

## Research Article

# An Improvement of the Hotelling $T^2$ Statistic in Monitoring Multivariate Quality Characteristics

Ashkan Shabbak<sup>1,2</sup> and Habshah Midi<sup>2,3</sup>

<sup>1</sup> Statistical Research and Training Center (SRTC), 1433873487 Tehran, Iran

<sup>2</sup> Laboratory of Computational Statistics and Operation Research, Institute for Mathematical Research, University Putra Malaysia, 43400 Serdang, Malaysia

<sup>3</sup> Mathematics Department, Faculty of Science, University Putra Malaysia, 43400 Serdang, Malaysia

Correspondence should be addressed to Ashkan Shabbak, ashkan@inspem.upm.edu.my

Received 4 November 2011; Revised 3 February 2012; Accepted 3 February 2012

Academic Editor: Hung Nguyen-Xuan

Copyright © 2012 A. Shabbak and H. Midi. This is an open access article distributed under the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

The Hotelling  $T^2$  statistic is the most popular statistic used in multivariate control charts to monitor multiple qualities. However, this statistic is easily affected by the existence of more than one outlier in the data set. To rectify this problem, robust control charts, which are based on the minimum volume ellipsoid and the minimum covariance determinant, have been proposed. Most researchers assess the performance of multivariate control charts based on the number of signals without paying much attention to whether those signals are really outliers. With due respect, we propose to evaluate control charts not only based on the number of detected outliers but also with respect to their correct positions. In this paper, an Upper Control Limit based on the median and the median absolute deviation is also proposed. The results of this study signify that the proposed Upper Control Limit improves the detection of correct outliers but that it suffers from a swamping effect when the positions of outliers are not taken into consideration. Finally, a robust control chart based on the diagnostic robust generalised potential procedure is introduced to remedy this drawback.

## 1. Introduction

In statistical quality control, a process changes into an out-of-control situation when outliers appear in two different ways, namely, outliers that are randomly distributed within a data set and outliers that sequentially occur after a specific observation during a specific period of time in the data set. The former and the latter situations are referred to as scatter outliers and a sustained step shift, respectively, [1–4].

The detection of correct outliers in phase I of the monitoring scheme is crucial. If outliers are not correctly detected, the result leads to model misspecification and to incorrect results during phase II [5]. The Hotelling  $T^2$  statistic, which was first introduced by Hotelling

in [6], is the most popular statistic used in multivariate control charts to monitor multiple quality characteristics [7–11].

Vargas [12] demonstrated that a  $T^2$  statistic based on the usual classical estimators fails to detect multiple scatter outliers for individual observations ( $n = 1$ ), although this statistic is effective in the presence of a small number of outliers. It is now evident that the Hotelling  $T^2$  statistic, which is based on the usual classical sample mean vector and variance-covariance matrix, is easily affected by the existence of more than one outlier in the HDS (Historical Data Set). In addition, the  $T^2$  statistic suffers from a masking or swamping effect [13, 14]. Sullivan and Woodall [10] showed that the  $T^2$  statistic based on the usual sample variance-covariance matrix for individual observations is not only less effective in detecting scatter outliers in the HDS but also poor in sustained step shifts in the mean vector.

Robust methods for multivariate data, based on the MVE and MCD, have been widely used in regression contexts for diagnosing influential observations and high leverage points and outliers [14, 15] but have only recently been applied to multivariate quality control process applications. Vargas [12] employed the minimum volume ellipsoid (MVE) and the minimum covariance determinant (MCD) as two robust estimates of location and dispersion [16], instead of the usual classical sample mean vector and covariance matrix in the Hotelling  $T^2$  statistic. The application of robust control charts for individual observations based on the MVE and MCD has also been discussed extensively by Jensen et al. [17]. It is worth mentioning that there is no guarantee that the mathematical distribution of the  $T^2$  statistic is preserved by replacing the location and scale estimators with robust versions. To remedy this problem, Vargas [12] and Jensen et al. [17] used an empirical distribution of the robust  $T^2$  statistic for calculating the empirical upper control limits (UCLs) of their proposed robust control charts.

We have seen the application of MVE and MCD methods in the development of robust control charts. Other methods such as outlier identification in high dimensions [18] and some proposed multivariate outlier detection techniques [19] may also be considered in the control process applications. Our main aim in this paper is to propose a robust multivariate control chart based on the diagnostic robust generalised potential (DRGP), which was initiated by Habshah et al. [20]. This is the first attempt to introduce another robust control chart as an alternative to the two existing robust MVE-based and MCD-based control charts. Hence, our focus in this paper is only limited to the DRGP-based control chart and compare its performance with the two preceding charts. We do not wish to compare these charts with other charts based on other methods mentioned in Wilcox [19] and Filzmoser et al. [18].

Vargas [12] and Jensen et al. [17] evaluated the performance of each control chart based on the number of detected outliers, regardless of whether they came from a correct outlier positions. In other words, their work is only devoted to the detection of outliers regardless of whether the outliers are true (correct outlier position) or false outliers. Their work has motivated us to consider the identification of correct outliers when evaluating different control charts. In this regard, we introduce another empirical method for calculating the UCLs of the robust  $T^2$  statistic based on the median and the MAD of the estimators.

## 2. Diagnostics for the Identification of High Leverage Points

High leverage points are a type of influential observation that is substantially different for one or more predictor variables [21, 22]. It is now evident that high leverage points are responsible for leading to model misspecifications and misleading results [23–25].

There are some methods in the literature for identifying high leverage points in linear regression models [14, 22, 25–28]. The  $i$ th diagonal elements of hat matrix  $H$  are referred to as leverage points and are denoted by

$$h_{ii} = \mathbf{v}_i^T (V^T V) \mathbf{v}_i, \quad i = 1, 2, \dots, m, \quad (2.1)$$

where  $V$  is an  $m \times k$  matrix of predictor variables of regression model. Hoaglin and Welsch [29] considered observations to be high leverage points when  $h_{ii}$  are greater than  $2k/m$ . Hadi [30] introduced another measure to diagnose high leverage points, which is known as “potential measures.” According to Hadi [30], the  $i$ th potential is defined as follows:

$$p_{ii} = \mathbf{v}_i^T (V_{(i)}^T V_{(i)})^{-1} \mathbf{v}_i, \quad (2.2)$$

where  $V_{(i)}$  is the data matrix with its  $i$ th row deleted. By using simple matrix algebra, it is easy to obtain a relationship between the potentials and the diagonal elements of  $H$ , as follows:

$$p_{ii} = \frac{h_{ii}}{1 - h_{ii}}. \quad (2.3)$$

Hadi [30] suggested a confidence bound cutoff point for  $p_{ii}$  as follows:

$$\text{Median}(p_{ii}) + c\text{MAD}(p_{ii}), \quad (2.4)$$

where  $\text{MAD}(p_{ii}) = \text{Median}\{|p_{ii} - \text{Median}(p_{ii})|\} / 0.6745$  and  $c$  is a constant that is chosen between 2 or 3, as appropriate. Robust version of the Mahalanobis distance is also being used to identify high leverage points [14, 18, 20]. Habshah et al. [20] pointed out that although these robust diagnostic techniques can rectify the masking problem, they are affected by the swamping effect, which is not desirable either.

To remedy this problem, Habshah et al. [20] proposed a unified approach, which is called the diagnostic robust generalised potential (DRGP) method which accommodates both the diagnostic and robust approaches together. The robust approach is utilised to detect the suspected high leverage points and then diagnostic approach is utilised to confirm our suspicion. The DRGP partitions the data into two sets. The first set consists of suspicious cases, which are deleted from the original observations (denoted by  $D$ ), and the second set contains the remaining data (denoted by  $R$ ). It is clear that if  $d$  is the number of cases that includes  $D$ , then the  $R$  set contains  $m - d$  observations and  $d < m - k$ . Without loss of generality, we assume that  $d$  cases are placed in the last  $d$  rows of  $Y$  and  $V$  so that the hat matrix is partitioned and that the  $i$ th deletion leverage is defined as follows:

$$h_{ii}^{(-D)} = \mathbf{v}_i^T (V_R^T V_R)^{-1} \mathbf{v}_i, \quad (2.5)$$

where  $V_R$  indicates the remaining observation matrix [31]. By considering (2.2), (2.3) and (2.5), the generalised potential is defined as follows:

$$p_{ii}^* = \begin{cases} \frac{h_{ii}^{(-D)}}{1 - h_{ii}^{(-D)}} & \text{for } i \in R \\ h_{ii}^{(-D)} & \text{for } i \in D. \end{cases} \quad (2.6)$$

Similar to the potential values and with regard to (2.4), the cutoff point for  $p_{ii}^*$  is

$$\text{Median}(p_{ii}^*) + c\text{MAD}(p_{ii}^*). \quad (2.7)$$

It is worth mentioning that the DRGP employs the robust Mahalanobis distance in the first step to detect high leverage points as preliminary suspicious observations; these points are placed in the  $D$  set. Next, in the second step, only the cases that are greater than (2.7) are reported as the final detections. In the next section, the DRGP is employed to effectively detect outliers in multivariate quality control charts for individual observations.

### 3. Multivariate Robust $T^2$ Control Charts

Suppose that there is an HDS in the phase I monitoring scheme that consists of  $m$  time-ordered observation vectors of dimension  $p$ , which are observed independently, where  $p$  is the number of quality characteristics that are measured ( $p < m$ ). It is assumed that each vector comes from a  $p$ -variate normal distribution. Thus, if  $X_i \in \mathbb{R}^p$  is a vector in the HDS for the  $i$ th time period,  $X_i \sim N_p(\mu, \Sigma)$ , where  $\mu$  and  $\Sigma$  are the population mean vector and the variance-covariance matrix, respectively. As mentioned earlier, the Hotelling  $T^2$  statistic is used to detect outliers in multivariate control charts. The general form of this statistic is

$$T_i^2 = (X_i - \mu)^T \sum_{i=1}^{-1} (X_i - \mu), \quad i = 1, 2, \dots, m. \quad (3.1)$$

Because the parameters in (3.1) are usually unknown, the usual sample mean and variance-covariance matrix are used as the classical estimations of  $\mu$  and  $\Sigma$ . In practice, these variables are expressed by

$$\begin{aligned} \bar{X} &= \frac{1}{m} \sum_{i=1}^m X_i, \\ S &= \frac{1}{m-1} \sum_{i=1}^m (X_i - \bar{X})(X_i - \bar{X})^T. \end{aligned} \quad (3.2)$$

In phase I, the parameters are retrospectively estimated based on the current HDS; as a result, the vector  $X_i$  is not independent of the estimators  $\bar{X}$  and  $S$ . In this situation,

the statistical distribution of (3.1) is given as

$$T^2 \sim \left[ \frac{(m-1)^2}{m} \right] B\left(\frac{p}{2}, \frac{m-p-1}{2}\right), \quad (3.3)$$

where  $B(p/2, (m-p-1)/2)$  represents a beta distribution with parameters  $p/2$  and  $(m-p-1)/2$  [8, 11]. From (3.3), the upper control limit (UCL) of  $T^2$  is  $((m-1)^2/m)B(\alpha, p/2, (m-p-1)/2)$ , where  $\alpha$  is the probability of a false alarm for each point plotted on the control chart and  $B(\alpha, p/2, (m-p-1)/2)$  is the  $\alpha$ th upper quantile of the beta distribution with parameters  $p/2$  and  $(m-p-1)/2$ .

The lower control limit (LCL) is often set to zero [8, 32]. It should be noted that the aforementioned UCLs are exact when applied to a single point in phase I, whereas phase I is a retrospective analysis of all observations. Therefore, the values of  $\alpha$  cannot be applied to a set of points. In this situation, if all of the statistics were distributed independently, then the overall probability of a false alarm would be

$$\alpha' = 1 - (1 - \alpha)^m, \quad (3.4)$$

where  $\alpha$  indicates the probability of a false alarm, which is assigned for each observation plotted on the control chart in a subgroup of size  $m$ . In practice, it is reasonable to determine the UCL by simulation to give a specified overall false alarm [5, 33, 34]. Hereafter, in this paper, we still refer to the overall false,  $\alpha'$  as  $\alpha$  for simplicity.

The use of (3.1) is not effective in the presence of multiple outliers, so robust alternatives are proposed. In this regard, two of the recently proposed robust alternative approaches to  $T^2$  are based on the MVE and the MCD estimators, which will be denoted by  $T_{\text{mve}}^2$  and  $T_{\text{mcd}}^2$ , respectively, and are defined as follows:

$$T_{\text{mve},i}^2 = \left( X_i - \bar{X}_{\text{mve}} \right)^T S_{\text{mve}}^{-1} \left( X_i - \bar{X}_{\text{mve}} \right)^2, \quad i = 1, 2, \dots, m, \quad (3.5)$$

$$T_{\text{mcd},i}^2 = \left( X_i - \bar{X}_{\text{mcd}} \right)^T S_{\text{mcd}}^{-1} \left( X_i - \bar{X}_{\text{mcd}} \right)^2, \quad i = 1, 2, \dots, m, \quad (3.6)$$

where  $\bar{X}_{\text{mve}}$  and  $\bar{X}_{\text{mcd}}$  are the robust estimations of the sample mean and  $S_{\text{mve}}^{-1}$  and  $S_{\text{mcd}}^{-1}$  are the corresponding estimators of the sample variance-covariance matrix. As previously mentioned in Section 1, due to the unknown distribution of  $T_{\text{mve}}^2$  and  $T_{\text{mcd}}^2$ , empirical methods are used to determine the UCLs. The empirical simulated UCLs for (3.5) or (3.6) are usually determined by finding the  $\alpha$ th upper quantile of the empirical distribution of the corresponding statistic. For this situation, Jensen et al. [17] and Vargas [12] defined  $\alpha$  as the overall false alarm. Following the idea of Habshah et al. [20], another empirical UCL is proposed as follows:

$$\text{Median}\left(T_{\text{mve/mcd},i}^2\right) + c\text{MAD}\left(T_{\text{mve/mcd},i}^2\right), \quad (3.7)$$

where MAD is the median absolute deviation of either  $T_{\text{mve},i}^2$  or the  $T_{\text{mcd},i}^2$ , as defined in (2.4). For simplicity, the first empirical UCL is referred to as Empr, and the proposed UCL will be denoted by Med-Mad. As will be shown later, (3.7) tends to declare too many observations as

outliers, by detecting outliers regardless of whether the detected outliers are true (in a correct outlier position) or false, even though (3.7) has a better performance in detecting real outliers at their correct positions in the HDS than Empr UCL. In this regard, we propose to apply the DRGP based on the MVE and MCD with the Med-Mad UCLs to reduce the number of undue signals caused by the detection of outliers irrespective of their correct positions and, at the same time, to effectively detect correct outliers.

Vargas [12] and Jensen et al. [17] employed the probability of signals to evaluate and compare control charts based on  $T_{mve}^2$  and  $T_{mcd}^2$ . Their work is only based on the number of detections, without considering whether those detections are correct outliers. It is worth mentioning that the probability of signals cannot be properly judged if there is a swamping effect in the monitoring scheme. In the following section, we will present our proposed control scheme and explain how it can detect real outliers at their correct positions.

#### 4. Simulation Study

In this section, a Monte Carlo simulation study is carried out to assess the performance of the control schemes discussed previously. The simulation is designed based on three subsamples, each of size  $m = 30, 50,$  and  $100,$  with a number of characteristics  $p = 2, 3, 5,$  and  $10.$  Let us assume that the in-control process is a  $p$ -variate normal distribution with mean vector  $\mu_0$  and covariance matrix  $\Sigma.$

The simulated empirical UCLs are obtained by generating 5000 in-control data sets for each combination of  $m$  and  $p.$  Due to the affine equivariant property of the  $T^2$  statistics, these limits are applicable to any values of  $\mu$  and  $\Sigma.$  Many researchers, such as Jensen et al. [17] and Vargas [12], have determined the Empr UCLs by calculating all of the  $T^2$  statistics for each observation in generated data sets of each subgroup of size  $m$  and recording the maximum value of the  $T^2$  statistics. Subsequently, the upper  $\alpha$ th percentile of the 5000 recorded maximum values of the  $T^2$  statistics is declared as the Empr UCL. In this manner, they defined  $\alpha$  as the overall false alarm. In this paper, we consider the probability of the overall false alarm, which is equal to  $\alpha = 0.05.$

To make the overall false alarm of the control charts based on the Med-Mad UCLs equivalent to Empr UCLs, which is equal to  $\alpha = 0.05,$  the following steps are considered. The Med-Mad UCLs are attained based on all  $m \times 5000$  values of the simulated  $T^2$  statistics. It is worth mentioning here that, to keep the overall false alarm at  $\alpha = 0.05,$  the probability of a false alarm for each observation in each subgroup of size  $m$  is calculated based on (3.4). Hence, to have the same overall false alarm, the values for  $c$  in (3.7) and for  $m = 30, 50,$  and  $100$  must be chosen as  $\Phi_{(0.002)}^{-1} \simeq 2.88, \Phi_{(0.001)}^{-1} \simeq 3,$  and  $\Phi_{(0.001)}^{-1} \simeq 3,$  respectively. The Empr and Med-Mad UCLs of  $T^2, T_{mve}^2,$  and  $T_{mcd}^2$  for each combination of  $m$  and  $p$  are presented in Table 1.

Then, a contaminated data set of size  $m$  in the  $p$  dimension is generated for different values of the noncentrality parameter (ncp). The out-of-control process is a  $p$ -variate normal distribution with the same covariance matrix but with a shifted mean vector of  $\mu_1.$  Thus, the variation here remains stable. The magnitude of the shift is measured by a scalar defined as follows:

$$(\mu_1 - \mu_0)^T \Sigma^{-1} (\mu_1 - \mu_0). \quad (4.1)$$

**Table 1:** The simulated UCLs for all of the  $T^2$  statistics.

$p$	$m$	Empirical upper control limit (UCL)					
		Empr			Med-Mad		
		$T^2$	$T_{mve}^2$	$T_{mcd}^2$	$T^2$	$T_{mve}^2$	$T_{mcd}^2$
2	30	10.561	24.351	58.812	5.501	6.080	8.796
	50	12.302	21.042	34.615	5.666	6.289	7.967
	100	14.124	20.478	24.441	5.680	6.314	7.145
3	30	12.202	29.028	108.530	7.995	8.606	13.204
	50	13.976	25.145	50.430	8.266	8.883	11.985
	100	16.503	23.011	30.371	8.278	8.967	10.425
5	30	14.950	49.371	319.625	11.969	12.780	20.501
	50	17.581	33.792	111.725	12.482	13.160	19.823
	100	20.054	28.448	41.639	12.639	13.366	16.142
10	30	20.033	80.693	1555.711	19.541	19.812	25.827
	50	24.250	63.427	393.625	20.894	21.852	33.893
	100	27.964	44.287	84.655	21.470	22.287	29.604

**Table 2:** The number of correctly detected outliers and the number of detected outliers for  $p = 2$  and  $m = 30$ .

$\epsilon$	ncp	Empr UCL			Med-Mad UCL		
		$T^2$	$T_{mve}^2$	$T_{mcd}^2$	$T^2$	$T_{mve}^2$	$T_{mcd}^2$
(5%) 2 outliers	5	0 (0)	0 (0)	0 (0)	1 (2)	1 (2)	1 (4)
	15	0 (0)	1 (1)	0 (0)	1 (2)	2 (3)	2 (5)
	25	1 (1)	1 (1)	1 (1)	2 (2)	2 (3)	2 (5)
	35	1 (1)	2 (2)	1 (1)	2 (2)	2 (3)	2 (5)
	45	1 (1)	2 (2)	1 (1)	2 (2)	2 (3)	2 (5)
	55	1 (1)	2 (2)	2 (2)	2 (2)	2 (3)	2 (5)
(10%) 3 outliers	5	0 (0)	0 (0)	0 (0)	1 (2)	1 (2)	1 (4)
	15	0 (0)	1 (1)	0 (0)	2 (2)	2 (4)	3 (5)
	25	0 (0)	2 (2)	1 (1)	2 (3)	3 (4)	3 (6)
	35	0 (0)	2 (2)	1 (1)	2 (3)	3 (4)	3 (6)
	45	0 (0)	3 (3)	2 (2)	3 (3)	3 (4)	3 (6)
	55	0 (0)	3 (3)	2 (2)	3 (3)	3 (4)	3 (6)
(15%) 5 outliers	5	0 (0)	0 (0)	0 (0)	1 (2)	1 (2)	2 (4)
	15	0 (0)	1 (1)	0 (0)	1 (2)	3 (4)	4 (6)
	25	0 (0)	2 (2)	1 (1)	2 (2)	4 (5)	5 (7)
	35	0 (0)	4 (4)	2 (2)	2 (2)	5 (5)	5 (7)
	45	0 (0)	4 (4)	3 (3)	2 (2)	5 (5)	5 (7)
	55	0 (0)	5 (5)	4 (4)	2 (2)	5 (5)	5 (7)
(20%) 6 outliers	5	0 (0)	0 (0)	0 (0)	1 (1)	1 (3)	2 (5)
	15	0 (0)	1 (1)	1 (1)	1 (2)	3 (4)	5 (7)
	25	0 (0)	3 (3)	1 (1)	1 (2)	5 (6)	6 (7)
	35	0 (0)	4 (4)	2 (2)	1 (2)	6 (6)	6 (7)
	45	0 (0)	5 (5)	3 (3)	1 (2)	6 (7)	6 (7)
	55	0 (0)	6 (6)	4 (4)	1 (2)	6 (7)	6 (8)

This measure is called the noncentrality parameter and is hereafter referred to as ncp. Four outlier percentage levels are considered, which are denoted by  $\epsilon = 5\%$ ,  $10\%$ ,  $15\%$ , and



**Table 3:** The number of correctly detected outliers and the number of detected outliers for  $p = 2$  and  $m = 50$ .

$\epsilon$	ncp	Empr UCL			Med-Mad UCL		
		$T^2$	$T_{mve}^2$	$T_{mcd}^2$	$T^2$	$T_{mve}^2$	$T_{mcd}^2$
(5%) 3 outliers	5	0 (0)	0 (0)	0 (0)	1 (3)	1 (4)	1 (5)
	15	0 (0)	1 (1)	1 (1)	2 (4)	3 (5)	3 (6)
	25	1 (1)	2 (2)	2 (2)	3 (4)	3 (5)	3 (6)
	35	1 (1)	3 (3)	2 (2)	3 (4)	3 (6)	3 (7)
	45	1 (2)	3 (3)	3 (3)	3 (4)	3 (6)	3 (7)
	55	2 (2)	3 (3)	3 (3)	3 (4)	3 (6)	3 (7)
(10%) 5 outliers	5	0 (0)	0 (0)	0 (0)	1 (3)	2 (4)	2 (5)
	15	0 (0)	1 (1)	1 (1)	3 (4)	4 (6)	4 (7)
	25	0 (0)	3 (3)	3 (3)	4 (4)	5 (7)	5 (8)
	35	0 (0)	5 (5)	4 (4)	4 (5)	5 (7)	5 (8)
	45	0 (0)	5 (5)	5 (5)	4 (5)	5 (7)	5 (8)
	55	0 (0)	5 (5)	5 (5)	4 (5)	5 (7)	5 (8)
(15%) 8 outliers	5	0 (0)	0 (0)	0 (0)	2 (3)	2 (4)	3 (5)
	15	0 (0)	2 (2)	1 (1)	2 (3)	6 (7)	7 (9)
	25	0 (0)	5 (5)	4 (4)	3 (4)	8 (9)	8 (10)
	35	0 (0)	7 (7)	6 (6)	3 (4)	8 (10)	8 (10)
	45	0 (0)	8 (8)	7 (7)	3 (4)	8 (10)	8 (10)
	55	0 (0)	8 (8)	8 (8)	3 (4)	8 (10)	8 (10)
(20%) 10 outliers	5	0 (0)	0 (0)	0 (0)	1 (3)	2 (3)	3 (5)
	15	0 (0)	2 (2)	1 (1)	2 (3)	7 (8)	8 (10)
	25	0 (0)	5 (5)	4 (4)	2 (3)	9 (11)	10 (11)
	35	0 (0)	8 (8)	7 (7)	2 (3)	10 (11)	10 (12)
	45	0 (0)	10 (10)	9 (9)	2 (3)	10 (12)	10 (12)
	55	0 (0)	10 (10)	10 (10)	2 (3)	10 (12)	10 (12)

20%. It is clear from (4.1) that the severity of the shift only depends on the values of  $\mu_1$ . Hence, without loss of generality, it can be assumed that  $\mu_0$  is a zero vector and  $\Sigma$  is a  $p \times p$  identity matrix.

The control charts are assessed based on the proposed criterion, which are based on the number of true detected outliers with regard to the correct position of the generated out-of-control observations in the data set. The number of outliers detected without regard to their positions is also presented for comparison.

Repeating this process 5000 times, the number of detections and the number of correctly detected outliers with correct positions are recorded for each replication. The numbers of detected outliers are determined by comparing each of the  $T^2$  values with the respective UCLs given in Table 1. Each detected outlier is checked by its position in the data set to determine whether it can truly be generated from the intentional simulated contaminated points in the data set. The true or correctly detected outliers refer to the outliers detected at the correct position. The number of detected outliers simply indicates the outliers that have been detected irrespective of their correct position. The average number of detections over 5000 iterations is presented for  $p = 2$ ,  $m = 30$ , 50, and 100 in Tables 2, 3, and 4. It is important to note that the presented values were rounded up to two digits. The values in parentheses represent the number of outliers detected regardless of their correct positions.



**Table 4:** The number of correctly detected outliers and the number of detected outliers for  $p = 2$  and  $m = 100$ .

$\epsilon$	ncp	Empr UCL			Med-Mad UCL		
		$T^2$	$T^2_{mve}$	$T^2_{mcd}$	$T^2$	$T^2_{mve}$	$T^2_{mcd}$
(5%) 5 outliers	5	0 (0)	0 (0)	0 (0)	2 (6)	2 (7)	2 (8)
	15	1 (1)	2 (2)	2 (2)	4 (7)	5 (10)	5 (10)
	25	1 (1)	4 (4)	4 (4)	5 (7)	5 (10)	5 (11)
	35	2 (2)	5 (5)	5 (5)	5 (7)	5 (10)	5 (11)
	45	2 (2)	5 (5)	5 (5)	5 (7)	5 (10)	5 (11)
	55	3 (3)	5 (5)	5 (5)	5 (7)	5 (10)	5 (11)
(10%) 10 outliers	5	0 (0)	0 (0)	0 (0)	3 (6)	4 (8)	4 (9)
	15	0 (0)	3 (3)	3 (3)	6 (8)	9 (13)	9 (13)
	25	0 (0)	7 (7)	7 (7)	7 (9)	10 (14)	10 (15)
	35	0 (0)	9 (9)	9 (9)	8 (10)	10 (15)	10 (15)
	45	0 (0)	10 (10)	10 (10)	8 (10)	10 (15)	10 (15)
	55	0 (0)	10 (10)	10 (10)	9 (11)	10 (15)	10 (15)
(15%) 16 outliers	5	0 (0)	0 (0)	0 (0)	3 (6)	4 (8)	5 (9)
	15	0 (0)	3 (3)	3 (3)	5 (7)	13 (16)	14 (17)
	25	0 (0)	11 (11)	10 (10)	6 (8)	16 (19)	16 (19)
	35	0 (0)	15 (15)	14 (14)	6 (8)	16 (20)	16 (20)
	45	0 (0)	16 (16)	16 (16)	7 (8)	16 (20)	16 (20)
	55	0 (0)	16 (16)	16 (16)	7 (9)	16 (20)	16 (20)
(20%) 20 outliers	5	0 (0)	0 (0)	0 (0)	3 (6)	5 (7)	5 (8)
	15	0 (0)	3 (3)	3 (3)	4 (6)	15 (18)	16 (19)
	25	0 (0)	12 (12)	12 (12)	4 (6)	20 (23)	20 (23)
	35	0 (0)	18 (18)	17 (17)	5 (6)	20 (23)	20 (23)
	45	0 (0)	20 (20)	19 (19)	5 (6)	20 (24)	20 (23)
	55	0 (0)	20 (20)	20 (20)	4 (6)	20 (24)	20 (23)

**Table 5:** The number of correctly detected outliers and the number of outliers detected using the DRGP when  $p = 2$ .

ncp	$m = 30$				$m = 100$			
	$\epsilon = 5\%$		$\epsilon = 20\%$		$\epsilon = 5\%$		$\epsilon = 20\%$	
	$T^2_{mve}$	$T^2_{mcd}$	$T^2_{mve}$	$T^2_{mcd}$	$T^2_{mve}$	$T^2_{mcd}$	$T^2_{mve}$	$T^2_{mcd}$
5	1 (3)	1 (3)	1 (2)	2 (2)	2 (6)	2 (6)	5 (7)	5 (6)
15	2 (3)	2 (3)	3 (4)	5 (5)	5 (9)	5 (9)	15 (15)	16 (14)
25	2 (3)	2 (3)	5 (6)	5 (5)	5 (8)	5 (8)	20 (20)	20 (20)
35	2 (3)	2 (3)	6 (6)	6 (6)	5 (8)	5 (8)	20 (21)	20 (21)
45	2 (3)	2 (3)	6 (6)	6 (6)	5 (8)	5 (8)	20 (21)	20 (21)
55	2 (3)	2 (3)	6 (6)	6 (6)	5 (8)	5 (8)	20 (21)	20 (21)

Let us first focus on the results of  $T^2$ ,  $T^2_{mve}$ , and  $T^2_{mcd}$  obtained by using the Empr and Med-Mad UCLs. Following these results, we will see later why the DRGP approach is proposed.

As can be seen from these tables, the UCLs based on the Med-Mad approach for all three control charts have better performance in detecting the real outliers at their correct positions, compared to Empr UCL. The classical  $T^2$  chart, based on Empr UCL, performs very

**Table 6:** The correct detection rate for  $m = 50$  and  $p = 2$ .

$\epsilon$	Detection method		Correct detection rate					
			ncp = 5	ncp = 15	ncp = 25	ncp = 35	ncp = 45	ncp = 55
(5%) 3 outliers	DRGP	$T_{mve}^2$	0.38	0.66	0.69	0.68	0.67	0.68
		$T_{mcd}^2$	0.37	0.65	0.69	0.69	0.68	0.69
	Without DRGP	$T_{mve}^2$	0.36	0.59	0.60	0.60	0.59	0.59
		$T_{mcd}^2$	0.31	0.52	0.53	0.53	0.53	0.54
(10%) 5 outliers	DRGP	$T_{mve}^2$	0.49	0.78	0.81	0.81	0.81	0.80
		$T_{mcd}^2$	0.49	0.78	0.82	0.82	0.82	0.80
	Without DRGP	$T_{mve}^2$	0.47	0.71	0.72	0.72	0.72	0.71
		$T_{mcd}^2$	0.42	0.66	0.67	0.68	0.68	0.68
(15%) 8 outliers	DRGP	$T_{mve}^2$	0.55	0.85	0.91	0.91	0.92	0.91
		$T_{mcd}^2$	0.56	0.87	0.91	0.91	0.92	0.91
	Without DRGP	$T_{mve}^2$	0.54	0.8	0.83	0.83	0.83	0.83
		$T_{mcd}^2$	0.52	0.78	0.82	0.81	0.81	0.81
(20%) 10 outliers	DRGP	$T_{mve}^2$	0.54	0.81	0.93	0.95	0.95	0.95
		$T_{mcd}^2$	0.55	0.85	0.94	0.95	0.95	0.95
	Without DRGP	$T_{mve}^2$	0.57	0.81	0.87	0.87	0.87	0.86
		$T_{mcd}^2$	0.56	0.84	0.87	0.86	0.87	0.86

poorly. It can be seen that the Med-Mad UCLs are more reliable in detecting the correctly detected outliers for the robust control charts compared to Empr UCLs, particularly when the percentage of outliers increases. However, the results signify that both  $T_{mve}^2$  and  $T_{mcd}^2$  based on the Med-Mad UCLs suffer from a swamping effect due to the detection of more outliers without regard to their correct positions.

For example, at 10%, regarding outliers in the HDS with  $m = 100$ ,  $p = 2$ , and  $ncp = 25$  and both  $T_{mve}^2$  and  $T_{mcd}^2$  based on the Med-Mad UCLs, the methods detect exactly 10 outliers at the correct positions, but they also detect 15 outliers irrespective of their correct positions. It is interesting to note that the performance of the robust control charts based on the Empr UCLs is reasonably close to that of the robust control charts based on the Med-Mad UCLs for very large values of  $ncp$ , such as  $ncp \geq 45$ . The  $T_{mve}^2$  and  $T_{mcd}^2$  parameters, based on the Med-Mad UCLs, are equally good at detecting correct outliers. Nonetheless, with increasing subgroup size,  $T_{mcd}^2$  based on the Med-Mad UCL is slightly better than the  $T_{mve}^2$  based on Med-Mad, particularly for a large percentage of outliers.

We can see that although the classical  $T^2$  Med-Mad-based method is better than the  $T^2$  Empr-based method, it detects a smaller number of exact outliers as the percentage of outliers increases. These results are consistent with other values of  $p$  but are not reported here due to space constraints. The findings of Tables 2 to 4 seem to suggest that the Med-Mad UCLs are more reliable than Empr UCLs in detecting correct outliers. However, the Med-Mad UCLs detected more outliers irrespective of their correct positions due to swamping effects. In other words, we have shown that the robust control charts based on the Med-Mad UCLs effectively detect the number of correct outliers, but they overdetected outliers irrespective of their correct positions. As such, we need to employ control charts that can reduce such undue detections. In this regard, we suggest applying the DRGP procedure discussed in Section 2. The same simulation procedure was then carried out, and the DRGP was applied to the data sets.

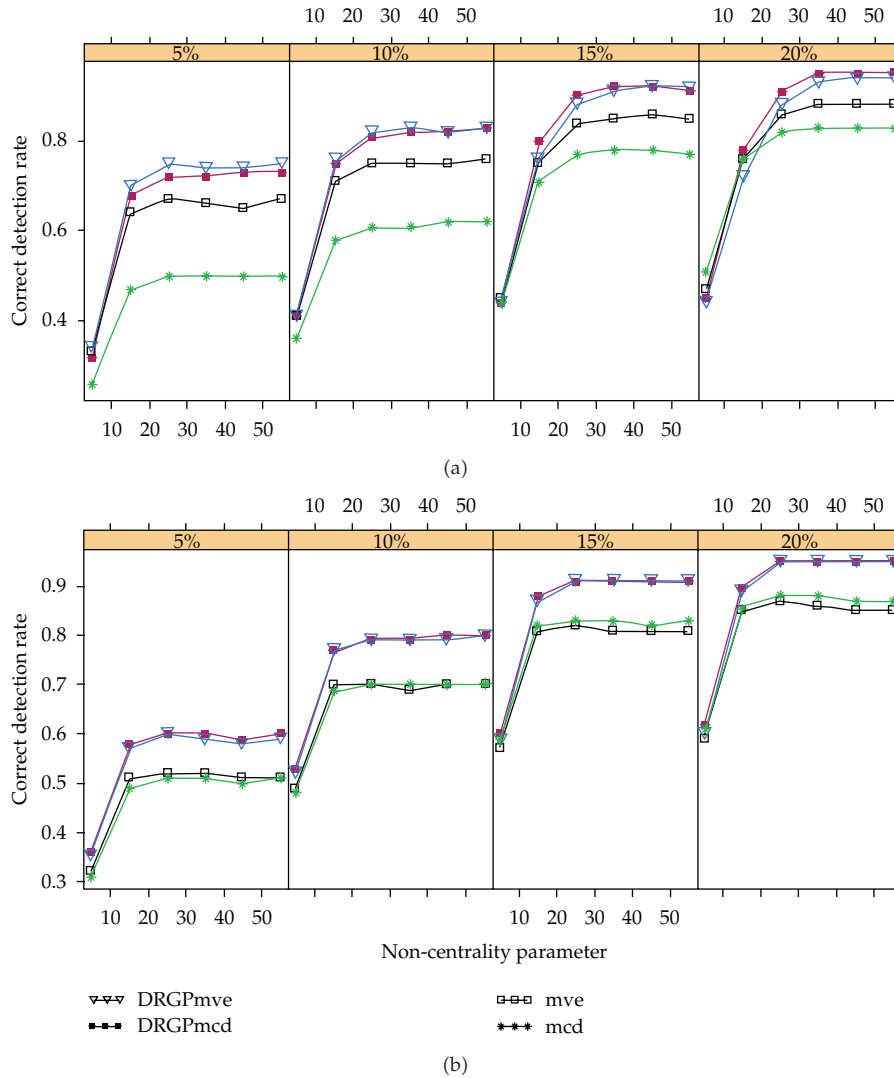


Figure 1: The correct detection rate for different outlier percentages for  $p = 2$ :  $m = 30$  (a) and  $m = 100$  (b).

The results obtained by using the DRGP approach are exhibited in Table 5. Due to space limitations, the results are presented only for  $m = 30, 100$  and  $\epsilon = 5\%, 20\%$ . As can be seen from Table 5, there is a steady decrease in the number of undue observations when the DRGP approach is applied. For example, as shown in Table 4, the total number of detections by the Med-Mad UCL with  $ncp = 55$  is 24 for the MVE and MCD, while it decreases to 21 in Table 5. To simplify the presentation of results, the proportion of correctly detected outliers to detected outliers is calculated and referred to as the correct detection rate. The values of the correct detection rates for  $m = 50$  and  $p = 2$  are shown in Table 6. The results indicate that the DRGP approach provides higher correct detection rates compared to the other methods.

However, these results were not very encouraging for small shifts ( $ncp = 5$ ). It can be seen from Table 6 that there is a gradual and steady rise in the correct detection rate with increasing outlier percentage. The results are similar for other values of  $m$  and  $p$ , which are

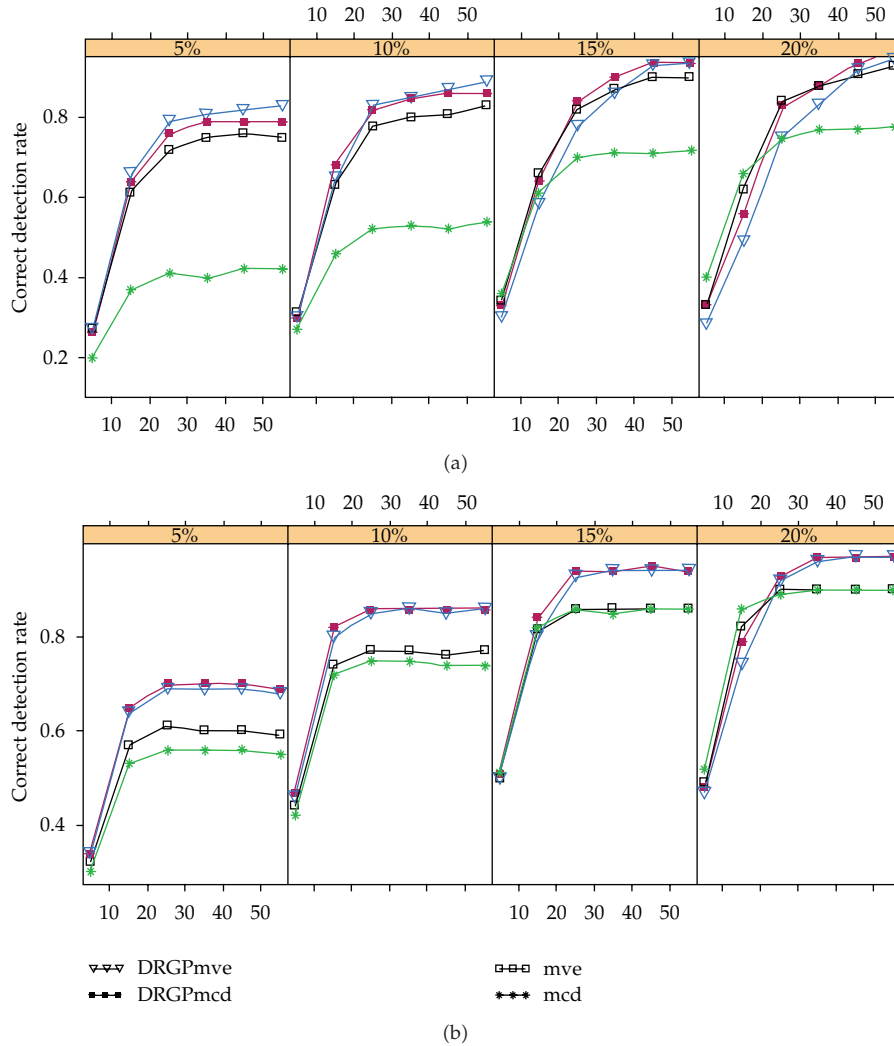


Figure 2: The correct detection rate for different outlier percentages for  $p = 3$  :  $m = 30$  (a) and  $m = 100$  (b).

not tabulated due to space limitations. For more clarification, the correct detection rates for various values of  $m$  and  $p$  are plotted in Figures 1 and 2. These figures confirm that the DRGP approach gives a higher correct detection rate.

### 5. Numerical Example

In this section, a numerical example is introduced to assess the performance of our method. This is a bivariate data set which is taken from Shewhart [35]. It presents the measurements of the depth of sapwood and the depth of penetration of creosote in telephone poles. The subgroup size of the original dataset is 10 and we only focus on the first column of the data set based on 20 subgroups. The  $T^2$ ,  $T_{mve}^2$ ,  $T_{mcd}^2$ , and the DRGP statistics were then applied to the data set. The Empr UCLs and the Med-Mad UCLs for  $T^2$ ,  $T_{mve}^2$ , and  $T_{mcd}^2$  statistics

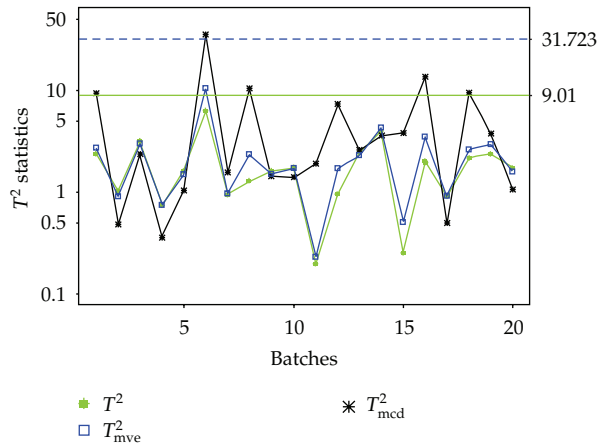


Figure 3: Control charts for the creosoting Telephone Poles Data, based on the Empr UCLs.

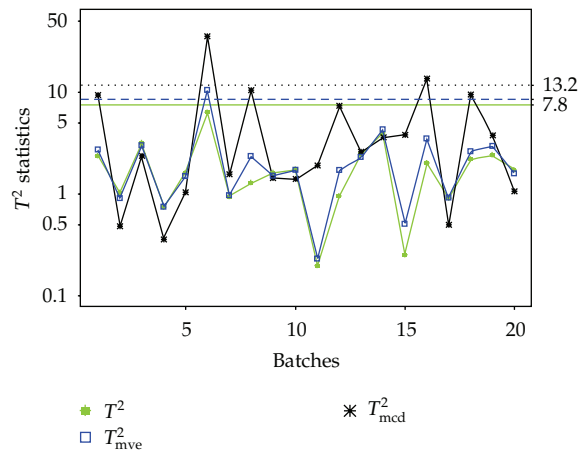


Figure 4: Control charts for the creosoting Telephone Poles Data, based on the Med-Mad UCLs.

are 9.010, 31.723.738, and 113.140 and 7.010, 5.910, and 9.461, respectively. These UCLs are calculated using the simulation, as discussed in Section 4, for  $\alpha = 0.05$ .

Figure 3 shows the different  $T^2$  control charts for Empr UCLs. As can be seen from this graph, none of the control charts based on the Empr UCL is able to detect any outlier in the data set. On the other hand, robust  $T^2$  statistics based on the Med-Mad UCLs can identify severed outliers (Figure 4).

The points 6 and 16 are detected by  $T^2_{mve}$  as outliers, and  $T^2_{mcd}$  identified the cases of 1, 6, 8, 16, and 21. The robust control charts based on the DRGP approach detected one observation as outlier which is observation number 6. The results are not included here due to space limitations.

## 6. Conclusions

Most research studies evaluate the performance of robust multivariate  $T^2_{mve}$  and  $T^2_{mcd}$  control charts based on the number of outliers detected without regard to the correct position of

outliers in the data set. The detection of real outliers (outliers at the correct position in the data set) is crucial to avoid making wrong inferences. This study has shown that although the  $T_{mve}^2$  and the  $T_{mcd}^2$  based on the Med-Mad UCLs are effective in the detection of correct outliers, they have the tendency to declare undue observations as outliers (irrespective of whether they are true or false outliers) due to a swamping effect. In this respect, in the evaluation of robust control charts, we suggest not only a consideration of the number of outlier detections but also a consideration of the correct position of the detected outliers.

Our findings also suggest that the proposed Med-Mad UCLs have better performance than the commonly used Empr UCLs in detecting outliers with regard to the correct position of outliers, especially for higher proportions of outliers. The practical finding of this paper is that the robust control chart based on the DRGP with the proposed Med-Mad UCLs gives credible performance.

The limitations of our study are that inferences or conclusions are only confined to the detection of multiple outliers for individual observations, scatter outlier situations, and moderate dimensional data sets ( $p \leq 10$ ).

## References

- [1] A. Shabbak and H. Midi, "Robust multivariate control charts to detect small shifts in mean," *Mathematical Problems in Engineering*, vol. 2011, Article ID 923463, 2011.
- [2] M. Habshah, A. Shabbak, B. A. Talib, and M. N. Hassan, "Multivariate control chart based on robust estimator," *Journal of Quality Measurement and Analysis*, vol. 5, no. 1, pp. 17–33, 2009.
- [3] J. H. Sullivan, "Detection of multiple change points from clustering individual observations," *Journal of Quality Technology*, vol. 34, no. 4, pp. 371–383, 2002.
- [4] J. A. N. Vargas and J. C. Lagos, "Comparison of multivariate control charts for process dispersion," *Quality Engineering*, vol. 19, no. 3, pp. 191–196, 2007.
- [5] J. D. Williams, W. H. Woodall, J. B. Birch, and J. H. Sullivan, "Distribution of hotelling's  $T^2$  statistic based on the successive differences estimator," *Journal of Quality Technology*, vol. 38, no. 3, pp. 217–229, 2006.
- [6] H. Hotelling, "The generalization of Student's ratio," *The Annals of Mathematical Statistics*, vol. 2, no. 3, pp. 360–378, 1931.
- [7] F. B. Alt and N. D. Smith, "17 multivariate process control," *Handbook of Statistics*, vol. 7, pp. 333–351, 1988.
- [8] R. L. Mason and J. C. Young, *Multivariate Statistical Process Control with Industrial Applications*, Society for Industrial Mathematics, Philadelphia, Pa, USA, 2002.
- [9] A. Shabbak, H. Midi, and M. N. Hassan, "The performance of robust multivariate statistical control charts based on different cutoff-points with sustained shift in mean," *Journal of Applied Sciences*, vol. 11, no. 1, pp. 56–65, 2011.
- [10] J. H. Sullivan and W. H. Woodall, "A comparison of multivariate control charts for individual observations," *Journal of Quality Technology*, vol. 28, no. 4, pp. 398–408, 1996.
- [11] N. D. Tracy, J. C. Young, and R. L. Mason, "Multivariate control charts for individual observations," *Journal of Quality Technology*, vol. 24, no. 2, pp. 88–95, 1992.
- [12] J. A. N. Vargas, "Robust estimation in multivariate control charts for individual observations," *Journal of Quality Technology*, vol. 35, no. 4, pp. 367–376, 2003.
- [13] P. J. Rousseeuw and B. C. van Zomeren, "Unmasking multivariate outliers and leverage points," *Journal of the American Statistical Association*, vol. 85, no. 411, pp. 633–639, 1990.
- [14] P. J. Rousseeuw and A. M. Leroy, *Robust Regression and Outlier Detection*, Wiley-IEEE, New York, NY, USA, 2nd edition, 2003.
- [15] R. A. Maronna, R. D. Martin, and V. J. Yohai, *Robust Statistics, Theory and Methods*, John Wiley and Sons, New York, NY, USA, 2006.
- [16] P. J. Rousseeuw, "Least median of squares regression," *Journal of the American Statistical Association*, vol. 79, no. 388, pp. 871–880, 1984.

- [17] W. A. Jensen, J. B. Birch, and W. H. Woodall, "High breakdown estimation methods for phase I multivariate control charts," *Quality and Reliability Engineering International*, vol. 23, no. 5, pp. 615–629, 2007.
- [18] P. Filzmoser, R. Maronna, and M. Werner, "Outlier identification in high dimensions," *Computational Statistics and Data Analysis*, vol. 52, no. 3, pp. 1694–1711, 2008.
- [19] R. R. Wilcox, "Some small-sample properties of some recently proposed multivariate outlier detection techniques," *Journal of Statistical Computation and Simulation*, vol. 78, no. 8, pp. 701–712, 2008.
- [20] M. Habshah, M. R. Norazan, and A. H. M. R. Imon, "The performance of diagnostic-robust generalized potentials for the identification of multiple high leverage points in linear regression," *Journal of Applied Statistics*, vol. 36, no. 5, pp. 507–520, 2009.
- [21] J. F. Hair Jr., R. E. Anderson, R. L. Tatham, and W. C. Black, *Multivariate Data Analysis: With Readings*, Prentice-Hall, Upper Saddle River, NJ, USA, 1995.
- [22] M. H. Kutner, C. J. Nachtsheim, and J. Neter, *Applied Linear Regression Models*, McGraw-Hill, New York, NY, USA, 5th edition, 2005.
- [23] A. Imon, "Identification of high leverage points in logistic regression," *Pakistan Journal of Statistics*, vol. 22, no. 2, pp. 147–156, 2006.
- [24] D. Peña and V. J. Yohai, "The detection of influential subsets in linear regression by using an influence matrix," *Journal of the Royal Statistical Society B*, vol. 57, no. 1, pp. 145–156, 1995.
- [25] T. P. Ryan, *Modern Regression Methods*, Wiley-Interscience, New York, NY, USA, 2nd edition, 2008.
- [26] S. Chatterjee and A. S. Hadi, *Regression Analysis by Example*, John Wiley and Sons, New York, NY, USA, 2006.
- [27] A. S. Hadi and J. S. Simonoff, "Procedures for the identification of multiple outliers in linear models," *Journal of the American Statistical Association*, vol. 88, no. 424, pp. 1264–1272, 1993.
- [28] A. Imon and M. M. Ali, "Simultaneous identification of multiple outliers and high leverage points in linear regression," *Journal of Korean*, vol. 16, pp. 429–444, 2005.
- [29] D. C. Hoaglin and R. E. Welsch, "The hat matrix in regression and ANOVA," *American Statistician*, vol. 32, no. 1, pp. 17–22, 1978.
- [30] A. S. Hadi, "A new measure of overall potential influence in linear regression," *Computational Statistics and Data Analysis*, vol. 14, no. 1, pp. 1–27, 1992.
- [31] A. Imon, "Identifying multiple high leverage points in linear regression," *Journal of Statistical Studies*, vol. 3, pp. 207–218, 2002.
- [32] D. C. Montgomery, *Introduction to Statistical Quality Control*, John Wiley & Sons, New York, NY, USA, 5th edition, 2005.
- [33] T. P. Ryan, *Statistical Methods for Quality Improvement*, John Wiley and Sons, New York, NY, USA, 2nd edition, 2000.
- [34] J. H. Sullivan and W. H. Woodall, "Adapting control charts for the preliminary analysis of multivariate observations," *Communications in Statistics B*, vol. 27, no. 4, pp. 953–979, 1998.
- [35] W. A. Shewhart, *Economic Control of Quality of Manufactured Product*, D. Van Nostrand, Princeton, NJ, USA, 1931.





# Hindawi

Submit your manuscripts at  
<http://www.hindawi.com>

