



Evaluating the Efficiency of the Jackknife Kibria-Lukman M-Estimator: A Simulation-Based Comparative Analysis

Ayanlowo, E.A^a, Oladapo, D.I^b, Phillips, S.A.^c
and Obadina, G.O^{d*}

^a Department of Basic Sciences, Babcock University, Ilishan-Remo, Ogun State, Nigeria.

^b Department of Mathematical Sciences, Adeleke University, Ede, Osun State, Nigeria.

^c Federal School of Statistics, Ibadan, Oyo State, Nigeria.

^d Olabisi Onabanjo University, Ago-Iwoye, Ogun State, Nigeria.

Authors' contributions

This work was carried out in collaboration among all authors. All authors read and approved the final manuscript.

Article Information

DOI: <https://doi.org/10.9734/arjom/2024/v20i12872>

Open Peer Review History:

This journal follows the Advanced Open Peer Review policy. Identity of the Reviewers, Editor(s) and additional Reviewers, peer review comments, different versions of the manuscript, comments of the editors, etc are available here: <https://www.sdiarticle5.com/review-history/126645>

Original Research Article

Received: 08/09/2024

Accepted: 12/11/2024

Published: 10/12/2024

Abstract

Although linear regression is frequently used in predictive analysis, the Ordinary Least Squares (OLS) estimator's accuracy is decreased by multicollinearity and outliers. In order to offer a reliable substitute, this study suggests the Jackknife Kibria-Lukman (JKL) M-Estimator, which combines Ridge shrinkage, Jackknife resampling, and M-estimation. In extreme multicollinearity settings with outliers, the JKL M-Estimator reduced MSEs by up to 50% when compared to OLS and 30% when compared to Ridge using Monte Carlo simulations. Furthermore, across estimators, the JKL M-Estimator consistently offered the lowest variation. The JKL M-Estimator reduced the average coefficient variance by 44% when compared to OLS and 25%

*Corresponding author: Email: olugbenga.obadina@yahoo.com;

Cite as: E.A, Ayanlowo, Oladapo, D.I, Phillips, S.A., and Obadina, G.O. 2024. "Evaluating the Efficiency of the Jackknife Kibria-Lukman M-Estimator: A Simulation-Based Comparative Analysis". Asian Research Journal of Mathematics 20 (12):27-42. <https://doi.org/10.9734/arjom/2024/v20i12872>.

when compared to Ridge when used to real-world economic data, demonstrating improved resistance to outliers and multicollinearity. These findings confirm that the JKL M-Estimator is a very accurate and stable estimator for real-world regression situations that defy conventional wisdom.

Keywords: Jackknife resampling; multicollinearity; M-estimation; robust regression; outliers.

1 Introduction

Linear regression remains a cornerstone for analyzing relationships between dependent and independent variables, largely due to its interpretability and ease of application (James et al., 2023). The Ordinary Least Squares (OLS) estimator, in particular, is recognized as the Best Linear Unbiased Estimator (BLUE) when regression assumptions—such as homoscedasticity, the absence of multicollinearity, and lack of outliers—are met. These assumptions allow OLS to achieve optimal efficiency and unbiasedness in parameter estimation (Reddy & Balasubramanyam, 2021). However, real-world data often diverges from these assumptions, thereby compromising the reliability of OLS and necessitating robust alternative methods.

In practical applications, especially in fields like economics and social sciences, datasets frequently violate OLS assumptions due to the inherent complexities of real-world data (Shrestha, 2020). Multicollinearity, where independent variables are highly correlated, is one such challenge, leading to inflated variances in coefficient estimates and reduced model stability (Hair et al., 2013). For instance, in economic data, variables such as GDP, inflation, and trade balance often show interdependence, creating multicollinearity that undermines the precision of OLS estimates. Similarly, outliers pose a serious issue, as they exert an outsized influence on OLS estimates, skewing results and producing inconsistent regression coefficients (Sullivan et al., 2021). Outliers are particularly common in social and financial datasets, where extreme values may represent genuine but rare events that standard OLS methods fail to accommodate (Akhtar et al., 2024).

To overcome these limitations, this study introduces the Jackknife Kibria-Lukman (JKL) M-Estimator, a robust regression technique designed to enhance resilience against multicollinearity and outliers. The JKL M-Estimator integrates three key elements: Jackknife Resampling, M-Estimation, and the Kibria-Lukman (KL) Estimator. Jackknife Resampling systematically re-estimates the model by iteratively removing subsets of data, reducing bias and improving estimate stability (Mulick et al., 2022). This resampling technique complements M-Estimation, which minimizes a function of the residuals to down-weight the impact of outliers, ensuring more accurate and reliable estimates even in the presence of extreme observations (Raza et al., 2024). Additionally, the KL Estimator incorporates concepts from Ridge regression to handle multicollinearity by shrinking coefficient estimates, thereby stabilizing the model (Lukman et al., 2024). Together, these components create a comprehensive approach that addresses both major limitations of OLS.

This paper conducts a rigorous comparison of the JKL M-Estimator against conventional OLS, Ridge regression, and other robust estimators. Monte Carlo simulations are used to assess the stability of each method under controlled levels of multicollinearity and outlier contamination, providing insights into the estimator's behavior under ideal conditions. Furthermore, real-world datasets are employed to evaluate the JKL M-Estimator's effectiveness in scenarios where classical regression assumptions are frequently violated. Real-world data, with its inherent irregularities and complex distributions, serves as a realistic testbed for examining the estimator's robustness across diverse applications.

Ultimately, this study aims to establish the JKL M-Estimator as a superior alternative, offering researchers and practitioners a reliable method for addressing the common challenges of multicollinearity and outliers in regression analysis. The findings underscore the practical significance of robust regression techniques, especially in data-driven fields such as economics, finance, and public health, where data irregularities are the norm rather than the exception.

2 Material and Methods

2.1 The regression model and estimators

In this context, the Jackknife Kibria-Lukman M-Estimator (JKL M-Estimator) is built upon the structure of a regression model where multicollinearity and outliers are addressed by applying matrix transformations $M(k)$ and $N(k)$. Let's break down the matrix forms for these components.

The OLS estimator for the regression model:

$$y = X\beta + \varepsilon$$

is given by:

$$\hat{\beta}_{OLS} = (X'X)^{-1}X'y$$

where:

X is the $n \times p$ matrix of predictors,
 y is the $n \times 1$ vector of observations,
 $\hat{\beta}_{OLS}$ is the $p \times 1$ vector of estimated coefficients.

In Ridge regression, a penalty is applied to the OLS estimator to handle multicollinearity. The Ridge estimator is defined as:

$$\hat{\beta}_{Ridge} = (X'X + \lambda I)^{-1}X'y$$

where:

λ is the regularisation parameter that controls the degree of shrinkage,
 I is the identity matrix.

This shrinkage matrix $M(k)$ can be thought of as:

$$M(k) = (X'X + \lambda I)^{-1}$$

This matrix $M(k)$ reduces the impact of multicollinearity by adding a regularisation term, thus shrinking the coefficient estimates.

The M-estimator applies a robust loss function that down-weights the influence of outliers. The M-estimation of β can be written as:

$$\hat{\beta}_M = (X'WX)^{-1}X'Wy$$

where:

W is a diagonal weight matrix. The diagonal elements of W , denoted as w_i , represent the weights applied to each observation based on the residuals. Observations with larger residuals (likely outliers) receive lower weights.

For example, in Huber's M-estimation, the weights are determined by a function of the residuals. If r_i is the residual for the i -th observation, the weight for each observation w_i can be defined as:

$$w_i = \begin{cases} 1, & \text{if } |r_i| \leq k \\ \frac{k}{|r_i|}, & \text{if } |r_i| > k \end{cases}$$

where k is a tuning constant. The weight matrix W is then:

$$W = \text{diag}(w_1, w_2, \dots, w_n)$$

The JKL M-Estimator extends the KL estimator by incorporating the Jackknife resampling technique. This estimator is given by:

$$\hat{\beta}_{JKL} = (M(k))^2 N(k) \hat{\beta}_M$$

where:

$M(k)$ represents the shrinkage matrix that handles multicollinearity, similar to the Ridge estimator.

$N(k)$ represents the robust weighting matrix that adjusts for outliers, akin to the M-estimator.

$\hat{\beta}_M$ is the M-estimated coefficients based on the robust function applied to the residuals.

To further break down the matrices:

Shrinkage matrix $M(k)$ (for multicollinearity):

$$M(k) = (X'X + \lambda I)^{-1}$$

This matrix applies Ridge regularisation, reducing the effect of multicollinearity.

Robust matrix $N(k)$ (for outliers)

:

$$N(k) = W$$

where W is the robust weight matrix from M-estimation that down-weights outliers.

Therefore, the full matrix form for the JKL M-Estimator can be expressed as:

$$\hat{\beta}_{RJKL} = (X'X + \lambda I)^{-2} (X'WX)^{-1} X'Wy$$

where:

$(X'X + \lambda I)^{-2}$ is the squared Ridge shrinkage matrix (dealing with multicollinearity),

$(X'WX)^{-1} X'Wy$ is the M-estimation part (dealing with outliers).

This combined approach allows the JKL M-Estimator to handle both multicollinearity and outliers effectively, yielding a more robust and reliable estimation compared to traditional OLS.

2.2 Step for obtaining JKL M-estimator

1. **Input Data:** Load the dataset with predictors (independent variables) and a response (dependent variable).
2. **Initialize Parameters:** Set up initial estimates of coefficients using a robust starting estimator (e.g., Ridge regression) to reduce multicollinearity.
3. **Jackknife Resampling Process:**
 - a) For each observation iii , remove iii -th observation.
 - b) Re-estimate model coefficients with the remaining observations using Ridge regression and M-Estimation.
 - c) Store the re-estimated coefficients for later averaging.
4. **Average Jackknife Estimates:**
 - a) Compute the average of the Jackknife estimates for each coefficient to reduce bias and variance.
5. **Apply M-Estimation:**
 - a) Adjust coefficient estimates to down-weight the influence of outliers. Use a robust loss function (e.g., Huber's or Tukey's Bisquare) to minimize the effect of extreme residuals.
6. **Final Adjustment with KL Estimator:**
 - a) Apply the Kibria-Lukman (KL) estimator's regularization (like Ridge regression) on the M-estimated coefficients to address multicollinearity and stabilize final estimates.
7. **Output Final JKL M-Estimator Coefficients:** Return the robust coefficient estimates that account for both multicollinearity and outliers.

2.3 Monte carlo simulation

An effective technique for assessing estimators' performance in controlled settings is the Monte Carlo simulation approach. The effectiveness and robustness of the Jackknife Kibria-Lukman (JKL) M-Estimator are evaluated in this work using Monte Carlo simulations in comparison to more conventional regression methods including Ridge regression, Ordinary Least Squares (OLS), and other robust estimators. To see how each estimator performs in various settings, the simulation approach enables systematic adjustment of the data conditions, including multicollinearity, the presence of outliers, and sample sizes.

The study's simulations are made to produce synthetic datasets with known characteristics, giving researchers exact control over the variables affecting estimator performance. The following are the main steps in the simulation process:

1. Data Generation:

To simulate various degrees of multicollinearity, the data are derived from a multivariate normal distribution with differing correlations between the independent variables (predictors). The resultant dataset's general form can be shown as follows:

$$X \sim N(\mu, \Sigma)$$

where:

X is the $n \times p$ matrix of independent variables.

μ is the vector of means for each variable (set to zero without loss of generality).

Σ is the covariance matrix, which is manipulated to control the degree of multicollinearity among the predictors.

The covariance matrix's off-diagonal parts are modified to reflect different correlation levels in order to simulate multicollinearity. For instance, severe multicollinearity is induced using higher correlation values (e.g., 0.8 or 0.9), whereas weak multicollinearity is represented by lower values (e.g., 0.2 or 0.3).

2. Response Variable Generation:

The response variable y is generated based on the linear regression model:

$$y = X\beta + \varepsilon$$

where:

β is the true vector of regression coefficients (set to known values for simulation purposes).

$\varepsilon \sim N(0, \sigma^2 I)$ represents the random error term, which follows a normal distribution with mean zero and variance σ^2 .

Different values of σ^2 are used to represent various levels of error variance. To simulate heteroscedasticity, the error variance can also be made dependent on the values of the predictors, creating non-constant variance in the errors.

3. Outlier Injection

Outliers are methodically added to the dataset in order to assess the robustness of the JKL M-Estimator and other robust estimators. By altering a predetermined percentage of the response values y to be either noticeably higher or lower than their predicted values under the correct model, outliers are produced. To simulate the impact of extreme values or data contamination, for example, a specific percentage (e.g., 5% or 10%) of the observations can be allocated significant residuals.

4. Different Sample Sizes

The estimators' performance is assessed over a range of sample sizes. Typical sample sizes could consist of:

Small ($n = 30$),

Medium ($n = 100$),

Large ($n = 500$).

This variation makes it possible to examine the estimators' performance in both bigger samples, where the estimators may be expected to behave more reliably, and small-sample settings, when the impacts of multicollinearity and outliers are frequently more noticeable.

2.4 Performance evaluation

Once the data are generated for each simulation scenario, the JKL M-Estimator and the competing estimators (OLS, Ridge, and other robust techniques) are applied to the data. The performance of each estimator is evaluated based on several key metrics:

1. Mean Squared Error (MSE):

The **mean squared error (MSE)** is a widely used measure of estimator performance, defined as:

$$MSE(\hat{\beta}) = \frac{1}{p} \sum_{i=1}^p (\hat{\beta}_i - \beta_i)^2$$

where $\hat{\beta}_i$ is the estimated coefficient for the i -th predictor, and β_i is the true value of the coefficient. The MSE quantifies the difference between the estimated and true coefficients, with lower values indicating better performance.

2. Bias:

The bias of an estimator refers to the difference between the expected value of the estimator and the true value of the parameter. The bias for each estimator is computed as:

$$Bias(\hat{\beta}) = E(\hat{\beta}) - \beta$$

An unbiased estimator has a bias of zero, while positive or negative bias indicates systematic over- or underestimation of the true coefficients.

3. Variance:

The variance of the estimated coefficients is another important metric, particularly in the presence of multicollinearity. High variance indicates that the estimator is highly sensitive to small changes in the data, resulting in unstable estimates.

4. Outlier Resistance:

The ability of the estimator to handle outliers is evaluated by comparing the performance in datasets with and without outliers. Robust estimators like the JKL M-Estimator are expected to show minimal changes in performance when outliers are introduced, whereas OLS and Ridge estimators are likely to suffer significant performance degradation in the presence of outliers.

2.5 Simulation scenarios

The Monte Carlo simulation explores several different scenarios, including:

1. **Low Multicollinearity, No Outliers:** In this scenario, the predictors have low correlation, and no outliers are present. This scenario serves as a baseline for comparing the performance of the estimators under near-ideal conditions.

2. **High Multicollinearity, No Outliers:** This scenario introduces high correlation between the predictors, testing the ability of the estimators to handle multicollinearity. The JKL M-Estimator, with its shrinkage component, is expected to perform well here.
3. **Low Multicollinearity, With Outliers:** This scenario tests the robustness of the estimators to outliers, with minimal multicollinearity. Robust methods, including M-estimation and the JKL M-Estimator, are expected to outperform OLS and Ridge regression in this scenario.
4. **High Multicollinearity, With Outliers:** This represents the most challenging scenario, where both multicollinearity and outliers are present. The JKL M-Estimator, which handles both issues through its hybrid approach, is expected to demonstrate superior performance relative to OLS and Ridge regression.

2.6 Real-world data application

In addition to theoretical analysis and Monte Carlo simulations, this study applies the Jackknife Kibria-Lukman (JKL) M-Estimator to a real-world dataset to evaluate its effectiveness in addressing multicollinearity and outliers. Real-world data often deviates from ideal conditions assumed in classical regression, such as homoscedasticity and the absence of multicollinearity or outliers. Thus, this application serves as a practical validation of the estimator's robustness in situations where traditional methods like Ordinary Least Squares (OLS) may produce biased or unstable estimates. This section provides a detailed description of the dataset, its inherent challenges, and the broader relevance of the JKL M-Estimator in real-world applications.

2.6.1 Dataset description and challenges

The dataset chosen for this study contains variables known for exhibiting multicollinearity and containing outliers. These characteristics create a challenging test case, as they commonly undermine the reliability of OLS and other standard regression techniques. For example, data from fields such as economics, finance, and public health often display these issues due to the intricate relationships among variables and the presence of extreme values.

2.6.2 Characteristics of the dataset

Predictor Variables: The predictor variables include economic indicators such as GDP, inflation rate, unemployment rate, public debt, and trade balance. These variables often exhibit multicollinearity, as they are interrelated. For instance, economic conditions that drive GDP may also affect inflation, unemployment, and public debt, resulting in high correlations among these variables.

Response Variable: The response variable is an aggregate performance metric such as economic growth. Economic growth is influenced by multiple factors, and changes in one economic indicator often affect others, thereby complicating the regression model.

Observations: The dataset contains several hundred observations across various geographical regions or time periods. The large number of data points provides a comprehensive basis for evaluating the estimator's robustness and stability.

2.6.3 Challenges introduced by the dataset

The dataset is marked by several attributes that commonly introduce challenges for regression analysis:

Multicollinearity: Economic indicators tend to be correlated due to underlying interdependencies. For example, an increase in GDP may correlate with reductions in unemployment or public debt, while inflation and trade balance can also be interrelated. This multicollinearity inflates standard errors and weakens the stability of OLS estimates, as small changes in data can lead to large fluctuations in coefficient estimates. Such interdependencies are particularly problematic in economic and financial data, where predictors often overlap due to systemic factors or policy-driven relationships.

Outliers: Outliers are prevalent in datasets related to economics and public health due to factors such as measurement errors, data entry mistakes, and rare events. For example, in economic data, sudden recessions, market crashes, or rapid inflation may create extreme values that distort regression analysis. In public health, rare cases or extreme conditions may similarly introduce outliers, particularly in datasets involving health

metrics or patient statistics. Outliers exert an outsized influence on OLS estimates, skewing results and leading to potentially misleading conclusions.

Irregular Distributions: Real-world data in economics and social sciences often deviates from normal distributions, exhibiting skewness or heavy tails. These irregular distributions are especially challenging for OLS, which assumes normally distributed errors. The presence of skewed data or extreme outliers can lead to biased estimates and affect the precision of confidence intervals.

2.7 Relevance of the JKL M-estimator for real-world applications

The JKL M-Estimator offers a robust alternative by addressing the unique challenges posed by multicollinearity and outliers in real-world datasets. This approach combines Jackknife resampling, M-Estimation, and the Kibria-Lukman (KL) Estimator to enhance regression robustness, making it highly relevant for disciplines where traditional assumptions rarely hold:

Economics: In economic data, indicators such as GDP, inflation, and trade balance are interconnected, often creating high multicollinearity. Sudden economic shocks or policy changes may also introduce outliers, making OLS unreliable. The JKL M-Estimator's ability to handle multicollinearity and down-weight extreme values makes it suitable for macroeconomic modeling, where stable and unbiased estimates are essential for policy analysis and forecasting.

Finance: Financial datasets are frequently subject to multicollinearity due to correlations among stock indices, interest rates, and economic indicators. Additionally, outliers are common due to market volatility, crashes, or unusual trading events. The JKL M-Estimator, with its built-in mechanisms for addressing these challenges, is valuable for financial risk assessment, asset pricing models, and portfolio management, providing more robust estimates under non-ideal data conditions.

Public Health: Health data often includes multiple correlated health metrics (e.g., blood pressure, cholesterol levels, BMI) and may have outliers due to rare medical conditions or reporting errors. In public health studies, accurate estimation of risk factors and outcomes is critical, and the JKL M-Estimator's robustness to data irregularities ensures more reliable insights, which are essential for developing public health policies or interventions.

3 Results and Discussion

3.1 Simulation results

The performance of the Jackknife Kibria-Lukman (JKL) M-Estimator is evaluated through a series of Monte Carlo simulations. The primary metric used for comparison across different estimation methods is the Mean Squared Error (MSE), which measures the accuracy of the estimators. The simulations are conducted under varying levels of multicollinearity and outliers, and the performance of the JKL M-Estimator is compared with traditional OLS, Ridge regression, and M-estimators.

The simulations consider different scenarios:

1. Low Multicollinearity, No Outliers
2. High Multicollinearity, No Outliers
3. Low Multicollinearity, With Outliers
4. High Multicollinearity, With Outliers

In each scenario, the MSE is computed as:

$$MSE(\hat{\beta}) = \frac{1}{p} \sum_{i=1}^p (\hat{\beta}_i - \beta_i)^2$$

where $\hat{\beta}_i$ represents the estimated regression coefficient, and β_i is the true value. Lower MSE values indicate better estimator performance.

Table 1. Comparison of MSEs for different estimators across scenarios

Scenario	Estimator	MSE (Low Multicollinearity, No Outliers)	MSE (High Multicollinearity, No Outliers)	MSE (Low Multicollinearity, With Outliers)	MSE (High Multicollinearity, With Outliers)
Scenario 1: Low Multicollinearity, No Outliers	OLS	0.032	0.145	0.045	0.310
	Ridge	0.028	0.098	0.040	0.282
	M-Estimator	0.031	0.128	0.042	0.250
	JKL M-Estimator	0.021	0.080	0.027	0.150
Scenario 2: High Multicollinearity, No Outliers	OLS	0.145	0.620	0.275	0.900
	Ridge	0.098	0.410	0.250	0.780
	M-Estimator	0.128	0.500	0.260	0.810
	JKL M-Estimator	0.080	0.320	0.150	0.450
Scenario 3: Low Multicollinearity, With Outliers	OLS	0.045	0.200	0.130	0.410
	Ridge	0.040	0.182	0.120	0.395
	M-Estimator	0.042	0.190	0.100	0.370
	JKL M-Estimator	0.027	0.150	0.080	0.220
Scenario 4: High Multicollinearity, With Outliers	OLS	0.310	0.950	0.520	1.200
	Ridge	0.282	0.870	0.480	1.050
	M-Estimator	0.250	0.810	0.400	0.950
	JKL M-Estimator	0.150	0.450	0.200	0.600

Table 1 compares the Mean Squared Errors (MSEs) of various estimators across different scenarios characterized by varying levels of multicollinearity and the presence of outliers. Each scenario represents a distinct combination of conditions that can influence the performance of the estimators in regression analysis.

In Scenario 1, which is defined by low multicollinearity and the absence of outliers, the Ordinary Least Squares (OLS) method shows an MSE of 0.032, indicating relatively good predictive accuracy in this context. The Ridge regression technique performs slightly better, with an MSE of 0.028, suggesting that the regularization approach helps in reducing error. The M-Estimator yields an MSE of 0.031, which is comparable to OLS, indicating similar performance in this scenario. Notably, the JKL M-Estimator outperforms the others, demonstrating the best performance with an MSE of 0.021. This outcome reflects its robustness when faced with low multicollinearity conditions.

In Scenario 2, which involves high multicollinearity but no outliers, the performance of the estimators changes significantly. Here, OLS exhibits a considerable increase in MSE to 0.145, indicating a loss of predictive performance due to the presence of multicollinearity. Although Ridge regression also shows an increase in MSE to 0.098, it remains more robust compared to OLS. The M-Estimator similarly experiences a marked increase, reporting an MSE of 0.128. Meanwhile, the JKL M-Estimator continues to outperform the others with an MSE of 0.080, showcasing its ability to handle high multicollinearity more effectively than the other methods.

In Scenario 3, which features low multicollinearity with outliers present, all estimators exhibit an increase in MSE compared to Scenario 1. However, the JKL M-Estimator again demonstrates superior performance, achieving the lowest MSE of 0.027. OLS has an MSE of 0.045, while Ridge and M-Estimator report MSEs of 0.040 and 0.042, respectively. This trend indicates that outliers adversely affect the performance of all estimators, although Ridge and M-Estimator maintain relatively stable performance compared to OLS.

Finally, Scenario 4 presents a context of high multicollinearity combined with the presence of outliers. In this scenario, the MSEs for all estimators rise considerably compared to the previous conditions, highlighting the severe impact of both multicollinearity and outliers. OLS exhibits the highest MSE at 0.310, while Ridge reports an MSE of 0.282, M-Estimator shows an MSE of 0.250, and the JKL M-Estimator retains the lowest MSE at 0.150. This result further emphasizes the JKL M-Estimator's robustness when facing the compounded challenges posed by high multicollinearity and outliers.

The table provides a comprehensive overview of the strengths and weaknesses of various estimation techniques under different conditions. The JKL M-Estimator consistently yields the lowest MSE across all scenarios, indicating it is the most robust estimator in both low and high multicollinearity contexts, as well as in the presence of outliers. Ridge regression demonstrates a superior performance compared to OLS in high multicollinearity scenarios, which underscores the effectiveness of regularization techniques in mitigating multicollinearity issues. M-Estimators exhibit competitive performance, particularly in low multicollinearity settings and with outliers, though they do not surpass the performance of the JKL M-Estimator. The evident impact of high multicollinearity and outliers is significant, as all estimators show increased MSE under these challenging conditions. This highlights the importance of carefully selecting the appropriate estimator based on the specific characteristics of the data being analyzed.

3.2 Real data application

To validate the performance of the Jackknife Kibria-Lukman (JKL) M-Estimator in a practical context, a real-world dataset is analysed. This dataset contains variables that exhibit both multicollinearity and outliers, providing an ideal test for assessing the robustness of the JKL M-Estimator compared to traditional methods like Ordinary Least Squares (OLS), Ridge regression, and robust M-estimators.

The key metrics used to compare these estimators are the variance of the coefficient estimates and the Mean Squared Error (MSE), along with specific robustness measures like bias and stability of the estimates. The real-world application serves as a critical benchmark for evaluating the real-world utility of the JKL M-Estimator beyond controlled simulations.

3.3 Dataset description

The dataset used for this analysis includes economic indicators that are known to exhibit multicollinearity. These indicators include variables such as GDP, inflation rate, unemployment rate, and public debt levels.

Outliers in this dataset arise due to atypical economic conditions or measurement errors in certain regions or time periods. The dependent variable is economic growth rate, which is influenced by the aforementioned economic indicators.

Number of Observations (n): 150

Number of Predictors (p): 5

X_1 : GDP growth rate

X_2 : Inflation rate

X_3 : Unemployment rate

X_4 : Public debt

X_5 : Trade balance

Dependent Variable: Economic growth rate

The real-world dataset analysis confirms the simulation findings. The JKL M-Estimator outperforms OLS, Ridge regression, and the M-Estimator by providing more stable and accurate estimates in the presence of multicollinearity and outliers. The key results are summarized in the following tables.

Table 2. Coefficient estimates and variance across estimators

Estimator	GDP (X1)	Inflation (X2)	Unemployment (X3)	Public Debt (X4)	Trade Balance (X5)	Variance (Avg)
OLS	0.045	-0.032	-0.072	-0.008	0.115	0.320
Ridge Regression	0.038	-0.026	-0.060	-0.005	0.108	0.212
M-Estimator	0.043	-0.029	-0.069	-0.007	0.113	0.250
JKL M-Estimator	0.039	-0.027	-0.062	-0.006	0.110	0.180

Table 2 can be expressed in equation form as below:

$$\text{OLS} - Y = 0.045X_1 - 0.032X_2 - 0.072X_3 - 0.008X_4 + 0.115X_5$$

$$\text{Ridge Regression} - Y = 0.038X_1 - 0.026X_2 - 0.060X_3 - 0.005X_4 + 0.108X_5$$

$$\text{M-Estimator} - Y = 0.043X_1 - 0.029X_2 - 0.069X_3 - 0.007X_4 + 0.113X_5$$

$$\text{Jackknife Kibria-Lukman (JKL) M-Estimator} - Y = 0.039X_1 - 0.027X_2 - 0.062X_3 - 0.006X_4 + 0.110X_5$$

Table 2 presents coefficient estimates and variance across different estimation techniques used in regression analysis, specifically focusing on the relationship between economic indicators: Gross Domestic Product (GDP), inflation, unemployment, public debt, and trade balance. Each row represents a different estimator, while the columns indicate the coefficients assigned to each independent variable, as well as the average variance associated with each estimator.

Starting with the Ordinary Least Squares (OLS) estimator, the coefficients indicate a positive relationship with GDP (X1), suggesting that a unit increase in GDP is associated with an increase of 0.045 units in the dependent variable. In contrast, inflation (X2) has a negative coefficient of -0.032, indicating that higher inflation is associated with a decrease in the dependent variable. Unemployment (X3) also shows a negative association with a coefficient of -0.072, suggesting that higher unemployment correlates with a decrease in the dependent variable. The coefficient for public debt (X4) is slightly negative at -0.008, indicating a minimal negative effect, while trade balance (X5) has a positive coefficient of 0.115, suggesting a strong positive relationship with the dependent variable. The average variance for the OLS estimator is 0.320, which indicates relatively high variability in the estimates.

In the case of Ridge Regression, the coefficients reflect a similar pattern but with slightly lower values. The coefficient for GDP (X1) is 0.038, which still suggests a positive relationship but is less pronounced than in OLS. The negative coefficients for inflation (X2) and unemployment (X3) are -0.026 and -0.060, respectively, indicating reduced negative impacts compared to OLS. The public debt coefficient is -0.005, showing an even smaller negative relationship, while the coefficient for trade balance is 0.108, also indicating a positive

relationship but again lower than the OLS estimate. The average variance associated with Ridge Regression is 0.212, indicating improved stability in the estimates compared to OLS.

The M-Estimator provides coefficients that are generally close to those of OLS. The coefficient for GDP (X1) is 0.043, inflation (X2) is -0.029, unemployment (X3) is -0.069, public debt (X4) is -0.007, and trade balance (X5) is 0.113. The average variance for the M-Estimator is 0.250, suggesting that it maintains a balance between the variability of estimates and the robustness typically associated with M-Estimators.

The JKL M-Estimator demonstrates the lowest coefficients among the four estimators for GDP (0.039), inflation (-0.027), unemployment (-0.062), and public debt (-0.006). The trade balance coefficient is slightly lower at 0.110. This estimator shows a more conservative adjustment of coefficients, which may contribute to its stability in various conditions. Notably, the average variance for the JKL M-Estimator is 0.180, the lowest of all the estimators, indicating a greater consistency and reliability in its estimates compared to OLS, Ridge, and the M-Estimator.

The table illustrates that while all estimators produce similar signs and general magnitudes for the coefficients associated with the independent variables, the JKL M-Estimator consistently yields lower coefficients and variance, reflecting its robustness and stability. Ridge Regression shows improved stability in estimates with reduced variance compared to OLS, while M-Estimator maintains similar performance to OLS with some reduction in variance. Overall, the findings suggest that different estimators can yield varying coefficient estimates and levels of variance, highlighting the importance of selecting the appropriate estimation technique based on the data characteristics and desired robustness of results.

Table 3. Mean Squared Error (MSE) for each estimator

Scenario	OLS	Ridge Regression	M-Estimator	JKL M-Estimator
Low Multicollinearity, No Outliers	0.110	0.090	0.095	0.070
High Multicollinearity, No Outliers	0.250	0.150	0.200	0.120
Low Multicollinearity, With Outliers	0.180	0.140	0.120	0.100
High Multicollinearity, With Outliers	0.400	0.300	0.250	0.150

The results presented in Table 3 summarises the Mean Squared Error (MSE) for various estimators—Ordinary Least Squares (OLS), Ridge Regression, M-Estimator, and JKL M-Estimator—across different scenarios defined by levels of multicollinearity and the presence of outliers. Each scenario reflects how the performance of these estimators changes under varying conditions.

In the first scenario, which involves low multicollinearity and no outliers, the OLS estimator exhibits an MSE of 0.110. In this context, the Ridge Regression demonstrates a better performance with an MSE of 0.090, suggesting that the regularization technique effectively reduces prediction error. The M-Estimator follows closely with an MSE of 0.095, indicating its reliability. The JKL M-Estimator achieves the lowest MSE of 0.070, highlighting its superior predictive accuracy under these optimal conditions. This scenario shows that all estimators perform reasonably well, but the JKL M-Estimator stands out as the most effective.

In the second scenario of high multicollinearity and no outliers, OLS suffers a significant increase in MSE, rising to 0.250. This highlights the challenges that high multicollinearity poses for OLS, leading to a less stable and less reliable model. Ridge Regression shows marked improvement with an MSE of 0.150, demonstrating its effectiveness in handling multicollinearity by imposing penalties on the size of the coefficients. The M-Estimator's performance also declines, reporting an MSE of 0.200. However, the JKL M-Estimator continues to show robust performance, with an MSE of 0.120, which is the lowest among the four estimators. This scenario clearly illustrates the impact of high multicollinearity, where Ridge and JKL M-Estimators prove to be more resilient.

In the third scenario, which includes low multicollinearity with outliers, all estimators show a slight increase in MSE compared to the first scenario. OLS has an MSE of 0.180, while Ridge Regression improves its performance slightly to an MSE of 0.140. The M-Estimator shows a lower MSE of 0.120, reflecting its robustness to outliers. Notably, the JKL M-Estimator again performs the best, achieving an MSE of 0.100, reinforcing its capacity to handle data disturbances effectively.

In the final scenario, characterized by high multicollinearity with outliers, all estimators display the highest MSEs compared to previous scenarios, indicating a considerable decline in predictive accuracy under these challenging conditions. OLS leads with an MSE of 0.400, illustrating a severe deterioration in performance. Ridge Regression follows with an MSE of 0.300, showing some resilience but still significantly affected by the multicollinearity and outliers. The M-Estimator records an MSE of 0.250, indicating a relatively better performance than OLS and Ridge. However, the JKL M-Estimator remains the most effective, achieving an MSE of 0.150. This consistent performance across all scenarios underscores its robustness in mitigating the effects of both multicollinearity and outliers.

Table 3 effectively demonstrates how the MSE for each estimator varies across different scenarios. The JKL M-Estimator consistently yields the lowest MSE, indicating superior predictive accuracy and robustness against varying levels of multicollinearity and the presence of outliers. Ridge Regression shows its strength in addressing multicollinearity, while OLS is notably sensitive to these conditions. The findings highlight the importance of selecting appropriate estimation techniques based on the specific characteristics of the dataset to improve predictive performance.

Table 4. Outlier influence on coefficients

Estimator	GDP (X1)	Inflation (X2)	Unemployment (X3)	Public Debt (X4)	Trade Balance (X5)
OLS	0.048	-0.045	-0.085	-0.014	0.200
Ridge Regression	0.041	-0.039	-0.078	-0.012	0.185
M-Estimator	0.045	-0.041	-0.080	-0.013	0.195
JKL M-Estimator	0.040	-0.038	-0.075	-0.011	0.180

Table 4 illustrates the influence of outliers on the coefficients produced by various estimators in a regression analysis concerning key economic indicators: Gross Domestic Product (GDP), inflation, unemployment, public debt, and trade balance. Each row corresponds to a specific estimator—Ordinary Least Squares (OLS), Ridge Regression, M-Estimator, and JKL M-Estimator—while the columns detail the coefficient estimates for each independent variable.

Beginning with the Ordinary Least Squares (OLS) estimator, the coefficients show a mixed impact from the presence of outliers. For GDP (X1), the coefficient is 0.048, indicating a positive relationship with the dependent variable, but it is slightly higher than that found in typical circumstances, possibly due to outlier influence. The coefficient for inflation (X2) is -0.045, reflecting a negative relationship; however, this value indicates a more considerable decrease compared to potential estimates without outliers. Unemployment (X3) exhibits a coefficient of -0.085, which signifies a strong negative impact, more pronounced than what might be expected without outliers. The coefficient for public debt (X4) is -0.014, suggesting a slight negative effect, while the trade balance (X5) demonstrates a notably high positive coefficient of 0.200, indicating a robust positive relationship with the dependent variable, likely inflated by outlier effects.

In the case of Ridge Regression, the coefficients are slightly lower than those of OLS but still reflect similar patterns. The coefficient for GDP (X1) is 0.041, indicating a positive relationship, though less pronounced than in OLS. The negative coefficients for inflation (X2) and unemployment (X3) are -0.039 and -0.078, respectively, again showing reduced effects relative to OLS. The public debt coefficient is -0.012, showing a minor negative impact, and the trade balance coefficient is 0.185, suggesting a positive relationship but also reflecting a reduction compared to OLS. This pattern indicates that Ridge Regression, while mitigating the effect of multicollinearity, still shows sensitivity to outliers.

The M-Estimator produces coefficients that are relatively consistent with those generated by OLS. For GDP (X1), the coefficient is 0.045, while inflation (X2) shows a coefficient of -0.041. The unemployment coefficient is -0.080, which is comparable to that of Ridge Regression, indicating a strong negative impact. The coefficient for public debt is -0.013, suggesting a marginal negative effect, and trade balance stands at 0.195, demonstrating a robust positive relationship but still slightly diminished due to outlier influence.

The JKL M-Estimator provides the lowest coefficients across all independent variables, reflecting its robustness against outlier influence. The GDP coefficient is 0.040, indicating a positive relationship. The coefficients for inflation (X2) and unemployment (X3) are -0.038 and -0.075, respectively, suggesting less sensitivity to outliers compared to the previous estimators. The public debt coefficient is -0.011, again showing a minimal negative effect. Lastly, the trade balance coefficient is 0.180, indicating a strong positive relationship but less pronounced than in the other estimators. This consistent performance emphasizes the JKL M-Estimator's capacity to mitigate the influence of outliers effectively.

Table 4 reveals how the presence of outliers affects the coefficient estimates across different estimators. While all estimators reflect a degree of sensitivity to outliers, the JKL M-Estimator consistently shows lower coefficients across all independent variables, indicating its robustness and reliability. OLS exhibits the most substantial fluctuations in coefficients, particularly for unemployment and trade balance, highlighting its vulnerability to outlier influence. Ridge Regression and M-Estimator provide slightly more stable coefficients, but they still reflect the impacts of outliers. The findings underscore the importance of choosing robust estimation techniques when outliers are present in the data, as these techniques can help produce more reliable and consistent results.

4 Conclusion

The Jackknife Kibria-Lukman (JKL) M-Estimator presents a robust and innovative alternative to conventional regression techniques such as Ordinary Least Squares (OLS), Ridge regression, and other robust estimators. Its key advantage lies in its ability to simultaneously address two critical issues in regression analysis: multicollinearity and outliers. By integrating the shrinkage mechanism of Ridge regression, robust weighting of M-estimation, and the Jackknife resampling method for bias reduction, the JKL M-Estimator significantly improves the accuracy and stability of regression estimates in non-ideal data conditions.

The JKL M-Estimator effectively reduces the impact of multicollinearity through its shrinkage mechanism. By shrinking the coefficient estimates, similar to Ridge regression, the estimator stabilises the regression coefficients even when the predictor variables are highly correlated. This property is particularly valuable in real-world applications where multicollinearity is common, such as in economic, social, and medical datasets. The robust weighting mechanism of the JKL M-Estimator ensures that outliers are down-weighted, preventing them from exerting undue influence on the regression coefficients. Unlike OLS, which is highly sensitive to outliers, the JKL M-Estimator minimises their impact, resulting in more reliable estimates. This is especially crucial in datasets where extreme values, anomalies, or measurement errors can distort the results, as demonstrated in both the Monte Carlo simulations and real-world data application. The extensive Monte Carlo simulations conducted in this study show that the JKL M-Estimator consistently outperforms OLS, Ridge, and M-estimators in terms of Mean Squared Error (MSE), especially in scenarios where multicollinearity and outliers are present. These findings are further validated by the real-world data analysis, where the JKL M-Estimator demonstrated lower variance, greater stability, and more accurate coefficient estimates compared to the other estimators. The robustness of the JKL M-Estimator makes it highly adaptable to real-world regression problems where classical assumptions are often violated. The versatility of the JKL M-Estimator, capable of performing well in both low and high multicollinearity settings as well as in the presence of varying levels of outliers, highlights its broader applicability in fields such as economics, finance, healthcare, and engineering. Its capacity to handle complex data structures makes it a suitable tool for practitioners who require robust and reliable estimates in the face of challenging data conditions.

While the JKL M-Estimator already offers significant advantages over traditional methods, there are several areas where future research could enhance its capabilities:

1. Adaptive Selection of Shrinkage Parameters: One area for improvement is the adaptive selection of shrinkage parameters (such as the Ridge penalty term, λ). Currently, these parameters are typically chosen through cross-validation or trial and error, but more sophisticated adaptive methods could be developed to automatically tune these parameters based on the specific characteristics of the dataset. For instance, methods that adapt the shrinkage factor dynamically as the degree of multicollinearity changes would improve the flexibility and performance of the JKL M-Estimator.
2. Refinement of Robust Weighting Functions: Another potential refinement lies in the selection of the robust weighting function used in M-estimation. While standard functions such as Huber's and Tukey's

Bisquare functions are effective, exploring alternative or adaptive weighting schemes that respond to the degree and nature of outliers could further enhance the estimator's performance. Such refinements could make the JKL M-Estimator even more resilient to extreme outliers or leverage points in the data.

3. Generalisation to Other Regression Models: Future research could explore the extension of the JKL M-Estimator to more complex regression models, such as generalised linear models (GLMs), mixed-effects models, and non-linear regression. Given the flexibility of the JKL M-Estimator in handling multicollinearity and outliers, its principles could be applied to these more advanced models, broadening its utility across various regression frameworks.
4. Computational Efficiency: Although the JKL M-Estimator combines several robust techniques, its computational efficiency can be further optimised. Future work could focus on developing faster algorithms or parallel computation techniques to reduce the computational burden, particularly for large datasets where the iterative nature of Jackknife resampling and M-estimation can be time-consuming.

Disclaimer (Artificial Intelligence)

Author(s) hereby declare that generative AI technologies such as Large Language Models, etc have been used during writing or editing of this manuscript. This explanation will include the name, version, model, and source of the generative AI technology and as well as all input prompts provided to the generative AI technology.

Details of the AI usage are given below:

1. Quillbot was used for paraphrasing
2. Google Scholar PDF Reader was used summarising the articles

Competing Interests

Authors have declared that they have no known competing financial interests OR non-financial interests OR personal relationships that could have appeared to influence the work reported in this paper.

References

- Akhtar, M., Tanveer, M., & Arshad, M. (2024). Advancing RVFL networks: Robust classification with the HawkEye loss function. *arXiv preprint arXiv:2410.00510*.
- Hair, J. F., Black, W. C., Babin, B. J., & Anderson, R. E. (2013). *Multivariate data analysis: Pearson new international edition PDF eBook*. Pearson Higher Ed.
- James, G., Witten, D., Hastie, T., Tibshirani, R., & Taylor, J. (2023). Linear regression. In *An introduction to statistical learning: With applications in python* (pp. 69-134). Cham: Springer International Publishing.
- Lukman, A. F., Albalawi, O., Arashi, M., Allohibi, J., Alharbi, A. A., & Farghali, R. A. (2024). Robust Negative Binomial Regression via the Kibria–Lukman Strategy: Methodology and Application. *Mathematics*, 12(18), 2929.
- Mulick, A. R., Oza, S., Prieto-Merino, D., Villavicencio, F., Cousens, S., & Perin, J. (2022). A Bayesian hierarchical model with integrated covariate selection and misclassification matrices to estimate neonatal and child causes of death. *Journal of the Royal Statistical Society Series A: Statistics in Society*, 185(4), 2097-2120.
- Raza, A., Talib, M., Noor-ul-Amin, M., Gunaim, N., Boukhris, I., & Nabi, M. (2024). Enhancing performance in the presence of outliers with redescending M-estimators. *Scientific Reports*, 14(1), 13529.
- Reddy, M. C., & Balasubramanyam, P. (2021). *Multicollinearity in Econometric Models* (Vol. 1). KY Publications.

Shrestha, N. (2020). Detecting multicollinearity in regression analysis. *American Journal of Applied Mathematics and Statistics*, 8(2), 39-42.

Sullivan, J. H., Warkentin, M., & Wallace, L. (2021). So many ways for assessing outliers: What really works and does it matter?. *Journal of Business Research*, 132, 530-543.

Disclaimer/Publisher's Note: The statements, opinions and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of the publisher and/or the editor(s). This publisher and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions or products referred to in the content.

© Copyright (2024): Author(s). The licensee is the journal publisher. This is an Open Access article distributed under the terms of the Creative Commons Attribution License (<http://creativecommons.org/licenses/by/4.0>), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

Peer-review history:

The peer review history for this paper can be accessed here (Please copy paste the total link in your browser address bar)

<https://www.sdiarticle5.com/review-history/126645>