

## Illustrating The Invariance Property of Hotelling's $T^2$ Statistic Using Three Variables

<sup>1</sup>Nura Muhammad, <sup>2</sup>Ambi Polycarp Nagwai ,

<sup>1</sup>Lecturer, Department of Computer and Statistics, Kano State polytechnic, Kano, Nigeria<sup>1</sup>

<sup>2</sup>Professor, Department of Statistics, University of Ewe, Ghana<sup>2</sup>

**ABSTRACT:** The Hotelling's  $T^2$  Chart is used in Process monitoring when a quick detection of shifts in the mean vector is desired. When there are shifts in a multivariate Control Charts, it is clearly shows that special cause variation is present in the Process, but the major drawback of Hotelling's  $T^2$  is inability to identify which variable(s) is/are the source of the signals. Hence effort must be made to identify which variable(s) is/are responsible for the out- of- Control situation. In this article, we employ Mason, Young and Tracy (MYT) approach in identifying the variables for the signals.

**KEYWORDS:** Quality Control, Multivariate Statistical Process Control, Hotelling's  $T^2$  Statistic.

### I. INTRODUCTION:

Nowadays, One of the most powerful tools in Quality Control is the Statistical Control Chart developed in the 1920s by Walter Shewharts, the Control Chart found wide spread use during World War II and has been employed, with various modifications, ever since. Multivariate Statistical Process Control (SPC) using Hotelling's  $T^2$  statistic is usually employed to detect shifts. However, Hotelling's  $T^2$  Control Chart has a shortcoming as it can't figure out the causes of the change. Thus, decomposition of  $T^2$  is recommended and aims at paving a way of identifying the variable(s) significantly contributing to an out- of- Control signals.

**Multivariate Statistical Process Control:** The Shewhart control charts have been widely applied in a variety of industries because it is very simple to implement and the information generated from the Shewhart control charts is also easy for plant staff to understand. However, monitoring each process variable with separate ShewhartControl chart ignores the correlation between variables and does not fully reflect the real process situation. Nowadays, the process industry has become more complex than it was in the past and inevitably that number of process variables need to be monitored has increased dramatically. It's very often, these variables are multivariate in nature and using Shewhart control charts becomes insufficient.

### II. HOTELLING'S $T^2$ STATISTIC

Hotelling H. (2) can be viewed as the originator of multivariate control charts. Hotelling Proposed a concept of generalized distance between new observations to its sample mean. We first illustrate how this method works with a Multivariate case. Assuming these  $X_1, X_2$  and  $X_3$  are distributed according to the Multivariate normal distribution. Let represent a  $p$  dimensional vector of measurements made on a Process. Assuming that when the process is in control, the  $X_i$  are independent and follow a multivariate normal distribution with mean vector  $\mu$  and covariance matrix  $\Sigma$ , normally  $\mu$  and  $\Sigma$  are unknown, but we can use estimated from a historical data set with  $n$  observations.

**Phase I and Phase II :** The application of Hotelling's  $T^2$  statistic shall be categorized into two phases. Phase I tests whether the preliminary process was in control and phase II tests whether the future observation remains in-control (Alt.), (1). Phase I operation refers to the construction of in-control data set. Same idea as Shewhart control chart, control limits are estimated from a period of in-control data. To obtain this in-control data, the raw data set needs to be purged. For instance, the outliers need to be removed and the missing data needs to be substituted with an estimate. During phase I Operation, Hotelling's  $T^2$  statistic is calculated for each measurement and compared to the control limit, which will follows Chi-square distribution (according to Richard, A.J. & Dean, W.W.), (5)

$$T^2 = (X - \bar{X})S^{-1}(X - \bar{X}) \sim X_{\alpha, p}^2 \quad (\text{Chi - square distribution}) \quad (1)$$

Also other research shows that the control limit follows Beta distribution (Mason, Young & Tracy,). (3)

$$T^2 = (X - \bar{X})S^{-1}(X - \bar{X}) \sim \frac{(n-1)^2}{2} \beta\left(\alpha, \frac{p}{2}, \frac{n-p-1}{2}\right) \quad (2)$$

$n$ : number of preliminary observations Both control limits will be approximate when the number of observations is large. The control limit based on Chi-square distribution is established on the assumption that are true values  $\mu$  and  $\Sigma$ , which is just an approximate situation (Mason, Young & Tracy,) (3). Beta distribution is more precise and is a recommendable choice. After purging the raw data with Hotelling's  $T^2$  statistic, the in-control data set is ready for monitoring future observations which is termed as phase II operation. The control limit for determining future observation is different from the one in phase I. It follows an F distribution with  $p$  and  $(n-p)$  degrees of freedom. The idea of using Hotelling's  $T^2$  statistic in phase I and phase II is the same. Each measurement is examined whether it is out-of-control by checking if it deviates extraordinarily from its sample mean. It should be reminded to choose the correct upper control limit on different purposes. The  $T^2$  statistic calculated from all the observation will be plotted in a chart against observation serious and compared to the upper control limit. It should be noticed that there is no center line and the lower control limit is set to zero, because the meaning of  $T^2$  statistic is a generalized distance between the observation and its sample mean.

### III. CASE DEMONSTRATION

Starting by selecting any one of the  $p$  variables, followed by any of the  $(p-1)$  remaining Variables to condition on the first selected variable, the next step is to select any of the remaining  $(p-2)$  variables to condition on the first two selected variables. Iterating the same procedure will generate all the decomposition equations which give the same overall  $T^2$  statistic. For instance, taking the first  $p$  variables, assuming we select  $p=1$  is selected (remaining 2nd, and 3rd variables). The decomposition results are presented as follows;

$$T_1^2 + T_{2.1}^2 + T_{3.12}^2$$

$$T_1^2 + T_{3.1}^2 + T_{2.13}^2$$

The next step is to keep the selected  $p$  variable and the  $(p-1)$  remaining variables conditioned on the  $p$  selected variable constant (that is, the first and the second decomposition terms in the above decompositions). Then, any of the remaining  $(p-2)$  variables not used in the first decomposition to be condition on the first two kept variables is selected. Computing in the same procedure will generate all the decomposition equations which produce the same overall  $T^2$  statistic. Hence, the complete  $T^2$  decomposition results using three variables ( $p = 3$ ) are presented in equation (5). Hotelling's  $T^2$  statistic helps us to determine whether a measurement is in-control or not.

**Description of the methodology used:** To reach the main purpose of this research, the data were generated from Wapco plc. The data is then used to generate a real data set collected at every point where there is a variation; three variables were monitored for twenty five observations were recorded.

**The MYT Decomposition :** The  $T^2$  statistic can be decomposed into P-orthogonal components (Mason, Tracy and Young), (3) for instance, if you have P-variables to decompose, the procedure is as follows:

$$T^2 = T_1^2 + T_{2.1}^2 + \dots + T_{p.1,2,\dots,p-1}^2 \quad (3)$$

The first term is an unconditional  $T^2$ , decomposing it as the first variable of the

$$T_1^2 = \frac{(X_1 - \bar{X}_1)^2}{S_1^2} = \frac{(-18)^2}{54} = 6 \quad (4)$$

Where  $X_1$  and  $S_1$  is the mean and standard deviation of the variable  $X_1$  respectively.  $T_j^2$  Will follow an F-Distribution which can be used as upper Control Limit

$$UCL(x_j) = \frac{n+1}{n} F_{\alpha, 1, n-1},$$

Taking a case of three variables as an example, it can be decomposed as,

$$T^2 = \begin{cases} T_1^2 + T_{2.1}^2 + T_{3.12}^2 \\ T_1^2 + T_{3.1}^2 + T_{2.13}^2 \\ T_2^2 + T_{1.2}^2 + T_{3.12}^2 \\ T_2^2 + T_{3.2}^2 + T_{1.23}^2 \\ T_3^2 + T_{1.3}^2 + T_{2.13}^2 \\ T_3^2 + T_{2.3}^2 + T_{1.23}^2 \end{cases} \quad (5)$$

These decompositions give the same overall  $T^2$  statistic as stated by Mason, *et al.*(4).

**Model for The  $T^2$  Decomposition Using Three Variables :** Starting by selecting any one of the  $p$  variables, followed by any of the  $(p-1)$  remaining Variables to condition on the first selected variable, the next step is to select any of the remaining  $(p-2)$  variables to condition on the first two selected variables. Iterating the same procedure will generate all the decomposition equations which give the same overall  $T^2$  statistic. For instance, taking the first  $p$  variables, assuming we select  $p=1$  is selected (remaining 2nd and 3rd). The decomposition results are presented as follows;

$$T_1^2 + T_{2.1}^2 + T_{3.12}^2$$

$$T_1^2 + T_{3.1}^2 + T_{2.13}^2$$

The next step is to keep the selected  $p$  variable and the  $(p-1)$  remaining variables conditioned on the  $p$  selected variable constant (that is, the first and the second decomposition terms in the above decompositions). Then, any of the remaining  $(p-2)$  variables not used in the first decomposition to be condition on the first two kept variables is selected. These yield the following results:

$$T_2^2 + T_{1.2}^2 + T_{3.12}^2$$

$$T_2^2 + T_{3.2}^2 + T_{1.23}^2$$

$$T_3^2 + T_{1.3}^2 + T_{2.13}^2$$

$$T_3^2 + T_{2.3}^2 + T_{1.23}^2$$

The first terms are unconditional while others are conditional. The importance of the above result is that it allows one to examine the  $T^2$  statistic from many different perspectives. From the decomposition, it shows that an increase in the number of variables will lead to an increase in the number of terms. These make computation become troublesome.

**Table 1: Unique Decomposition Terms (cited from Mason *et al.*, 1997a)**

No. of variables ( $p$ )	No. of possible terms $p! \times p$	No. of unique terms $p \times 2^{(p-1)}$
2	4	4
3	18	12
4	96	32
5	600	80
10	36288000	5120

#### IV. ILLUSTRATION:

In this segment, we intend to demonstrate by way of example, how an assignable signal could be detected, the data that was used for the purpose of the research was collected from a Portland Cement Company in Lagos, Nigeria. The data consists of temperature from a boiler machine in which twenty five observations were taken under different temperatures. The mean and covariance matrix given below was obtained from the historical data set.

$$s = \begin{pmatrix} 54 & 0.958 & 20.583 \\ 0.958 & 4.84 & 2.963 \\ 20.583 & 2.963 & 22.993 \end{pmatrix}, \quad X = \begin{pmatrix} 525 \\ 513.56 \\ 538.92 \end{pmatrix}, \quad \text{and } X' = (507 \quad 516 \quad 527)$$

In order to illustrate the invariance property of  $T^2$ , we look at

$$T^2 = (X - \bar{X})' S^{-1} (X - \bar{X}) = 10.85$$

Therefore, to get the contribution of each component, the  $T^2$  above is decomposed as shown below:

$$\begin{aligned} T_1^2 + T_{2,1}^2 + T_{3,12}^2 &= 6 + 0.63 + 4.22 = 10.85 \\ T_1^2 + T_{3,1}^2 + T_{2,13}^2 &= 6 + 1.54 + 3.31 = 10.85 \\ T_2^2 + T_{1,2}^2 + T_{3,12}^2 &= 1.44 + 5.19 + 4.22 = 10.85 \\ T_2^2 + T_{3,2}^2 + T_{1,23}^2 &= 1.44 + 6.83 + 2.58 = 10.85 \\ T_3^2 + T_{1,3}^2 + T_{2,13}^2 &= 6.18 + 1.36 + 3.31 = 10.85 \\ T_3^2 + T_{2,3}^2 + T_{1,23}^2 &= 6.18 + 2.09 + 2.58 = 10.85 \end{aligned}$$

## V. CONCLUSION

Multivariate process control had gained ground in statistical process control and the needs to improve the MSPC in a production process become the order of the Day. A well-known MYT model was adopted in this article to decompose the Hotelling's  $T^2$  statistic into orthogonal components, each of which states the contribution of individual process variable. Also, we were able to demonstrate the invariance property of  $T^2$  statistic, where the decomposition had the same overall results regardless of the ordering of the observation vector.

## REFERENCES:

1. Alt, F.B., (1985), Multivariate quality control, *Encyclopedia of the statistical sciences*, Vol.6, Pp.111-122
2. Hotelling, H., (1931). The generalization of student's ratio. *The Annals of Mathematical Statistics* v.2, pp.360-378
3. Mason, R.L., Tracy, N.D. & Young, J.C., (1992), Multivariate control charts for individual observations, *Journal of Quality Technology*, Vol.24, No.2, pp.88-95
4. Mason, R.L., Tracy, N.D., and Young, J.C. (1997a). A practical Application for Interpreting Multivariate T2 Control Chart Signals, *Journal of Quality Technology*, 29(4): 396-501
5. Richard, A.J. & Dean, W.W., (2002). *Applied multivariate statistical analysis*, N.J.: Prentice Hall

## APPENDIX I

Boiler Temperature Data

N0.of observation	X <sub>1</sub>	X <sub>2</sub>	X <sub>3</sub>
1	507	516	527
2	512	513	533
3	520	512	537
4	520	514	538
5	530	515	542
6	528	516	542
7	522	513	537
8	527	509	537
9	533	514	528
10	530	512	528
11	530	512	541
12	527	513	541
13	529	514	542
14	522	509	539
15	532	515	545
16	531	514	543
17	535	514	542
18	516	515	537
19	514	510	532
20	536	512	540
21	522	514	540
22	520	514	540
23	526	517	546
24	527	514	543
25	529	518	544
Total	13125	12839	13473
Average	525	513.36	538.92