

## Exact distribution of Hadi's ( $H^2$ ) influence measure and identification of potential outliers

G.S.David Sam Jayakumar\*<sup>†</sup>  and A. Sulthan<sup>‡</sup> 

### Abstract

This paper proposed an exact distribution of Hadi's influence measure that can be used to evaluate the potential outliers in a linear multiple regression analysis. The authors explored a relationship between the measure in terms of two independent F-ratios and they derived density function of the measure in a complicated series expression form with Gauss hyper-geometric function. Moreover, the first two moments of the distribution are derived in terms of Beta function and the authors computed the critical points of Hadi's measures at 5% and 1% significance level for different sample sizes and varying no. of predictors. Finally, the numerical example shows the identification of the potential outliers and the results extracted from the proposed approaches are more scientific, systematic and its exactness outperforms the Hadi's traditional approach.

**Keywords:** Hadi's measure, Potential outliers, Series expression form, Gauss hyper-geometric function, Moments, Beta function, Critical points.

*Mathematics Subject Classification (2010):* 62H10

*Received :* 14.12.2017 *Accepted :* 28.06.2018 *Doi :* 10.15672/HJMS.2018.614

---

\*Assistant Professor, Jamal Institute of Management, Trichy 20, India Email: samjaya77@gmail.com

<sup>†</sup>Corresponding Author.

<sup>‡</sup>Research Scholar, Jamal Institute of Management, Trichy 20, India Email: contact@iamsulthan.in

## 1. Introduction and Related work

While fitting a regression model it is well understood that not all observations in a dataset play an equal role. Some observations have more impact than the others. Those observations which significantly influence the results of a regression analysis are called influential observations. Andrews and Pregibon [2] highlighted the need to find the outliers that matter. This means not all outliers need to be harmful in that they have an undue influence on the estimation of the parameters in the regression model. Hence, examining the residuals alone might not lead us to the detection of aberrant or unusual observations. Thus, other ways for finding influential observations are needed. Hoaglin and Welsch [9] discussed the importance of the projection matrix in linear regression, where the projection matrix is the matrix that projects onto the regression space. They argued that the diagonal elements of the projection matrix are important ingredients in influence analysis. The diagonal elements are referred to as leverages since they can be thought of as the amount of leverage concerning the response value on the corresponding predicted response value. Perhaps the most well-known influence measure was proposed by Cook [6], referred as Cook's distance. It is an influence measure used for assessing the influence of the observations on the estimated parameter vector in the linear regression model. It is widely used by practitioners for detecting influential observations. There are other influence measures to use in the linear regression analysis for assessing the influence of the observations on various results of the regression analysis. Such as, Andrews and Pregibon [2] derived a measure of the influence of an observation on the estimated parameters. This measure the AP statistic is based on the change in volume of confidence ellipsoids with and without a particular observation. Moreover, Belsley et al. [3] suggested an influence measure for assessing the influence of an observation on the variance of the estimated parameters in the linear regression model, known as COVRATIO. Besides the influence measures mentioned here there exist much more, see e.g. Chatterjee and Hadi [5] and Hadi [8] for excellent overviews of influence measures. Graphical investigation of data is a powerful tool in exploratory analysis. It can be used to examine relationships between variables and discover observations deviating from other. Hence, influential observations can also be detected using graphical tools. Mosteller and Tukey [11] introduced the added variable plot, which is used for graphically detecting observations that have a large influence on the parameter estimates. For details concerning the added variable plot, such as construction and properties, see, Belsley et al. [3], where the plot is referred to as the partial regression leverage plot, and Cook and Weisberg [7]. Other results on graphical tools in influence analysis are provided by Johnson and McCulloch [10]. It is important to note that the graphical tools used in influence analysis are not conclusive, but rather suggestive. From the previous discussions, we can see that the 1970's and the 1980's were the decades when most research results on influence analysis in linear regression came to see the light. However, influence analysis in linear regression is still an active research area. Nurunnabi et al. [12] proposed a modification of Cook's distance. This modification enables the identification of multiple influential observations. Furthermore, Beyaztas and Alin [4] used a combined Bootstrap and Jackknife algorithm to detect influential observations. In applied data analysis, there is an increasing availability of data sets containing a large number of variables. When such data is in the hands of the researcher sparse regression can be implemented, which is another field of research active today. In sparse regression, a penalty term on the regression parameters is added which shrinks the number of parameters. Common approaches to estimate the parameter in the sparse regression are, however, sensitive to influential observations and new methods are needed. Alfons et al. [1] and Park et al. [13] proposed robust

estimation methods, where influential observations are not harmful to the resulting estimates. Considering the above reviews, the authors proposed the exact distribution of Hadi's influence function ( $H_i^2$ ) which exactly identifies the influential data points and it is discussed in the subsequent sections.

## 2. Relationship between Hadi's ( $H^2$ ) influence measure and F-ratios

The multiple linear regression models with random error is given by

$$(2.1) \quad Y = X\beta + e$$

where  $Y$  is the matrix of the dependent variable,  $\beta$  is the vector of beta coefficients or partial regression coefficients and  $e$  is the residual followed normal distribution

$N(0, \sigma_e^2 I_n)$ . From (2.1), statisticians concentrate and give importance to the error diagnostics such as outlier detection, identification of leverage points and evaluation of influential observations. Several error diagnostics techniques exist in the literature proposed by statisticians, but Hadi's ( $H_i^2$ ) influence measure is the interesting technique based on the simple fact that potentially influential observations are outliers in X-space, Y-space or both. The general form of the Hadi's influence measure of the  $i^{th}$  observation is given by

$$(2.2) \quad H_i^2 = \frac{(p+1) \hat{e}_i^2}{(1-h_{ii}) \left( \hat{e}^t \hat{e} - \hat{e}_i^2 \right)} + \frac{h_{ii}}{1-h_{ii}}$$

Where  $\hat{e}_i^2$  is the vector of squared estimated residuals,  $p$  is the no. of predictors,  $\hat{e}^t \hat{e}$  is the sum of the squared estimated residuals and  $h_{ii}$  is the hat values of  $i^{th}$  observation or diagonal elements of the hat matrix ( $H = X(X'X)^{-1}X'$ ). This diagnostic measure is the sum of two components each of which has an interpretation. A large value for the first term indicates that the model has a poor fit (a large prediction error) and a large value for the second term indicates the presence of an outlier in the X-space. Similarly, Hadi pointed this diagnostic measure possess several desirable properties and it is also supplemented by a graphical display which shows the source of influence. He suggested, ( $H_i^2$ ) for observations more than a cut-off of  $E(H_i^2) + c\sqrt{V(H_i^2)}$  which is treated as a potential outlier. Hadi's influence measure can also be written in an alternative form as

$$(2.3) \quad H_i^2 = \frac{p+1}{(1-h_{ii}) \left( \left( \frac{\hat{e}^t \hat{e}}{\hat{e}_i^2} \right) - 1 \right)} + \frac{h_{ii}}{1-h_{ii}}$$

It is known the unbiased estimate of the true error variance is  $s^2 = \hat{e}^t \hat{e} / n - p - 1$  and substitute  $\hat{e}^t \hat{e} = s^2(n-p-1)$  in (2.3) to get

$$(2.4) \quad H_i^2 = \frac{p+1}{\frac{\hat{e}_i^2}{s^2} - (1-h_{ii})} + \frac{h_{ii}}{1-h_{ii}}$$

Rewrite (2.4), in terms of the internally studentized residual ( $r_i$ ) which is equal to  $\hat{e}_i / s\sqrt{1-h_{ii}}$  and it is given as

$$(2.5) \quad H_i^2 = \frac{p+1}{((n-p-1)/r_i^2) - (1-h_{ii})} + \frac{h_{ii}}{1-h_{ii}}$$

Though Hadi's influence measure is scientific and the yardstick used to detect the influential observation is not scientific and the authors believe it is based on the rule of thumb approach. Because  $(H_i^2)$  is non-normally distributed and the usage of mean and variance in the cut-off  $(E(H_i^2) + c\sqrt{V(H_i^2)})$  is meaningless and illogic. Secondly, when using the cut-off, it is not recommended by the author to use a specific and fixed value for the constant  $(c)$ . Finally, the usage of plots and graphs to identify the potential outliers and sources of influence leads to imprecision and ambiguity. In order to overcome this rule of thumb approach of identifying the influential observations, authors proposed the exact distribution for Hadi's influence measure and established a scientific yardstick to scrutinize the exact influential observations. For this, authors utilize the relationship among the Hadi's  $(H_i^2)$ , internally studentized residual  $(r_i)$  and hat elements  $(h_{ii})$ . The terms  $(r_i)$  and  $(h_{ii})$  are independent because the computation of  $(r_i)$  involves the error term  $e_i \sim N(0, \sigma_e^2)$  and  $h_{ii}$  values involves the set of predictors  $(H = X(X'X)^{-1}X')$ . Therefore, from the property of least squares  $E(eX) = 0$ , so  $r_i$  and  $h_{ii}$  are also uncorrelated and independent. Using this assumption, authors first determine the distribution of  $(r_i)$  based on the relationship given by Weisberg (1980) as

$$(2.6) \quad t_i = r_i \sqrt{\frac{n-p-2}{(n-p-1) - r_i^2}} \sim t_{(n-p-2)}$$

From (2.6) it follows student's  $t$ -distribution with  $(n-p-2)$  degrees of freedom and it can be written in terms of the F-ratio as

$$r_i^2 = \frac{(n-p-1)t_i^2}{(n-p-2) + t_i^2}$$

$$(2.7) \quad r_i^2 = \frac{(n-p-1)F_{i(1, n-p-2)}}{(n-p-2) + F_{i(1, n-p-2)}}$$

From (2.7), if  $t_i$  follows student's  $t$ -distribution with  $(n-p-2)$  degrees of freedom, then  $t_i^2$  follows  $F_{(1, n-p-2)}$  distribution with  $(1, n-p-2)$  degrees of freedom. Similarly, authors identify the distribution of  $h_{ii}$  based on the relationship proposed by Belsey et al [3] and they showed when the set of predictors is multivariate normal with  $(\mu_X, \Sigma_X)$ , then

$$(2.8) \quad \frac{(n-p)(h_{ii} - 1/n)}{(p-1)(1-h_{ii})} \sim F_{(p-1, n-p)}$$

From (2.8) it follows F-distribution with  $(p-1, n-p)$  degrees of freedom and it can be written in an alternative form as

$$(2.9) \quad E(H_i^2) = (p+1)(\phi_1(p, n)) + \frac{n}{n-1}(\phi_2(p, n)) - 1$$

In order to derive the exact distribution of  $(H_i^2)$ , substitute (2.7) and (2.9) in (2.5), authors get the Hadi's  $(H_i^2)$  measure in terms of the two independent F-ratios with  $(1, n-p-2)$  and  $(p-1, n-p)$  degrees of freedom respectively and the relationship is given as

$$(2.10) \quad H_i^2 = \frac{p+1}{\frac{(n-p-2) + F_{i(1, n-p-2)}}{F_{i(1, n-p-2)}} - \frac{(n-1)/n}{1 + ((p-1)/(n-p))F_{i(p-1, n-p)}}} + \frac{((p-1)/(n-p))F_{i(p-1, n-p)} + 1/n}{(n-1)/n}$$

$$(2.11) \quad H_i^2 = \frac{p+1}{1 + (n-p-2)F_{i(n-p-2,1)} - \frac{(n-1)/n}{1+((p-1)/(n-p))F_{i(p-1,n-p)}}} + \frac{(1 + ((p-1)/(n-p))F_{i(p-1,n-p)}) - (n-1)/n}{(n-1)/n}$$

From (2.11), it can be further simplified and  $(H_i^2)$  is expressed in terms of two independent beta variables namely  $\theta_{1i}$  and  $\theta_{2i}$  of the first kind by using the following facts

$$(2.12) \quad \frac{1}{1 + (n-p-2)F_{i(n-p-2,1)}} = \theta_{1i} \sim \beta_1 \left( \frac{1}{2}, \frac{n-p-2}{2} \right)$$

$$(2.13) \quad \frac{1}{1 + ((p-1)/(n-p))F_{i(p-1,n-p)}} = \theta_{2i} \sim \beta_1 \left( \frac{n-p}{2}, \frac{p-1}{2} \right)$$

Then, without loss of generality (2.11) can be written as

$$(2.14) \quad H_i^2 = \frac{p+1}{(1/\theta_{1i}) - ((n-1)/n)\theta_{2i}} + \frac{(1/\theta_{2i}) - (n-1)/n}{(n-1)/n}$$

$$(2.15) \quad H_i^2 = \frac{(p+1)\theta_{1i}}{1 - ((n-1)/n)\theta_{1i}\theta_{2i}} + (n/n - 1)(1/\theta_{2i}) - 1$$

From (2.15), the authors showed the Hadi's influence measure in terms of  $\theta_{1i} \sim \beta_1 \left( \frac{1}{2}, \frac{n-p-2}{2} \right)$  and  $\theta_{2i} \sim \beta_1 \left( \frac{n-p}{2}, \frac{p-1}{2} \right)$  which followed beta distribution of first kind with two shape parameters  $p$  and  $n$  respectively. To avoid complexity further, the relationship from (2.15) modified as

$$(2.16) \quad \frac{n-1}{n} (1 + H_i^2) = \left( \frac{1+p((n-1)/n)\theta_{1i}\theta_{2i}}{1 - ((n-1)/n)\theta_{1i}\theta_{2i}} \right) (1/\theta_{2i}) = \psi_i$$

Based on the identified relationship from (2.16), the authors derived the distribution of the Hadi's  $(H_i^2)$  and it is discussed in the next section.

### 3. Exact Distribution of Hadi's $(H_i^2)$

Using the technique of two-dimensional Jacobian of transformation, the joint probability density function of the two Beta variables of kind-1 namely  $\theta_{1i}, \theta_{2i}$  were transformed into density function of  $\psi_i$  and it is given as

$$(3.1) \quad f(\psi_i, u_i) = f(\theta_{1i}, \theta_{2i}) |J|$$

From (3.1), It is known that  $\theta_{1i}$  and  $\theta_{2i}$  are independent then rewrite (3.1) as

$$(3.2) \quad f(\psi_i, u_i) = f(\theta_{1i}) f(\theta_{2i}) |J|$$

Using the change of variable technique, substitute  $\theta_{2i} = u_i$  in (2.16) to get

$$(3.3) \quad \theta_{1i} = \frac{\psi_i u_i - 1}{((n-1)/n) u_i (p + \psi_i u_i)}$$

Then partially differentiate (3.3) and compute the Jacobian determinant in (3.2) as

$$(3.4) \quad f(\psi_i, u_i) = f(\theta_{1i}) f(\theta_{2i}) \left| \frac{\partial(\theta_{1i}, \theta_{2i})}{\partial(\psi_i, u_i)} \right|$$

$$(3.5) \quad f(\psi_i, u_i) = f(\theta_{1i}) f(\theta_{2i}) \left| \begin{array}{cc} \frac{\partial \theta_{1i}}{\partial \psi_i} & \frac{\partial \theta_{1i}}{\partial u_i} \\ \frac{\partial \theta_{2i}}{\partial \psi_i} & \frac{\partial \theta_{2i}}{\partial u_i} \end{array} \right|$$

From (3.5), It is known that the  $\theta_{1i}$  and  $\theta_{2i}$  are independent, then the density function of the joint distribution of  $\theta_{1i}$  and  $\theta_{2i}$  is given as

$$(3.6) \quad f(\theta_{1i}, \theta_{2i}) = \frac{1}{B\left(\frac{1}{2}, \frac{n-p-2}{2}\right)} \theta_{1i}^{\frac{1}{2}-1} (1-\theta_{1i})^{\frac{n-p-2}{2}-1} \\ \times \frac{1}{B\left(\frac{n-p}{2}, \frac{p-1}{2}\right)} \theta_{2i}^{\frac{n-p}{2}-1} (1-\theta_{2i})^{\frac{p-1}{2}-1}$$

where  $0 \leq \theta_{1i}, \theta_{2i} \leq 1, n, p > 0$   
and

$$(3.7) \quad \left| \begin{array}{cc} \frac{\partial \theta_{1i}}{\partial \psi_i} & \frac{\partial \theta_{1i}}{\partial u_i} \\ \frac{\partial \theta_{2i}}{\partial \psi_i} & \frac{\partial \theta_{2i}}{\partial u_i} \end{array} \right| = \left| \begin{array}{cc} \frac{p+1}{((n-1)/n)(p+\psi_i u_i)} & \frac{p-\psi_i^2 u_i^2 + 2\psi_i u_i}{((n-1)/n)u_i^2(p+\psi_i u_i)^2} \\ 0 & 1 \end{array} \right| = \frac{p+1}{((n-1)/n)(p+\psi_i u_i)^2}$$

Then substitute (3.6) and (3.7) in (3.5) in terms of the substitution of  $u_i$ , to get the joint distribution of  $\psi_i$  and  $u_i$  as

$$(3.8) \quad f(\psi_i, u_i) = \\ \frac{1}{B\left(\frac{1}{2}, \frac{n-p-2}{2}\right)} \left( \frac{\psi_i u_i - 1}{\left(\frac{n-1}{n}\right) u_i (p + \psi_i u_i)} \right)^{\frac{1}{2}-1} \left( 1 - \frac{\psi_i u_i - 1}{\left(\frac{n-1}{n}\right) u_i (p + \psi_i u_i)} \right)^{\frac{n-p-2}{2}-1} \\ \times \frac{1}{B\left(\frac{1}{2}, \frac{n-p-2}{2}\right)} u_i^{\frac{n-p-2}{2}-1} (1-u_i)^{\frac{1}{2}-1} \times |J|$$

where  $\frac{n-1}{n} \leq \psi_i < \infty, 0 \leq u_i \leq 1$  and  $|J| = \frac{p+1}{((n-1)/n)(p+\psi_i u_i)^2}$

Rearrange (3.8) and integrate with respect to  $u_i$ , to get the marginal distribution of  $\psi_i$  as

$$(3.9) \quad f(\psi_i; p, n) = \alpha(p, n) \left( \sum_{r=0}^{\frac{n-p-2}{2}-1} \sum_{s=0}^{r+\frac{1}{2}-1} \binom{\frac{n-p-2}{2}-1}{r} \left(r + \frac{1}{2} - 1\right) \left(\frac{n-1}{n}\right)^{r+\frac{1}{2}-1} p^{-\left(r+\frac{3}{2}\right)} \right) \\ \times (-1)^{2r+s+\frac{1}{2}-1} \psi_i^s \int_0^1 u_i^{\frac{n-p-3}{2}+s-r-1} (1-u_i)^{\frac{p-1}{2}-1} \left(\frac{1}{1+(\psi_i/p)u_i}\right)^{r+\frac{3}{2}} du_i$$

where  $\frac{n-1}{n} \leq \psi_i < \infty$  and  $\alpha(p, n) = \frac{p+1}{\left(\frac{n-1}{n}\right)B\left(\frac{1}{2}, \frac{n-p-2}{2}\right)B\left(\frac{n-p}{2}, \frac{p-1}{2}\right)}$

It is known, from (3.9)

$$(3.10) \quad \int_0^1 u_i^{\frac{n-p-3}{2}+s-r-1} (1-u_i)^{\frac{p-1}{2}-1} \left(\frac{1}{1+(\psi_i/p)u_i}\right)^{r+\frac{3}{2}} du_i = \\ \frac{\Gamma\left(\frac{p-1}{2}\right)}{(\psi_i/p)\Gamma\left(r+\frac{3}{2}\right)} (\Omega_1(\psi_i; p, n, r, s) + \Omega_2(\psi_i; p, n, r, s))$$

where

$$\Omega_1(\psi_i; p, n, r, s) = (\psi_i/p)^{-\left(\frac{n-p-3}{2}+s-r\right)+1} \frac{\Gamma\left(\frac{n-p-3}{2}+s-r\right)\Gamma\left(2r+\frac{3}{2}-\left(\frac{n-p-3}{2}+s\right)\right)}{\Gamma\left(\frac{p-1}{2}\right)}$$

$$\begin{aligned}
 & {}_2F_1\left(\frac{n-p-3}{2} + s - r, \frac{p+1}{2}; 1 - \left(r + \frac{3}{2}\right) + \frac{n-p-3}{2} + s - r; \frac{1}{\psi_i/p}\right) \\
 \Omega_2(\psi_i; p, n, r, s) &= (\psi_i/p)^{-(r+\frac{3}{2})+1} \frac{\Gamma\left(r + \frac{3}{2}\right)\Gamma\left(\frac{n-p}{2} + s - 2r\right)}{\Gamma\left(\frac{p-1}{2}\right)} \\
 & {}_2F_1\left(r + \frac{3}{2}, 1 - \left(\frac{n-7}{2} + s - 2r\right); 1 + r + \frac{3}{2} - \left(\frac{n-p-3}{2} + s - r\right); -\frac{1}{\psi_i/p}\right)
 \end{aligned}$$

Then substitute (3.10) in (3.9) and arrange the terms, to get the density function of  $\psi_i$  in the series expression form as

$$\begin{aligned}
 (3.11) \quad f(\psi_i; p, n) &= \lambda(p, n) \left( \sum_{r=0}^{\frac{n-p-2}{2}-1} \sum_{s=0}^{r+\frac{1}{2}-1} \left(\frac{n-p-2}{2} - 1\right) \left(r + \frac{1}{s} - 1\right) \left(\frac{n}{n-1}\right)^{r+\frac{1}{2}-1} p^{-(r+\frac{3}{2})} \right. \\
 & \quad \left. \times (-1)^{2r+s+\frac{1}{2}-1} \psi_i^{s-1} \frac{1}{\Gamma\left(r+\frac{3}{2}\right)} (\Omega_1(\psi_i; p, n, r, s) + \Omega_2(\psi_i; p, n, r, s)) \right)
 \end{aligned}$$

where,  $(n-1)/n \leq \psi_i < \infty, n, p > 0, n > p$  and  $\lambda(p, n) = (p+1)\alpha(p, n) = \frac{(p+1)\Gamma\left(\frac{p-1}{2}\right)}{\left(\frac{n-1}{n}\right)B\left(\frac{1}{2}, \frac{n-p-2}{2}\right)B\left(\frac{n-p}{2}, \frac{p-1}{2}\right)}$  is the normalizing constant as a function of  $p$  and  $n$ . In order to derive the density function of Hadi's measure, the authors again utilize the relationship between  $\psi_i$  and  $(H_i^2)$ . It is known  $H_i^2 > 0$ , then  $(n-1)/n \leq \psi_i < \infty$ . Hence from (2.16), using one-dimensional Jacobian of transformation, the density function of  $(H_i^2)$  can be written as

$$(3.12) \quad f(H_i^2) = f(\psi_i) |J|$$

$$(3.13) \quad f(H_i^2) = f(\psi_i) \left| \frac{d\psi_i}{dH_i^2} \right|$$

Then substitute  $\psi_i = \frac{n-1}{n} (1 + H_i^2), \frac{d\psi_i}{dH_i^2} = \frac{n-1}{n}$  and (3.11) in (3.13), to get the final form of the density function of  $(H_i^2)$  as

$$\begin{aligned}
 (3.14) \quad f(H_i^2; p, n) &= \\
 \varphi(p, n) & \left( \sum_{r=0}^{\frac{n-p-2}{2}-1} \sum_{s=0}^{r+\frac{1}{2}-1} \left(\frac{n-p-2}{2} - 1\right) \left(r + \frac{1}{s} - 1\right) \left(\Gamma\left(r + \frac{3}{2}\right)\right)^{-1} \left(\frac{n}{n-1}\right)^{r-s+\frac{1}{2}} \right. \\
 & \quad \left. p^{-(r+\frac{3}{2})+1} (-1)^{2r+s+\frac{1}{2}-1} (1 + H_i^2)^{s-1} (\Omega_1(H_i^2; p, n, r, s) + \Omega_2(H_i^2; p, n, r, s)) \right)
 \end{aligned}$$

where,  $0 \leq H_i^2 < \infty, n, p > 0, n > p$  and  $\varphi(p, n) = \frac{n-1}{n} \lambda(p, n) = \frac{(p+1)\Gamma\left(\frac{p-1}{2}\right)}{B\left(\frac{1}{2}, \frac{n-p-2}{2}\right)B\left(\frac{n-p}{2}, \frac{p-1}{2}\right)}$

$$\begin{aligned}
 \Omega_1(H_i^2; p, n, r, s) &= \frac{\Gamma\left(\frac{n-p-3}{2} + s - r\right)\Gamma\left(2r + \frac{3}{2} - \left(\frac{n-p-3}{2} + s\right)\right)}{\Gamma\left(\frac{p-1}{2}\right)} \left(\frac{n-1}{np} (1 + H_i^2)\right)^{-\left(\frac{n-p-3}{2} + s - r\right)+1} \\
 & \quad \times {}_2F_1\left(\frac{n-p-3}{2} + s - r, \frac{p+1}{2}; 1 - \left(r + \frac{3}{2}\right) + \frac{n-p-3}{2} + s - r; -\frac{n-1}{np} \frac{1}{(1+H_i^2)}\right)
 \end{aligned}$$

$$\begin{aligned}
 \Omega_2(H_i^2; p, n, r, s) &= \frac{\Gamma\left(r + \frac{3}{2}\right)\Gamma\left(\frac{n-p}{2} + s - 2r\right)}{\Gamma\left(\frac{n-7}{2} + s - 2r\right)} \left(\frac{n-1}{np} (1 + H_i^2)\right)^{-\left(r+\frac{3}{2}\right)+1} \\
 & \quad \times {}_2F_1\left(r + \frac{3}{2}, 1 - \left(\frac{n-7}{2} + s - 2r\right); 1 + r + \frac{3}{2} - \left(\frac{n-p-3}{2} + s - r\right); -\frac{n-1}{np} \frac{1}{(1+H_i^2)}\right)
 \end{aligned}$$

From (3.14), it is the density function of Hadi's  $(H_i^2)$  influence measure which involves the following such as  $\Omega_1(H_i^2; p, n, r, s), \Omega_2(H_i^2; p, n, r, s)$  are the auxiliary functions,  ${}_2F_1$  is the Gauss hypergeometric function and the normalizing constant  $\lambda(p, n)$  comprised of Beta and Gamma functions  $(B\left(\frac{1}{2}, \frac{n-p-2}{2}\right), B\left(\frac{p-1}{2}, \frac{n-p}{2}\right), \Gamma\left(\frac{p-1}{2}\right))$  with two shape parameters  $(p, n), n$  is the sample size and  $p$  is the no. of predictors used in a multiple linear

regression model. In order to know the location and dispersion of Hadi's ( $H_i^2$ ), the authors derived the first two moments in terms of mean, variance from and it is shown as follows. Using (2.15), rewrite in terms of series expression form as

$$(3.15) \quad H_i^2 = (p+1) \left( \sum_{k=0}^{\infty} \left( \frac{n-1}{n} \right)^k \theta_{1i}^{k+1} \theta_{2i}^k \right) + \frac{n}{n-1} \left( \frac{1}{\theta_{2i}} \right) - 1$$

Now take expectation and substitute the moments of two independent beta variables  $\theta_{1i}$  and  $\theta_{2i}$  of kind-1, to get the first moment of ( $H_i^2$ ) as

$$E(H_i^2) = (p+1) \left( \sum_{k=0}^{\infty} \left( \frac{n-1}{n} \right)^k E(\theta_{1i}^{k+1}) E(\theta_{2i}^k) \right) + \frac{n}{n-1} E\left(\frac{1}{\theta_{2i}}\right) - 1$$

$$(3.16) \quad E(H_i^2) = (p+1) (\phi_1(p, n)) + \frac{n}{n-1} (\phi_2(p, n)) - 1$$

Where  $\phi_1(p, n) = \sum_{k=0}^{\infty} \left( \frac{n-1}{n} \right)^k \left( \frac{B(k+\frac{3}{2}, \frac{n-p-2}{2})}{B(\frac{1}{2}, \frac{n-p-2}{2})} \right) \left( \frac{B(\frac{n-p}{2}+k, \frac{p-1}{2})}{B(\frac{n-p}{2}, \frac{p-1}{2})} \right)$

$\phi_2(p, n) = \frac{B(\frac{n-p-2}{2}, \frac{p-1}{2})}{B(\frac{n-p}{2}, \frac{p-1}{2})}$  and  $B$  is the beta function respectively.

From (2.15), rewrite and square both sides, then take expectation, to get the second moment of ( $H_i^2$ ) as

$$(H_i^2 + 1)^2 = (p+1)^2 \theta_{1i}^2 \left( \frac{1}{1 - \left( \frac{n-1}{n} \right) \theta_{1i} \theta_{2i}} \right)^2 + \left( \frac{n}{n-1} \right)^2 \left( \frac{1}{\theta_{2i}} \right)^2$$

$$+ \frac{2n(p+1)}{n-1} \left( \frac{\theta_{1i}}{\theta_{2i}} \right) \frac{1}{1 - \left( \frac{n-1}{n} \right) \theta_{1i} \theta_{2i}}$$

$$(H_i^2 + 1)^2 = (p+1)^2 \left( \sum_{k=0}^{\infty} (k+1) \left( \frac{n-1}{n} \right)^k \theta_{1i}^{k+2} \theta_{2i}^k \right) + \left( \frac{n}{n-1} \right)^2 \left( \frac{1}{\theta_{2i}} \right)^2$$

$$+ \frac{2n(p+1)}{n-1} \left( \sum_{k=0}^{\infty} \left( \frac{n-1}{n} \right)^k \theta_{1i}^{k+1} \theta_{2i}^{k-1} \right)$$

$$E(H_i^2 + 1)^2 = (p+1)^2 \left( \sum_{k=0}^{\infty} (k+1) \left( \frac{n-1}{n} \right)^k E(\theta_{1i}^{k+2}) E(\theta_{2i}^k) \right) + \left( \frac{n}{n-1} \right)^2 E\left(\frac{1}{\theta_{2i}}\right)^2$$

$$+ \frac{2n(p+1)}{n-1} \left( \sum_{k=0}^{\infty} \left( \frac{n-1}{n} \right)^k E(\theta_{1i}^{k+1}) E(\theta_{2i}^{k-1}) \right)$$

$$(3.17) \quad E(H_i^2)^2 = (p+1)^2 \left( \sum_{k=0}^{\infty} (k+1) \left( \frac{n-1}{n} \right)^k E(\theta_{1i}^{k+2}) E(\theta_{2i}^k) \right) + \left( \frac{n}{n-1} \right)^2 E\left(\frac{1}{\theta_{2i}}\right)^2$$

$$+ \frac{2n(p+1)}{n-1} \left( \sum_{k=0}^{\infty} \left( \frac{n-1}{n} \right)^k E(\theta_{1i}^{k+1}) E(\theta_{2i}^{k-1}) \right) - 2E(H_i^2) - 1$$

Therefore, It is known

$$(3.18) \quad V(H_i^2) = E(H_i^2)^2 - (E(H_i^2))^2$$

Then substitute (3.16) and (3.17) in (3.18), to get the variance of ( $H_i^2$ ) as

$$V(H_i^2) = (p+1)^2 (\Phi_1(p, n)) + \left( \frac{n}{n-1} \right)^2 (\Phi_2(p, n)) + \frac{2n(p+1)}{n-1} (\Phi_3(p, n))$$

where

$$\Phi_1(p, n) = \left( \sum_{k=0}^{\infty} (k+1) \left( \frac{n-1}{n} \right)^k \frac{B(k+\frac{5}{2}, \frac{n-p-2}{2}) B(\frac{n-p}{2}+k, \frac{p-1}{2})}{B(\frac{1}{2}, \frac{n-p-2}{2}) B(\frac{n-p}{2}, \frac{p-1}{2})} \right)$$

$$- \left( \sum_{k=0}^{\infty} \left( \frac{n-1}{n} \right)^k \frac{B(k+\frac{3}{2}, \frac{n-p-2}{2}) B(\frac{n-p}{2}+k, \frac{p-1}{2})}{B(\frac{1}{2}, \frac{n-p-2}{2}) B(\frac{n-p}{2}, \frac{p-1}{2})} \right)^2$$



$$\Phi_2(p, n) = \frac{B\left(\frac{n-p-4}{2}, \frac{p-1}{2}\right)}{B\left(\frac{n-p}{2}, \frac{p-1}{2}\right)} - \left(\frac{B\left(\frac{n-p-2}{2}, \frac{p-1}{2}\right)}{B\left(\frac{n-p}{2}, \frac{p-1}{2}\right)}\right)^2$$

$$\Phi_3(p, n) = \left(\sum_{k=0}^{\infty} \binom{n-1}{n}^k \frac{B\left(k+\frac{3}{2}, \frac{n-p-2}{2}\right)B\left(\frac{n-p-2}{2}+k, \frac{p-1}{2}\right)}{B\left(\frac{1}{2}, \frac{n-p-2}{2}\right)B\left(\frac{n-p}{2}, \frac{p-1}{2}\right)}\right) - \left(\frac{B\left(\frac{n-p-2}{2}, \frac{p-1}{2}\right)}{B\left(\frac{n-p}{2}, \frac{p-1}{2}\right)}\right) \left(\sum_{k=0}^{\infty} \binom{n-1}{n}^k \frac{B\left(k+\frac{3}{2}, \frac{n-p-2}{2}\right)B\left(\frac{n-p}{2}+k, \frac{p-1}{2}\right)}{B\left(\frac{1}{2}, \frac{n-p-2}{2}\right)B\left(\frac{n-p}{2}, \frac{p-1}{2}\right)}\right)$$

By using the mean and variance of Hadi’s measure from (3.16) and (3.18), the authors established the upper control limit of  $(H_i^2)$  for different combination of  $(p, n)$  by using (3.20). Therefore

$$(3.19) \quad UCL(H_i^2) = E(H_i^2) + \sqrt{V(H_i^2)}$$

$$(3.20) \quad UCL(H_i^2) = (p+1)\phi_1(p, n) + \frac{n}{n-1}\phi_2(p, n) - 1 + \sqrt{(p+1)^2\Phi_1(p, n) + \left(\frac{n}{n-1}\right)^2\Phi_2(p, n) + \frac{2n(p+1)}{n-1}\Phi_3(p, n)}$$

By using (3.19), as a first approach, the authors utilize the upper control limit as a cut-off to identify the influential observation in linear multiple regression models. The computed  $(H_i^2)$  of any observation is greater than upper control limit, then the observation is said to be influential and it may be a potential outlier. As a second approach, the authors adopted the test of significance approach of evaluating and identifying the influential observations in a sample. The approach is to derive the critical points of the Hadi’s  $(H_i^2)$  measure by using the following relationship from (2.10) is given as

$$(3.21) \quad H_{i(p,n)}^2(\alpha) = \frac{p+1}{\frac{(n-p-2)+F_{i(1,n-p-2)}(\alpha)}{F_{i(1,n-p-2)}}} - \frac{(n-1)/n}{1+((p-1)/(n-p))F_{i(p-1,n-p)}(\alpha)} + \frac{((p-1)/(n-p))F_{i(p-1,n-p)}(\alpha) + 1/n}{(n-1)/n}$$

From (3.21) for a different combination of values of  $(p, n)$  and for the significance probability  $p(H_i^2 > H_{i(p,n)}^2(\alpha)) = \alpha$ , authors computed the critical points of Hadi’s  $(H_i^2)$  measure. By using the critical points, it is possible to test the significance of the influential observation computed from a multiple linear regression model. The following table-1 visualizes the upper control limit of the Hadi’s  $(H_i^2)$  measure computed from (3.20) and tables 2,3 exhibits the significant percentage points of the distribution of Hadi’s  $(H_i^2)$  measure for varying sample size  $(n)$  and no.of predictors  $(p)$  at 5% and 1% significance  $(\alpha)$ .



**Table 3.** Significant two-tail percentage points of Hadi's

<i>n</i>	<i>p</i>									
	1	2	3	4	5	6	7	8	9	10
3	6.5000	-	-	-	-	-	-	-	-	-
4	.3338	69.0471	-	-	-	-	-	-	-	-
5	.2502	14.4658	128.0323	-	-	-	-	-	-	-
6	.2001	6.5596	24.8542	183.7271	-	-	-	-	-	-
7	.1668	3.9604	10.6671	34.5340	237.7744	-	-	-	-	-
8	.1429	2.7611	6.2112	14.4528	43.8928	290.8798	-	-	-	-
9	.1251	2.0933	4.2220	8.2658	18.0998	53.0713	343.3963	-	-	-
10	.1112	1.6749	3.1419	5.5444	10.2376	21.6699	62.1368	395.5203	-	-
11	.1000	1.3909	2.4786	4.0844	6.8090	12.1642	25.1921	71.1262	447.3694	-
12	.0909	1.1867	2.0356	3.1965	4.9825	8.0420	14.0624	28.6819	80.0619	499.0186
13	.0834	1.0334	1.7213	2.6083	3.8785	5.8565	9.2553	15.9414	32.1486	88.9580
14	.0770	.9143	1.4879	2.1939	3.1509	4.5409	6.7154	10.4553	17.8068	35.5982
15	.0715	.8193	1.3084	1.8881	2.6406	3.6769	5.1910	7.5641	11.6457	19.6621
16	.0667	.7419	1.1664	1.6541	2.2654	3.0729	4.1926	5.8327	8.4054	12.8289
17	.0625	.6776	1.0514	1.4698	1.9793	2.6300	3.4961	4.7011	6.4685	9.2412
18	.0588	.6234	.9566	1.3213	1.7548	2.2932	2.9866	3.9132	5.2046	7.0997
19	.0556	.5771	.8771	1.1992	1.5742	2.0294	2.5999	3.3378	4.3259	5.7042
20	.0527	.5372	.8096	1.0972	1.4262	1.8178	2.2975	2.9016	3.6849	4.7352
30	.0345	.3166	.4551	.5880	.7260	.8740	1.0358	1.2151	1.4161	1.6437
40	.0256	.2241	.3156	.3998	.4838	.5703	.6608	.7567	.8591	.9692
60	.0170	.1413	.1953	.2432	.2892	.3349	.3811	.4283	.4767	.5268
80	.0127	.1032	.1413	.1746	.2060	.2368	.2673	.2980	.3290	.3606
100	.0101	.0812	.1107	.1361	.1600	.1830	.2058	.2284	.2510	.2738
120	.0084	.0670	.0910	.1116	.1307	.1491	.1672	.1851	.2029	.2207
∞	0	0	0	0	0	0	0	0	0	0

### 4. Numerical Results and Discussion

To evaluate the potential outliers based on Hadi's influence measure of the *i*th observation in a regression model in this section the authors showed the results of a numerical study. For this, the authors fitted Step-wise linear regression models with a different set of predictors in a Brand equity study. The study comprised of 18 different attributes about a car brand. The Step-wise regression results reveal 4 nested models were extracted from the regression procedure. For each model, the Hadi's ( $H_2$ ) were computed, and a comparison of proposed approaches I and II with the Hadi's traditional approach of identifying the potential outliers are visualized in the following tables.

**Table 4.** Identification of Potential Outliers, Comparative results of Hadi's approach and proposed approach-I

Model	<i>p</i>	Hadi's traditional approach								Proposed approach-I	
		*Cut-off $(H_i^2) = (E(H_i^2) + c\sqrt{V(H_i^2)})$								**Cut-off $(H_i^2) = (E(H_i^2) + \sqrt{V(H_i^2)})$	
		<i>c</i> = 1	<i>n</i> > <b>A</b>	<i>c</i> = 2	<i>n</i> > <b>A</b>	<i>c</i> = 3	<i>n</i> > <b>A</b>	<i>c</i> = 4	<i>n</i> > <b>A</b>		
1	1	0.04271	17	0.07393	14	0.10514	7	0.13636	5	0.021590	31
2	2	0.06726	15	0.11531	12	0.16337	7	0.21142	5	0.035821	29
3	3	0.09110	15	0.15518	12	0.21926	5	0.28334	4	0.048238	33
4	4	0.11565	13	0.19635	12	0.27706	5	0.35776	4	0.061802	31

*p*-no. of predictors *n*=275 \*A-Cut-off  $(H_i^2) = (E(H_i^2) + c\sqrt{V(H_i^2)})$  \*\*B- Cut-off  $(H_i^2) = (E(H_i^2) + \sqrt{V(H_i^2)})$ -refer(3.20)

**Table 5.** Identification of Potential Outliers, Comparative results of Hadi's approach and proposed approach-II

Model	p	Hadi's traditional approach								Proposed approach-II			
		*Cut-off $(H_i^2) = (E(H_i^2) + c\sqrt{V(H_i^2)})$								**Critical( $H_i^2$ ) at 5% level		**Critical( $H_i^2$ ) at 1% level	
		c = 1	n > A	c = 2	n > A	c = 3	n > A	c = 4	n > A	n > B	n > B	n > B	n > B
1	1	0.04271	17	0.07393	14	0.10514	7	0.13636	5	0.003679	78	0.003651	78
2	2	0.06726	15	0.11531	12	0.16337	7	0.21142	5	0.017942	50	0.028388	35
3	3	0.09110	15	0.15518	12	0.21926	5	0.28334	4	0.026061	58	0.038219	43
4	4	0.11565	13	0.19635	12	0.27706	5	0.35776	4	0.033031	57	0.046480	47

*p-no. of predictors*  $n=275$  \*A-Cut-off  $(H_i^2)$  \*\*B-Critical  $(H_i^2)$

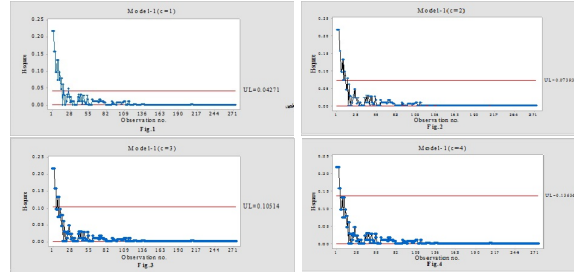
**Table 6.** Identification of Potential Outliers, Comparative results of Proposed approach I and II

Model	p	Proposed approach-I		Proposed approach-II			
		*Cut-off $(H_i^2) = (E(H_i^2) + c\sqrt{V(H_i^2)})$		*Critical( $H_i^2$ ) at 5% level		*Critical( $H_i^2$ ) at 1% level	
		n > A	n > A	n > B	n > B	n > B	n > B
1	1	0.021590	31	0.003679	78	0.003651	78
2	2	0.035821	29	0.017942	50	0.028388	35
3	3	0.048238	33	0.026061	58	0.038219	43
4	4	0.061802	31	0.033031	57	0.046480	47

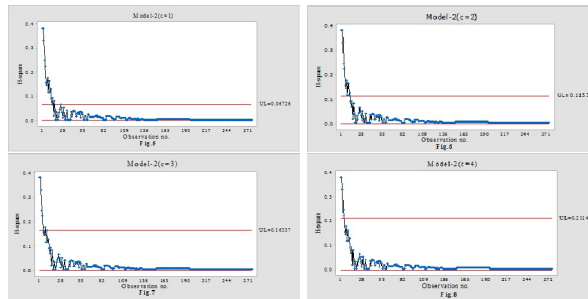
*p-no. of predictors*  $n=275$  \*A-Cut-off  $(H_i^2)$  \*\*B-Critical  $(H_i^2)$

Table 4 and 5 visualizes the comparative results of Hadi's traditional approach of evaluating the potential outliers with the proposed approached 1 and 2. Under Hadi's traditional approach, 4 nested multiple regression models are evaluated and the cut-offs' for different  $c$  values are shown in the table. As far as the fitted model-1 is a concern, the computed Hadi's influence measure for 17, 14, 7 and 5 observations were above the cut-off value and hence these observations are said to be potential outliers. Similarly, model-2 is concern 15,12,7 and 5 observations are finalized as potential outliers, in the same manner, in model-3, the calculated Hadi's influence measure for 15,12,5 and 4 observations was above the cut-off and hence these observations are said to be the potential outliers. Moreover, in model-4, 13, 12, 5 and 4 observations are treated as potential outliers because these observations exceeding the Cut-off. Under Hadi's approach at what value of  $c$ , an analyst can identify the potential outliers in the fitted models? For this question, the proposed approach-I has the answer. Under proposed approach-I, the cut-off was scientifically determined and in model-1, the calculated value of Hadi's influential measure for 31 observations are above the cut-off and in model-2 29 observations, in model-3, 33 observations and in model-4, 31 observations are exceeding the scientifically determined cut-off. Hence these observations are treated as potential outliers. Under the proposed approach-II, the authors utilize the test of significance approach to identify the potential outliers. As far as the model-1 is a concern, the computed values of Hadi's influential measure for 78 observations are greater than the critical Hadi's  $H_2$  value at 5% significance level. Similarly, model-2, model-3, and model-4 are also evaluated and the authors identified 50, 58, 57 observations are potential outliers at 5% significance level. Likewise, 78, 35, 43 and 47 observations are treated as potential outliers at 1 % significance level for model-1, model-2, model-3 and model-4 respectively. Finally, among the three approaches to evaluate the outliers, the proposed approach-II is systematic and scientific when compared to Hadi's traditional approach and proposed approach-I, because the proposed approach II identified more number of outliers at different significance level and the cut-off critical Hadi's  $(H_i^2)$  value is also scientifically determined from the distribution of Hadi's influence measure. Hence the authors observed, the proposed approach-2 outperforms the Hadi's traditional approach

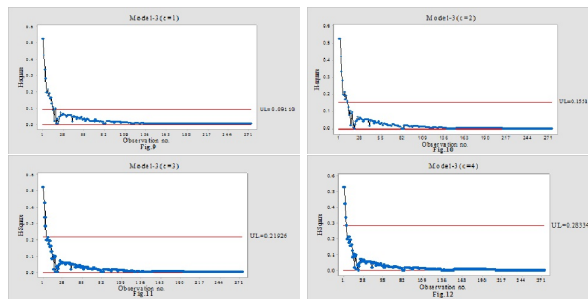
and it will be the better when you compared it with the proposed approach-I. Finally, the comparative results emphasize the superiority of proposed approaches over the traditional approach and it is visualized through the graphical display from the following control charts.



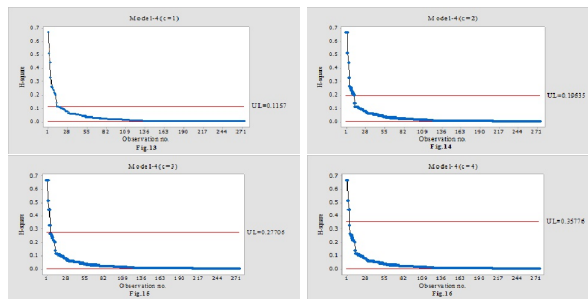
**Figure 1.** Control chart for fitted Model-1 shows the Identification of potential outliers based on Hadis approach



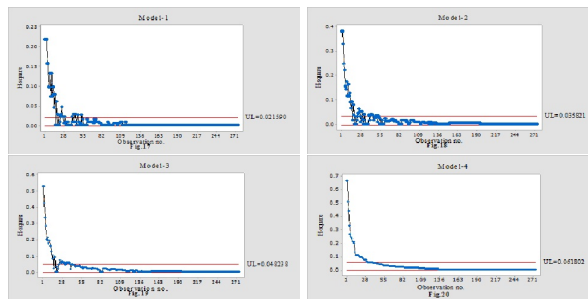
**Figure 2.** Control chart for fitted Model-2 shows the Identification of potential outliers based on Hadis approach



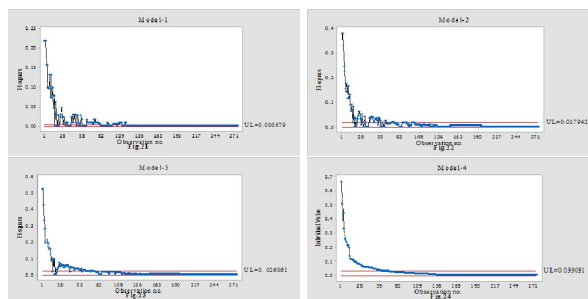
**Figure 3.** Control chart for fitted Model-3 shows the Identification of potential outliers based on Hadis approach



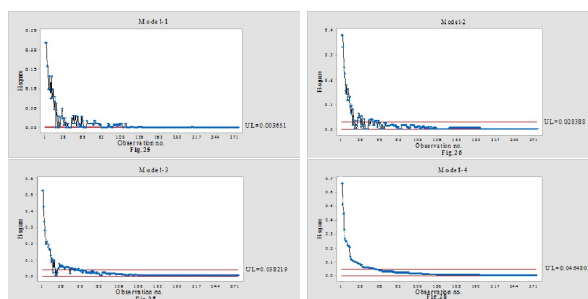
**Figure 4.** Control chart for fitted Model-4 shows the Identification of potential outliers based on Hadis approach



**Figure 5.** Control chart for each fitted model shows the Identification of potential outliers based on Proposed approach-I



**Figure 6.** Control chart for each fitted model shows the Identification of potential outliers at 5% level based on proposed approach-II



**Figure 7.** Control chart for each fitted model shows the Identification of potential outliers at 1% level based on proposed approach-II

## 5. Conclusion

From the previous sections, the authors proposed a scientific approach that is based on the test of significance for Hadi's influence measure to evaluate the potential outliers in a multiple linear regression model. At first, the exact distribution of the Hadi's ( $H_i^2$ ) was derived and the authors visualized the density function of  $H_i^2$  in terms of complicated series expression form in terms of Gauss hypergeometric function and with two shape parameters namely  $p$  and  $n$ . Moreover, the authors computed the critical percentage points of ( $H_i^2$ ) at 5 %, 1% level of significance and it is utilized to evaluate the potential outliers. Finally, the proposed approach II is more systematic and scientific because it is based on the test of significance and the results were superior when compared it with Hadi's traditional approach and proposed approach-I. Hence, the authors conclude, the proposed approach-II overrides the use of traditional approach, proposed approach-I and also it outperforms the traditional Hadi's approach in identifying more potential outliers in multiple regression models. Though Hadi's measure is used in the applied statistics for many years but authors found the absence of this technique in statistical software, limits the application of this efficient technique in the research. So the authors recommend the software developers and computational data analyst to include this valuable and pragmatic method in academic and commercial software in near future. Similarly, the authors believe that the scientific approach introduced in this study made Hadi's method a more significant tool in outlier detection as well as to the frequent users of linear multiple regression analysis.

## References

- [1] Alfons, A., Croux, C. & Gelper, S. *Sparse least trimmed squares regression for analyzing high-dimensional large data sets*. The Annals of Applied Statistics, 7, 226-248, 2013.
- [2] Andrews, D. F., & Pregibon, D. *Finding the outliers that matter*. Journal of the Royal Statistical Society. Series B (Methodological), 85-93, 1978.
- [3] Belsley, D. A., & Kuh, E. Welsch., RE. *Regression Diagnostics: Identifying Influential Data and Sources of Collinearity*. Wiley Series in Probability and Mathematical Statistics, 1980.
- [4] Beyaztas, U., & Alin, A. *Sufficient jackknife-after-bootstrap method for detection of influential observations in linear regression models*. Statistical Papers, 55(4), 1001-1018, 2014.
- [5] Chatterjee, S., & Hadi, A. S. *Influential observations, high leverage points, and outliers in linear regression*. Statistical Science, 379-393, 1986.
- [6] Cook, R. D. *Detection of influential observation in linear regression*. Technometrics, 19(1), 15-18, 1977.

- [7] Cook, R. D., & Weisberg, S. *Criticism and influence analysis in regression*. Sociological methodology, 13(3), 313-361, 1982.
- [8] Hadi, A. S. *Identifying multiple outliers in multivariate data*. Journal of the Royal Statistical Society. Series B (Methodological), 761-771, 1992.
- [9] Hoaglin, D. C., & Welsch, R. E. *The hat matrix in regression and ANOVA*. The American Statistician, 32(1), 17-22, 1978.
- [10] Johnson, B. W., & McCulloch, R. E. *AddedVariable Plots in Linear Regression*. Technometrics, 29(4), 427-433, 1987.
- [11] Mosteller, F., & Tukey, J. W. *Data analysis and regression: a second course in statistics*. Addison-Wesley Series in Behavioral Science: Quantitative Methods, 1977.
- [12] Nurunnabi, A. A. M., Hadi, A. S., & Imon, A. H. M. R. *Procedures for the identification of multiple influential observations in linear regression*. Journal of Applied Statistics, 41(6), 1315-1331, 2014.
- [13] Park, H., Sakaori, F. & Konishi, S. *Robust sparse regression and tuning parameter selection via the efficient bootstrap information criteria*. Journal of Statistical Computation and Simulation, 84, 1596-1607, 2014.