


Robust Multivariate Regression Based on Shrinkage S_n Estimator

Lakshmi R 

Department of Statistics and Data Science
Assistant Professor
Christ (Deemed to be University)
Bengaluru, India

Sajesh T A 

Department of Statistics
Associate Professor
St. Thomas College (Autonomous)
Thrissur, India

Abstract

The primary goal of multivariate regression analysis is to estimate model parameters. However, when the data includes outliers or extreme observations, the maximum likelihood estimator may not be suitable for estimation. Therefore, it's essential to identify a parameter estimation method that remains relatively unaffected by minor changes in the data. In this paper, we present a robust approach to multivariate regression that relies on the robust estimation of the joint location and scatter matrix for both the explanatory and response variables. We make use of shrinkage based robust location and scatter matrix proposed by Lakshmi and Sajesh (2025). Through simulations, we explore the finite-sample performance and robustness of the estimator. To improve efficiency, we suggest a reweighted estimator selected from multiple reweighting options. We demonstrate that the multivariate regression estimator possesses the equivariance properties. The proposed estimator achieves a balance of high robustness and efficiency in estimation. Proposed estimator's efficacy is illustrated using no: of benchmark dataset.

Keywords: multivariate regression, shrinkage S_n estimator.

1. Introduction

In statistical modeling, regression analysis is a method used to estimate the relationships between variables. It involves a range of techniques for modeling and analyzing multiple variables, with an emphasis on the relationship between dependent (response) variables and independent (predictor) variables. More precisely, regression analysis helps to understand how the expected value of the dependent variables shifts when any of the independent variables are altered. As discussed, the multivariate regression model allows us to examine the influence of several variables on one or more dependent variables within the same model.

Let $\mathbf{x} = (x_1, x_2, \dots, x_p)^t$ be a p - dimensional predictor and a q - dimensional response $\mathbf{y} = (y_1, y_2, \dots, y_p)^t$. Consider a multivariate regression model $\mathbf{y} = \mathbf{B}^t \mathbf{x} + \alpha + \varepsilon$, \mathbf{B} is $p \times q$ slope matrix, α is q -dimensional intercept vector, and $\varepsilon = (\varepsilon_1, \dots, \varepsilon_q)^t$ denotes the *i.i.d.* error term with mean zero and $\text{cov}(\varepsilon) = \Sigma_\varepsilon$, is a $q \times q$ positive definite matrix. Let μ denote the

location of joint variables (\mathbf{x}, \mathbf{y}) and Σ denotes their scatter matrix. Partition matrix of μ and Σ be as follows: $\mu = \begin{pmatrix} \mu_x \\ \mu_y \end{pmatrix}$ and $\Sigma = \begin{pmatrix} \Sigma_{xx} & \Sigma_{xy} \\ \Sigma_{yx} & \Sigma_{yy} \end{pmatrix}$. Conventional maximum likelihood estimates of μ and Σ , empirical mean $\hat{\mu}$ and empirical covariance matrix $\hat{\Sigma}$ are often used as estimates. The components of the resulting estimators $\hat{\mu}$ and $\hat{\Sigma}$ are utilized in the least squares equations (Johnson and Wichern 2002) as follows: It is widely recognized that classical multiple regression is highly sensitive to outliers in the data (see, Suhail, Chand, and Kibria (2021), Lukman, Farghali, Kibria, and Oluyemi (2023), Wasim, Zaman, Ahmad, and Kibria (2025), Alghamdi, Hammad, Golam Kibria, Abd-Elmougod, Sapkota, and Gemeay (2025), Yasmin and Kibria (2025)). This issue is equally prevalent in the context of multivariate regression, since classical $\hat{\mu}$ and $\hat{\Sigma}$ are sensitive in the presence of anomalies. To address this issue, one can substitute the classical estimates of location and scatter with highly robust estimates that are less sensitive to outliers, enabling a more robust analysis. Maronna and Yohai (1997) provide an overview of robust multivariate regression algorithms in the context of simultaneous equations models. Although Koenker and Portnoy (1990) presented a M-type approach, their estimator lacked affine equivariance. Rousseeuw, Van Aelst, Van Driessen, and Gulló (2004) proposed estimator based on the robust estimate of the location and dispersion of the joint distribution of the (\mathbf{x}, \mathbf{y}) variables using Minimum Covariance Determinant, which is computationally time consuming and expensive. Ollila, Oja, and Hettmansperger (2002) and Ollila, Oja, and Koivunen (2003) introduced multivariate regression based on rank covariance matrices and evaluated their properties in their paper. The proposed estimator possess good efficiency, but fails to perform in usual sense of robustness. A multivariate regression extension of the least trimmed squares estimator (MLTS) was investigated by Agulló, Croux, and Van Aelst (2008), in which slope matrix is obtained as minimum of the determinant of the robust MCD scatter matrix of the residuals. Van Aelst and Willems (2005) introduced S-estimators in multivariate regression similar to that of MLTS. Both LTS and S based multivariate regression estimators performs with least squares as initial fit. Sajana and Sajesh (2018) introduced multivariate regression estimator based on Comedian covariance and empirically developed the properties of the proposed estimator, also checked the performance through simulation study. The estimator performs better in terms of robustness and efficiency with respect to other estimators compared in their paper.

In this study we propose to make use of reweighted, Shrinkage based S_n covariance estimator and Shrinkage L_1 median as robust alternatives to the classical estimates of Σ and μ respectively, which developed by Lakshmi and Sajesh (2025), multivariate regression estimator. We make use of the proposed estimator here and compare it with classical MLE, Rousseeuw *et al.* (2004), Orthogonalized Gnanadesikan - Kettenring (OGK) estimates of Maronna and Zamar (2002) based multivariate regression estimates, Sajana and Sajesh (2018) and Kunjunni and Abraham (2022) based multivariate regression estimates. Unlike MCD-based regression and S-estimators, which rely on high-breakdown but computationally intensive covariance estimation and often suffer from efficiency loss under moderate contamination, the proposed method is built on a multivariate extension of the S_n scale estimator combined with shrinkage regularization. This avoids subsampling or iterative re-weighting schemes inherent to MCD and S-procedures, leading to improved numerical stability and scalability in moderate-to-high dimensions. In contrast to OGK-based approaches, which depend on pairwise scale estimates and orthogonalization that may propagate outlier effects across components, the proposed estimator directly incorporates robust scale information within a shrinkage framework. Moreover, the present approach extends S_n to a multivariate shrinkage setting, explicitly addressing multicollinearity and high-dimensional structure. As a result, the proposed method achieves a favorable balance between robustness, efficiency, and regularization that is not simultaneously attained by existing approaches.

Next section of this article describes our proposed multivariate regression estimator and following section shows the simulation study on robustness and efficiency property of the proposed estimator. We also prove the equivariance property empirically. We have shown the

real - world application and lastly gives the conclusion of our work.

2. Shrinkage S_n multivariate regression

The principle behind shrinkage estimation lies in the notion of “shrinking” an estimator \hat{E} towards a target estimator \hat{T} , which serves to effectively diminish estimation errors. Utilizing a shrinkage estimator offers a significant benefit by balancing bias and variance. Cabana, Lillo, and Laniado (2021) proposed the shrinkage estimator of L_1 - median as a robust alternative to the location. And the shrinkage estimator based on L_1 -median is defined as:

$$\hat{\boldsymbol{\mu}}_{Sh} = (1 - \eta)\hat{\boldsymbol{\mu}}_{MM} + \eta\nu_{\boldsymbol{\mu}}\mathbf{e}, \quad (1)$$

where $\nu_{\boldsymbol{\mu}}\mathbf{e}$ is the shrinkage target matrix, \mathbf{e} is a vector of ones with p - dimension and $\hat{\boldsymbol{\mu}}_{MM}$ is the L_1 -median from the samples. Scaling factor $\nu_{\boldsymbol{\mu}}$ and the shrinkage intensity η should be such that, they minimize the expected quadratic loss. S_n covariance of two random variables X and Y be:

$$S_n(X, Y) = 1.4304(\text{med}_i[\text{med}_{j \neq i}\{(x_i - x_j)(y_i - y_j)\}]).$$

Let \mathbf{X} be $n \times p$ matrix with sample size n , number of variables p , and $\mathbf{X}_j (j = 1, 2, \dots, p)$ be the column of the matrix. The covariance matrix of \mathbf{X} based on S_n would be: $\hat{S}_n = S_n(\mathbf{X}_i, \mathbf{X}_j)$. The influence function of S_n covariance used in shrinkage S_n matrix is bounded (Croux, Rousseeuw, and Hössjer 1994). Breakdown of univariate median is nearly 50% because out of n points, if $[(n - 1)/2]$ points changed, median remains bounded. The S_n covariance used for obtaining Shrinkage S_n is nothing but two folded repeated median which is bounded and has an asymptotic breakdown of 50% by the following theorem improved from Siegel (1982). The repeated median estimator S_n will remain bounded whenever more than $[(n - 1)/2]$ points of n observations are confined fixed while the remaining points are arbitrarily moved. By established results for repeated median estimators (Siegel 1982) and S_n covariance estimators, this functional possesses a bounded influence function and an asymptotic breakdown point of 50%. Ledoit and Wolf (2004) and DeMiguel, Martin-Utrera, and Nogales (2013) shown shrinkage covariance estimation technique itself a consistent estimation technique too. Consequently, the regression estimator constructed from the joint Shrinkage S_n -based scatter matrix of the explanatory and response variables inherits these robustness properties. Moreover, the subsequent reweighting step, based on robust residuals from the initial S_n -based fit, improves statistical efficiency while preserving robustness (Lopuhaa and Rousseeuw 1991; Lopuhaä 1999).

The robust version of covariance matrix based on Shrinkage S_n would be:

$$\hat{\boldsymbol{\Sigma}}_{Sh} = (1 - \eta)\hat{E} + \eta\hat{T}, \quad \text{where} \quad \hat{E} = \hat{S}_n. \quad (2)$$

Lakshmi and Sajesh (2025) proposed the shrinkage-based S_n covariance matrix as a robust alternative to the traditional covariance estimator. The shrinkage estimator defined in (2), where T is the shrinkage target matrix, the shrinkage intensity η is estimated from the data, for which theoretical justification is provided in the Supplementary Material. In practice, the shrinkage intensity η is computed once using the closed-form expression derived in the Supplementary Material and is directly used in the construction of the shrinkage-based S_n covariance estimator. In this paper we utilize the above defined location estimate and covariance matrix estimate and there by propose a reweighted regression estimator.

Consider $\mathbf{z} = (\mathbf{x}, \mathbf{y})$, the joint variable with location and covariance matrix $\boldsymbol{\mu}$, $\boldsymbol{\Sigma}$ respectively. The associated squared Mahalanobis distance for each observation $\mathbf{z}_i, i = 1, 2, 3, \dots, n$, based on initial estimates of Shrinkage L_1 median $\hat{\boldsymbol{\mu}}_{Sh}$ and Shrinkage S_n covariance matrix $\hat{\boldsymbol{\Sigma}}_{Sh}$ be:

$$\text{RD}^2(\mathbf{z}_i) = (\mathbf{z}_i - \hat{\boldsymbol{\mu}}_{Sh})^t \hat{\boldsymbol{\Sigma}}_{Sh}^{-1} (\mathbf{z}_i - \hat{\boldsymbol{\mu}}_{Sh}). \quad (3)$$

The weight function based on above defined robust Mahalanobis distance be $w_i = w(\text{RD}^2(\mathbf{z}_i))$, where a weight of 1 be assigned to the observations (\mathbf{z}_i) with a Mahalanobis distance less than $\frac{\chi_{0.95,p}^2 \times \text{median}(\text{RD}^2(\mathbf{z}_i))}{\chi_{0.5,p}^2}$. The cutoff used for the robust distance-based weights follows the same construction as that proposed by Maronna and Zamar (2002) and Lakshmi and Sajesh (2025) for robust distance measures. Thus the reweighted shrinkage location and S_n covariance matrix be defined as:

$$\hat{\boldsymbol{\mu}}^1 = \frac{\sum_{i=1}^n w_i \mathbf{z}_i}{\sum_{i=1}^n w_i}, \quad \hat{\boldsymbol{\Sigma}}^1 = \frac{\sum_{i=1}^n w_i (\mathbf{z}_i - \hat{\boldsymbol{\mu}}^1)(\mathbf{z}_i - \hat{\boldsymbol{\mu}}^1)^t}{\sum_{i=1}^n w_i}. \quad (4)$$

It makes sense to use weights in regression analyses based on the residuals from the original fit (Rousseeuw and Leroy 1987). Let the residuals based on weighted regression estimates (WR) be:

$$\mathbf{r}_i = \mathbf{y}_i - \hat{\mathbf{B}}^t \mathbf{x}_i - \hat{\boldsymbol{\alpha}}. \quad (5)$$

Now we again reweight the WR estimator based on the residuals defined above. Consider the weights $w_{\mathbf{r}_i} = w(\text{RD}^2(\mathbf{r}_i))$ which assign a weight 1 to the residuals (\mathbf{r}_i) with Mahalanobis distance less than $\chi_{1,0.99}^2$, where $\text{RD}^2(\mathbf{r}_i)$ be the Mahalanobis distance of WR residuals, defined as $\text{RD}^2(\mathbf{r}_i) = \mathbf{r}_i^t (\hat{\boldsymbol{\Sigma}}_\epsilon)^{-1} \mathbf{r}_i$. The final reweighted regression estimators be:

$$\mathbf{T}^R = \left(\sum_{i=1}^n w_{\mathbf{r}_i} \mathbf{u}_i \mathbf{u}_i^t \right)^{-1} \sum_{i=1}^n w_{\mathbf{r}_i} \mathbf{y}_i \mathbf{u}_i, \quad (6)$$

and

$$\hat{\boldsymbol{\Sigma}}_\epsilon^R = \frac{\sum_{i=1}^n w_{\mathbf{r}_i} (\mathbf{r}_i^R)_i (\mathbf{r}_i^R)_i^t}{\sum_{i=1}^n w_{\mathbf{r}_i}}, \quad (7)$$

where $\mathbf{T}^R = ((\hat{\mathbf{B}}^R)^t, \hat{\boldsymbol{\alpha}}^R)^t$, $\mathbf{u}_i = (\mathbf{x}_i^t, 1)^t$ and $(\mathbf{r}_i^R)_i = \mathbf{y}_i - (\hat{\mathbf{B}}^R)^t \mathbf{x}_i - \hat{\boldsymbol{\alpha}}^R$ and \mathbf{T}^R be the proposed reweighted shrinkage multivariate regression estimator (SS_n). The superscript R implies the weights were based on initial regression. The robustness of these reweighted regression estimators is derived from the properties of the initial regression estimators. It is important to note that the weights now depend solely on the magnitude of the residual distance $w_{\mathbf{r}_i}$. Unlike the initial estimates, good leverage points are no longer down weighted.

Algorithm 1: Shrinkage \mathbf{S}_n -based reweighted multivariate regression (SS_n)

1. Form the joint observations $\mathbf{z} = (\mathbf{x}, \mathbf{y})$.
2. Compute the $\hat{\boldsymbol{\mu}}_{Sh}$ and covariance $\hat{\boldsymbol{\Sigma}}_{Sh}$.
3. Compute the shrinkage S_n covariance estimator $\hat{\boldsymbol{\Sigma}}_{Sh} = (1 - \eta)\hat{\boldsymbol{\Sigma}}_{S_n} + \eta T$, where the shrinkage intensity η is obtained from the closed-form expression derived in the Supplementary Material.
4. Compute robust squared Mahalanobis distances

$$\text{RD}^2(\mathbf{z}_i) = (\mathbf{z}_i - \hat{\boldsymbol{\mu}}_{Sh})^t \hat{\boldsymbol{\Sigma}}_{Sh}^{-1} (\mathbf{z}_i - \hat{\boldsymbol{\mu}}_{Sh}). \quad (8)$$

5. Define weights $w_i = w(\text{RD}^2(\mathbf{z}_i))$ based on $\text{RD}^2(\mathbf{z}_i)$.
6. Compute weighted location and scatter estimators

$$\hat{\boldsymbol{\mu}}^1 = \frac{\sum_{i=1}^n w_i \mathbf{z}_i}{\sum_{i=1}^n w_i}, \quad \hat{\boldsymbol{\Sigma}}^1 = \frac{\sum_{i=1}^n w_i (\mathbf{z}_i - \hat{\boldsymbol{\mu}}^1)(\mathbf{z}_i - \hat{\boldsymbol{\mu}}^1)^t}{\sum_{i=1}^n w_i}. \quad (9)$$

7. Compute residuals

$$\mathbf{r}_i = \mathbf{y}_i - \hat{\mathbf{B}}^t \mathbf{x}_i - \hat{\boldsymbol{\alpha}}. \quad (10)$$

8. Consider $w\mathbf{r}_i = w(\text{RD}^2(\mathbf{r}_i))$ which assign a weight 1 to the residuals (\mathbf{r}_i) with $\text{RD}^2(\mathbf{r}_i)$ less than $\chi_{1,0.99}^2$.

9. Compute the final reweighted regression estimator $\mathbf{T}^R = ((\hat{\mathbf{B}}^R)^t, \hat{\boldsymbol{\alpha}}^R)^t$.

Note. The shrinkage intensity η is computed once and is not updated iteratively.

3. Efficiency

To assess the efficiency of the proposed regression estimator, we conducted the following simulation study. For different sample sizes n and various choices of p and q , we generated r datasets of size n from the multivariate standard Gaussian distribution $N(0, I_{p+q})$, where $\mathbf{B} = 0$ and $\boldsymbol{\alpha} = 0$. For each data set ($l = 1, 2, \dots, r$), we conduct the reweighted regression estimation as defined in Section 2 based on shrinkage estimators, yielding $p \times q$ slope estimate $\hat{\mathbf{B}}^{(l)}$, intercept estimate $\hat{\boldsymbol{\alpha}}^{(l)}$ and the covariance matrix estimate $\hat{\boldsymbol{\Sigma}}_{\boldsymbol{\varepsilon}}^{(l)}$ of the errors.

The variance estimate of slope coefficient is obtained as:

$$\text{var}(\hat{\mathbf{B}}_{jk}^{(l)}) = n\text{var}(\hat{\mathbf{B}}_{jk}^{(l)}), \quad \text{where } j = 1, 2, \dots, p \quad \text{and} \quad k = 1, 2, \dots, q. \quad (11)$$

Then the corresponding finite sample efficiency of the slope estimate is defined as $1/\text{ave}_{j,k}(\text{var}(\hat{\mathbf{B}}_{jk}))$. Similarly, we calculate the finite-sample efficiency of the intercept vector. To assess the accuracy of the error scatter matrix, we use the standardized variance (Bickel and Lehmann 2012) of the elements of the error covariance matrix, defined as follows:

$$\text{st.var}(\hat{\boldsymbol{\Sigma}}_{\boldsymbol{\varepsilon}jk}) = \frac{n\text{var}_l((\hat{\boldsymbol{\Sigma}}_{\boldsymbol{\varepsilon}}^{(l)})_{jk})}{[\text{ave}_l \text{ave}_j((\hat{\boldsymbol{\Sigma}}_{\boldsymbol{\varepsilon}}^{(l)})_{jj})]^2} \quad \text{for } j = 1, \dots, q \quad \text{and} \quad k = 1, \dots, q.$$

The overall finite-sample efficiency of the off-diagonal elements is then defined as follows $1/\text{ave}_{j \neq k}(\text{st.var}((\hat{\boldsymbol{\Sigma}}_{\boldsymbol{\varepsilon}})_{jk}))$. And the finite sample efficiency of diagonal elements is given by $2/\text{ave}_j(\text{st.var}((\hat{\boldsymbol{\Sigma}}_{\boldsymbol{\varepsilon}})_{jj}))$. The Table 1 and Tables 24, 25, 26 in Supplementary material show the results. It is very clear, under non outlier scenario MLE is performing well. And our proposed estimator stands next to MLE showing better performance than other robust estimators by possessing least MSE, bias.

4. Robustness

We performed simulations to investigate the finite-sample robustness with dataset containing outliers. A vertical outlier is a point $(\mathbf{x}_i, \mathbf{y}_i)$ whose \mathbf{x}_i is not outlying but does not follow the linear trend of the majority of the data. A point $(\mathbf{x}_i, \mathbf{y}_i)$ where \mathbf{x}_i is an outlier is referred to as a leverage point. If this $(\mathbf{x}_i, \mathbf{y}_i)$ deviates from the pattern of the majority, it is termed a bad leverage point. Conversely, if it aligns with the majority pattern and does not negatively impact the fit, it is considered a good leverage point. Regression estimators often fail when confronted with vertical outliers or bad leverage points. In this study we created datasets that include both types of outliers for the evaluation of estimators. For a given sample size n , we generate $r = 1000$ datasets from multivariate standard Normal distribution with mean 0 and identity matrix I_{p+q} as variance covariance matrix. For the evaluation purpose, we consider 10%, 20% and 40% contamination in datasets. Keeping \mathbf{x}_i , the q response variables are taken from multivariate normal distribution $N(2\sqrt{\chi_{p+q,0.99}^2}, I_q)$. This produces vertical outliers, because here only response variables are outlying. We also replaced the data with

Table 1: Efficiency table, $n = 100$

Estimators	Slope	Intercept	Σ_{diag}	$\Sigma_{offdiag}$
$p = 4, q = 4$				
SS_n	1.133559	0.987338	2.26897	1.166734
MLE	1.098142	1.062179	2.179757	1.118315
Comedian	1.208254	1.036421	2.626631	1.28434
MCD	1.855456	1.611511	4.231406	2.018595
OGK	1.251197	1.134508	2.785136	1.343533
S_n	1.055214	1.006972	2.244428	1.14198
$p = 4, q = 8$				
SS_n	1.069516	1.05058	2.171636	1.082605
MLE	1.079249	1.052644	2.194896	1.087536
Comedian	1.216045	1.096154	2.454006	1.268528
MCD	1.781371	1.614493	3.718948	1.850268
OGK	1.253904	1.080404	2.717749	1.311943
S_n	1.074581	1.0278	2.206211	1.093251
$p = 8, q = 4$				
SS_n	1.187912	1.027722	2.450559	1.203182
MLE	1.131539	1.114378	2.46635	1.156902
Comedian	1.27453	1.158034	2.739073	1.334702
MCD	2.173524	1.856234	4.767762	2.382701
OGK	1.324617	1.135834	2.935003	1.443301
S_n	1.132202	1.067909	2.397812	1.122268
$p = 10, q = 10$				
SS_n	1.197286	1.11769	2.353268	1.186694
MLE	1.146476	1.141881	2.296645	1.135159
Comedian	1.318311	1.204311	2.674198	1.319547
MCD	2.350712	2.063462	4.494091	2.282524
OGK	1.414083	1.240724	2.915869	1.421498
S_n	1.261643	1.165052	2.627408	1.261015

bad leverage points for which the p independent variables are generated according to multivariate normal distribution $N(2\sqrt{\chi_{p,0.99}^2}, I_p)$ and the q dependent variables are generated from multivariate normal $N(2\sqrt{\chi_{q,0.99}^2}, I_q)$. Apart from the mentioned distribution, we also considered observations from multivariate normal $N(2\sqrt{\chi_{p+q,0.99}^2}, 0.1 \times I_q)$, $N(2\sqrt{\chi_{q,0.99}^2}, 0.1 \times I_q)$, $N(2\sqrt{\chi_{p,0.99}^2}, 0.1 \times I_p)$ respectively for vertical outliers and bad leverage points. Rather than using observations from different distributions, which would make it easier for the estimators to detect outliers and perform, we opted to use the same multivariate normal distribution with varying parameters. As in efficiency study, we generate each dataset ($l = 1, 2, \dots, r$), compute the slope matrix estimate $\hat{\mathbf{B}}^{(l)}$, the intercept estimate $\hat{\boldsymbol{\alpha}}^{(l)}$ and the covariance matrix estimate $\hat{\boldsymbol{\Sigma}}_{\boldsymbol{\varepsilon}}^{(l)}$ of the errors. In order to measure robustness, we make use of bias and mean squared error (MSE). As in the case of uni variate component, the bias and MSE of the slope are defined as:

$$\text{bias}(\hat{\mathbf{B}}) = \sqrt{\text{ave}_{j,k}(\text{bias}(\hat{\mathbf{B}}_{jk})^2)}$$

and

$$\text{MSE}(\hat{\mathbf{B}}) = \text{ave}_{j,k}(\text{MSE}(\hat{\mathbf{B}}_{jk})).$$

Similarly bias, MSE for the intercept $\hat{\boldsymbol{\alpha}}$, for diagonal and off diagonal elements of $\hat{\boldsymbol{\Sigma}}_{\boldsymbol{\varepsilon}}$ are calculated. Only one table of robustness property for dimension $(p, q) = (4, 4)$ is given in Table 2, remaining tables (more than 22 tables of combination of dimension and sample size) of all dimension and sample sizes considered in the study for robustness property provided in Supplementary material. The overall results from the supplementary tables show that irrespective of dimension and sample size our proposed estimator shows least MSE, bias than other robust estimators compared in the study. For higher contamination, especially for 20%, 40% also, our proposed estimator exhibits better performance in terms of MSE, bias, i.e., SS_n based multivariate regression estimator possesses the least MSE, bias than other compared estimators. This gave un an empirical evidence on robustness of proposed regression estimator.

Generalized approaches to regression, scale, affine equivariance, and the robustness of multiple regression estimators were introduced by Rousseeuw and Leroy (1987). Regression equivariance implies that if a linear function of the explanatory variables is added to the responses, the coefficients of that linear function are similarly added to the estimator. The \mathbf{y} - equivariance of the estimator means that a linear transformation of the response variables results in the estimator being transformed in the same way. \mathbf{x} - equivariance indicates that if the predictor variables undergo a linear transformation, the estimator will transform correspondingly. The three equivariance properties are proven for our proposed estimator empirically by considering contaminated scenarios same as for checking robustness, additionally 0% contamination too considered. MSE values are considered for the evaluation. Results are tabulated in Tables 3 and 4 given below.

Under moderate to severe contamination, the proposed shrinkage SS_n -based reweighted estimator consistently exhibits lower mean squared error and improved stability compared to classical estimators and several established robust alternatives. This advantage is particularly pronounced for larger dimensions and higher contamination levels, where methods such as MCD, OGK, and Comedian often suffer from increased variability or breakdown.

In contrast, under clean or near-clean data settings, classical estimators based on maximum likelihood tend to be more efficient, as expected, and may outperform robust procedures including the proposed method. This behavior reflects the usual robustness–efficiency trade-off and indicates that the SS_n estimator is not designed to replace the MLE in ideal conditions, but rather to provide reliable performance in the presence of contamination.

Table 2: Robustness table - contaminants from $N\left(2\sqrt{\chi_{p+q,0.99}^2}, 0.1 \times I_q\right)$ with delta = 10%, $p = 4, q = 4, \text{diag} = 0.1$

	Slope MSE	Intercept MSE	Diagonal element mse	Non-diagonal elements mse	Slope bias	Intercept bias	Diagonal element bias	Non-diagonal element bias
<i>n</i> = 50								
S_n	0.027062	0.024811	0.869399	0.02140	0.003809	0.001816	0.908377	-0.00498
MCD	0.056792	0.029701	0.805927	0.02885	-0.00122	-7.10e-5	0.880357	-0.00533
Comedian	0.027862	0.024672	0.852523	0.02171	0.002757	0.000244	0.898504	-0.00392
OGK	0.032135	0.025488	0.85953	0.02285	0.000946	0.00393	0.872867	-0.00075
SS_n	0.027537	0.02368	0.80369	0.01998	0.00056	5.10e-3	0.87283	-0.00295
<i>n</i> = 100								
S_n	0.012804	0.011781	0.903699	0.011783	-0.00017	-0.00153	0.938678	0.006895
MCD	0.015236	0.01145	0.880362	0.01302	-0.00177	0.000362	0.950676	0.008832
Comedian	0.012505	0.011895	0.902279	0.011021	-0.00057	-0.00269	0.93753	0.00686
OGK	0.014642	0.013542	0.877704	0.012663	0.000848	0.000761	0.926497	0.008697
SS_n	0.01248	0.0113	0.864418	0.010806	0.000346	0.002428	0.917793	0.006029
<i>n</i> = 200								
S_n	0.005973	0.005682	0.926016	0.005729	-0.000800	0.000141	0.956202	0.003477
MCD	0.006437	0.00573	0.946856	0.006063	0.000298	0.000251	0.913761	0.003293
Comedian	0.005955	0.005523	0.923031	0.005722	-0.000590	-0.000910	0.954904	0.00261
OGK	0.013524	0.011409	0.91934	0.011265	0.001151	0.003127	0.891392	0.00972
SS_n	0.00591	0.005506	0.898654	0.005729	0.000355	-0.00044	0.941665	0.002543
<i>n</i> = 500								
S_n	0.002309	0.002166	0.93471	0.002198	4.65e-5	-0.00025	0.964442	0.001258
MCD	0.002404	0.00217	0.925404	0.002333	-1.90e-5	0.000597	0.969063	0.002317
Comedian	0.00231	0.002195	0.933463	0.002302	-0.00019	0.001192	0.963769	0.001277
OGK	0.002426	0.002143	0.938785	0.002358	-4.90e-5	0.000273	0.96283	0.000895
SS_n	0.00266	0.00242	0.92209	0.00258	-0.00025	0.00093	0.95751	0.00021

Table 3: \mathbf{y} -affine equivariance table

	$p = 4, q = 4$		$p = 6, q = 10$		$p = 10, q = 6$		$p = 10, q = 10$	
	Slope	Intercept	Slope	Intercept	Slope	Intercept	Slope	Intercept
delta = 0%								
<i>n</i> = 50	3.08e-5	8.80e-5	1.67e-5	6.62e-5	1.60e-5	9.53e-5	1.19e-5	8.59e-5
<i>n</i> = 100	7.67e-6	2.21e-5	3.68e-6	1.48e-5	2.63e-6	1.92e-5	2.61e-6	1.69e-5
<i>n</i> = 200	1.71e-6	5.71e-6	7.03e-7	3.42e-6	5.37e-7	4.13e-6	4.25e-7	3.60e-6
<i>n</i> = 500	3.24e-7	1.13e-6	9.88e-8	5.39e-7	1.17e-7	9.51e-7	6.41e-8	5.16e-7
delta = 10%								
<i>n</i> = 50	4.07e-5	9.28e-5	2.21e-5	6.86e-5	1.30e-5	1.15e-4	1.21e-5	8.00e-5
<i>n</i> = 100	1.16e-5	3.22e-5	3.64e-6	1.23e-5	3.10e-6	2.13e-5	2.29e-6	1.59e-5
<i>n</i> = 200	4.63e-6	1.69e-5	8.22e-7	3.53e-6	6.20e-7	4.77e-6	4.89e-7	3.59e-6
<i>n</i> = 500	6.00e-6	6.27e-6	1.40e-7	6.60e-7	2.14e-7	6.08e-6	9.14e-8	6.73e-7
delta = 20%								
<i>n</i> = 50	1.04e-4	0.002708	4.94e-6	1.46e-5	5.47e-6	0.00021	4.03e-6	2.70e-5
<i>n</i> = 100	4.96e-5	0.002438	1.31e-6	4.29e-6	2.74e-6	0.00026	9.84e-7	5.26e-6
<i>n</i> = 200	6.42e-5	0.003637	3.18e-7	1.15e-6	1.28e-6	1.93e-4	1.55e-7	1.21e-6
<i>n</i> = 500	8.49e-5	0.004147	6.71e-8	2.38e-7	1.88e-6	9.19e-4	3.54e-8	4.32e-6

Table 4: \mathbf{x} -affine equivariance table

	$p = 4, q = 4$		$p = 6, q = 10$		$p = 10, q = 6$		$p = 10, q = 10$	
	Slope	Intercept	Slope	Intercept	Slope	Intercept	Slope	Intercept
delta = 0%								
<i>n</i> = 50	1.33e-1	4.11e-4	1.45e-1	1.94e-4	1.75e-1	5.19e-4	1.29e-1	2.91e-4
<i>n</i> = 100	1.06e-1	8.88e-5	1.01e-1	5.39e-5	1.14e-1	1.02e-4	9.18e-2	8.18e-5
<i>n</i> = 200	6.33e-2	2.79e-5	5.79e-2	1.40e-5	8.21e-2	2.52e-5	5.60e-2	1.85e-5
<i>n</i> = 500	6.53e-2	4.73e-6	4.32e-2	2.27e-6	3.37e-2	3.25e-6	3.72e-2	2.63e-6
delta = 10%								
<i>n</i> = 50	2.10e-1	4.25e-4	1.43e-1	2.22e-4	2.04e-1	5.88e-4	1.42e-1	2.93e-4
<i>n</i> = 100	8.75e-2	3.93e-5	9.87e-2	6.58e-5	1.42e-1	1.99e-4	1.05e-1	1.03e-4
<i>n</i> = 200	1.06e-1	4.51e-5	7.96e-2	1.90e-5	8.38e-2	5.43e-5	5.19e-2	3.13e-5
<i>n</i> = 500	9.93e-2	1.04e+0	2.78e-2	3.18e-6	5.63e-2	9.95e-6	4.14e-2	4.98e-6
delta = 20%								
<i>n</i> = 50	0.26603	0.00015	1.57e-1	1.23e-4	2.25e-1	0.00024	1.85e-1	1.57e-4
<i>n</i> = 100	1.46e-1	6.29e-5	1.15e-1	2.84e-5	1.59e-1	9.06e-5	1.07e-1	5.22e-5
<i>n</i> = 200	9.45e-2	9.08e-6	5.58e-2	9.44e-6	1.13e-1	2.91e-5	5.48e-2	1.37e-5
<i>n</i> = 500	7.71e-2	1.28e-6	5.89e-2	1.17e-6	5.67e-2	2.28e-6	2.79e-2	2.34e-6

Overall, the simulation results indicate that the proposed method offers a favorable balance between robustness and efficiency, performing competitively in clean settings while providing substantial gains in accuracy and stability when contamination is present. Detailed numerical results supporting these conclusions are reported in the Supplementary Material.

5. Real life example

5.1. Pulpfibre data

We examine a dataset from [Lee, Roy, Hong, and Whiting \(1993\)](#) that includes measurements of various properties of pulp fibers as well as the characteristics of the paper produced from these fibers. The objective is to explore the relationships between the properties of pulp fibers and the characteristics of the resulting paper. The dataset comprises $n = 62$ measurements of four specific pulp fiber attributes: arithmetic fiber length, long fiber fraction, fine fiber fraction, and zero span tensile. The dataset includes measurements of four paper properties: breaking length, elastic modulus, stress at failure, and burst strength. Our objective is to predict the four paper properties using the four fiber characteristics. To achieve this, we initially applied classical multivariate regression to the dataset. Figure 1 of classical estimator is displayed along with other compared estimators. This diagnostic plot integrates information on regression outliers and leverage points, providing a more comprehensive view than analyzing each distance individually. For classical estimator, we can see that observations 51, 52 and 56 are identified as vertical outliers. Conversely, while some observations are recognized as leverage points (with observations 60 and 61 being the most prominent), they are not classified as regression outliers. To validate the results from classical multivariate regression, Rousseeuw applied univariate robust LTS regression to each response individually, using the same regressors. The results indicate that the univariate LTS regressions identify observations 51, 52, 56, and 61 as outliers. The results shows the necessity of finding outliers and substantiate the results using robust estimators. However, using univariate robust methods on each response variable individually does detects outliers only along the coordinate directions of the responses and fails to identify outliers that may be masked within these directions. Therefore, it is generally preferable to use a robust multivariate regression estimation method that can detect all outliers and is efficient both statistically and computationally. From a regression perspective, the presence of vertical outliers and leverage points has a direct impact on the estimated relationships between fibre properties and paper characteristics. The proposed SS_n estimator downweights these influential observations, resulting in more stable coefficient estimates that better reflect the dominant trend in the majority of the data.

Diagnostic plot allows us to categorize data points as regular observations, vertical outliers, good leverage points, or bad leverage points. Additionally, it helps us determine whether a point is an extreme outlier or just a borderline case. For plotting we consider the horizontal cutoff line as the cut off of respective used Mahalanobis distance and vertical cut off as $\sqrt{\chi_{q,0.99}^2}$. The choice the vertical line cutoff is usually between $\sqrt{\chi_{q,0.975}^2}$ and $\sqrt{\chi_{q,0.99}^2}$. Here we chose later as our cutoff line. Rousseeuw in his paper explored the sample and found the observation (59 – 62) were produced from fir wood. Most of the outlying observations were obtained from different pulping process. That is observation 62 obtained from chemi - thermomechanical process. Likewise observations 22, 46 - 48, 58 - 61 are obtained from different pulping processes. Evidently from the diagnostic plot, we can say our proposed estimator capable of detecting all these observations as outliers which validate the usefulness of our estimator.

The outlying observations detected by the robust methods correspond to samples produced using different pulping processes, such as chemi-thermomechanical processing and fir wood sources, which are known to generate fundamentally different fibre characteristics. Treating these observations on equal footing with the bulk of the data, as in the classical analysis, may

therefore distort the inferred relationships between fibre properties and paper strength.

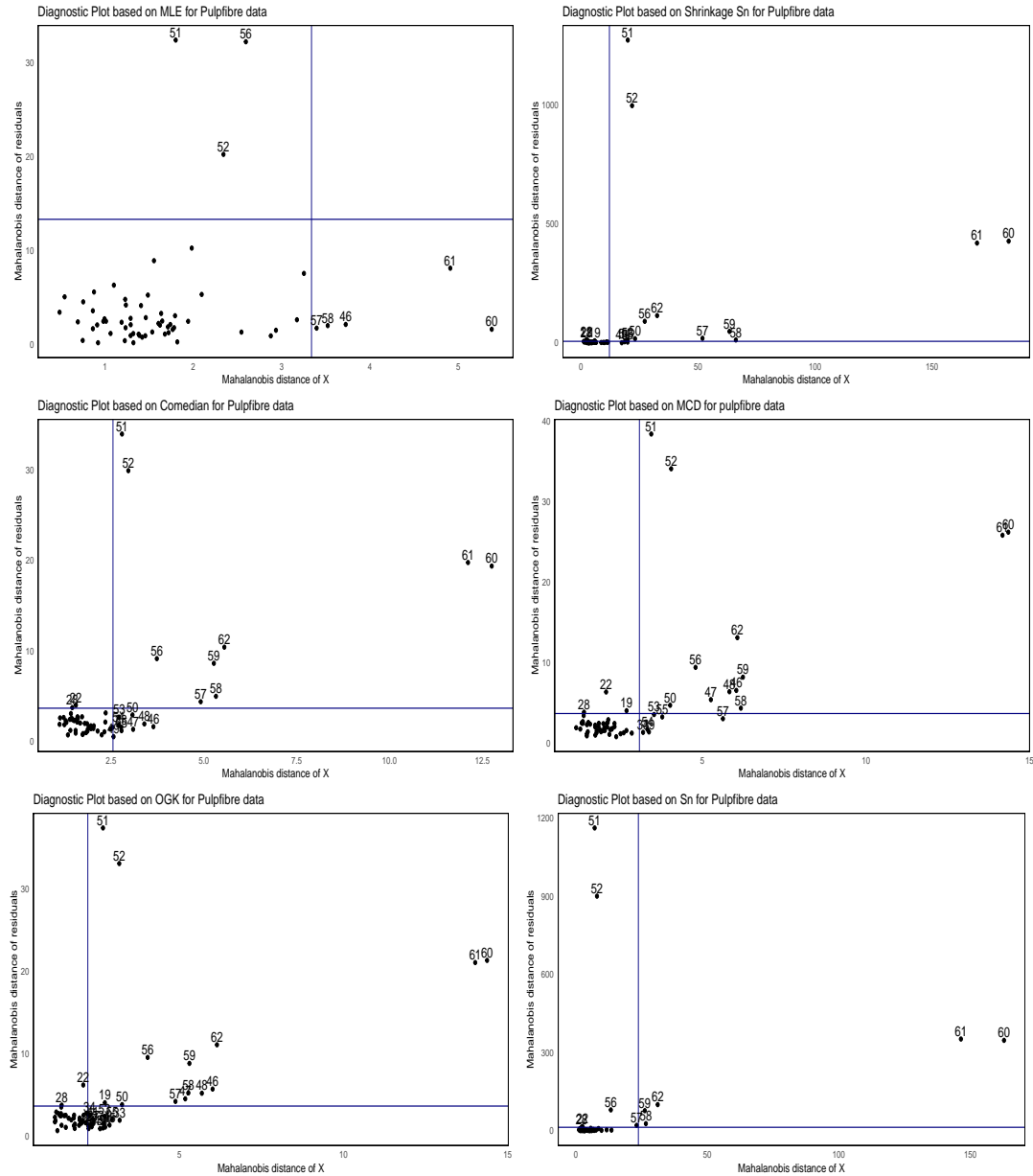


Figure 1: Diagnostic plot based on different estimators for Pulpfibre data

Table 5: Pulpfibre data

	Slope MSE	Intercept MSE	Slope Bias	Intercept Bias
MCD	561.7581	2158.558	9.37032	-40.2621
OGK	532.1857	2073.755	9.020945	-38.7037
Comedian	634.0527	2230.64	10.72183	-40.7999
S_n	634.0527	2230.64	10.72183	-40.7999
SS_n	534.7054	2076.603	10.58657	-43.9223
MLE	847.0833	2786.09	12.30624	-45.6509

5.2. School data

This data consists of $n = 70$ observations on school sites in US (Charnes, Cooper, and Rhodes 1981; Roelant, Van Aelst, and Croux 2009). Data contains three response variable and 5 explanatory variables. The response variables are: total reading score measured by the Metropolitan Achievement Test, total mathematics score measured by the Metropolitan Achievement Test and the Coopersmith self-esteem inventory. The independent variables are: education level of mother, highest occupation of a family member, number of parental visits to the school, parent counseling concerning school-related topics and the number of teachers at the school. The diagnostic plots are given in Figure 2. Cutoff lines are considered same as above mentioned. In the school data example, classical multivariate regression identifies only a limited number of atypical observations, resulting in coefficient estimates that are sensitive to a few high-leverage schools. Robust methods, including the proposed SS_n estimator, detect several additional outlying and leverage observations, which substantially alters the fitted relationships between parental involvement variables and student performance outcomes.

Classical estimator detects only one observation as vertical outliers and 6 observations as good leverage points. But robust estimators detect the observations 33, 35, 44, 59 as outliers, more than two moderate large good leverage points (1, 10, 50, 54, 66, 67...). Thus it is always better to make use of robust technique and our proposed estimator performs well in detecting the vertical outliers, leverage points.

Also, we have tabulated the intercept, slope MSE, Bias of these two examples and are given below. The table values shows the inconsistent performance of classical MLE in these datasets. And our proposed estimator exhibit outstanding performance with minimum MSE than other estimators, indicating that the influence of atypical schools with extreme characteristics has been mitigated. This is important from an interpretive standpoint, as policy conclusions regarding the effect of parental visits, counseling, and school size on academic outcomes can differ depending on whether such atypical institutions dominate the analysis.

Table 6: School Data

	Slope MSE	Intercept MSE	Slope Bias	Intercept Bias
MCD	3.8891	6.7543	0.7086	2.7644
OGK	3.8172	6.6208	0.7015	1.3937
Comedian	3.8073	6.6095	0.6528	2.2576
S_n	3.2752	6.7108	0.6532	2.2301
SS_n	3.0669	6.6013	0.6427	2.2106
MLE	3.222	0.0487	0.6837	-0.181

6. Conclusion

Classical regression estimation method is highly sensitive to the presence of outliers in the dataset, which can significantly affect their robustness. Therefore, alternative methods capable of detecting and withstanding outliers are necessary to ensure reliable results, even when outliers are present. Several works like Singer and Sen (1985) and Koenker and Portnoy (1990) introduced robust regression estimation methods using M estimators. Ollila *et al.* (2002) and Ollila *et al.* (2003) proposed estimators based on affine equivariant signs and ranks. Rousseeuw *et al.* (2004) proposed regression based on reweighted MCD estimator. Roelant *et al.* (2009) introduced non reweighted multivariate regression based on generalised S estimators. Sajana and Sajesh (2018) introduced non reweighted multivariate regression based Comedian estimator. In this paper we develop a two step reweighted multivariate regression SS_n based on estimator proposed by Lakshmi and Sajesh (2025). We evaluated the performance of our proposed multivariate regression estimator and compared with other existing estimators. We evaluated robustness, efficiency and affine equivariance through the metrics Slope, Intercept

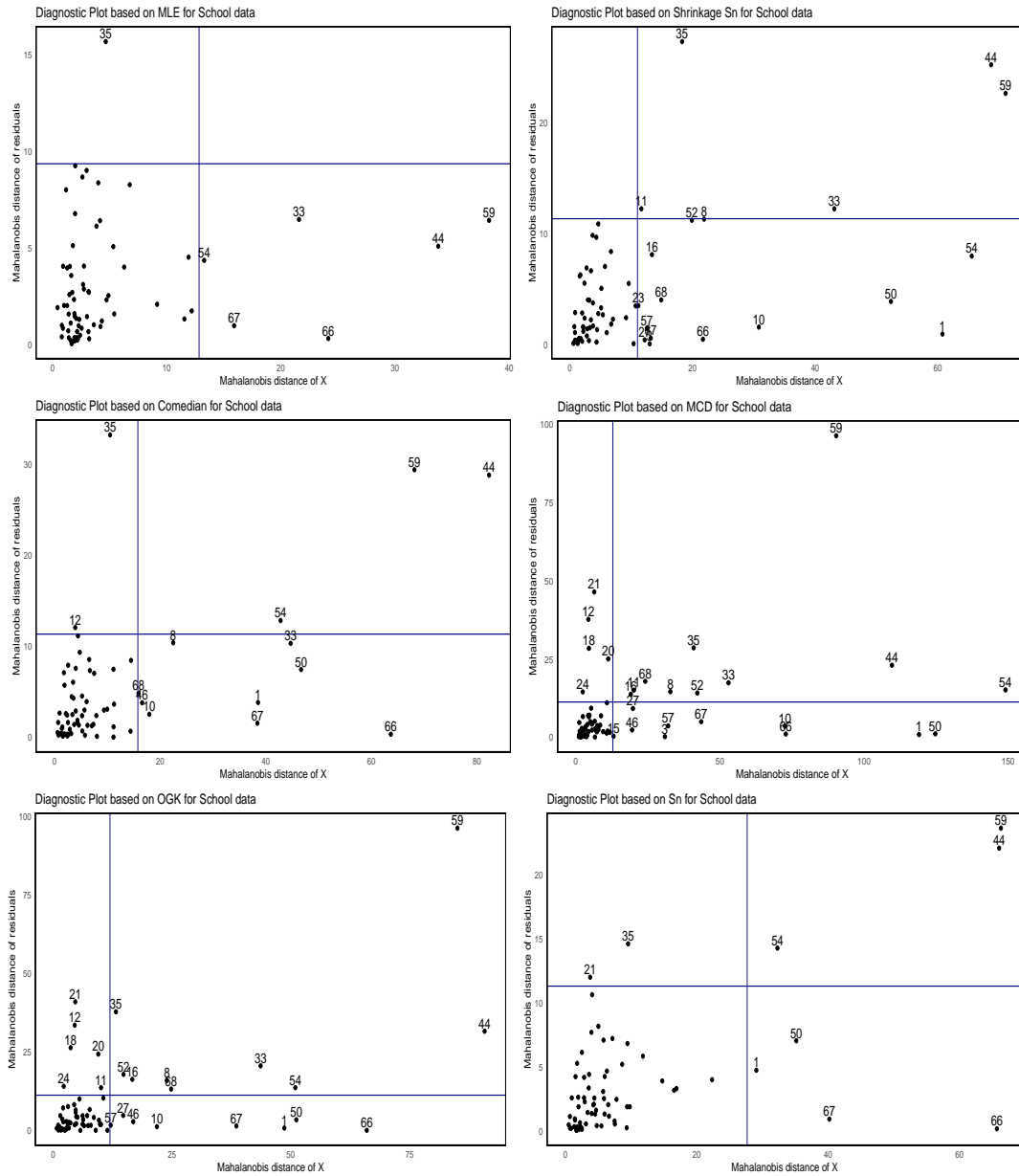


Figure 2: Diagnostic plot based on different estimators for School data

bias and MSE. To enhance efficiency, we made use of reweighting schemes and found that the best results are achieved by using Shrinkage S_n -based robust distances. These distances are used to create a reweighted estimator of location and scatter, which then serves as the foundation for the initial regression. The robust residuals from this initial regression are then used to determine the weights for the final regression. This two step reweightedness gave us finely performing regression estimator SS_n with empirically high asymptotic efficiency. For comparing the robustness, we used contaminated simulated datasets. Contamination schemes are finely chosen to explore the performance of all the estimators. The results shows our proposed estimator outperforms other estimators in terms of MSE and bias validating high robustness empirically. We empirically verified the affine equivariance, and the results provided positive confirmation of our proposed estimator's equivariance property. We checked the performance of the estimator in two real life data sets and made use of Diagnostic plot for evaluation purpose. The plots enforce our proposed estimation method is capable of detecting outliers in the dataset. The MSE, bias values shows how classical estimator provide derogatory values in the presence of outliers and shows the capability of our proposed estimator. Thus our study assures using our proposed SS_n multivariate regression estimator in datasets with multiple outliers for multivariate regression estimation.

7. Disclosure statement

No potential conflict of interest was reported by the author(s).

8. Data availability

The availability of Data supporting the findings of this study are mentioned in the respective section and references.

References

- Agulló J, Croux C, Van Aelst S (2008). "The Multivariate Least Trimmed Squares Estimator." *Journal of Multivariate Analysis*, **99**(3), 311–338.
- Alghamdi FM, Hammad AT, Golam Kibria BM, Abd-Elmougod GA, Sapkota LP, Gemeay AM (2025). "On Robust and Non-Robust Modified Liu Estimation in Poisson Regression Model with Multicollinearity and Outliers." *International Journal of Uncertainty, Fuzziness and Knowledge-Based Systems*, **33**(06), 787–823. doi:10.1142/S0218488525500266.
- Bickel PJ, Lehmann EL (2012). "Descriptive Statistics for Nonparametric Models. III. Dispersion." In *Selected Works of EL Lehmann*, pp. 499–518. Springer. doi:10.1007/978-1-4614-1412-4_44.
- Cabana E, Lillo RE, Laniado H (2021). "Multivariate Outlier Detection Based on a Robust Mahalanobis Distance with Shrinkage Estimators." *Statistical Papers*, **62**, 1583–1609. doi:10.1007/s00362-019-01148-1.
- Charnes A, Cooper WW, Rhodes E (1981). "Evaluating Program and Managerial Efficiency: An Application of Data Envelopment Analysis to Program Follow Through." *Management Science*, **27**(6), 668–697.
- Croux C, Rousseeuw PJ, Hössjer O (1994). "Generalized S-Estimators." *Journal of the American Statistical Association*, **89**(428), 1271–1281. doi:10.2307/2290990.

- DeMiguel V, Martin-Utrera A, Nogales FJ (2013). “Size Matters: Optimal Calibration of Shrinkage Estimators for Portfolio Selection.” *Journal of Banking & Finance*, **37**(8), 3018–3034.
- Johnson RA, Wichern DW (2002). *Applied Multivariate Statistical Analysis*. Pearson Prentice Hall: Upper Saddle River, NJ. ISBN 978-0130925534.
- Koenker R, Portnoy S (1990). “M Estimation of Multivariate Regressions.” *Journal of the American Statistical Association*, **85**(412), 1060–1068. doi:10.2307/2289602.
- Kunjunni SO, Abraham ST (2022). “Multidimensional Outlier Detection and Robust Estimation Using S_n Covariance.” *Communications in Statistics-Simulation and Computation*, **51**(7), 3912–3922. doi:10.1080/03610918.2020.1725820.
- Lakshmi R, Sajesh TA (2025). “A Robust Distance-Based Approach for Detecting Multidimensional Outliers.” *Journal of Applied Statistics*, **52**(6), 1278–1298. doi:10.1080/02664763.2024.2422403.
- Ledoit O, Wolf M (2004). “A Well-Conditioned Estimator for Large-Dimensional Covariance Matrices.” *Journal of Multivariate Analysis*, **88**(2), 365–411. doi:10.1016/S0047-259X(03)00096-4.
- Lee J, Roy DN, Hong M, Whiting P (1993). “Relationships Between Properties of Pulp-Fibre and Paper.” In *Products in Papermaking*, pp. 159–182. doi:10.15376/frc.1993.1.159.
- Lopuhaä HP (1999). “Asymptotics of Reweighted Estimators of Multivariate Location and Scatter.” *Annals of Statistics*, **27**(5), 1638–1665. doi:10.1214/aos/1017939145.
- Lopuhaa HP, Rousseeuw PJ (1991). “Breakdown Points of Affine Equivariant Estimators of Multivariate Location and Covariance Matrices.” *Annals of Statistics*, **19**(1), 229–248. doi:10.1214/aos/1176347978.
- Lukman AF, Farghali RA, Kibria BMG, Oluyemi OA (2023). “Robust-Stein Estimator for Overcoming Outliers and Multicollinearity.” *Scientific Reports*, **13**, 9066. doi:10.1038/s41598-023-36053-z.
- Maronna RA, Yohai VJ (1997). “Robust Estimation in Simultaneous Equations Models.” *Journal of Statistical Planning and Inference*, **57**(2), 233–244.
- Maronna RA, Zamar RH (2002). “Robust Estimates of Location and Dispersion for High-Dimensional Datasets.” *Technometrics*, **44**(4), 307–317. doi:10.1198/004017002188618509.
- Ollila E, Oja H, Hettmansperger TP (2002). “Estimates of Regression Coefficients Based on the Sign Covariance Matrix.” *Journal of the Royal Statistical Society Series B: Statistical Methodology*, **64**(3), 447–466.
- Ollila E, Oja H, Koivunen V (2003). “Estimates of Regression Coefficients Based on Lift Rank Covariance Matrix.” *Journal of the American Statistical Association*, **98**(461), 90–98. doi:10.1198/016214503388619120.
- Roelant E, Van Aelst S, Croux C (2009). “Multivariate Generalized S-Estimators.” *Journal of Multivariate Analysis*, **100**(5), 876–887.
- Rousseeuw PJ, Leroy AM (1987). *Robust Regression and Outlier Detection*. John Wiley & Sons. doi:10.1002/0471725382.
- Rousseeuw PJ, Van Aelst S, Van Driessen K, Gulló JA (2004). “Robust Multivariate Regression.” *Technometrics*, **46**(3), 293–305. doi:10.1198/004017004000000329.

- Sajana OK, Sajesh TA (2018). “Empirical Robust Multivariate Regression Parameter Estimation Using Median Approach.” *International Journal of Scientific Research in Mathematical and Statistical Sciences*, **5**(5), 65–71. doi:10.26438/ijstrmss/v5i5.6571.
- Siegel AF (1982). “Robust Regression Using Repeated Medians.” *Biometrika*, **69**(1), 242–244. doi:10.2307/2335877.
- Singer JM, Sen PK (1985). “M-Methods in Multivariate Linear Models.” *Journal of Multivariate Analysis*, **17**(2), 168–184.
- Suhail M, Chand S, Kibria BMG (2021). “Quantile-Based Robust Ridge M-Estimator for Linear Regression Model in Presence of Multicollinearity and Outliers.” *Communications in Statistics-Simulation and Computation*, **50**(11), 3194–3206. doi:10.1080/03610918.2019.1621339.
- Van Aelst S, Willems G (2005). “Multivariate Regression S-Estimators for Robust Estimation and Inference.” *Statistica Sinica*, **15**(4), 981–1001.
- Wasim D, Zaman Q, Ahmad M, Kibria BMG (2025). “Mitigating Multicollinearity and Outliers in Regression: Comparison of Some New and Old Robust Ridge M-Estimators.” *Journal of Statistical Computation and Simulation*, **95**(16), 3526–3547. doi:10.1080/00949655.2025.2538110.
- Yasmin N, Kibria BMG (2025). “Performance of Some Improved Estimators and their Robust Versions in Presence of Multicollinearity and Outliers.” *Sankhya B, The Indian Journal of Statistics*, **87**, 173–219. doi:10.1007/s13571-025-00352-4.

Affiliation:

Lakshmi R
Assistant Professor
Department of Statistics and Data Science
Christ (Deemed to be University)
Bengaluru, India
E-mail: lakshmi.nss19@gmail.com