

## Enhanced weighted least squares regression: A robust approach for managing outliers and heteroscedasticity

Ayanlowo E. A<sup>1</sup>, Oladapo D.I<sup>2</sup>, Odeyemi, A. S<sup>3</sup> and Obadina, G.O<sup>4,\*</sup>

<sup>1</sup> Department of Basic Sciences, Babcock University, Ilishan-Remo, Ogun State, Nigeria.

<sup>2</sup> Department of Mathematical Sciences, Adeleke University, Ede, Osun State, Nigeria.

<sup>3</sup> Department of Statistics University of Fort Hare Alice, Eastern Cape, South Africa.

<sup>4</sup> Department of Statistics, Olabisi Onabanjo University, Ago-Iwoye, Ogun State, Nigeria.

International Journal of Science and Technology Research Archive, 2024, 07(02), 097-106

Publication history: Received on 04 October 2024; revised on 20 December 2024; accepted on 23 December 2024

Article DOI: <https://doi.org/10.53771/ijstra.2024.7.2.0064>

### Abstract

The Ordinary Least Squares (OLS) approach performs best in the best-case scenario, but when the homoscedasticity and outlier-absence assumptions are broken, it performs noticeably worse. With an emphasis on Robust Weighted Least Squares (RWLS), M-estimation, and Least Trimmed Squares (LTS), this paper assesses robust alternatives to OLS. OLS was demonstrated to be extremely sensitive to outliers using a dataset on state-level crime rates in the US; the Mean Squared Error (MSE) increased from 5.480 (original data) to 12.580 (with outliers). On the other hand, RWLS using Tukey's Bisquare function produced the most consistent coefficient estimates and the smallest MSE, increasing from 4.520 to 5.340. In comparison to OLS, M-estimation and LTS also demonstrated lower MSE and higher resilience. The findings show how well robust approaches—in particular, RWLS with Tukey's Bisquare—address heteroscedasticity and outliers, which makes them essential for practical regression analysis.

**Keywords:** Heteroscedasticity; Outliers; Robust regression; M-estimation; Tukey's Bisquare

### 1 Introduction

The Ordinary Least Squares (OLS) estimator is frequently employed in regression analysis because, under certain important assumptions, it is straightforward, computationally efficient, and has unbiased characteristics (Reddy & Henze, 2023). To ensure that OLS performs at its best, homoscedasticity—the state in which the variance of the errors stays constant across all levels of the independent variables—is one of these requirements (Knaub Jr, 2021). In addition to being the Best Linear Unbiased Estimator (BLUE) under homoscedasticity, OLS generates effective estimates with the least standard errors (Bhatti, et al., 2023).

Heteroscedasticity results from the frequent violation of the homoscedasticity assumption in real-world data. When the variance of the errors fluctuates across various levels of the independent variables rather than remaining constant, this is known as heteroscedasticity. The standard errors of the OLS estimator are inflated by this variance, rendering the conclusions from hypothesis tests untrustworthy. As a result, OLS-based p-values and confidence intervals lose their credibility and may provide inaccurate findings (Ruan, 2024).

OLS estimation is significantly hampered by outliers in addition to heteroscedasticity. Extreme observations that significantly depart from the data's general pattern are known as outliers. OLS is especially sensitive to significant data deviations because it depends on reducing the sum of squared residuals; outliers have a disproportionate impact on the

\* Corresponding author: Obadina, G.O

regression coefficients. Therefore, even a small number of outliers might result in skewed and inconsistent parameter estimations, which will skew the regression analysis's findings (Boukerche, Zheng & Alfandi, 2020)

Robust regression methods that can manage outliers and heteroscedastic errors are obviously needed in light of these constraints (Kim & Li, 2023). This study's main objective is to investigate solid OLS substitutes that offer better accurate predictions even when there are anomalies in the data. Within the framework of Robust Weighted Least Squares (RWLS), the study specifically focusses on two robust regression techniques:

**M-Estimation:** By adding a weighting scheme that lessens the influence of outliers, this method modifies the OLS minimisation procedure. By applying a robust loss function and down-weighting big residuals, M-estimators lessen the impact of extreme observations instead of minimising the sum of squared residuals. Thus, M-estimation improves the robustness of the parameter estimations by providing a more robust method for dealing with outliers.

Another reliable regression technique that guards against outliers is Least Trimmed Squares (LTS), which trims a predetermined percentage of the highest residuals prior to model fitting. LTS reduces the effect of outliers and guarantees that the regression model captures the data's central tendency by concentrating on the subset of data with the fewest residuals. When a sizable percentage of the dataset is made up of outliers, this method is quite helpful.

Robust Weighted Least Squares (RWLS) is based on the combination of these robust techniques with Weighted Least Squares (WLS). By giving each observation a variable weight, WLS—a generalisation of OLS—addresses heteroscedasticity by giving observations with larger error variance less weight. In order to ensure that observations with more variance have less of an impact on the final parameter estimations, WLS corrects for heteroscedasticity by adding weights inversely proportionate to the estimated variance of the errors.

RWLS is a potent method that simultaneously handles heteroscedasticity and outliers by combining M-estimation and LTS with WLS. Even in cases where the assumptions of constant error variance and the lack of outliers are broken, RWLS offers more accurate parameter estimations by combining robust approaches with weighted regression.

### *Research Objectives*

The particular goals of this research are:

- To assess how well Least Trimmed Squares (LTS) and M-Estimation perform when heteroscedasticity and outliers are present.
- To illustrate how, under less-than-ideal data circumstances, Robust Weighted Least Squares (RWLS) can increase the dependability of regression estimations.
- To contrast RWLS with conventional OLS and WLS, emphasising the benefits of robust approaches in producing precise, objective estimates in the face of difficult data circumstances.

This study intends to offer workable strategies for managing data inconsistencies that commonly occur in applied regression analysis through the use of robust methodologies, producing more reliable and understandable findings. RWLS is a useful tool for real-world regression problems because it combines the advantages of M-estimation and LTS within the WLS framework, providing a thorough method for addressing the twin concerns of outliers and heteroscedasticity.

## **2 Methodology**

This study examines linear regression models under the assumption of heteroscedastic errors. The general regression model is represented as:

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

where  $y$  is the vector of observed dependent values,  $X$  represents the matrix of independent variables,  $\beta$  is the vector of regression parameters, and  $\varepsilon$  denotes the error terms. When errors exhibit heteroscedasticity, the OLS estimator becomes unreliable, necessitating a weighted least squares approach.

## 2.1 Weighted Least Squares (WLS)

The Weighted Least Squares (WLS) method is an extension of the Ordinary Least Squares (OLS) estimator, specifically designed to address heteroscedasticity in regression analysis. Heteroscedasticity, the violation of the constant variance assumption, can result in inefficient OLS estimates, leading to biased standard errors and unreliable statistical inferences. WLS corrects for this by applying a weighting scheme to each observation, ensuring that observations with higher error variance exert less influence on the estimated coefficients.

In the general linear regression model:

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

where:

$y$  is the  $n \times 1$  vector of observed dependent variable values,  
 $X$  is the  $n \times p$  matrix of predictors (independent variables),  
 $\beta$  is the  $p \times 1$  vector of regression coefficients to be estimated,  
 $\varepsilon$  is the  $n \times 1$  vector of error terms.

In the presence of heteroscedasticity, the error terms  $\varepsilon$  have non-constant variance, denoted as:

$$\text{Var}(\varepsilon) = \sigma_i^2$$

where  $\sigma_i^2$  varies for each observation  $i$ . In contrast to OLS, which assumes that all error terms have equal variance, WLS introduces a set of weights to account for these varying error variances. The WLS estimator,  $\hat{\beta}_{WLS}$ , is expressed as:

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

where:

$W$  is an  $n \times n$  diagonal matrix containing the weights assigned to each observation, with the weights typically being the inverse of the estimated error variances, i.e.,  $W = \text{diag}(w_1, w_2, \dots, w_n)$  where  $w_i = \frac{1}{\sigma_i^2}$ .

By applying weights inversely proportional to the error variances, WLS ensures that observations with larger variances have less influence on the coefficient estimates. This is particularly useful when dealing with heteroscedasticity, where the error variance is not constant across all observations. The WLS estimator minimises the weighted sum of squared residuals, rather than the unweighted sum as in OLS, thereby adjusting for the varying levels of error variance:

$$\hat{\beta}_{WLS} = \arg \min_{\beta} \sum_{i=1}^n w_i (y_i - X_i\beta)^2$$

## 2.2 M-Estimation

M-estimation is a robust regression technique developed to reduce the impact of outliers by down-weighting extreme residuals. Unlike the Ordinary Least Squares (OLS) method, which minimises the sum of squared residuals, M-estimation minimises a more flexible loss function,  $\rho(e)$ , that assigns different weights to residuals based on their magnitude. The result is a regression model that is more resistant to the influence of outliers, making it a preferred method when the data contain deviations from normality or anomalies that can distort the parameter estimates.

### 2.2.1 Concept of M-Estimation

In OLS, the objective is to minimise the sum of squared residuals,  $e_i^2$ , for all observations  $i = 1, 2, \dots, n$ , where  $e_i = y_i - X_i\beta$  is the residual for the  $i$ -th observation. The OLS loss function is given by:

$$\sum_{i=1}^n e_i^2$$

However, this quadratic loss function places equal weight on all residuals, causing large residuals (outliers) to exert disproportionate influence on the parameter estimates. To address this, M-estimation introduces a more robust

objective function that minimises a general function of the residuals,  $\rho(e)$ , rather than their squared values. The M-estimator solves the following optimisation problem:

$$\hat{\beta}_M = \arg \min_{\beta} \sum_{i=1}^n \rho(e_i)$$

where  $\rho(e_i)$  is a function designed to reduce the influence of large residuals. By choosing an appropriate  $\rho(e)$ , M-estimation diminishes the impact of outliers while maintaining efficient estimation for non-outlying data points.

### 2.2.2 Common Choices of $\rho(e)$

The choice of the function  $\rho(e)$  is crucial for determining how aggressively outliers are down-weighted. Several robust loss functions have been proposed, each with different properties in handling residuals:

- **Huber's Function:** Huber's loss function is one of the most widely used in M-estimation. It behaves quadratically for small residuals, like OLS, but linearly for large residuals, limiting their influence. The function is defined as:

$$\rho(e_i) = \begin{cases} \frac{1}{2}e_i^2 & \text{if } |e_i| \leq k \\ k\left(|e_i| - \frac{1}{2}k\right) & \text{if } |e_i| > k \end{cases}$$

where  $k$  is a tuning constant that controls the threshold at which the loss function transitions from quadratic to linear. For residuals smaller than  $k$ , the function behaves like the traditional squared loss, and for residuals larger than  $k$ , it becomes linear, reducing the impact of outliers.

- **Tukey's Bisquare Function:** Tukey's bisquare function (also called Tukey's biweight) is more aggressive in down-weighting outliers than Huber's function. It provides nearly no influence for very large residuals while treating small residuals similarly to OLS. Tukey's function is defined as:

$$\rho(e_i) = \begin{cases} k^2 \left[ 1 - \left( 1 - \left( \frac{e_i}{k} \right)^2 \right)^3 \right] & \text{if } |e_i| \leq k \\ k^2 & \text{if } |e_i| > k \end{cases}$$

For residuals smaller than  $k$ , Tukey's function behaves smoothly, but for residuals larger than  $k$ , it effectively limits the impact by setting the influence close to zero. This makes Tukey's bisquare function particularly effective in datasets with a large number of extreme outliers.

- **Other Loss Functions:** Other robust loss functions used in M-estimation include Hampel's three-part function and Cauchy's function, each designed to limit the influence of outliers while maintaining efficiency for central data points. These functions provide different trade-offs between robustness and efficiency, depending on the nature of the data.

### 2.2.3 Iterative Weighted Least Squares (IWLS) in M-Estimation

M-estimation is typically computed using an Iteratively Reweighted Least Squares (IRLS) procedure. In this approach, an initial estimate of the regression coefficients is obtained (usually via OLS), and then the residuals are computed. Based on these residuals, the weight for each observation is updated at each iteration, with larger residuals receiving lower weights, until the estimates converge. The iterative procedure can be summarised as follows:

- **Initial Estimate:** Start with an initial estimate of the regression coefficients,  $\hat{\beta}^{(0)}$ , typically obtained through OLS.
- **Compute Residuals:** At each iteration, calculate the residuals  $e_i^{(k)} = y_i - X_i \hat{\beta}^{(k)}$ , where  $k$  represents the iteration number.
- **Update Weights:** Using the residuals, update the weights  $w_i$  for each observation based on the chosen robust loss function. For example, for Huber's function, the weight at iteration  $k$  is:

$$w_i^{(k)} = \begin{cases} 1 & \text{if } |e_i| \leq k \\ \frac{k}{|e_i^{(k)}|} & \text{if } |e_i| > k \end{cases}$$

- **Recompute Coefficients:** Using the updated weights, recompute the coefficient estimates using the weighted least squares formula:

$$\hat{\beta}^{(k+1)} = (X'W(k)X)^{-1}X'W(k)y$$

where  $W^{(k)}$  is the diagonal matrix of updated weights at iteration  $k$ .

- **Convergence:** Repeat steps 2–4 until the coefficient estimates converge, i.e., the change in  $\hat{\beta}^{(k)}$  between iterations is sufficiently small.

### 2.3 Least Trimmed Squares (LTS)

Least Trimmed Squares (LTS) is a robust regression technique designed to limit the influence of outliers by selectively excluding a certain portion of the largest residuals from the estimation process. Unlike the Ordinary Least Squares (OLS) estimator, which minimises the sum of squared residuals across all observations, LTS aims to minimise the sum of the smallest squared residuals. By focusing on a subset of the data that is less likely to be contaminated by outliers, LTS ensures that the resulting parameter estimates are not distorted by extreme values, making it highly resistant to outliers.

#### 2.3.1 Concept of LTS Estimation

In standard OLS regression, the goal is to minimise the total sum of squared residuals:

$$\sum_{i=1}^n e_i^2 = \sum_{i=1}^n (y_i - X_i\beta)^2$$

where:

$e_i = y_i - X_i\beta$  is the residual for the  $i$ -th observation,

$y_i$  is the observed value of the dependent variable,

$X_i\beta$  is the predicted value based on the model.

However, OLS is highly sensitive to outliers because large residuals (from outlying observations) contribute disproportionately to the sum of squared residuals. LTS addresses this issue by trimming the largest residuals and only using a subset of the smallest residuals in the estimation process.

The LTS estimator is defined as:

$$\hat{\beta}_{LTS} = \arg \min_{\beta} \sum_{i=1}^h e_{(i)}^2$$

where:

$e_{(i)}$  are the ordered residuals, meaning the residuals are ranked from smallest to largest (in absolute terms), and only the smallest residuals are included in the sum,

$h$  is the number of residuals used in the minimisation, which is less than or equal to the total number of observations  $n$ .

Thus, LTS minimises the sum of the smallest  $h$  squared residuals, excluding the largest residuals that are most likely to be influenced by outliers. By trimming the largest residuals, LTS effectively reduces the impact of outliers on the parameter estimates, making it a robust alternative to OLS.

### 2.3.2 Trimming Parameter $h$

The choice of  $h$  is crucial in LTS estimation because it determines how many observations are included in the regression. Specifically,  $h$  must be selected to balance robustness against outliers and the efficiency of the estimator. The trimming parameter  $h$  is typically chosen as a fraction of the total sample size  $n$ , with the following general guidelines:

- If  $h = n$ , no trimming is applied, and the LTS estimator reduces to OLS (no outliers are excluded).
- If  $h$  is too small, too many observations are excluded, leading to inefficiency because useful data points may be discarded.
- A typical recommendation is to choose  $h$  as approximately  $[n/2] + (p + 1)/2$  where  $p$  is the number of predictors, ensuring that a majority of the data is used while still excluding potential outliers.

For example, if there are  $n = 100$  observations and  $p = 5$  predictors, a reasonable choice for  $h$  would be around 57 to 60. This choice ensures that enough observations are included to produce efficient estimates, while the outliers with the largest residuals are excluded from the calculation.

**Table 1** Comparison of LTS and OLS

Property	OLS	LTS
Sensitivity to Outliers	High sensitivity; outliers significantly affect estimates	Highly resistant; outliers are trimmed and do not influence estimates
Efficiency	High efficiency when no outliers are present	Lower efficiency than OLS due to trimming, but much more robust in the presence of outliers
Breakdown Point	Low breakdown point; a few outliers can severely distort the model	High breakdown point; up to 50% of data can be outliers without significant distortion
Handling Leverage Points	Sensitive to high leverage points	Sensitive to high leverage points, although residual trimming helps

## 3 Numerical Example

This section presents a numerical example using data on state-wide crime rates in the United States (1993) to illustrate the performance of different regression techniques: Ordinary Least Squares (OLS), Robust Weighted Least Squares (RWLS), M-estimation, and Least Trimmed Squares (LTS). The dataset contains heteroscedasticity and outliers, making it an ideal test case for robust estimation methods. The dataset includes six predictor variables (e.g., population density, unemployment rate, income inequality, etc.) and one response variable (crime rate). To further test the robustness of the estimators, outliers were deliberately introduced by modifying the first and 25th observations.

Predictor Variables:

- X1: Population density (persons per square mile)
- X2: Unemployment rate (%)
- X3: Income inequality (Gini index)
- X4: Percentage of population below the poverty line
- X5: Median household income (in USD)
- X6: Urban population (%)

Response Variable:

- Y: Crime rate (number of crimes per 100,000 population)

Procedure:

We applied OLS, M-estimation, LTS, and RWLS (using both Huber's and Tukey's Bisquare functions) to the original dataset, as well as the modified dataset containing two outliers in the first and 25th observations. The outliers were created by artificially inflating these values by 200%, simulating extreme deviations from the central trend. The

performance of the estimators was evaluated using the Mean Squared Error (MSE), parameter estimates, standard errors, and residual plots to compare the robustness of each method in the presence of heteroscedasticity and outliers.

**Table 2** Parameter Estimates, Standard Errors, and MSE for Different Estimators (Original Data)

Estimator	X1 Coef.	X2 Coef.	X3 Coef.	X4 Coef.	X5 Coef.	X6 Coef.	MSE
OLS	0.450	1.320	2.200	0.630	-0.150	0.800	5.480
RWLS (Huber)	0.425	1.290	2.100	0.610	-0.140	0.770	4.780
RWLS (Tukey)	0.415	1.275	2.080	0.590	-0.135	0.750	4.520
M-estimation	0.430	1.300	2.150	0.620	-0.145	0.780	4.950
LTS	0.440	1.310	2.180	0.625	-0.148	0.795	5.000

Table 2 provides an insightful comparison of parameter estimates, standard errors, and Mean Squared Error (MSE) for different regression methods applied to the original dataset. This allows for evaluating how effectively each method handles potential irregularities, such as heteroscedasticity and outliers, which can influence the accuracy of regression models.

The Ordinary Least Squares (OLS) method, often used as the default regression approach, yielded the highest MSE of 5.480. This high error indicates that OLS is less efficient compared to the robust methods, highlighting its sensitivity to outliers and heteroscedasticity. The parameter estimates for each predictor, while reasonable, are vulnerable to distortion from data irregularities. OLS assumes constant variance and no extreme outliers, so in practice, it tends to perform poorly when these assumptions are violated.

In contrast, the Robust Weighted Least Squares (RWLS) method using Huber's function significantly improved the model's performance, reducing the MSE to 4.780. The Huber function adjusts the impact of extreme residuals by down-weighting outliers, resulting in more reliable parameter estimates than OLS. The coefficient estimates for all predictors were slightly smaller than those obtained through OLS, reflecting a more stable model with less influence from extreme data points.

The best overall performance was achieved by RWLS using Tukey's Bisquare function, which delivered the lowest MSE at 4.520. Tukey's Bisquare function further minimized the influence of outliers compared to Huber's method, resulting in even more stable coefficient estimates. The slight reductions in the coefficients, along with the lowest error rate, demonstrate Tukey's robustness and suitability for datasets prone to outliers and heteroscedastic errors. This makes it the preferred method for obtaining accurate regression estimates in the presence of data irregularities.

M-estimation also showed a reduction in MSE (4.950) compared to OLS, indicating moderate robustness to outliers. However, the MSE was slightly higher than that of RWLS, particularly with Tukey's function. The parameter estimates were closer to those of OLS but were adjusted to mitigate the impact of outliers. M-estimation reduces the influence of extreme values but is not as aggressive as Tukey's Bisquare function, leading to a modest improvement in reliability.

Finally, the Least Trimmed Squares (LTS) method, which trims the most extreme residuals to exclude outliers from the regression analysis, showed an MSE of 5.000. While LTS provided more reliable estimates than OLS by protecting against outliers, it was less efficient than the RWLS methods. The coefficient estimates were close to those of OLS, and while LTS effectively reduces the influence of outliers, it is not as effective as Tukey's Bisquare function in minimizing overall error.

The results show that RWLS with Tukey's Bisquare function provided the most reliable parameter estimates with the lowest MSE, demonstrating its superiority in addressing both heteroscedasticity and outliers. RWLS with Huber's function and M-estimation also improved performance compared to OLS, with lower MSEs and more stable estimates, although Tukey's function consistently performed better. LTS provided reliable estimates by trimming outliers but was slightly less efficient than the RWLS methods. Overall, the robust methods significantly outperformed OLS, with RWLS using Tukey's Bisquare emerging as the most effective solution for real-world data characterized by outliers and non-constant error variance.

**Table 3** Parameter Estimates, Standard Errors, and MSE for Different Estimators (Modified Data with Outliers)

Estimator	X1 Coef.	X2 Coef.	X3 Coef.	X4 Coef.	X5 Coef.	X6 Coef.	MSE
OLS	1.100	2.500	4.100	1.500	-0.300	1.200	12.580
RWLS (Huber)	0.480	1.350	2.250	0.650	-0.170	0.830	5.600
RWLS (Tukey)	0.470	1.320	2.210	0.640	-0.160	0.810	5.340
M-estimation	0.490	1.360	2.280	0.670	-0.180	0.840	5.880
LTS	0.500	1.370	2.300	0.680	-0.190	0.850	6.000

Table 3 presents the parameter estimates, standard errors, and Mean Squared Error (MSE) for different estimators applied to a modified dataset containing deliberately introduced outliers. This comparison reveals how each method handles outliers, offering insight into the robustness of OLS, RWLS (Huber), RWLS (Tukey), M-estimation, and LTS.

The Ordinary Least Squares (OLS) method performed poorly in the presence of outliers, as evidenced by the inflated coefficient estimates and significantly increased MSE of 12.580. The outliers had a drastic impact on the OLS estimates, causing overestimation of the coefficients, particularly for X1, X2, and X3. For example, the coefficient for X1 increased from 0.450 (in the original dataset) to 1.100, and for X3, it rose to 4.100, which indicates that OLS is extremely sensitive to extreme values. This significant increase in MSE suggests that OLS is unreliable when outliers are present, as the model fits these extreme observations rather than capturing the true relationship in the majority of the data.

In contrast, RWLS (Huber) demonstrated greater robustness, with a significantly lower MSE of 5.600. The parameter estimates for X1, X2, and X3 were much closer to their values in the original dataset, indicating that Huber's method effectively down-weighted the influence of outliers. For example, the X1 coefficient was reduced to 0.480, far lower than the inflated OLS estimate, reflecting the stability of Huber's approach. This reduction in MSE, coupled with more reasonable coefficient estimates, shows that Huber's function successfully mitigates the effects of outliers while maintaining model accuracy.

RWLS (Tukey's Bisquare function) provided the best overall performance in the presence of outliers, with the smallest MSE of 5.340. Tukey's function is more aggressive in reducing the influence of outliers than Huber's, which is reflected in the even more stable parameter estimates. For instance, the X1 coefficient was reduced to 0.470, and the MSE was slightly lower than Huber's. This demonstrates that Tukey's Bisquare function is the most effective technique in handling extreme values, providing the most reliable and robust parameter estimates while minimizing the error.

M-estimation performed moderately well, with an MSE of 5.880, which is lower than OLS but slightly higher than both RWLS methods. The parameter estimates were reasonably close to their true values, but the method was not as efficient as Tukey's Bisquare in limiting the impact of outliers. For instance, the X1 coefficient was 0.490, slightly larger than the RWLS estimates, and the MSE was higher, indicating that while M-estimation is robust, it is not as aggressive in controlling for outliers as RWLS with Tukey's function.

Least Trimmed Squares (LTS) showed a moderate increase in MSE to 6.000. While LTS trimmed the largest residuals to reduce the effect of outliers, its performance was slightly less efficient than the RWLS methods. The X1 coefficient was 0.500, close to the values produced by the other robust methods, but the higher MSE suggests that LTS was not as effective as Tukey's Bisquare in minimizing overall error. Nonetheless, LTS still performed much better than OLS in controlling the influence of outliers.

RWLS with Tukey's Bisquare function emerged as the most robust and reliable estimator for this dataset, effectively controlling for outliers and delivering the lowest error, while OLS showed significant vulnerability to outliers, underscoring the importance of using robust methods in data contaminated with extreme values.

---

## 4 Results

This section presents the results of the analysis comparing the performance of Ordinary Least Squares (OLS), Robust Weighted Least Squares (RWLS) using Huber and Tukey's Bisquare functions, and Least Trimmed Squares (LTS). The analysis was performed on both the original dataset and a modified version with introduced outliers. The objective was



to observe how the presence of outliers affects parameter estimates, standard errors, and t-values across different estimation techniques.

**Table 4** Parameter Estimates, Standard Errors, and t-values for Original and Modified Datasets

Estimator	Coefficient	Standard Error (Original)	t-value (Original)	Standard Error (Modified)	t-value (Modified)
OLS	0.500	0.120	4.17	0.310	1.61
RWLS (Huber)	0.480	0.115	4.17	0.135	3.56
RWLS (Tukey)	0.475	0.110	4.32	0.125	3.80
LTS	0.490	0.113	4.34	0.140	3.50

Table 4 compares the parameter estimates, standard errors, and t-values for various estimators, applied to both the original and modified datasets (with outliers). This comparison provides insight into how each method manages the data's irregularities, particularly the impact of outliers on the reliability of the estimates.

Starting with Ordinary Least Squares (OLS), the method produced a coefficient estimate of 0.500 in the original dataset, with a standard error of 0.120 and a t-value of 4.17. These values suggest that, under ideal conditions, OLS performs reasonably well, showing a significant relationship between the predictor and the response variable. However, when outliers were introduced in the modified dataset, OLS's performance drastically declined. The standard error increased substantially to 0.310, causing the t-value to drop sharply to 1.61, rendering the coefficient statistically insignificant. This drastic shift highlights OLS's sensitivity to outliers, as the outliers inflated the error and weakened the significance of the coefficient estimate. The increase in the standard error and decline in t-value demonstrate that OLS is unreliable in the presence of extreme values.

In contrast, RWLS (Huber) showed greater robustness. In the original dataset, RWLS with Huber's function produced a coefficient estimate of 0.480, with a standard error of 0.115 and a t-value of 4.17, closely resembling OLS. However, when applied to the modified dataset, Huber's method demonstrated its ability to limit the influence of outliers. The standard error increased only slightly to 0.135, and the t-value remained relatively high at 3.56. This indicates that Huber's function effectively down-weighted the extreme observations, maintaining the significance of the coefficient. The smaller increase in the standard error, compared to OLS, reflects Huber's robustness, as it mitigated the outlier impact while still producing reliable estimates.

RWLS (Tukey's Bisquare) function consistently performed the best. In the original dataset, Tukey's method yielded a coefficient estimate of 0.475, with the lowest standard error of 0.110 and the highest t-value of 4.32. This result suggests that Tukey's Bisquare function provides the most precise estimates under normal conditions. In the modified dataset, Tukey's function maintained its superiority, producing a standard error of 0.125 and a t-value of 3.80. The small increase in standard error and the relatively high t-value reflect Tukey's ability to aggressively down-weight outliers, leading to the most reliable and robust estimates among all the methods tested. The minimal change in the t-value indicates that Tukey's function is particularly effective at preserving the significance of the coefficient even in the presence of data irregularities.

For Least Trimmed Squares (LTS), the performance was generally robust, though slightly less efficient than RWLS. In the original dataset, LTS produced a coefficient estimate of 0.490, with a standard error of 0.113 and a t-value of 4.34, closely aligning with RWLS methods. However, when outliers were introduced, the standard error increased to 0.140, and the t-value dropped to 3.50. While this performance was better than OLS, it did not match the efficiency of RWLS, particularly Tukey's function. The larger increase in standard error indicates that while LTS mitigates the effects of outliers, it is not as effective as the RWLS methods in reducing the overall error.

OLS demonstrated significant vulnerability to outliers, leading to inflated standard errors and reduced t-values in the modified dataset, highlighting its unreliability in the presence of extreme values. RWLS with Huber's function provided more stability, maintaining lower standard errors and higher t-values, indicating better resistance to outliers. RWLS with Tukey's Bisquare function was the most robust, consistently providing the lowest standard errors and highest t-values, both in the original and modified datasets. LTS also improved robustness compared to OLS, though it was slightly

less efficient than RWLS methods. Overall, RWLS with Tukey's Bisquare function offered the most reliable and precise estimates, demonstrating superior robustness in handling data with outliers.

---

## 5 Conclusion

This study highlights the limitations of Ordinary Least Squares (OLS) in regression analysis, particularly when the assumptions of homoscedasticity and the absence of outliers are violated. OLS is highly sensitive to these issues, resulting in inflated standard errors, biased estimates, and unreliable statistical inferences. To address these shortcomings, several robust methods were evaluated for their ability to provide more reliable parameter estimates.

M-estimation improves robustness by down-weighting outliers, resulting in more stable estimates, particularly in the presence of large residuals. Least Trimmed Squares (LTS) goes further by trimming the largest residuals, which makes it effective for datasets with a high proportion of outliers. Both methods offer more reliable alternatives to OLS. Robust Weighted Least Squares (RWLS), which combines the advantages of Weighted Least Squares (WLS) and robust techniques like M-estimation, was found to be the most effective. Specifically, RWLS using Tukey's Bisquare function provided the most reliable estimates, with the lowest standard errors and highest t-values, proving to be the best method for handling both heteroscedasticity and outliers.

The key takeaway from the study is the superiority of robust methods over OLS in scenarios involving data irregularities. Among these, RWLS with Tukey's Bisquare function emerged as the most reliable estimator. The study suggests that future research could explore integrating other robust techniques, applying these methods to non-linear regression models, and developing adaptive tuning methods for better performance in high-dimensional datasets.

In conclusion, robust regression techniques, especially RWLS with Tukey's Bisquare function, play an essential role in delivering accurate and dependable regression results in the presence of heteroscedasticity and outliers. These methods are critical as data complexity increases across various fields.

---

## Compliance with ethical standards

### *Disclosure of conflict of interest*

No conflict of interest to be disclosed.

---

## References

- [1] Bhatti, S. H., Khan, F. W., Irfan, M., & Raza, M. A. (2023). An effective approach towards efficient estimation of general linear model in case of heteroscedastic errors. *Communications in Statistics-Simulation and Computation*, 52(2), 392-403.
- [2] Boukerche, A., Zheng, L., & Alfandi, O. (2020). Outlier detection: Methods, models, and classification. *ACM Computing Surveys (CSUR)*, 53(3), 1-37.
- [3] Kim, J., & Li, J. C. H. (2023). Which robust regression technique is appropriate under violated assumptions? A simulation study. *Methodology*, 19(4), 323-347.
- [4] Knaub Jr, J. R. (2021). When Would Heteroscedasticity in Regression Occur?. *Pakistan Journal of Statistics*, 37(4).
- [5] Reddy, T. A., & Henze, G. P. (2023). Linear Regression Analysis Using Least Squares. In *Applied Data Analysis and Modeling for Energy Engineers and Scientists* (pp. 169-221). Cham: Springer International Publishing.
- [6] Ruan, Y. (2024). Exploring Multiple Regression Models: Key Concepts and Applications. *Science and Technology of Engineering, Chemistry and Environmental Protection*, 1(7).