
PROPOSED HOTELLING'S T^2 CONTROL CHARTS BASED ON LOCATION AND SCATTER MATRIX.

***Obafemi, Oluwafemi Samuel¹, Oyeyemi, Gafar Matanmi² and Bolarinwa, Folashade Adeola³**

¹Department of mathematics and Statistics, Federal Polytechnics Ado- Ekiti, Nigeria.

²Department of Statistics, University of Ilorin, Ilorin, Kwara State, Nigeria.

³Department of mathematics and Statistics, Federal Polytechnics Ado- Ekiti, Nigeria.

(corresponding author)

ABSTRACT

In many quality control settings, the product (process) under examination may have two or more correlated quality characteristics; hence, an appropriate approach is needed to simultaneously monitor all the quality characteristics. The Hotelling T^2 control chart based on the usual sample mean vector and variance - covariance matrix performs poorly, especially when there are multiple out of control points in the multivariate data set. Several alternative methods have been proposed, this includes methods based on the minimum volume ellipsoid (MVE) and the minimum covariance determinant (MCD) among many other methods. These control charts are powerful in detecting a reasonable number of outlying data. In this paper we propose a modified Hotelling T^2 control charts using the eigen-values obtained from scatter matrix/ variance-covariance matrix of the multivariate data. The methods were used on a real-life data set. The studies show that this method outperforms the classical Hotelling T^2 control charts and compete well with charts based on MVE and MCD for a small number of observations when the number of out of control points was increased

KEYWORDS: Hotelling T^2 , Control charts, Outliers, eigen-value, quality Characteristics, Mean shift.

1. INTRODUCTION

Statistical Process Control (SPC) field is closely-related to univariate outlier detection methods. It looks into the cases where the univariable stream of measure represents a stochastic process. Control charts are the most popular tools and techniques used in Statistical Process Control (SPC) to monitor the quality characteristics of products and services in organizations and industries. In many of these industrial processes, it is frequently required to monitor several quality characteristics at the same time. Such quality characteristics may includes weight, degree of hardness, thickness, width and length of a certain type of tablets (Liu, 1995). For the fact that the quality characteristics of these products are clearly correlated, the separate univariate control charts for monitoring such quality characteristics may not be efficient in signaling out-of-control points and changes in the overall quality of the products, therefore, it is desirable to have a control charts that can measure and monitor these characteristics simultaneously, multivariate control charts are the most appropriate tools applicable in such situations (Alt, 1985).

A subtle approach towards monitoring and improving quality is the statistical process control (SPC) charts that aim at quality improvement through reduction of variation. It is one of the primary techniques used in the controlling of product quality.

In many quality control settings, the product (process) under examination may have two or more correlated quality characteristics (variables); hence, an appropriate approach is needed to monitor all these characteristics simultaneously. This leads to the multivariate quality control problem, which is the subject of research by many quality control experts. As the objective of performing multivariate statistical process control is to monitor the process over time, in order to detect any unusual events that allow quality and process improvement, it is essential to track the cause of an out of control signal. However, as opposed to univariate control charts, the complexity of multivariate control charts and the cross-correlation among variables make the analysis of assignable causes of out-of-control signals difficult. This has been the basis for extensive research performed in the field of multivariate control chart since the 1940s, when Hotelling (1947) recognized that the quality of a product might depend on several correlated characteristics.

Outliers can heavily influence the estimation of scatter matrix and subsequently the parameters or statistics that are needed to be derived from it. Therefore, a robust estimate of scatter matrix that would not be affected by outliers is required to obtain valid results (Hubert and Engelen, 2007).

An outlier is an observation that deviates so much from other observations as to arouse suspicion that it was generated by different mechanism as defined in statistical quality control concepts (Hawkins, 1980).

Because of the importance of multivariate control charts in monitoring at least two correlated quality characteristics and to detects and remove outlying variable (quality characteristics), this research proposed an alternative hotelling T^2 control chart based on robust method of estimating location and scatter matrix and compared it efficiency with Hotelling T^2 using classical, Minimum Volume Ellipsoid (MVE) and Minimum Covariance Determinant (MCD)

Material and methodology

The classical estimators for μ and Σ are the empirical mean and covariance matrix respectively where

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i, \quad S_n = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})^T.$$

Both the empirical mean and variance-covariance are highly sensitive to outliers.

The Minimum Covariance Determinant (MCD) method developed by Rousseeuw, (1984) is a highly robust estimator of multivariate location and scatter matrix. The empirical covariance matrix $C = C(X)$ of X is the $p \times p$ matrix defined by;

$$C(X) = \frac{1}{h} \sum (x_i - t)(x_i - t)^T \text{ where } t_i = \frac{1}{h} \sum_{j=1}^h x_{ij} \text{ or in a matrix notation implying that}$$

$$C(X) = \frac{1}{h} \sum_{i=1}^h x_i x_i^T - t t^T$$

Minimum volume ellipsoid (MVE) estimator is first proposed by Rousseeuw (1984), and it has been studied extensively for non-control charts settings and often used as multivariate outlier’s detector. This method involves, drawing a random sub-sample of $(p + 1)$ different observation, indexed by

$$J_j = \{x_{j1}, x_{j2}, \dots, x_{jp+1}\}.$$

For this sub-sample, the mean and covariance matrix are computed as;

$$\bar{x}_J = \frac{1}{p+1} \sum_{i=1}^{p+1} x_{ji} \text{ and } S_J = \frac{1}{p} \sum_{i=1}^{p+1} (x_{ji} - \bar{x})(x_{ji} - \bar{x})^T$$

The corresponding ellipsoid is then increased or decreased to contain exactly h

Observations. Usually, $h = \frac{n + p + 1}{2}$,

The MVE estimators is given $\bar{x}_{MVE} = \bar{x}_{J^*}$ and $S_{MVE} = c_{n,p}^2 (\chi_{p,0.5}^2)^{-1} m_J^2 S_{J^*}$

where $c_{n,p}^2$ is a correction factor for small sample, and $\chi_{p,0.5}^2$ is the median of the chi-squared distribution with p degrees of freedom? The mahalanobis distances base on this estimator, proved to be very effective in detecting several outliers in a multivariate point cloud Rousseeuw and Zomeren (1990).

Obafemi and Oyeyemi (2018) proposed a robust method of estimating the location (mean) and scatter (variance-covariance) matrix for a multivariate data set in the presence of outliers using eigen-values.

Given y_1, y_2, \dots, y_p for multivariate normal, i.e $Y_p \sim N_p(\mu, \Sigma)$ where Σ is positive definite. The proposed method of estimating the parameter μ and Σ is focusing more on the eigen-values of variance covariance matrix. Given a p -dimensional multivariate normal data $Y_{p \times m}$ with m observation $\{y_i\}_{i=1}^m$, the interest here is to obtain subsets of $\{y_i\}_{i=1}^m$ of size $k = p+1$ that satisfy some criteria stated below:

The minimum of the arithmetic mean of eigen-value; minimum of the harmonic mean of the eigen-value; and minimum of the Geometric mean of the Eigen-value obtained from the classical variance-covariance

matrix.

A sample of size k from m is therefore drawn that will give C_{p+1}^m possible subsets of size p + 1. The

variance-covariance matrix Σ_j is therefore estimated as $\Sigma_j = \frac{1}{p+1} (y_j - \bar{y}_j)(y_j - \bar{y}_j)^T$.

For each of the p x p matrix Σ_j , the eigen-values $\lambda_{j1}, \lambda_{j2}, \dots, \lambda_{jp}$ are obtained. For each set of eigen-values, the following are obtained; the Arithmetic mean (A), the harmonic mean (H) and the geometric mean (G) of the eigen-values,

The objective here is to obtain data points in which the eigen-values of its variance-covariance matrix will satisfy at least two of the following criteria; Minimum Arithmetic mean (A_{min}); Minimum Harmonic mean (H_{min}); and Minimum Geometric mean (G_{min}), taking into consideration when the variance covariance matrix is from uncorrelated variables and also dependent variable (correlated variables).

The resulting covariance matrix will be inflated or deflated to accommodate good data point $h = \frac{(n+p+1)}{2}$. within the observed data.

The classical mean and variance-covariance matrix of the h points is the proposed robust estimate of the vector of means and scatter matrix given as;

$$\bar{y}_{proposed} = \frac{1}{h} \sum_{i \in J} y_i \text{ and } S_{proposed} = \frac{1}{h-1} \sum_{i \in h} (y_i - \bar{y}_{proposed})(y_i - \bar{y}_{proposed})^T \cdot (\chi_{j,0.025}^2)^{\frac{1}{p}}$$

where $(\chi_{p,0.025}^2)^{\frac{1}{p}}$ is a correcting factor with p as the number of variables $h = \frac{(n+p+1)}{2}$.

hence the proposed Hotelling T^2 is obtained with the following parameter

$$T^2_{proposed} = (y - \bar{Y}_{proposed})^T S^{-1} (y - \bar{Y}_{proposed})^T$$

HOTELLING'S T² CONTROL CHART

The χ^2 statistic is used when a sample is drawn out of a population and the real population parameters are known. Ghare and Torgerson (1968) developed a bivariate control chart based on the χ^2 statistic. This control chart uses a graphical implementation with an elliptical in control region in the two-dimensional XY plane. A more general case of this control charts is the charts commonly referred to as χ^2 control chart. This differs from Ghare and Torgerson's control charts in two ways; first, the χ^2 control

chart is not limited to the bivariate case. Second, the implementation involves plotting χ^2 statistic over time and comparing the statistic to the determined critical value. This loses the separation of quality characteristic value that the bivariate method possesses. On the other hand, in the case where the covariance matrix is not known, then the Hotelling T^2 control chart, Hotelling (1947) becomes appropriate. And the statistics is as shown in equation 1.0 below

$$;T^2 = n(X - \bar{X})^T S^{-1} (X - \bar{X}) \dots\dots\dots (1.0)$$

The Hotelling T^2 chart is analogous to the Shewart X-bar chart when parameters are unknown. The lower control limit for a T^2 chart is zero (0) and the upper control limit is given by;

$UCL = \frac{p(m-1)(n-1)}{mn-m-p+1} F_{\alpha, p, mn-m-p+1}$, LCL = 0, Where $F_{\alpha, p, mn-m-p+1}$ represents the 100 $(1-\alpha)^{th}$ percentile of the F-distribution with degree of freedom p and $mn-m-p+1$ Following Jeson *et al.* (2007) and Vargas (2003) a robust minimum volume ellipsoid (MVE) as alternative to Hottelling T^2 statistics is defined as;

$T^2_{MVEi} = (x_i - \bar{X}_{MVE})^T S^{-1}_{MVE} (x_i - \bar{X}_{MVE})$, where \bar{X}_{MVE} and S_{MVE} are the estimates of mean vector and scatter matrix obtained through MVE respectively. And for the robust minimum covariance determinant (MCD) an alternative to Hottelling T^2 statistics is defined as;

$T^2_{MCDi} = (x_i - \bar{X}_{MCD})^T S^{-1}_{MCD} (x_i - \bar{X}_{MCD})$, where \bar{X}_{MCD} and S_{MCD} are the estimates of mean vector and scatter matrix obtained through MCD location respectively. The statistical software R is used to calculate MVE and MCD estimates based on a genetic Algorithm.

Determining the upper control limits for the control charts

The upper control limits were determined from 1000 simulation such that all the methods considered had overall false alarm probability of 0.05. The limits were obtained by generating 1000 data set for $n=30$ and $p=2$, The Hotelling T^2 statistic, T_i^2 were computed for $i = 1, 2, \dots, n$. The maximum value was recorded and the 95th percentile of the maximum value of Hotelling T^2 for $j = 1, 2, \dots, 1000$ was taken to be the upper control limits for the control chart. The values obtained are 9.02, 16.29, 16.29 and 15.42 for the normally distributed variables for Classical, MCD, MVE and Proposed methods, respectively

REAL LIFE DATA ILUSTRATION

We consider Tablet NIR spectral data, (2011) source from http://www.idrc-chambersburg.org/shootout_2002.htm. The data are spectra measured in the transmittance mode, of 460 (row) pharmaceutical tablets, 650 columns; the first two variables of 30 samples were considered and were

reproduced. The two variables are used in the construction of Hotelling T^2 control charts using the classical, the MCD, MVE and the proposed method the resulted charts are then compared.

The sample location and scattered matrix for the classical methods of table 1.0 are

$$\bar{X}_{classical} = \begin{bmatrix} 3.2773 \\ 3.2698 \end{bmatrix}, S_{Classical} = \begin{bmatrix} 0.009628 & 0.008058 \\ 0.008058 & 0.007781 \end{bmatrix},$$

The location and scattered matrix for the two existing robust method using R is given below:

Location and scatter matrix for MCD

$$\bar{X}_{MCD} = \begin{bmatrix} 3.2538 \\ 3.2502 \end{bmatrix}, S_{MCD} = \begin{bmatrix} 0.003885 & 0.003872 \\ 0.003872 & 0.003863 \end{bmatrix}$$

Location and scatter matrix for MVE

$$\bar{X}_{MVE} = \begin{bmatrix} 3.2538 \\ 3.2502 \end{bmatrix}, S_{MVE} = \begin{bmatrix} 0.003885 & 0.003872 \\ 0.003872 & 0.003863 \end{bmatrix}$$

The location and scattered matrix using the proposed robust method are given as follows;

$$\bar{X}_{Proposed} = \begin{bmatrix} 3.2602 \\ 3.2559 \end{bmatrix}, S_{Proposed} = \begin{bmatrix} 0.004917 & 0.004958 \\ 0.004958 & 0.005006 \end{bmatrix}$$

Columns 4,5,6 of Tables 1.0 and 1.1 show the values of Hotelling T^2 statistics, $T_i^2_{normal}$, $T_i^2_{MCD}$, $T_i^2_{MVE}$, and $T_i^2_{proposed}$, based on the Classical, MCD, MVE and proposed methods respectively.

Comparing the values obtained using the four statistics to their respective upper control limits which are 9.02, 16.29, 16.29 and 15.42 for the Classical, MCD, MVE and the proposed methods respectively. It is observed that all the methods signal the 3rd and the 20th observation as out of control points.

Finally, the observations 13, 25 and 30 are modified to (3.422, 3.4625), (3.4635, 3.5666) and (3.4333, 3.3944) respectively making the number of outlier equal to 6 (20%).

Table 1.1 gives the statistics of the multivariate control chats using the four different methods. The classical method is able to signal two observation as out-of-control while the other three methods; the

MCD, MVE and the proposed methods signal observations 3, 9, 13, 20, 25 and 30 as out of control which are all the outliers introduced.

Table 1.0 Data set and Hotelling T² statistic using, Classical, MCD, MVE and the proposed robust method' when there are 2 (6.7%) outliers

s/n	x1	x2	Classical	Mcd	Mve	Proposed
1	3.20889	3.20453	0.56479	6.35371	6.35371	3.61919
2	3.29773	3.29764	0.09571	0.25329	0.25329	0.25799
3	3.50000	3.46750	12.1864	20.7310	20.7310	17.2273
4	3.21821	3.21233	0.40312	9.40144	9.40144	5.64698
5	3.23531	3.23330	0.20071	1.93614	1.93614	1.09164
6	3.26813	3.26277	0.04464	6.94894	6.94894	4.84454
7	3.25249	3.25280	0.06737	0.16839	0.16839	0.10300
8	3.13896	3.13358	2.5502	12.7421	12.7421	7.02711
9	3.26651	3.25985	0.06543	10.2216	10.2216	7.11745
10	3.27590	3.27253	0.02753	3.12863	3.12863	2.25511
11	3.30798	3.30790	0.19005	0.39341	0.39341	0.39962
12	3.16344	3.16169	1.78471	4.68587	4.68587	2.81861
13	3.24515	3.23931	0.11851	8.47972	8.47972	5.48205
14	3.21516	3.20936	0.44963	9.31775	9.31775	5.54898
15	3.25219	3.24832	0.06565	4.19761	4.19761	2.69514
16	3.33588	3.33082	0.77107	6.76416	6.76416	6.03419
17	3.23255	3.23170	0.24541	1.00624	1.00624	0.55904
18	3.26281	3.25498	0.08736	13.74195	13.74195	9.50382
19	3.22073	3.21799	0.38389	3.22881	3.22881	1.80564
20	3.50403	3.65556	27.28189	42.4731	42.4731	29.21725
21	3.36075	3.36281	1.12098	2.41735	2.41735	1.57932
22	3.12058	3.11619	3.36419	12.0089	12.0089	6.76837
23	3.30934	3.30583	0.27433	3.48457	3.48457	2.98266
24	3.32884	3.32623	0.54822	2.63227	2.63227	2.52697
25	3.40090	3.3947	2.87449	11.9588	11.9588	11.6791
26	3.27128	3.26555	0.05384	7.77759	7.77759	5.49070
27	3.23634	3.23345	0.18544	2.99495	2.99495	1.75115
28	3.31354	3.30856	0.37153	6.17724	6.17724	5.14467
29	3.2983	3.29601	0.13190	1.76869	1.76869	1.50533
30	3.17548	3.17632	1.49387	2.50627	2.50627	2.07431

The bold observations are the outlying values

Hotelling T^2 control charts using, Classical, MCD, MVE and the proposed robust method' when there are 2 (6.7%) outliers

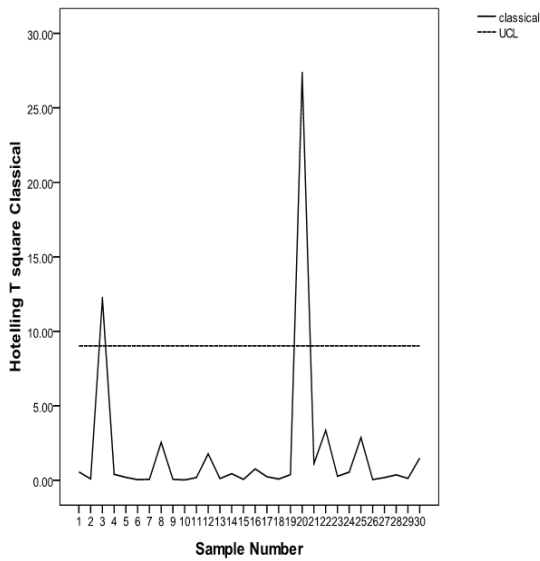


Figure 1.0a: Classical control charts with 2 (6.7%) outliers

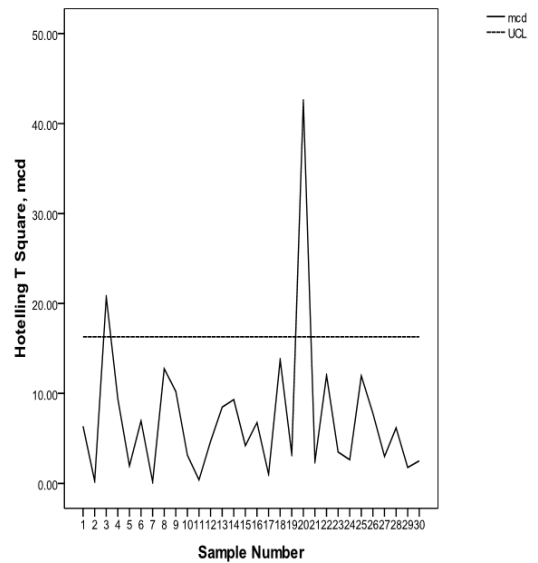


Figure 1.0b: MCD control charts with 2 (6.7%) outliers

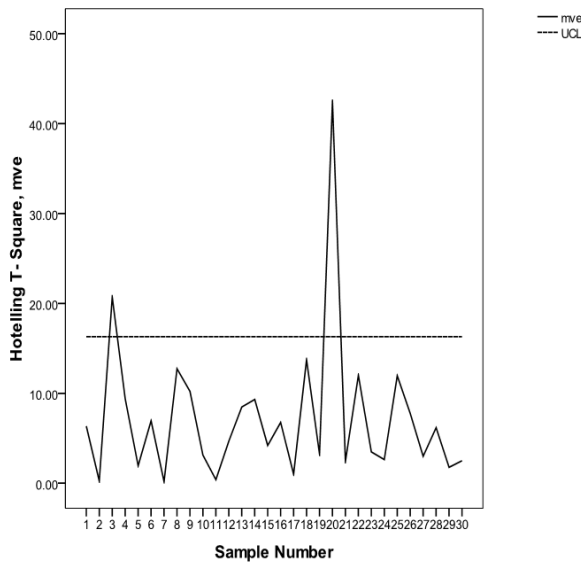


Figure 1.0c: MVE control charts with 2 (6.7%) outliers

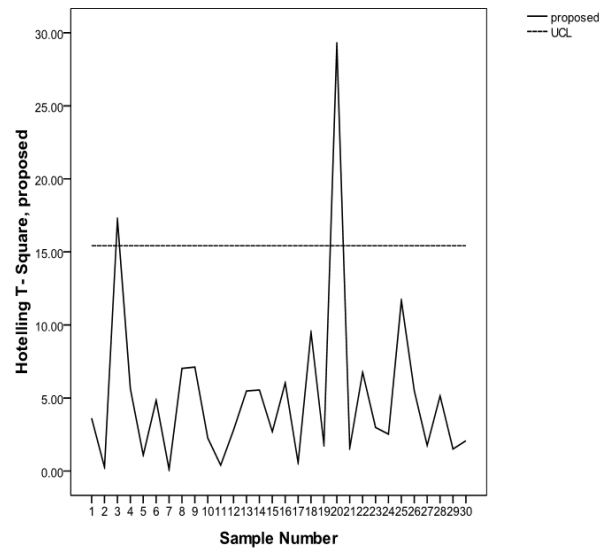


Figure 1.0d: Proposed control charts with 2 (6.7%) outliers

Table 1.1 Data set and Hotelling T2 statistic using, Classical, MCD, MVE and the proposed robust method' when there are 6 (20%) outliers

s/n	x1	x2	Classical	Mcd	Mve	Proposed
1	3.20889	3.20453	0.75571	2.91142	2.91142	1.18736
2	3.29773	3.29764	0.01107	5.24814	5.24814	4.49214
3	3.5225	3.41335	7.77481	219.751	219.751	182.923
4	3.21821	3.21233	0.61239	1.74652	1.74652	0.66373
5	3.23531	3.2333	0.36679	4.18953	4.18953	2.63072
6	3.26813	3.26277	0.08441	0.25506	0.25506	0.09424
7	3.25249	3.2528	0.19385	7.84084	7.84084	5.93276
8	3.13896	3.13358	2.43108	7.4816	7.4816	2.76751
9	3.46000	3.57220	9.35292	258.1855	258.1855	224.208
10	3.27590	3.27253	0.04467	1.05768	1.05768	0.75617
11	3.30798	3.30790	0.02471	5.05811	5.05811	4.41573
12	3.16344	3.16169	1.71346	9.98416	9.98416	5.02348
13	3.42220	3.46250	2.69891	37.97514	37.97514	33.3535
14	3.21516	3.20936	0.65971	1.89981	1.89981	0.71401
15	3.25219	3.24832	0.19643	1.31389	1.31389	0.69981

16	3.33588	3.33082	0.14126	0.40476	0.40476	0.15291
17	3.23255	3.23170	0.39904	6.41066	6.41066	4.32018
18	3.26281	3.25498	0.12552	1.19052	1.19052	1.00893
19	3.22073	3.21799	0.56153	3.93491	3.93491	2.13322
20	3.50403	3.31560	15.3282	664.272	664.272	560.759
21	3.36075	3.36281	0.43494	9.55558	9.55558	8.48215
22	3.12058	3.11619	3.01583	11.02737	11.0273	3.93144
23	3.30934	3.30583	0.01572	0.57012	0.57012	5.08003
24	3.32884	3.32623	0.10184	1.31494	1.31494	1.14402
25	3.4635	3.56660	8.55970	219.1118	219.111	190.541
26	3.27128	3.26555	0.06807	0.19248	0.19248	0.0845
27	3.23634	3.23345	0.35608	2.92295	2.92295	1.67394
28	3.31354	3.30856	0.02459	0.07242	0.07242	0.03546
29	3.2983	3.29601	0.00385	1.74877	1.74877	1.50015
30	3.4333	3.39444	1.94277	24.28449	24.28449	19.6555

The bold observations are the outlying values

Hotelling T² control charts using Classical, MCD, MVE and the proposed robust method' when there are 6 (20%) outliers

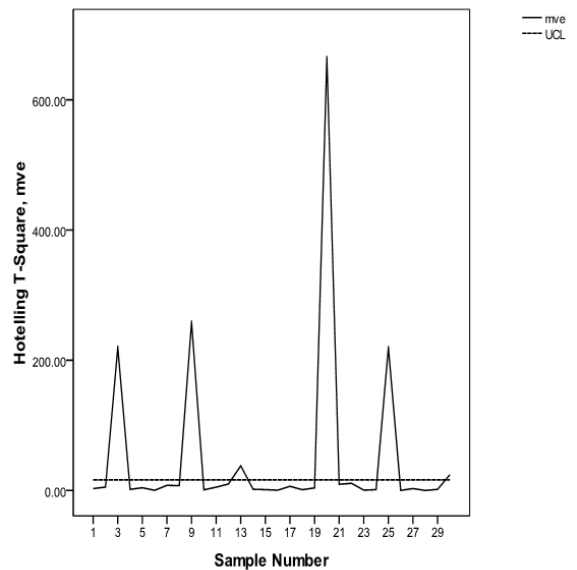
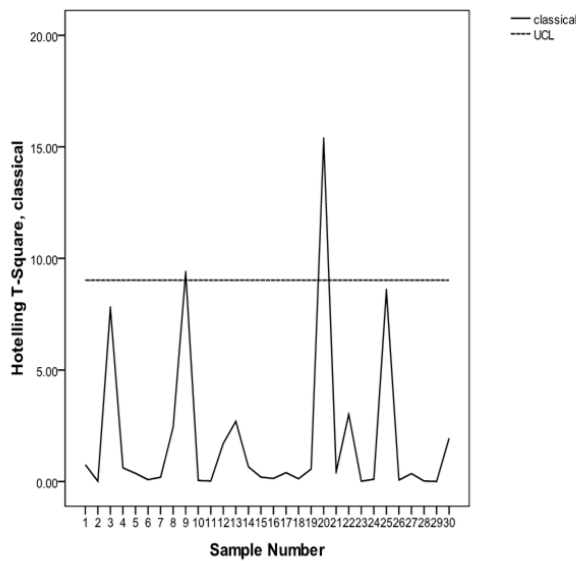


Figure 1.1a: Classical control charts with 6 (20%) outliers

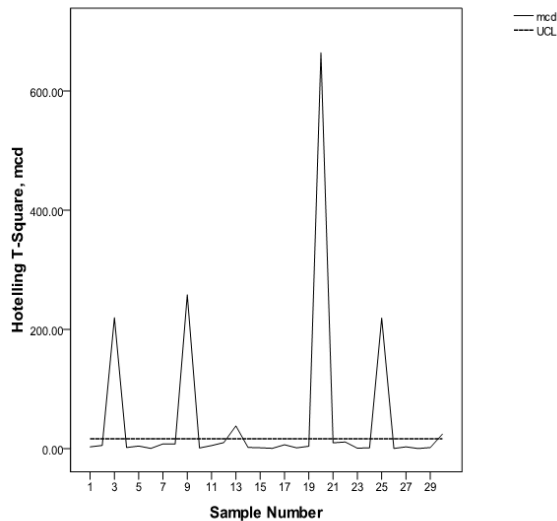


Figure 1.1a: Classical control charts with 6 (20%) outliers

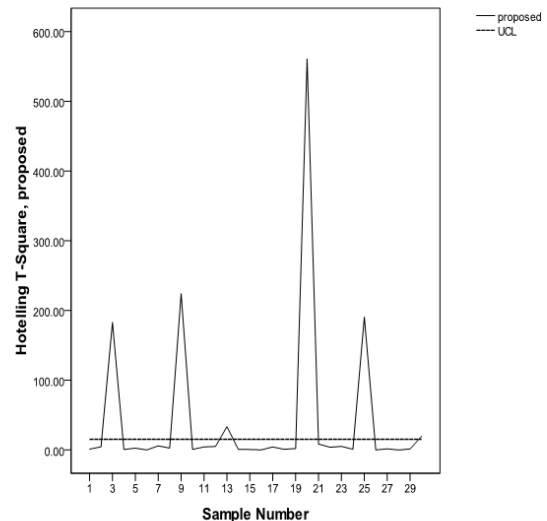


Figure 1.1b: MCD charts with 6 (20%) outliers

Figure 1.1d: Proposed control charts with 6 (20%) outliers

FINDINGS, CONCLUSION AND RECOMMENDATION

All the methods were employed in obtaining the Hotelling’s T^2 statistic which was used in obtaining the control chart. The basis for comparison was the proportions of times each chart correctly detect the presence of outliers by signaling an out of control.

The four methods were employed in the construction of Hotelling T^2 control charts using real life data, of spectra measured in the transmittance mode of pharmaceutical tablets with outliers introduced to the data arbitrarily at different magnitude. All the methods were able to signal an out of control when there were two outlying points in the multivariate data set. However, when the number of outliers was increased arbitrarily to 20%, the classical method only signaled two points as being out of control while the other three robust methods (MCD, MVE and Proposed) signaled all the outlying point (6 points) introduced. The points were above the upper control limits.

The proposed hotelling T^2 control charts compete favourably well with the two known robust hotelling T^2 (MVE and MCD) in detecting a shift in the mean by signaling an out of control when there are few or many outliers in the multivariate data set.

The proposed method performed better than the classical methods when there is higher percentage of outliers in the data set. In conclusion, the classical Hotelling T^2 is only efficient when there are no or very few outliers in the multivariate data set while the other methods studied in this work are not only efficient when there are presence of few outliers in the data set but also efficient in the presence of multiple outliers. However, the proposed robust method performed well and very efficient in the two extreme cases. The efficiencies of the classical and the existing and widely used methods (MVE and MCD) of estimation are combined by the proposed robust method. It follows that, if there is no information regarding the number of outlier in the multivariate data set as far as the analyst is concerned, it is recommended that the proposed robust method of Hotelling T^2 control charts can be employed, since it works well with the other robust methods.

REFERENCES

- Alt, F.B. (1995). Multivariate Quality control in S.Kotz and Johnson (Eds), The Eyclopedia of statistical Sciences, 6, 110-122. New York: John Wiley & Sons.
- Ghare, P. H. and Torgerson, P. E. (1968). The Multi-characteristics Control Chart. *Journal of Industrial Engineering*, 19, 269-272
- Hawkins, D.M. (1980). Identification of Outliers. New York, Chapman and Hall, LTD.
- Hotelling, H. (1974). Multivariate quality control in techniques of statistical analysis edited by Eisenhart, Hastay and wallis, McGraw-Hill, New York, NY.
- Hubert, M and Engelen, S. (2007). Fast Cross –validation of High- breakdown Resampling Algorithms for PCA. *Computational Statistics & data Analysis*, 51(10) 5013-5024.
- Jensen, W.A., Birch, J.B. and Woodall, W. H. (2000). High Breakdown Estimation Methods for Phase I Multivariate Control Charts. [WWW.Stat.vt.edu\(newsletter/ annal_report_statistics-2008.p.d.f\)](http://WWW.Stat.vt.edu/newsletter/annal_report_statistics-2008.p.d.f).
- Liu,R.V. (1995). Control Charts for Multivariate Processes. *Journal of the American Statistical Association*, 90 (432), 1380-1387
- Obafemi, O.S. and Oyeyemi, G.M. (2018). Alternative estimator for multivariate location and scatter matrix in the presence of outlier (*Annals computer science series*. 16th Toma 2nd fast)
- Rousseeuw, P. J. (1984). least modern of squares regression. *Journal of the American Statistical Association*, 79: 871-880

Rousseeuw P. J. and Van Zomeren B. C. (1990). unmasking multivariate outliers and leverage points. Journal of the American statistical Association. Vol. 85 (411), pp. 633-651.

Vargas, J.A. (2003). Robust estimation in multivariate control charts for individual observations. Journal of Quality Technology 35(4), pp. 367-376)