i\right) = Cov\left(\widehat{\alpha}_i\right) + b_i'b_i.$*

## 3.2 Comparison the AUALLE Over ALLE

In order to compare the AUALLE over ALLE, the following theorem explains the conditions that must be met to show the superiority of the AUALLE estimator over the ALLE estimator.

**Theorem 2.** *Under logistic regression model, when $m > \frac{1}{\lambda_i}$, the AUALLE is better than ALLE in the sense of MSEM if and only if $b_2'(D_1 + b_1'b_1)b_2 < 1$.*

*Proof.* We consider the MSEM difference of ALLE and AUALLE in order to show the superiority between them as follows:

$$\text{MMSE}\left(\widehat{\beta}_{\text{ALL}}\right) - \text{MMSE}\left(\widehat{\beta}_{\text{AUALL}}\right) = P\left(N_m\Lambda^{-1}N_m{}^t - R_m\Lambda^{-1}R_m{}^t\right)P^t + b_1b_1{}^t - b_2b_2{}^t = PD1P^t + b_1b_1{}^t - b_2b_2{}^t,$$

$$F = \left(I + \Lambda^{-1}\right)^{-1}\left(mI - \Lambda^{-1}\right),$$

$$N_m = (I + mI)\left(I + \Lambda^{-1}\right)^{-1},$$

$$R_m = \left[I - \left(I + \Lambda^{-1}\right)^{-2}\left(mI - \Lambda^{-1}\right)^2\right],$$

$$b_1 = \text{F}\beta \quad, \quad b_2 = (W_m - I)\beta \quad and$$

$$D1 = N_mS^{-1}N_m{}^t - R_mS^{-1}R_m{}^t.$$

Now we are starting for finding the conditions make $D_1$ a positive definite matrix (pd) namely; $D_1 > 0$. For that;

$$D_1 = N_mS^{-1}N_m' - R_mS^{-1}R_m{}^t = \text{diag}\left\{\frac{\lambda_i(1+m)^2}{(\lambda_i+1)^2} - \frac{(\lambda_i(1-m)+2)^2\lambda_i{}^2(1+m)^2}{\lambda_i(\lambda_i+1)^4}\right\}_{i=1}^{p+1}.$$

So, $D_1 > 0$ when,

$$\frac{\lambda_i(1+m)^2}{(\lambda_i+1)^2} > \frac{(\lambda_i(1-m)+2)^2\lambda_i(1+m)^2}{\lambda_i(\lambda_i+1)^4} \Rightarrow (\lambda_i+1)^2 > (\lambda_i(1-m)+2)^2 \text{ and then } m > \frac{1}{\lambda_i}.$$

Therefore, using Lemma 3, the proof is completed.

## 3.3 Superiority of the MAUALLE Over AUALLE

The following theorem shows the superiority of the MAUALLE over AUALLE by specifying the necessary conditions.

**Theorem 3.** *Under logistic regression model, when $m < \frac{1}{\lambda_i}$, $i = 1, \ldots, n$, the MAUALLE is better than AUALLE in the sense of MSEM if and only if $b_3'(D_2 + b_2'b_2)b_3 < 1$.*

*Proof.* The MSEM difference of MAUALLE and AUALLE is given as follows:

$$\Delta_1 = \text{MMSE}\left(\widehat{\beta}_{\text{AUALL}}\right) - \text{MMSE}\left(\widehat{\beta}_{\text{MAUALL}}\right) = D_2 + b_2b_2' - b_3b_3',$$

$$Where, \quad b_3 = (T_m - I)\beta \quad And, \qquad D_2 = W_mS^{-1}W_m' - T_mS^{-1}T_m'$$

We can rewrite $D_2$ in another form:

$$D_2 = W_mS^{-1}W_m' - T_mS^{-1}T_m' = W_mS^{-1}W_m' - W_mQ_mS^{-1}Q'_mW_m' = W_m\left\{S^{-1} - Q_mS^{-1}Q'_m\right\}W_m'.$$

Since $S^{-1} - Q_mS^{-1}Q'_m$ represent the difference between the variance of MLE and ALL estimators and as a result of (Theorem 1) from [1], the proof is completed.

The above theorems indicate that the proposed estimators are better than the other estimators under conditions. Also, the superiority of the estimators seems to depend on the unknown parameter $\beta$ and on the choice of the value of the biasing parameter $m$. For this reason and for practical purposes, we have to replace them by suitable estimates. Therefore,

we replace $\beta$ by MLE. Now, we have to estimate $m$ by using SMSE of AUALLE. The procedure is to minimize the SMSE of AUALLE by differentiate it with respect to $m$ as follows:

$$\frac{\partial \, \text{SMSE}\left(\widehat{\beta}_{\text{AUALL}}\right)}{\partial m} = \sum_{i=1}^{p+1} \frac{(m\lambda_i - 1)^3 \left(1 - \lambda_i \alpha_i^2\right) - (m\lambda_i - 1)(\lambda_i + 1)^2}{(\lambda_i + 1)^4}$$

Equate $\dfrac{\partial \, \text{SMSE}\left(\widehat{\beta}_{\text{AUALL}}\right)}{\partial m}$ to zero implies:

$$m\lambda_i^2 \left(m^2 - 1\right) - \lambda_i^2 \left(3m^2 + 2m - 1 - a\alpha_i^2\right) + 2\lambda_i = 0, \; i = 1, 2, \ldots, p+1$$

Where $a = m^3 \lambda_i^3 - 3m^2 \lambda_i^2 + 3m\lambda_i - 1$.

Due to the complexity of simplifying the equation with respect to $m$, therefore, it is recommended to utilize computer software to determine the optimal value of $m$ that can reduce the value of SMSE for AUALLE to minimum. Similarly, using the same approach can determine the optimal estimated value of $m$ for MAUALLE.

## 4 The Simulation Study

In this section, a simulation study is conducted to investigate and compare the accuracy of the new estimators AUALLE and MAUALLE with other exist estimators MLE and ALLE.

### 4.1 Algorithm

[12] The following method can be used to produce independent variables with different level of correlation:

$$x_{ij} = (1 - \rho^2)^{1/2} z_{ij} + \rho z_{ip} \qquad i = 1, 2, \ldots, n \quad j = 1, 2, \ldots, p \tag{25}$$

where $z_{ij}$ are numbers represent independent variables which are distributed as standard normal pseudo-random and $\rho$ suggested to be to be (0.80, 0.90, 0.95, and 0.99) as the correlation values between any two independent variables. As a limitation of this paper, the sample size n is considered to be 50, 100 and 200. In addition to that, the number of independent variables is set to $p = 4$ and $p = 8$ in order to obtain a clear vision of the performance of the new estimators. To analyze the dependent variable, the logistic regression model is employed. The values of pseudo random are derived based on the $\text{Be}(\pi_i)$ distribution, where:

$$\pi_i = \frac{\exp\left(x_i^t \beta\right)}{1 + \exp\left(x_i^t \beta\right)}$$

Following [15], $\beta$ is a vector and chosen to be the eigenvector corresponding to the largest eigenvalue of the matrix S such that $\beta' \beta = 1$. Further, we consider some selected values for $m$ (0.001, 0.005, 0.01, 0.05, 0.1, 0.5, 0.9). The simulation is repeated 10000 times and the estimated mean square error (MSE) values of the estimators are obtained using the following equation:

$$\text{MSE}(\beta^*) = \frac{1}{10000} \sum_{i=1}^{10000} (\beta_i^* - \beta)^T (\beta_i^* - \beta),$$

where $\beta_i^*$ is the obtained estimator by $i^{th}$ simulation. The computer software used for this purpose is R program.

### 4.2 Results of the Simulation Study

The MSE values of the estimators are reported in Tables 3-8 . In all cases, regardless of the sample size, degree of correlation and the number of explanatory variables, the performance of the proposed estimators was better than that of the rest of the estimators. On the other hand, it can be observed that the MLE estimator has the worst performance because it has the highest mean square error. In most cases, when the value of m is close to 1, the performance of the MAUALLE estimator is better than the AUALE estimator, while the ALLE estimator does not perform at the desired level in all cases compared to the performance of the proposed estimators. From Tables 3-8, the value of the mean square error is decreasing with the increase of sample size (n), while the effect of the number of explanatory variables on the performance of the estimators shows to be inversely related in terms of the value of the mean square error.

# 5 Application to Real Data

In this section, the myopia dataset examined by [16,5] is considered. The dataset is based on a study of myopia where it's from 618 of the subjects who had at least five years of follow up and were not myopic when they entered the study and includes 17 variables. However, following [16,5], only 5 variables as explanatory variables are used:
spherical equivalent refraction (SPHEQ), axial length (AL), anterior chamber depth (ACD), lens thickness (LT), vitreous chamber depth (VCD) which are all continuous variables of same scale (mm). The focus of analysis lies in the dependent variable, indicating the presence or absence of myopia, where myopia is represented by the numerical value 1 and its absence by 0. Furthermore, the data matrix X is centered and standardized so that $X^TX$ will be in the correlation form.
The IRLS algorithm is used to fit the logistic regression model. Estimated regression parameters and the scalar MSE values for MLE, ALLE, AULLE and MAUALLE estimators are given in Table 2 for different values of $m$.
According to Table 1, there is a high correlation between the explanatory variable's axial length and vitreous chamber depth (0.9419), and the condition number that is used as a measure of multicollinearity is calculated to be 393.3814, which means the existence of sever multicollinearity in the data set. Difference values of the biasing parameter $m$ have been selected randomly and for each value of $m$, the value of SMSE for MLE, ALL, AUALL and MAUALL are given in Table 2. The results in Table 2 detect that the proposed estimators AUALLE and MAUALLE outperform MLE and ALL estimators for all values of $0 < m < 1$.

On the other hand, it can be observed that the performance of MAUALL is better than AUALL for all $m$ values, which supports what was found in the simulation study.

|       | SPHEQ   | AL      | ACD     | LT      | VCD     |
|-------|---------|---------|---------|---------|---------|
| SPHEQ | 1.0000  | -0.3055 | -0.2388 | -0.0727 | -0.2471 |
| AL    | -0.3055 | 1.0000  | 0.4563  | -0.3289 | 0.9419  |
| ACD   | -0.2388 | 0.4563  | 1.0000  | -0.3393 | 0.1994  |
| LT    | -0.0727 | -0.3289 | -0.3393 | 1.0000  | -0.4516 |
| VCD   | -0.2471 | 0.9419  | 0.1994  | -0.4516 | 1.0000  |

**Table 1:** Correlation matrix of the data set

| $m$  | MLE      | ALL      | AUALL    | MAUALL   |
|------|----------|----------|----------|----------|
| 0.01 | 273.8544 | 34.04408 | 33.99584 | 33.08963 |
| 0.1  | 273.8544 | 34.03909 | 33.98851 | 33.08695 |
| 0.2  | 273.8544 | 34.03363 | 33.98069 | 33.08378 |
| 0.3  | 273.8544 | 34.02827 | 33.97323 | 33.08042 |
| 0.4  | 273.8544 | 34.02301 | 33.96615 | 33.07688 |
| 0.5  | 273.8544 | 34.01784 | 33.95939 | 33.07318 |
| 0.6  | 273.8544 | 34.01276 | 33.95295 | 33.06932 |
| 0.7  | 273.8544 | 34.00779 | 33.94684 | 33.06533 |
| 0.8  | 273.8544 | 34.0029  | 33.941   | 33.0612  |
| 0.9  | 273.8544 | 33.99812 | 33.93544 | 33.05697 |
| 0.99 | 273.8544 | 33.99389 | 33.9307  | 33.05306 |

**Table 2:** The SMSE values of different estimators of the data set

| m | $\rho = 0.80$ | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| | $p = 4$ | | | | $p = 8$ | | | |
| | ML | ALL | AUALL | MAUALL | ML | ALL | AUALL | MAUALL |
| 0.001 | 55.877 | 1.228 | 0.712 | 0.621 | 135.496 | 1.629 | 0.701 | 0.425 |
| 0.005 | 55.877 | 1.234 | 0.713 | 0.620 | 135.496 | 1.640 | 0.704 | 0.423 |
| 0.01 | 55.877 | 1.242 | 0.715 | 0.618 | 135.496 | 1.653 | 0.706 | 0.422 |
| 0.05 | 55.877 | 1.306 | 0.726 | 0.606 | 135.496 | 1.760 | 0.730 | 0.410 |
| 0.1 | 55.877 | 1.389 | 0.743 | 0.593 | 135.496 | 1.898 | 0.763 | 0.401 |
| 0.5 | 55.877 | 2.160 | 0.973 | 0.600 | 135.496 | 3.170 | 1.175 | 0.486 |
| 0.9 | 55.877 | 3.037 | 1.377 | 0.831 | 135.496 | 4.659 | 1.855 | 0.761 |
| m | $\rho = 0.90$ | | | | | | | |
| | $p = 4$ | | | | $p = 8$ | | | |
| | ML | ALL | AUALL | MAUALL | ML | ALL | AUALL | MAUALL |
| 0.001 | 105.641 | 0.953 | 0.620 | 0.589 | 260.435 | 1.060 | 0.529 | 0.388 |
| 0.005 | 105.641 | 0.957 | 0.621 | 0.588 | 260.435 | 1.066 | 0.530 | 0.386 |
| 0.01 | 105.641 | 0.962 | 0.622 | 0.586 | 260.435 | 1.073 | 0.531 | 0.385 |
| 0.05 | 105.641 | 1.006 | 0.627 | 0.573 | 260.435 | 1.137 | 0.542 | 0.374 |
| 0.1 | 105.641 | 1.063 | 0.637 | 0.560 | 260.435 | 1.219 | 0.560 | 0.364 |
| 0.5 | 105.641 | 1.584 | 0.803 | 0.578 | 260.435 | 1.958 | 0.827 | 0.454 |
| 0.9 | 105.641 | 2.140 | 1.127 | 0.821 | 260.435 | 2.786 | 1.321 | 0.693 |

**Table 3:** Estimated MSE of ML, ALL, AUALL and MAUALL for different values of m when n=200

| m | $\rho = 0.95$ | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| | $p = 4$ | | | | $p = 8$ | | | |
| | ML | ALL | AUALL | MAUALL | ML | ALL | AUALL | MAUALL |
| 0.001 | 205.992 | 0.791 | 0.568 | 0.571 | 513.531 | 0.738 | 0.435 | 0.370 |
| 0.005 | 205.992 | 0.795 | 0.568 | 0.570 | 513.531 | 0.742 | 0.436 | 0.369 |
| 0.01 | 205.992 | 0.798 | 0.568 | 0.568 | 513.531 | 0.747 | 0.436 | 0.368 |
| 0.05 | 205.992 | 0.830 | 0.571 | 0.556 | 513.531 | 0.785 | 0.441 | 0.357 |
| 0.1 | 205.992 | 0.872 | 0.577 | 0.542 | 513.531 | 0.834 | 0.450 | 0.348 |
| 0.5 | 205.992 | 1.242 | 0.705 | 0.565 | 513.531 | 1.260 | 0.637 | 0.441 |
| 0.9 | 205.992 | 1.601 | 0.984 | 0.816 | 513.531 | 1.693 | 1.030 | 0.666 |
| m | $\rho = 0.99$ | | | | | | | |
| | $p = 4$ | | | | $p = 8$ | | | |
| | ML | ALL | AUALL | MAUALL | ML | ALL | AUALL | MAUALL |
| 0.001 | 996.388 | 0.648 | 0.524 | 0.558 | 2521.382 | 0.459 | 0.358 | 0.359 |
| 0.005 | 996.388 | 0.650 | 0.524 | 0.556 | 2521.382 | 0.461 | 0.358 | 0.357 |
| 0.01 | 996.388 | 0.653 | 0.523 | 0.554 | 2521.382 | 0.463 | 0.358 | 0.356 |
| 0.05 | 996.388 | 0.674 | 0.523 | 0.542 | 2521.382 | 0.479 | 0.357 | 0.345 |
| 0.1 | 996.388 | 0.701 | 0.525 | 0.529 | 2521.382 | 0.499 | 0.359 | 0.337 |
| 0.5 | 996.388 | 0.932 | 0.620 | 0.556 | 2521.382 | 0.644 | 0.478 | 0.436 |
| 0.9 | 996.388 | 1.109 | 0.858 | 0.812 | 2521.382 | 0.713 | 0.786 | 0.654 |

**Table 4:** Estimated MSE of ML, ALL, AUALL and MAUALL for different values of m when n=200

| $m$ | $\rho = 0.80$ | | | | | | | |
| | $p = 4$ | | | | $p = 8$ | | | |
| | ML | ALL | AUALL | MAUALL | ML | ALL | AUALL | MAUALL |
| 0.001 | 62.328 | 1.131 | 0.683 | 0.625 | 158.344 | 1.436 | 0.657 | 0.443 |
| 0.005 | 62.328 | 1.137 | 0.684 | 0.623 | 158.344 | 1.444 | 0.659 | 0.441 |
| 0.01 | 62.328 | 1.144 | 0.685 | 0.621 | 158.344 | 1.456 | 0.661 | 0.439 |
| 0.05 | 62.328 | 1.201 | 0.693 | 0.607 | 158.344 | 1.547 | 0.677 | 0.424 |
| 0.1 | 62.328 | 1.276 | 0.705 | 0.591 | 158.344 | 1.667 | 0.701 | 0.408 |
| 0.5 | 62.328 | 1.984 | 0.894 | 0.559 | 158.344 | 2.786 | 1.021 | 0.422 |
| 0.9 | 62.328 | 2.810 | 1.240 | 0.730 | 158.344 | 4.120 | 1.569 | 0.636 |
| $m$ | $\rho = 0.90$ | | | | | | | |
| | $p = 4$ | | | | $p = 8$ | | | |
| | ML | ALL | AUALL | MAUALL | ML | ALL | AUALL | MAUALL |
| 0.001 | 120.204 | 0.889 | 0.601 | 0.593 | 304.308 | 0.958 | 0.511 | 0.406 |
| 0.005 | 120.204 | 0.893 | 0.602 | 0.592 | 304.308 | 0.964 | 0.511 | 0.404 |
| 0.01 | 120.204 | 0.898 | 0.602 | 0.590 | 304.308 | 0.971 | 0.512 | 0.402 |
| 0.05 | 120.204 | 0.938 | 0.605 | 0.575 | 304.308 | 1.026 | 0.519 | 0.387 |
| 0.1 | 120.204 | 0.990 | 0.612 | 0.559 | 304.308 | 1.098 | 0.530 | 0.371 |
| 0.5 | 120.204 | 1.477 | 0.744 | 0.537 | 304.308 | 1.765 | 0.731 | 0.398 |
| 0.9 | 120.204 | 2.017 | 1.021 | 0.724 | 304.308 | 2.531 | 1.130 | 0.603 |

**Table 5:** Estimated MSE of ML, ALL, AUALL and MAUALL for different values of m when n=100

| $m$ | $\rho = 0.95$ | | | | | | | |
| | $p = 4$ | | | | $p = 8$ | | | |
| | ML | ALL | AUALL | MAUALL | ML | ALL | AUALL | MAUALL |
| 0.001 | 236.190 | 0.747 | 0.555 | 0.576 | 596.129 | 0.684 | 0.432 | 0.390 |
| 0.005 | 236.190 | 0.750 | 0.555 | 0.575 | 596.129 | 0.688 | 0.432 | 0.389 |
| 0.01 | 236.190 | 0.754 | 0.555 | 0.573 | 596.129 | 0.692 | 0.432 | 0.386 |
| 0.05 | 236.190 | 0.783 | 0.556 | 0.558 | 596.129 | 0.725 | 0.434 | 0.372 |
| 0.1 | 236.190 | 0.821 | 0.558 | 0.542 | 596.129 | 0.769 | 0.438 | 0.357 |
| 0.5 | 236.190 | 1.175 | 0.659 | 0.527 | 596.129 | 1.163 | 0.572 | 0.393 |
| 0.9 | 236.190 | 1.536 | 0.898 | 0.727 | 596.129 | 1.578 | 0.887 | 0.597 |
| $m$ | $\rho = 0.99$ | | | | | | | |
| | $p = 4$ | | | | $p = 8$ | | | |
| | ML | ALL | AUALL | MAUALL | ML | ALL | AUALL | MAUALL |
| 0.001 | 1152.709 | 0.614 | 0.512 | 0.561 | 2913.353 | 0.436 | 0.362 | 0.376 |
| 0.005 | 1152.709 | 0.616 | 0.512 | 0.559 | 2913.353 | 0.437 | 0.361 | 0.374 |
| 0.01 | 1152.709 | 0.618 | 0.511 | 0.557 | 2913.353 | 0.439 | 0.361 | 0.372 |
| 0.05 | 1152.709 | 0.637 | 0.509 | 0.542 | 2913.353 | 0.453 | 0.357 | 0.357 |
| 0.1 | 1152.709 | 0.663 | 0.509 | 0.526 | 2913.353 | 0.471 | 0.355 | 0.343 |
| 0.5 | 1152.709 | 0.889 | 0.579 | 0.515 | 2913.353 | 0.616 | 0.429 | 0.385 |
| 0.9 | 1152.709 | 1.082 | 0.783 | 0.725 | 2913.353 | 0.708 | 0.668 | 0.586 |

**Table 6:** Estimated MSE of ML, ALL, AUALL and MAUALL for different values of m when n=100

| m | $\rho = 0.80$ | | | | | | | |
| | $p = 4$ | | | | $p = 8$ | | | |
| | ML | ALL | AUALL | MAUALL | ML | ALL | AUALL | MAUALL |
| 0.001 | 64.556 | 1.030 | 0.674 | 0.665 | 210.234 | 1.140 | 0.617 | 0.516 |
| 0.005 | 64.556 | 1.034 | 0.674 | 0.664 | 210.234 | 1.146 | 0.618 | 0.514 |
| 0.01 | 64.556 | 1.040 | 0.675 | 0.662 | 210.234 | 1.154 | 0.619 | 0.512 |
| 0.05 | 64.556 | 1.088 | 0.679 | 0.646 | 210.234 | 1.220 | 0.626 | 0.492 |
| 0.1 | 64.556 | 1.152 | 0.686 | 0.628 | 210.234 | 1.306 | 0.637 | 0.470 |
| 0.5 | 64.556 | 1.775 | 0.818 | 0.545 | 210.234 | 2.143 | 0.820 | 0.391 |
| 0.9 | 64.556 | 2.536 | 1.080 | 0.614 | 210.234 | 3.177 | 1.172 | 0.486 |
| m | $\rho = 0.90$ | | | | | | | |
| | $p = 4$ | | | | $p = 8$ | | | |
| | ML | ALL | AUALL | MAUALL | ML | ALL | AUALL | MAUALL |
| 0.001 | 125.662 | 0.813 | 0.595 | 0.629 | 407.191 | 0.774 | 0.494 | 0.466 |
| 0.005 | 125.662 | 0.816 | 0.595 | 0.627 | 407.191 | 0.778 | 0.494 | 0.464 |
| 0.01 | 125.662 | 0.820 | 0.595 | 0.625 | 407.191 | 0.783 | 0.494 | 0.462 |
| 0.05 | 125.662 | 0.853 | 0.595 | 0.608 | 407.191 | 0.821 | 0.494 | 0.442 |
| 0.1 | 125.662 | 0.898 | 0.597 | 0.589 | 407.191 | 0.873 | 0.495 | 0.419 |
| 0.5 | 125.662 | 1.331 | 0.679 | 0.511 | 407.191 | 1.375 | 0.590 | 0.349 |
| 0.9 | 125.662 | 1.844 | 0.883 | 0.600 | 407.191 | 1.980 | 0.833 | 0.455 |

**Table 7:** Estimated MSE of ML, ALL, AUALL and MAUALL for different values of m when n=50

| m | $\rho = 0.95$ | | | | | | | |
| | $p = 4$ | | | | $p = 8$ | | | |
| | ML | ALL | AUALL | MAUALL | ML | ALL | AUALL | MAUALL |
| 0.001 | 247.118 | 0.681 | 0.550 | 0.609 | 801.378 | 0.569 | 0.432 | 0.445 |
| 0.005 | 247.118 | 0.684 | 0.550 | 0.607 | 801.378 | 0.572 | 0.431 | 0.443 |
| 0.01 | 247.118 | 0.686 | 0.550 | 0.605 | 801.378 | 0.574 | 0.430 | 0.440 |
| 0.05 | 247.118 | 0.710 | 0.547 | 0.588 | 801.378 | 0.597 | 0.426 | 0.420 |
| 0.1 | 247.118 | 0.742 | 0.545 | 0.568 | 801.378 | 0.628 | 0.422 | 0.398 |
| 0.5 | 247.118 | 1.054 | 0.597 | 0.493 | 801.378 | 0.929 | 0.468 | 0.335 |
| 0.9 | 247.118 | 1.406 | 0.765 | 0.593 | 801.378 | 1.273 | 0.650 | 0.449 |
| m | $\rho = 0.99$ | | | | | | | |
| | $p = 4$ | | | | $p = 8$ | | | |
| | ML | ALL | AUALL | MAUALL | ML | ALL | AUALL | MAUALL |
| 0.001 | 1263.320 | 0.568 | 0.513 | 0.593 | 4117.243 | 0.389 | 0.380 | 0.430 |
| 0.005 | 1263.320 | 0.570 | 0.513 | 0.592 | 4117.243 | 0.390 | 0.379 | 0.428 |
| 0.01 | 1263.320 | 0.571 | 0.512 | 0.589 | 4117.243 | 0.391 | 0.378 | 0.425 |
| 0.05 | 1263.320 | 0.587 | 0.507 | 0.572 | 4117.243 | 0.399 | 0.369 | 0.405 |
| 0.1 | 1263.320 | 0.607 | 0.503 | 0.553 | 4117.243 | 0.411 | 0.361 | 0.383 |
| 0.5 | 1263.320 | 0.810 | 0.530 | 0.481 | 4117.243 | 0.530 | 0.362 | 0.324 |
| 0.9 | 1263.320 | 1.018 | 0.667 | 0.590 | 4117.243 | 0.638 | 0.487 | 0.439 |

**Table 8:** Estimated MSE of ML, ALL, AUALL and MAUALL for different values of m when n=50

## 6  Conclusion

In this paper, new estimators called almost unbiased logistic AL estimator (AUALL) and modified almost unbiased logistic AL estimator (MAUALLE) are proposed for logistic regression model when the multicollinearity problem exists. The superiority conditions for the proposed estimators with the existing estimators MLE, ALL are derived with respect to MSEM and SMSE criteria. Further, from the real data application and the Monte Carlo simulation study, it can be observed that the performance of MAUALL is better than AUALL for all $m$ values; where it has smaller SMSE than MLE, ALL, and AUALL when a high multicollinearity exists among the explanatory variables.

## Declarations

## References

[1] Alheety, M.I. and Gore, S. D. A new estimator in multiple linear regression model. *Model assisted statistics and applications*. **3** (2008), 187–200.

[2] Alheety, M.I. and Kibria, B. M. G. A new version of unbiased ridge regression estimator under the stochastic restricted linear regression model. *Communications in Statistics - Simulation and Computation*. **50(6)** (2021), 1589–1599.

[3] Alheety, M. I., Qasim, M., Kristofer, Mansson and Kibria, B. M. G. . Modifed almost unbiased two-parameter estimator for the Poisson regression model with an application to accident data. *SORT*. **45(2)** (2021) 1–22.

[4] Asar, Y. and Genc, A. . New Shrinkage Parameters for the Liu-type Logistic Estimators. *Communications in Statistics - Simulation and Computation*. **45(3)** (2016), 1094–1103.

[5] Asar, Y. , Some New Methods to Solve Multicollinearity in Logistic Regression. *Communications in Statistics - Simulation and Computation*. **46(4)** (2017), 2576–2586.

[6] C. Radhakrishna Rao., H. Toutenburg, Shalabh, C. H. , *Linear Models and Generalizations Least Squares and Alternatives*. 3rd edition Springer Berlin Heidelberg NewYork, (2008).

[7] Esra, E. and Kadri, U. A. A new Liu-type estimator in binary logistic regression models. *Communications in Statistics – Theory and Methods*. **51(13)**(2022), 4370–4394.

[8] Farebrother, R. W. Further results on the mean square error of ridge regression. *Journal of the Royal Statistical Society, Series B (Methodological)*, **38**(1976), 248-250.

[9] Inan, D. and Erdogan, B. E. Liu-Type Logistic Estimator. *Communications in Statistics - Simulation and Computation*. **35** (2013), 667–687.

[10] Jadhav, N. H., On linearized ridge logistic estimator in the presence of multicollinearity. Computational Statistics, **42** (2020), 1578–1586.

[11] Lukman, A.F, Emmanuel, A., Clement, O. A. et al., . A Modified Ridge – Type Logistic Estimator. *Iranian Journal of Science and Technology, Transactions A: Science,*(2020).

[12] McDonald, G. C. and Galarneau, D. I. . A Monte Carlo evaluation of some ridge type estimators. *Journal of the American Statistics Association*. **70(350)** (1975), 407–416.

[13] Nuriye S. and Deniz I. A new alternative estimation method for Liu-type logistic estimator via particle swarm optimization: an application to data of collapse of Turkish commercial banks during the Asian financial crisis. *Journal of Applied Statstics*. **48** (2021), 13–15.

[14] Nja, M.E., Ogoke, U.P. and Nduka, E.C. The Logistic Regression Model with a Modified Weight Function. *Journal of Statistical and Econometric Method*. **2** (2013), 161–171.

[15] Schaeffer, R. L., Roi, L. D., and Wolfe, R. A. A ridge logistic estimator. *Communications in Statistics: Theory and Methods*. **13** (1984), 99–113.

[16] Tyagi, G. and Chandra, S. A general restricted estimator in binary logistic regression in the presence of multicollinearity. *Brazilian Journal of Probability and Statistics*. **63(2)** (2022), 287–314.

[17] Varathan, N. and Wijekoon, P. Modified almost unbiased Liu estimator in logistic regression. *Communications in Statistics - Simulation and Computation*, DOI:10.1080/03610918.2019.1626888.(2019).

[18] Varathan, N. . An improved ridge type estimator for logistic regression. *Statistic in Transition new series*. **23(3)** (2022), 113–126.

[19] Wang, S. . *The Inequalities of Matrices*. (Hefei: The Education of Anhui Press), (1994).

[20] Xu, J. and Yang, H. More on the bias and variance comparisons of the restricted almost unbiased estimators. *Communications in Statistics - Theory and Methods*. **40(22)** (2011), 4053–4064.