

Model of Robust Regression with Parametric and Nonparametric Methods

Dr. Nadia H. AL – Noor* and Asmaa A. Mohammad**

*Department of Mathematics \ College of Science\ Al-Mustansiriya University-Iraq

**Department of Mathematics \ College of Science for Women, University of Baghdad-Iraq

Corresponding E-mail:nadialnoor@yahoo.com

Abstract

In the present work, we evaluate the performance of the classical parametric estimation method "ordinary least squares" with the classical nonparametric estimation methods, some robust estimation methods and two suggested methods for conditions in which varying degrees and directions of outliers are presented in the observed data. The study addresses the problem via computer simulation methods. In order to cover the effects of various situations of outliers on the simple linear regression model, samples were classified into four cases (no outliers, outliers in the X -direction, outliers in the Y -direction and outliers in the XY -direction) and the percentages of outliers are varied between 10%, 20% and 30%. The performances of estimators are evaluated in respect to their mean squares error and relative mean squares error.

Keywords: Simple Linear Regression model; Ordinary Least Squares Method; Nonparametric Regression; Robust Regression; Least Absolute Deviations Regression; M-Estimation Regression; Trimmed Least Squares Regression.

1. Introduction

The simple linear regression model is expressed as:

$$Y = \beta_0 + \beta_1 X + \varepsilon \quad (1)$$

Where: Y is called response variable or dependent variable; X is called predictor variable, regressor variable or independent variable, and ε is called prediction error or residual. The symbols β_0 and β_1 are called intercept and slope respectively which they represents the linear regression unknown parameters or coefficients.

The process of estimating the parameters of regression model is still one of important subjects despite of large number of papers and studies written in this subject which differ in techniques followed in the process of estimation. The ordinary least squares (OLS) method is the most popular classical parametric regression technique in statistics and it is often used to estimate the parameters of a model because of nice property and ease of computation. According to Gauss-Marcov theorem, the OLS estimators, in the class of unbiased linear estimators, have minimum variance i.e. they are best linear unbiased estimator (BLUE)[10]. Nonetheless, the OLS estimates are easily affected by the presence of outliers, "outliers are observations which are markedly different from the bulk of the data or from the pattern set by the majority of the observations. In a regression problem, observations corresponding to excessively large residuals are treated as outliers[18]", and will produce inaccurate estimates. The breakdown point of the OLS estimator is 0% which implies that it can be easily affected by a single outlier. So alternative methods such as nonparametric and robust methods should be put forward which are less affected by the outliers. However, most robust methods are relatively difficult and computationally complicated. As an alternative to OLS, least absolute deviations regression (LAD or L1) has been proposed by Boscovich in 1757, then Edgeworth in 1887. LAD regression is the first step toward a more robust regression [22][26]. The next direct step to obtain robust regression was the use of M-estimators. The class of M-estimators was defined by Huber (1964, 1968) for the location model and extended by him to the regression model in (1973) [12] as an alternative robust regression estimator to the least squares. This method based on the idea of replacing the squared residual in OLS by another symmetric function, ρ , of the residuals [13]. Rousseeuw and Yohai (1984) [24] introduced the Trimmed Least Squares (TLS) regression which is a highly robust method for fitting a linear regression model. The TLS estimator minimizes the sum of the (h) smallest squared residuals. Alma (2011) [1] compare some robust regression methods such that TLS and M-estimate against OLS regression estimation method in terms of the determination of coefficient. Bai (2012) [3] review various robust regression methods including "M-estimate and TLS estimate" and compare between them based on their robustness and efficiency through a simulation study where $n=20,100$. In other side, Theil (1950) [27] introduced a nonparametric procedure which is expected to perform well without regard to the distribution of the error terms. This procedure is based on ranks and uses the median as robust measures rather than using the mean as in OLS. Mood and Brown (1950) [19] proposed to estimate the intercept and slope simultaneously from two equations depending upon divide the observations for two groups according to the median of the variable

(X). Conover (1980) [5] calculate the estimate of the intercept by used the median of the response variables, estimated Thiel's slope and the median of the explanatory variables. Hussain and Sprent (1983) [14] presented a simulation study in which they compared the OLS regression estimator against the Theil pairwise median and weighted Theil estimators in a study using 100 replications per condition. Hussain and Sprent characterized the data modeled in their study as typical data patterns that might result from contamination due to outliers. Contaminated data sets were generated using a mixture model in which each error term is either a random observation from a unit normal distribution $[N(0,1)]$ or an observation from a normal distribution with a larger variance $[N(0, k^2), k > 1]$. Jajo (1989) [15] carried a simulation study to compare the estimators that obtained from (Thiel, Mood-Brown, M-estimation and Adaptive M-estimation) with the estimators that obtained from least squares of the simple linear regression model in the presence of outliers. Mutan (2004) [20] introduced a Monte Carlo simulation study to comparing regression techniques including (ordinary least squares, least absolute deviations, trimmed least squares, Theil and weighted Theil) for the simple linear regression model when the distribution of the error terms is Generalized Logistic. Meenai and Yasmeen (2008) [17] applied nonparametric regression methods to some real and simulated data.

In the present work, we evaluate the performance of the classical nonparametric estimation methods, some robust estimation methods "least absolute deviations, M-estimation and trimmed least squares" and two suggested methods "depending upon nonparametric and M-estimation" with the OLS estimation method for conditions in which varying degrees and directions of outliers are presented in the observed data. The study addresses the problem via computer simulation methods. In order to cover the effects of various situations of outliers on the simple linear regression model, samples were classified into four cases (no outliers, outliers in the X -direction, outliers in the Y -direction "error distributed as contaminated normal", and outliers in the XY -direction) and the percentages of outliers are varied between 10%, 20% and 30% . The performances of estimators are evaluated in respect to their mean squares error and relative mean squares error.

2. Classical Estimation Method for Regression Parameters [16]

The most well-known classical parametric method of estimating the regression parameters is to use a least square error (LSE) approach. The basic idea of ordinary least squares is to optimize the fit by minimizing the total sum of the squares of the errors (deviations) between the observed values y_i and the estimated values

$$\hat{\beta}_0 + \hat{\beta}_1 x_i : \sum_{i=1}^n r_i^2 = \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2 \quad (2)$$

where $\hat{\beta}_0$ and $\hat{\beta}_1$ are estimates of β_0 and β_1 , respectively. The least squares estimators of β_0 and β_1 , $\hat{\beta}_0$ and $\hat{\beta}_1$ are:

$$\hat{\beta}_1^{OLS} = \frac{\sum_{i=1}^n x_i y_i - (\sum_{i=1}^n x_i)(\sum_{i=1}^n y_i)/n}{\sum_{i=1}^n x_i^2 - (\sum_{i=1}^n x_i)^2/n} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} \quad (3)$$

$$\hat{\beta}_0^{OLS} = \bar{y} - \hat{\beta}_1 \bar{x} \quad (4)$$

Where: $\bar{y} = (1/n) \sum_{i=1}^n y_i$ and $\bar{x} = (1/n) \sum_{i=1}^n x_i$

3. Alternative Estimation Methods for Regression Parameters

3.1 Nonparametric Regression [5][11][14][15][21][27]

The OLS regression method described above assume normally distributed error terms in the regression model. In distinction, classical nonparametric methods to linear regression typically employ parameter estimation methods that are regarded as distribution free. Since nonparametric regression procedures are developed without relying on the assumption of normality of error distributions, the only presupposition behind such procedures is that the errors of prediction are independently and identically distributed (i.i.d.). Many nonparametric procedures are based on using the ranks of the observed data rather than the observed data themselves. The robust estimate of slope for nonparametric fitted line was first described by Theil (1950). He proposed two methods, namely, the complete and the incomplete method. Assumed that all the x_i 's are distinct, and lose no generality that the x_i 's are arranged in ascending order. The complete Theil slope estimate is computed by comparing each data pair to all others in a pairwise fashion. A data set of $n (X,Y)$ pairs will result in $N = \binom{n}{2} = \frac{n(n-1)}{2}$ pairwise comparisons. For each of these comparisons a slope $\Delta Y/\Delta X$ is computed. The median of all possible pairwise slopes is taken as the nonparametric Thiel's slope estimate, $\hat{\beta}_1^{Thiel}$, Where:

$$S_{ij} = \frac{\Delta Y}{\Delta X} = \frac{y_j - y_i}{x_j - x_i} \quad ; x_i \neq x_j, 1 \leq i < j \leq n \quad (5)$$

$$\hat{\beta}_1^{Thiel} = \text{median}(S_{ij}) ; 1 \leq i < j \leq n \quad (6)$$

For incomplete method, Theil suggested using only a subset of all S_{ij} , and took as estimator of β_1 the median of the subset ($S_{i,i+n^*}$); where:

$$S_{i,i+n^*} = \frac{y_{i+n^*} - y_i}{x_{i+n^*} - x_i} \quad ; i = 1, 2, \dots, n^* \quad (7)$$

If n is even then $n^* = n/2$. If n is odd, the observation with rank $(n+1)/2$ is not used. The incomplete Theil's slope estimator is:

$$\hat{\beta}_1^{Thiel^*} = \text{median}(S_{i,i+n^*}) ; i = 1, 2, \dots, n^* \quad (8)$$

For estimation the intercept parameter, Thiel's intercept estimate, $\hat{\beta}_0^{TH}$, is defined as:

$$\hat{\beta}_0^{TH} = \text{median}(y_i - \hat{\beta}_1^{TH} x_i) ; i = 1, 2, \dots, n \quad (9)$$

Where $\hat{\beta}_1^{TH}$ is the estimate of β_1 according to the complete or the incomplete Thiel's slope estimator.

Other estimators of intercept have been suggested. Conover suggested estimating β_0 by using the formula:

$$\hat{\beta}_0^{CON} = \text{median}(y_i) - \hat{\beta}_1^{TH} \cdot \text{median}(x_i) \quad (10)$$

This formula "Conover's estimator" assures that the fitted line goes through the point (X_{median}, Y_{median}) . This is analogous to OLS, where the fitted line always goes through the point (\bar{x}, \bar{y}) .

3.2 Robust Regression

Any robust method must be reasonably efficient when compared to the least squares estimators; if the underlying distribution of errors are independent normal, and substantially more efficient than least squares estimators, when there are outlying observations. There are various robust methods for estimation the regression parameters. The main focus of this subsection is to least absolute deviations regression, M-estimation and trimmed least squares regression which are the most popular robust regression coefficients with outliers.

3.2.1 Least Absolute Deviations Regression [4][8][20][25]

The least absolute deviations regression (LAD regression) is one of the principal alternatives to the ordinary least squares method when one seeks to estimate regression parameters.

The goal of the LAD regression is to provide a robust estimator which is minimized the sum of the absolute residuals.

$$\min \sum_{i=1}^n |r_i| \quad (11)$$

The LAD procedure was developed to reduce the influence of Y -outliers in the OLS. The Y -outliers have less impact on the LAD results, because it does not square the residuals, and then the outliers are not given as much weight as in OLS procedure. However, LAD regression estimator is just as vulnerable as least squares estimates to high leverage outliers (X -outliers). In fact, LAD estimate have low breakdown point (BP is $1/n$ or 0%). Although the concept of LAD is not more difficult than the concept of the OLS estimation, calculation of the LAD estimates is more troublesome. Since there are no exact formulas for LAD estimates, an algorithm is used. Birkes and Dodge (1993) explain this algorithm for the simple linear regression model. It is known that LAD regression line passes through two of the data points. Therefore, the algorithm begins with one of the data points, denoted by (x_0, y_0) , and tries to find the best line passing through it. The procedure for finding the best line among all lines passing through a given data point (x_0, y_0) is describe below.

For each data point (x_i, y_i) , the slope of the line passing through the two points (x_0, y_0) and (x_i, y_i) is calculated and it is equal to the $(y_i - y_0)/(x_i - x_0)$. If $x_i = x_0$ for some i , the slope is not defined. The data points are re-indexed in such a way that: $(y_1 - y_0)/(x_1 - x_0) \leq (y_2 - y_0)/(x_2 - x_0) \leq \dots \leq (y_n - y_0)/(x_n - x_0)$. Now, the searched point (x_j, y_j) is determined by the index j for which.

$$\left. \begin{aligned} |x_1 - x_0| + \dots + |x_{j-1} - x_0| &< \frac{1}{2}T \\ |x_1 - x_0| + \dots + |x_{j-1} - x_0| + |x_j - x_0| &> \frac{1}{2}T \end{aligned} \right\} \quad (12)$$

Where $T = \sum_{i=1}^n |x_i - x_0|$.

This conditions guarantee that $\hat{\beta}_1$ minimizes the quantity $\sum_{i=1}^n |(y_i - y_0) - \hat{\beta}_1(x_i - x_0)|$

Analogously to $\sum|r_i|$ for the regression lines passing through (x_0, y_0) . The $\hat{\beta}_0$ is computed in such a way that the regression line crosses (x_0, y_0) . So, the best line passing through (x_0, y_0) is the line $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$ where:

$$\hat{\beta}_1^{LAD} = \frac{y_j - y_0}{x_j - x_0} \quad (13)$$

$$\hat{\beta}_0^{LAD} = y_0 - \hat{\beta}_1^{LAD} x_0 \quad (14)$$

We can equally verify that it passes through the data point (x_j, y_j) . We just have to rename the point (x_j, y_j) by (x_1, y_1) and restart.

3.2.2 M-Estimation Regression [1][3][10][12]

The most common general method of robust regression is M-estimation, introduced by Huber (1973). The M in M-estimates stands for "maximum likelihood type". That is because M-estimation is a generalization of maximum likelihood estimates (MLE). The goal of M-estimation is minimized a sum of less rapidly increasing functions of the residuals, $\sum_{i=1}^n \rho\left(\frac{r_i}{s}\right)$ where s is an estimate of scale which can be estimated by using the formula:

$$s = \frac{\text{median}|r_i - \text{median}(r_i)|}{0.6745} \quad (15)$$

A reasonable ρ should satisfy the following properties: $\rho(r) \geq 0$; $\rho(r) = \rho(-r)$; $\rho(0) = 0$; $\rho(r_i) \geq \rho(r_j)$ for $|r_i| \geq |r_j|$

M-estimators are robust to outliers in the response variable with high efficiency. However, M-estimators are just as vulnerable as least squares estimates to high leverage outliers. In fact, the BP (breakdown point) of M-estimates is $1/n$ or 0%. Suppose simple linear regression model, the M-estimator minimizes the objective function:

$$\begin{aligned} \sum_{i=1}^n \rho\left(\frac{r_i}{s}\right) &= \sum_{i=1}^n \rho\left(\frac{y_i - \beta_0 - \beta_1 x_i}{s}\right) \\ &= \sum_{i=1}^n \rho\left(\frac{r_i(\beta)}{s}\right) = \sum_{i=1}^n \rho(u_i) \end{aligned} \quad (16)$$

Where $u_i = \frac{r_i(\beta)}{s}$ are called standardized residuals. Let $\psi(u) = \rho(u)$

Differentiating (16) with respect to β and setting the partial derivatives to zero, we get the normal equations:

$$\left. \begin{aligned} \sum_{i=1}^n \psi\left(\frac{r_i(\beta)}{s}\right) &= 0 \\ \sum_{i=1}^n \psi\left(\frac{r_i(\beta)}{s}\right) x_i &= 0 \end{aligned} \right\} \quad (17)$$

To solve (17) we define the weight function $W(x) = \frac{\psi(x)}{x}$; if $x \neq 0$ and $W(x) = \dot{\psi}(0)$; if $x = 0$. let $w_i = W(u_i)$.

Then equations (17) can be written as

$$\left. \begin{aligned} \sum_{i=1}^n w_i (y_i - \beta_0 - \beta_1 x_i) &= 0 \\ \sum_{i=1}^n w_i (y_i - \beta_0 - \beta_1 x_i) x_i &= 0 \end{aligned} \right\} \quad (18)$$

Solving the estimating equations¹ (18) is a weighted least squares problem, minimizing $\sum_{i=1}^n w_i^2 u_i^2$. The weights, however, depend upon the residuals, the residuals depend upon the estimated coefficients, and the estimated coefficients depend upon the weights. An iterative solution (called iteratively reweighted least squares) is therefore required. So, the solution of (18) can be found by iterating between w_i and β :

1. Select an initial estimates $\hat{\beta}_0^{(0)}$ and $\hat{\beta}_1^{(0)}$, such as the least squares estimates.
2. At each iteration t , calculate standardized residuals $u_i^{(t-1)}$ and associated weights $w_i^{(t-1)} = W(u_i^{(t-1)})$ from the previous iteration.

¹ Newton-Raphson and Iteratively Reweighted Least Squares (IRLS) are the two methods to solve the M-estimates nonlinear normal equations. IRLS is the most widely used in practice and we considered for this study.

3. Solve for new weighted least squares estimates $\hat{\beta}_1^{(t)}, \hat{\beta}_0^{(t)}$.

$$\hat{\beta}_1^{(t)} = \frac{(\sum_{i=1}^n w_i^{(t-1)})(\sum_{i=1}^n w_i^{(t-1)} x_i y_i) - (\sum_{i=1}^n w_i^{(t-1)} x_i)(\sum_{i=1}^n w_i^{(t-1)} y_i)}{(\sum_{i=1}^n w_i^{(t-1)})(\sum_{i=1}^n w_i^{(t-1)} x_i^2) - (\sum_{i=1}^n w_i^{(t-1)} x_i)^2} \quad (19)$$

$$\hat{\beta}_0^{(t)} = \frac{(\sum_{i=1}^n w_i^{(t-1)} x_i^2)(\sum_{i=1}^n w_i^{(t-1)} y_i) - (\sum_{i=1}^n w_i^{(t-1)} x_i)(\sum_{i=1}^n w_i^{(t-1)} x_i y_i)}{(\sum_{i=1}^n w_i^{(t-1)})(\sum_{i=1}^n w_i^{(t-1)} x_i^2) - (\sum_{i=1}^n w_i^{(t-1)} x_i)^2} \quad (20)$$

Also, we can find $\hat{\beta}_0^{(t)}$ as:

$$\hat{\beta}_0^{(t)} = \frac{\sum_{i=1}^n w_i^{(t-1)} y_i}{\sum_{i=1}^n w_i^{(t-1)}} - \hat{\beta}_1^{(t)} \frac{\sum_{i=1}^n w_i^{(t-1)} x_i}{\sum_{i=1}^n w_i^{(t-1)}} \quad (21)$$

4. Repeat step 2 and step 3 until the estimated coefficients converge. The iteration process continues until some convergence criterion is satisfied, $|\hat{\beta}^{(t)} - \hat{\beta}^{(t-1)}| \cong 0$.

Several choices of ρ have been proposed by various authors. Two of these are presented in table (1) together with the corresponding derivatives (ψ) and the resulting weights (w).

Table (1): Different ρ functions, together with the corresponding derivatives ψ and the resulting weights w

Type	$\rho(r_i)$	$\psi(r_i)$	$w(r_i)$
Huber	$\begin{cases} \frac{1}{2} r_i^2 & ; r_i \leq c \\ c(r_i - \frac{1}{2}c) & ; r_i > c \end{cases}$	$\begin{cases} r_i & ; r_i \leq c \\ c \text{ sign}(r_i) & ; r_i > c \end{cases}$	$\begin{cases} 1 & ; r_i \leq c \\ \frac{c}{ r_i } & ; r_i > c \end{cases}$
	$c = 1.345, 1.5, 1.7, 2.08$		
Welsch	$\frac{c^2}{2} \left(1 - e^{-\left(\frac{r_i}{c}\right)^2}\right); r_i < \infty$	$r_i e^{-\left(\frac{r_i}{c}\right)^2}; r_i < \infty$	$e^{-\left(\frac{r_i}{c}\right)^2}; r_i < \infty$
	$c = 2.4, 2.985$		

3.2.3 Trimmed Least Squares Regression [3][23][24]

Rousseeuw and Yohai (1984) proposed the trimmed least squares (TLS) estimator regression. Extending from the trimmed mean, TLS regression minimizes the h out of n ordered squared residuals. So, the objective function is minimize the sum of the smallest h of the squared residuals and is defined as:

$$\min \sum_{i=1}^h r_{(i)}^2 \quad (22)$$

where $r_{(i)}^2$ represents the i^{th} ordered squared residuals $r_{(1)}^2 \leq r_{(2)}^2 \leq \dots \leq r_{(n)}^2$ and h is called the trimming constant which has to satisfy $\frac{n}{2} < h < n$. This constant, h , determines the breakdown point of the TLS estimator. Using $h = [(n / 2) + 1]$ ensures that the estimator has a breakdown point equal to 50%. When $h = n$, TLS is exactly equivalent to OLS estimator whose breakdown point is 0%. Rousseeuw and Leroy (1987) recommended $h = [n (1 - \alpha) + 1]$ where α is the trimmed percentage. This estimator is attractive because can be selected to prevent some of the poor results other 50% breakdown estimator show. TLS can be fairly efficient if the number of trimmed observations is close to the number of outliers because OLS is used to estimate parameters from the remaining h observations.

4. Suggested Estimators

4.1 First Suggested Estimator: in this estimator, we suggest to modifying Thiel estimator (complete and incomplete method). Thiel suggest using the median as a robust estimator of location instead of the mean in OLS. So, we suggest using the Gastwirth's estimator instead of median in Thiel estimator in order to not exclude too much of the information from the regression. Gastwirth's location estimator is a weighted sum of three order statistics. It is based on median with two ordered observations and therefore it contains information regarding the sample more than the median. The formula Gastwirth's location estimator is [9]:

$$\text{GAS} = 0.3 x_{[\frac{n}{3}+1]} + 0.4 \text{ median} + 0.3 x_{(n-[\frac{n}{3}]})} \quad (23)$$

Where: $\left[\frac{n}{3} + 1 \right]$: The integer part of the real number $\left(\frac{n}{3} + 1 \right)$ and $\left[\frac{n}{3} \right]$: The integer part of the real number $\left(\frac{n}{3} \right)$.

4.2 Second Suggested Estimator: in this estimator, we suggest to use the following function as M-estimator which satisfies the proprieties of ρ function.

$$\rho(r_i) = \frac{c}{18} \log \left(1 + \left(\frac{3r_i}{c} \right)^2 \right) \quad ; |r_i| < \infty, c$$

$$= 9 \tag{24}$$

The ψ function will be as follow:

$$\psi(r_i) = \frac{r_i/c}{1 + \left(\frac{3r_i}{c} \right)^2} \quad ; |r_i| < \infty, c$$

$$= 9 \tag{25}$$

5. Simulation Study

In this section we introduced the simulation study which has been carried out to illustrate the robustness of the estimators under different cases. Simulation was used to compare the mean squares error (MSE) and relative mean squares error (RMSE) of the estimates of regression coefficients and model by using the ordinary least squares (OLS); least absolute deviation (LAD); nonparametric estimators contains "complete Thiel's estimator (CTH) and incomplete Thiel's estimator (ITH) with Conover's estimator for intercept"; suggested nonparametric estimator contains "complete Gastwirth's estimator (CGAS) and incomplete Gastwirth's estimator (IGAS) with Conover's estimator for intercept"; M-estimators "Huber's M-estimators with $c=1.345$ (H-M), Welsch's M-estimators with $c=2.4$ (W-M) and suggested M-estimator (SU-M)" and trimmed least squares (TLS) with proportion of trimmed (α) equal to (10%, 20%, 30% and 40%). The data sets are generated from the simple linear regression model as: $y_i = 1 + 3x_i + \varepsilon_i$ which means that the true value of regression parameters are $\beta_0 = 1$ and $\beta_1 = 3$. Since the parameters known, a detailed comparison can be made. The process was repeated 1000 times to obtain 1000 independent samples of Y and X of size n . The sample sizes varied from small (10), to medium (30) and large (50). In order to cover the effects of various situations on the regression coefficients and model, samples were classified into four cases, three of them where contaminated with outliers. In addition, three percentages of outliers (δ) were considered, $\delta = 10\%$, 20% and 30% . We treated with normal and contaminated normal distribution. The simulation programs were written using Visual Basic6 programming language.

Case (1) No-outliers "Normal Case":

- Generate errors, $\varepsilon_i \sim N(0,1)$; $i = 1, 2, \dots, n$.
- Generate the values of independent variable, $x_i \sim N(0,100)$; $i = 1, 2, \dots, n$.
- Compute the y_i values.

Case (2) X-outliers:

- Generate errors, $\varepsilon_i \sim N(0,1)$; $i = 1, 2, \dots, n$.
- Generate the values of independent variable with no X-outliers, $x_i \sim N(0,100)$; $i = 1, 2, \dots, n(1 - \delta)$.
- Generate ($n\delta$) of X-outliers for the values of independent variable, $x_i \sim N(100,100)$; $i = n(1 - \delta) + 1, n(1 - \delta) + 2, \dots, n$.
- Compute the y_i values.

Case (3) Y-outliers:

- Generate the values with no Y-outliers using errors, $\varepsilon_i \sim N(0,1)$; $i = 1, 2, \dots, n(1 - \delta)$.
- Generate the values with Y-outliers using errors, $\varepsilon_i \sim N(0,50)$; $i = n(1 - \delta) + 1, n(1 - \delta) + 2, \dots, n$.
- Generate the values of independent variable, $x_i \sim N(0,100)$; $i = 1, 2, \dots, n$.
- Compute the y_i values.

Case (4) XY-outliers:

- Generate the values with no Y-outliers using errors, $\varepsilon_i \sim N(0,1)$; $i = 1, 2, \dots, n(1 - \delta)$.
- Generate the values with Y-outliers using errors, $\varepsilon_i \sim N(0,50)$; $i = n(1 - \delta) + 1, n(1 - \delta) + 2, \dots, n$.
- Generate the values of independent variable with no X-outliers, $x_i \sim N(0,100)$; $i = 1, 2, \dots, n(1 - \delta)$.
- Generate ($n\delta$) of X-outliers for the values of independent variable, $x_i \sim N(100,100)$; $i = n(1 - \delta) + 1, n(1 - \delta) + 2, \dots, n$.
- Compute the y_i values.

For each case, random samples of size n were chosen and from each sample thus obtained, MSE and RMSE using OLS, LAD, CTH, ITH, CGAS, IGAS, H-M, W-M, SU-M, TLS10%, TLS 20%, TLS 30% and TLS 40% were found and compared. MSE can be a useful measure of the quality of parameter estimation and is computed as:

$$MSE(\hat{\beta}) = Var(\hat{\beta}) + [Bias(\hat{\beta})]^2 \quad (26)$$

$$Bias(\hat{\beta}) = \bar{\beta} - \beta; Var(\hat{\beta}) = \frac{\sum_{I=1}^R (\hat{\beta}(I) - \bar{\beta})^2}{R-1}; \bar{\beta} = \frac{\sum_{I=1}^R \hat{\beta}(I)}{R}; R = 1000$$

$$MSE(\hat{Y}) = \frac{\sum_{I=1}^R MSE(\hat{Y}(I))}{R} \quad (27)$$

$$MSE(\hat{Y}(I)) = \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{n-2}$$

A relative mean squares error has also been used as a measure of the quality of parameter estimation. We computed RMSE as:

$$RMSE(\hat{\beta}) = \frac{MSE(\hat{\beta}^{OLS}) - MSE(\hat{\beta}^{other\ method})}{MSE(\hat{\beta}^{OLS})} \quad (28)$$

$$RMSE(\hat{Y}) = \frac{MSE(\hat{Y}^{OLS}) - MSE(\hat{Y}^{other\ method})}{MSE(\hat{Y}^{OLS})} \quad (29)$$

The formulation (28) is useful for comparing estimator performance and is interpreted as a proportionate (or percent) change from baseline, using the OLS estimator MSE within a given data condition as a baseline value [21]. Positive values of RMSE refer to the proportional reduction in the MSE of a given estimator with respect to OLS estimation. Hence, RMSE is interpreted as a relative measure of performance above and beyond that of the OLS estimator.

6. Conclusions and Recommendations

Based on simulation results that have been shown in tables (2)...(5), the following conclusions could be reached:

Under ideal conditions (unit normal error distribution, no contamination) "normal case", table (2) note the following:

- ❖ OLS indicates the best performance (as expected) for all sample sizes. The decline in the performance of the rest of the estimation methods compare to the performance of ordinary least squares "which can be seen through the negative values for the RMSE" it is the only sacrifice paid by those methods in anticipation of the existence of outliers.

- ❖ Proposed method "SU-M" provided the second best performance of the estimates for all sample sizes, as well as provided a performance equal to the performance of OLS in estimating the slope with a sample size equal to 30 and 50, followed by the performance of both H-M and W-M respectively. Consequently, the method of M-estimations surpassed the performance of the alternative methods for OLS.

- ❖ In general, the MSE values of estimating the intercept are greater than the corresponding MSE values of estimating the slope. So, the results for intercept estimator need more consideration.

- ❖ The use of GAS estimator instead of median in Thiel method reduced inflation in MSE values of model as compared to OLS. From the value of RMSE we can see the reduction was between (22%-28%) for all sample sizes in complete method whereas was between (20%-26%) for $n = 30,50$ in incomplete method.

- ❖ As the sample size increases, the value of MSE decreases.

- ❖ LAD introduced better performance comparing with nonparametric estimators in estimating intercept and model.

Under contamination cases, tables (3), (4) and (5) note the following:

- ❖ Ordinary Least squares recorded a decline in performance when outlier exists while most of the other estimation methods are recorded good performances depending on the percentage and direction of contaminations.

- ❖ In general, TLS indicates the best performance for all sample sizes depending on the proportion of trimmed. TLS can be fairly efficient when the number of trimmed observations is close to the number of outliers because OLS is used to estimate parameters from the remaining h observations.

❖ The MSE values indicate that the degree of sensitivity of all methods, except the TLS in some situations, to the existence of outliers in Y -direction was small compared with the degree of sensitivity to the existence of outliers in the X -direction and XY -direction.

❖ LAD and M-estimators are very sensitive to the presence of outliers in X -direction and XY -direction. In addition, the negative values of RMSE of LAD and M-estimators in some results indicate that these methods are more affected by outliers comparing with OLS. Also, LAD estimators are more sensitive to outliers comparing with M-estimators especially for estimating intercept and model. So, LAD and M-estimators are not robust estimators against those directions, but they are robust estimators against outliers in Y -direction.

❖ Nonparametric estimators introduced better performance in the presence of outliers in X -direction and XY -direction comparing with OLS, LAD and M-estimators especially for estimating slope and model.

❖ Although the performance of nonparametric estimators are better than OLS in presence of outliers in X -direction and XY -direction, it seems less better in estimating intercept when we have no outliers, thus those estimators is not robust for estimation intercept according to criterion of a robust methods that is any robust method must be reasonably efficient when compared to the least squares estimators; if the underlying distribution of errors are independent normal, and substantially more efficient than least squares estimators, when there are outlying observations.

❖ The use of GAS estimator instead of median in Thiel method improves the performance of this method when outliers appear in Y -direction. Also this estimator improves the performance of this method in some cases when outliers appear in X -direction and XY -direction and the most improvements get when it is used in an incomplete method especially for estimating intercept and model with 10% percentage of contamination and for estimating slope and model with 30% percentage of contamination.

❖ In general, the MSE values decrease when the sample sizes increase while the MSE values increase as the proportion of contaminations "outliers" increases.

Now, after pointing to the conclusions that were obtained in the present work, the following Recommendations for future work are relevant:

❖ The poor performance of OLS estimators with the presence of outliers confirms our need for alternative methods. Therefore, before analyzing the data, we should first check the presence of outliers and then construct the necessary tests whether to see the underlying assumptions are satisfied. After that, we should conduct the appropriate estimation techniques.

❖ Choosing a nonparametric method, especially to estimate slope and model, or choosing a trimmed method when the outliers appear in X -direction or XY -direction.

❖ Choosing M-estimation and LAD method, or choosing a trimmed method when the outliers are appearing in Y -direction.

❖ When the outliers appear in X -direction or XY -direction, choose RMSE or mean absolute error (MAE) as criteria for comparing between methods to avoid dealing with the large values of MSE.

References

- [1] Alma, Ö. G. (2011). "Comparison of Robust Regression Methods in Linear Regression". *Int. J. Contemp. Math. Sciences*, Vol. 6, No. 9, pp. 409- 421.
- [2] Amphanthong, P. (2009). "Estimation of Regression Coefficients with Outliers". A Dissertation Submitted to the School of Applied Statistics, National Institute of Development Administration in Partial Fulfillment of the Requirements for the Degree of Doctor of Philosophy in Statistics.
- [3] Bai, X. (2012). "Robust Linear Regression". A report submitted in partial fulfillment of the requirements for the degree Master of Science Department of Statistics College of Arts and Sciences Kansas State University, Manhattan, Kansas.
- [4] Birkes, D. and Dodge, Y. (1993). "Alternative Methods of Regression". John Wiley & Sons Publication, New York.
- [5] Conover, W.L. (1980) "Practical nonparametric statistics". Second edition, John Wiley & Sons Publication, New York.
- [6] Dietz, E.J. (1986). "On estimating a slope and intercept in a nonparametric statistics course". Institute of Statistics Mimeograph Series No. 1689R, North Carolina State University.
- [7] Dietz, E. J. (1987) "A Comparison of Robust Estimators in Simple Linear Regression" *Communications in Statistics - Simulation and Computation*, Vol. 16, Issue 4, pp. 1209-1227.
- [8] Dodge, Y. (2008). "The Concise Encyclopedia of Statistics". Springer Science & Business Media.
- [9] Gastwirth, J. L. (1966). "On Robust Procedures". *J. Amer. Statist. Assn.*, Vol. 61, pp. 929-948.

- [10] Gulasirima, R. (2005). "Robustifying Regression Models". A Dissertation Presented to the School of Applied Statistics National Institute of Development Administration in Partial Fulfillment of the Requirements for the Degree of Doctor of Philosophy in Statistics.
- [11] Helsel D.R., and Hirsch, R.M. (2002). "Statistical methods in water resources - Hydrologic analysis and interpretation: Techniques of Water-Resources Investigations of the U.S. Geological Survey", chap. A3, book 4.
- [12] Huber, P. J. (1973). "Robust regression: Asymptotics, conjectures and Monte Carlo". *Ann. Stat.*, Vol. 1, pp. 799-821.
- [13] Huber, P. J. (1981). "Robust Statistics". John Wiley & Sons Publication, New York.
- [14] Hussain, S. S., and Sprent, P. (1983). Nonparametric Regression. *Journal of the Royal Statistical Society, Series A*, Vol. 146, pp. 182-191.
- [15] Jajo, N. K. (1989). "Robust estimators in linear regression model". A Thesis Submitted to the Second Education College\ Ibn Al-Haitham\ Baghdad University in partial fulfillment of the requirements for the degree of Master of Science in Mathematics.
- [16] Marques de Sá, J. P. (2007). "Applied Statistics Using SPSS, STATISTICA, MATLAB and R". Second Edition, Springer-Verlag Berlin Heidelberg, New York.
- [17] Meenai, Y. A. and Yasmeen, F. (2008). "Nonparametric Regression Analysis". *Proceedings of 8th Islamic Countries Conference of Statistical Sciences*, Vol. 12, PP. 415-424, Lahore-Pakistan.
- [18] Midi, H., Uraibi, H. S. and Talib, B. A. (2009). "Dynamic Robust Bootstrap Method Based on LTS Estimators". *European Journal of Scientific Research*, Vol.32, No.3, pp. 277-287
- [19] Mood, A. M. (1950). "Introduction to the theory of statistics". McGraw-Hill, New York.
- [20] Mutan, O. C. (2004). "Comparison of Regression Techniques Via Monte Carlo Simulation". A Thesis Submitted to the Graduate School of Natural and Applied Sciences of Middle East Technical University in partial fulfillment of the requirements for the degree of Master of Science in Statistics.
- [21] Nevitt, J. and Tam, H.P. (1998). "A comparison of robust and nonparametric estimators under the simple linear regression model: Multiple linear regression viewpoints", Vol. 25, pp. 54-69.
- [22] Ronchetti, E.M. (1987) "Statistical Data Analysis Based on the L1-Norm and Related Methods". North-Holland, Amsterdam.
- [23] Rousseeuw, P. J. and Leroy, A. M. (1987) "Robust Regression and Outlier Detection". John Wiley & Sons Publication, New York.
- [24] Rousseeuw, P.J. and Yohai, V. (1984). "Robust regression by means of S-estimators, *Lecture Notes in Statistics* No.26, pp. 256-272, Springer Verlag, New York.
- [25] Seber, G. A. and Lee, A. J. (2003). "Linear regression analysis". Second Edition, John Wiley & Sons Publication, New York.
- [26] Stigler, S. M. (1986). "The History of Statistics: The Measurement of Uncertainty before 1900". Harvard University Press, Cambridge.
- [27] Theil, H. (1950). "A rank - invariant method of linear and polynomial regression analysis". *Indagationes Mathematicae*, Vol. 12, pp. 85-91.

