

# A Control Chart Based on A Nonparametric Multivariate Change-Point Model

MARK D. HOLLAND

*Beckman Coulter, Inc., 1000 Lake Hazeltine Drive, Chaska, MN 55318*

DOUGLAS M. HAWKINS

*University of Minnesota School of Statistics, 313 Ford Hall, Minneapolis, MN 55455*

Phase-II statistical process control (SPC) procedures are designed to detect a change in distribution when a possibly never-ending stream of observations is collected. Several techniques have been proposed to detect a shift in location vector when each observation consists of multiple measurements. These procedures require the user to make assumptions about the distribution of the process readings, to assume that process parameters are known, or to collect a large training sample before monitoring the ongoing process for a change in distribution. We propose a nonparametric procedure for multivariate phase-II statistical process control designed to detect shifts in location vector that relaxes these requirements based on an approximately distribution free multivariate test statistic. This procedure may not be appropriate for some multivariate distributions with unusual dependence structure between vector components. A diagnostic tool that can be used if a historical sample of data is available is provided to assist user in determining if the proposed procedure is appropriate for a given application.

**Key Words:** Location Test; Phase-II Statistical Process Control.

## 1. Introduction

In phase-II statistical process control (SPC), a potentially never-ending stream of observations is collected and each time a new observation is obtained a SPC check is performed. This procedure continues until the SPC check signals that the process has gone out-of-control. We assume that variability in process readings comes from two main sources: *common cause* random variability that cannot be removed without making fundamental changes to the process, and *special cause* variability that can be identified and removed. The process is in statistical control when only common cause variability exists. The process may go out-of-control in two distinct

ways; an isolated special cause may affect one or a small number of process readings and then go away, while a sustained special cause continues until it is identified and fixed (Hawkins et al., 2003). A special cause may be a shift in mean, a change in variance, or the process readings may change distribution.

Many challenges arise when conducting SPC on multivariate process readings that do not exist in the univariate case. Checking the assumption of multivariate normality is difficult, and real world data is unlikely to follow the multivariate normal (MVN) distribution. In addition, transforming data to normality is often possible in one dimension, but transformations to multivariate normality are limited (Qiu, 2008). For univariate data, shifts in mean can only occur in two directions: up or down. However, for multivariate data, shifts in the mean vector can occur in any arbitrary direction. This causes difficulties when defining distribution-free multivariate SPC procedures. Distribution-free multivariate SPC pro-

---

Dr. Holland is a Biostatistician. He is a Member of ASQ. His email address is [holland.mark@gmail.com](mailto:holland.mark@gmail.com).

Dr. Hawkins is a Professor. He is a Fellow of ASQ. His email address is [dhawkins@umn.edu](mailto:dhawkins@umn.edu).

cedures are those based on a test statistic which has a null distribution which does not depend on the distribution of the original data vectors. Some nonparametric procedures are unable to detect shifts in all directions (Qiu and Hawkins, 2001). Others are able to detect shifts in all directions, but are designed to detect a specific type of shift, such as an upward or downward shift in a single component of the mean vector (Qiu and Hawkins, 2003).

The problem of multivariate SPC has been studied to some extent, and the proposed procedures require a variety of assumptions about the parameters of the in-control or out-of-control process. Crosier (1988) derived a multivariate generalization of the univariate CUSUM control chart that allows more rapid detection of sustained shifts in mean vector than the multivariate Shewart chart. Lowry et al. (1992) proposed a multivariate exponentially weighted moving average (MEWMA) control chart, which is an extension of the univariate exponentially weighted moving average procedure. Both procedures are designed for use when the data is distributed according to the multivariate normal distribution with fixed and known in-control mean and covariance matrix.

Zamba and Hawkins (2006) defined a change-point model for multivariate normal data that uses the generalized likelihood ratio test in a similar fashion to the univariate procedure of Hawkins et al. (2003). As in the univariate case, the change-point model alleviates the need for a large phase-I sample to estimate in-control parameters of the multivariate normal distribution, and the chart performs well across a wide range of shifts in mean vector.

Qiu and Hawkins (2001, 2003) proposed distribution-free CUSUM procedures that do not require the assumption of multivariate normality, but require knowledge of the in-control distribution of the antirank vectors. Qiu (2008) proposed a distribution-free CUSUM procedure that estimates the in-control distribution using log linear modeling. Because this procedure is distribution-free and a method for estimating the in-control distribution is specified, it is well suited for multivariate SPC use when little is known about the distribution of the process readings, but a large training sample is available.

In this paper, we propose a nonparametric multivariate phase-II SPC procedure designed to detect shifts in the location vector that does not require knowledge of the in-control or out-of-control distribution parameters and does not require a large histori-

cal sample of in-control data. This procedure is based on an approximately distribution free multivariate generalization of the Wilcoxon-Mann-Whitney test, and extends the nonparametric change-point model based on the univariate rank-sum test (Hawkins and Deng, 2010) to the multivariate setting. Although the underlying test statistic is approximately distribution free, this procedure may not be appropriate for some multivariate distributions with unusual dependence structure between vector components. A diagnostic is proposed to aid in determining if the proposed procedure is appropriate that can be used if a historical sample of data is available.

### Multivariate Nonparametric Location Test

Given a sample of  $p \times 1$  random vectors,  $\mathbf{x}_1, \dots, \mathbf{x}_n$ , define the multivariate change-point model as

$$\mathbf{x}_i \sim \begin{cases} F(\boldsymbol{\mu}) & \text{if } i \leq \tau \\ F(\boldsymbol{\mu} + \boldsymbol{\delta}) & \text{if } i > \tau, \end{cases}$$

where  $F$  represents the distribution function of a multivariate location family,  $\boldsymbol{\delta}$  is an arbitrary shift in location vector, and  $\tau$  is the index of the observation at which the change-point occurs. Then, the class of multivariate rank based tests proposed by Choi and Marden (1997) can be applied to test the hypothesis  $H_0 : \boldsymbol{\delta} = \mathbf{0}$  vs. the alternative  $H_1 : \boldsymbol{\delta} \neq \mathbf{0}$ , given a fixed value of  $\tau = k$ .

If  $p = 1$ , let  $r^{(i)}$  be the rank of  $x_i$  among the observations. To center the ranks, set  $R(x_i) = 2r^{(i)} - n - 1$ . Then,

$$R(x_i) = \sum_{j=1}^n \text{SIGN}(x_i - x_j),$$

where  $\text{SIGN}(x) = -1$  if  $x < 0$ ,  $0$  if  $x = 0$ , and  $1$  if  $x > 0$ . The Wilcoxon-Mann-Whitney test statistic for difference in location between the two samples  $\{x_1, \dots, x_k\}$  and  $\{x_{k+1}, \dots, x_n\}$  is

$$u_k = \sum_{i=1}^k R(x_i).$$

To define a multivariate generalization of this test statistic, we must choose a substitute for the SIGN function. Choi and Marden (1997) suggest using a kernel function  $\mathbf{h}(\mathbf{x}, \mathbf{y})$  such that

$$\mathbf{h}(\mathbf{x}, \mathbf{y}) = -\mathbf{h}(\mathbf{y}, \mathbf{x}).$$

Then,  $\mathbf{h}(\mathbf{x}, \mathbf{x}) = 0$ , and  $\mathbf{h}(\mathbf{x}, \mathbf{y})$  can be interpreted as a measure of the difference between  $\mathbf{x}$  and  $\mathbf{y}$ .

For  $1 \leq i \leq n$ , define the multivariate centered rank,

$$\mathbf{R}_n(\mathbf{x}_i) = \sum_{j=1}^n \mathbf{h}(\mathbf{x}_i, \mathbf{x}_j).$$

Choi and Marden (1997) recommend the use of a directional rank test, using the kernel function

$$\mathbf{h}(\mathbf{x}, \mathbf{y}) = \frac{\mathbf{x} - \mathbf{y}}{\|\mathbf{x} - \mathbf{y}\|}.$$

Then,  $\mathbf{R}_n(\mathbf{x}_i)$  is the average of the unit vectors pointing from each data point to  $\mathbf{x}_i$ , and is called the directional rank of  $\mathbf{x}_i$ .

To define the directional rank test statistic, we need notation that allows us to define within group rank vectors. For each possible change point,  $k = 1, \dots, n-1$ , define

$$\mathbf{R}_{n,k}^*(\mathbf{x}_i) = \sum_{j=k+1}^n \mathbf{h}(\mathbf{x}_i, \mathbf{x}_j)$$

and

$$\bar{\mathbf{r}}_n^{(k)} = \frac{1}{k} \sum_{i=1}^k \mathbf{R}_n(\mathbf{x}_i).$$

The pooled sample covariance matrix of the centered rank vectors,  $\hat{\Sigma}_{k,n}$ , is formed by independently estimating the covariance matrices of the left and right segments of data and combining them into a single estimate.

$$\begin{aligned} \tilde{\Sigma}_{k,n} = & \frac{n^2}{n-2} \left( \frac{1}{k^2} \sum_{i=1}^k \mathbf{R}_k(\mathbf{x}_i) \mathbf{R}_k(\mathbf{x}_i)' \right. \\ & \left. + \frac{1}{(n-k)^2} \sum_{i=k+1}^n \mathbf{R}_{n,k}^*(\mathbf{x}_i) \mathbf{R}_{n,k}^*(\mathbf{x}_i)' \right), \end{aligned}$$

The unpooled covariance matrix,  $\hat{\Sigma}_n$ , is the sample covariance matrix of all data vectors taken as a single sample.

$$\hat{\Sigma}_n = \frac{1}{n-1} \sum_{i=1}^n \mathbf{R}_n(\mathbf{x}_i) \mathbf{R}_n(\mathbf{x}_i)'.$$

Under the null hypothesis,

$$\frac{nk}{(n-k)} \bar{\mathbf{r}}_n^{(k)'} \tilde{\Sigma}_{k,n}^{-1} \bar{\mathbf{r}}_n^{(k)} \rightarrow \chi_p^2.$$

Choi and Marden (1997) note that the unpooled covariance estimator is a consistent estimator of the true covariance matrix of the rank vectors under the null hypothesis, so it will also yield a valid test statistic.

The directional rank statistic for testing for difference in location vector between the left segment of the data  $\{\mathbf{x}_1, \dots, \mathbf{x}_k\}$  and the right segment of the data  $\{\mathbf{x}_{k+1}, \dots, \mathbf{x}_n\}$  is

$$r_{k,n} = \bar{\mathbf{r}}_n^{(k)'} \hat{\Sigma}_{k,n}^{-1} \bar{\mathbf{r}}_n^{(k)}, \quad (1)$$

where  $\hat{\Sigma}_{k,n} = ((n-k)/nk)\hat{\Sigma}_n$ . One may suspect that defining the test statistic using the unpooled covariance matrix would result in a more powerful test to detect differences in location vector because, if  $\|\delta\| > 0$ , the pooled covariance matrix estimator would be polluted by the difference in location between the left and right segments of data. However, simulations showed that the power of the test is similar using either covariance matrix estimator Choi and Marden (1997). For computational efficiency, we use the unpooled covariance estimator (see appendix).

### Nonparametric Multivariate Phase-II SPC Procedure

In this section, we construct a nonparametric multivariate procedure for conducting phase-II SPC based on the test statistic defined in equation (1). The procedure consists of maximizing the test statistic across possible change-point values, excluding indices near the beginning and end of the sequence. We do not compute the test statistic for very short left or right segments of data because we observed through simulation that, for data of dimension  $p = 2, 5, 10$ , the distribution of  $r_{k,n}$  may be sensitive to the distribution of the original data vectors when  $k$  or  $n-k$  is small.

Define

$$r_{\max,c,n} = \max_{c < k < n-c} r_{k,n}.$$

The quarantine constant,  $c$ , determines the number of data points at either end of the sequence that are not considered in the search for a possible change-point. For large enough values of  $c$ , the distribution of  $r_{k,n}$  is approximately  $\chi_p^2$  for all  $c < k < n-c$ , so the procedure will be approximately distribution-free. However, choosing  $c$  too large will result in a procedure that cannot easily detect shifts that occur near the end of the data sequence. We assume that the change-point is not known beforehand, and we estimate the change-point using

$$\hat{\tau} = \arg \max_{c < k < n-c} r_{k,n}.$$

When observation  $\mathbf{x}_n$  is obtained, compute  $r_{\max,c,n}$  and if  $r_{\max,c,n} > h_{\alpha,p,c,n}$ , signal that a shift has oc-

curred. Otherwise, collect another observation and repeat.

The sequence of control limits  $\{h_{\alpha,p,c,n}\}$  can be chosen such that the conditional probability of a false alarm when observation  $n$  is collected is equal to  $\alpha$ , given that no previous false alarm has occurred (Hawkins et al., 2003). That is, we chose  $\{h_{\alpha,p,c,n}\}$  to satisfy

$$P[r_{\max,c,n} > h_{\alpha,p,c,n} \mid r_{\max,c,j} \leq h_{\alpha,p,c,n}; j < n] = \alpha.$$

As in the previous change-point models proposed for phase-II SPC, finding an analytical solution for the sequence of control limits is intractable (Hawkins et al., 2003; Zamba and Hawkins, 2006; Hawkins and Deng, 2010). Thus, we estimated the control limits using five million simulated sequences of uncorrelated MVN data. Because the simulation study to derive control limits can be quite tedious, an R package available on CRAN (Holland, 2013) contains tables of control limits for data of dimension  $p = 2, \dots, 10$  along with an implementation of the proposed change point model. Contents of the R package including the R code and data files are also provided as supplementary material available at <http://www.asq.org/pub/jqt/>.

The change-point model based on  $r_{k,n}$  does not depend on knowledge of any process parameters, so phase-II monitoring can begin without first collecting a phase-I training sample. However, we must collect  $n > p$  observations to ensure that the covariance matrix of directional rank values,  $\hat{\Sigma}_{k,n}$  is non-singular. Furthermore, the model requires  $n > 2c+1$  so at least one possible change-point is outside of the quarantined region. For control limit simulations, we began monitoring at observation  $n = \max(p + 10, 2c + 3)$ . When the quarantine constant does not require that  $n \geq p + 10$ , we still waited an extra ten observations before monitoring because in practice users will likely accumulate a few “learning” observations before process monitoring begins. A table of control limits for  $p = 5, c = 15$  for  $n \leq 500$  is provided in Table 1. The control limits increase in an approximately linear fashion for values of  $n > 100$ . Control limits can be extended beyond  $n = 500$  by fitting a linear regression model to the control limit values for  $n > 100$  and extrapolating based on the fitted linear model. A simulation study was performed to confirm that this extension of the control limits achieves the desired in control ARL for multivariate normal data with  $p = 2, 5, 10$  and  $\rho = 0, 0.9$  with a targeted in control ARL of 2000.

TABLE 1. Control Limits,  $h_{\alpha,p,c,n}$ , for Change-Point Model Using  $r_{k,n}$  with  $p = 5, c = 15$

$n$	In-Control ARL = $1/\alpha$				
	100	200	500	1000	2000
33	14.100	15.209	16.553	17.485	18.355
34	13.500	14.724	16.193	17.221	18.175
35	13.261	14.567	16.137	17.200	18.221
36	13.158	14.518	16.154	17.264	18.316
37	13.097	14.516	16.209	17.360	18.452
38	13.073	14.531	16.296	17.473	18.583
39	13.062	14.557	16.366	17.587	18.723
40	13.061	14.596	16.445	17.684	18.842
45	13.147	14.819	16.790	18.149	19.416
50	13.237	14.989	17.094	18.535	19.855
60	13.392	15.259	17.519	19.059	20.497
70	13.505	15.436	17.785	19.423	20.958
80	13.564	15.562	17.994	19.673	21.250
90	13.606	15.645	18.131	19.840	21.478
100	13.646	15.718	18.249	20.037	21.655
125	13.714	15.829	18.425	20.255	21.990
150	13.740	15.896	18.541	20.415	22.175
200	13.790	15.982	18.681	20.591	22.414
300	13.819	16.051	18.813	20.768	22.647
500	13.890	16.113	18.916	20.906	22.820

## Evaluation of Performance

In order to compare the performance of different phase-II SPC procedures, one typically compares the average run length (ARL) of the procedures. The ARL is defined as the number of observations collected before the first signal occurs. Given that the in-control ARL of each procedure is constrained to be greater than or equal to  $1/\alpha$ , procedures can be compared by their out-of-control ARL values, which should be small.

To evaluate the in-control performance of the proposed procedure, we applied the procedure to simulated multivariate normal data and data simulated according to the following non-normal multivariate distributions. These distributions were chosen to represent common non-normal distributions observed in practice. The multivariate gamma distributions have correlated components with marginal distributions that are bounded below and skewed. The multivariate Cauchy distribution has correlated components with heavy tailed marginals.

### Cherian and Ramabhadran's Multivariate Gamma Distribution (CR MVG) (Kotz et al., 2000)

Let  $y_0, y_1, \dots, y_p$  be independent gamma random variables with pdf's

$$p_{y_i}(y_i) = \frac{1}{\Gamma(\theta_i)} e^{-y_i} y_i^{\theta_i-1}, \quad y_i > 0, \theta_i > 0.$$

Define  $\mathbf{x} = (y_0 + y_1, y_0 + y_2, \dots, y_0 + y_p)^T$ . Then, the marginal distribution of each  $x_i$  is a univariate gamma distribution with shape parameter  $\theta_0 + \theta_i$ . From the construction of  $x_i$ , it can be shown that

$$\text{cov}(x_i, x_j) = \text{var}(y_0) = \theta_0.$$

Thus,

$$\rho = \text{corr}(x_i, x_j) = \frac{\theta_0}{\sqrt{(\theta_0 + \theta_i)(\theta_0 + \theta_j)}}. \quad (2)$$

The CR MVG distribution with equicorrelated components can be parameterized with  $\theta_0$  and  $\rho$  using equation (2). Thus for simulation results that follow, values of  $\theta_0$  and  $\rho$  are given for each simulation scenario. When we refer to the CR MVG distribution with zero correlation between components, we mean that each component is an independent univariate gamma random variable with shape parameter  $\theta_0$ .

### Multivariate Transformed Gamma Distribution

We define another multivariate gamma distribution by transforming each component of a multivariate normal random vector to follow the gamma distribution. The marginals are the same as the CR multivariate gamma distribution, but the dependence structure is determined by the underlying multivariate normal distribution.

Let  $\mathbf{y} \sim N_p(\mathbf{0}, \Sigma)$ . For  $i = 1, \dots, p$ , set  $u_i = \Phi(y_i/\sigma_i)$ , where  $\Phi$  is the CDF of a standard normal random variable and  $\sigma_i$  is the standard deviation of  $y_i$ . By the probability integral transform,  $u_i \sim \text{Uniform}(0, 1)$ . Let  $x_i = F_i^{-1}(u_i)$ , where  $F_i$  is the CDF of a  $\text{gamma}(\theta, 1)$  random variable. Then,  $x_i \sim \text{gamma}(\theta, 1)$ , and we say that  $\mathbf{x}$  follows the transformed gamma distribution. For all simulation studies, we chose  $\Sigma$  to have all diagonal elements equal to one and off-diagonal elements equal to  $\rho$ .

### Multivariate Cauchy Distribution

The multivariate Cauchy (MVC) distribution is equivalent to the multivariate T distribution with one degree of freedom. To define the multivariate T

distribution, let  $\mathbf{y} \sim N_p(\boldsymbol{\mu}, \Sigma)$  and let  $w \sim \chi_\nu^2$ . Define

$$\mathbf{x} = \frac{1}{\sqrt{w/\nu}} \mathbf{y}.$$

Then,  $\mathbf{x}$  follows the multivariate T distribution. If  $\nu = 1$ ,  $\mathbf{x}$  follows the multivariate Cauchy distribution. Again, we chose  $\Sigma$  to have all diagonal elements equal to one and off-diagonal elements equal to  $\rho$  for all simulation studies.

### Study of In-Control ARL

Because the distribution of  $r_{k,n}$  may depend on the distribution of the data vectors in small to moderately sized samples, we used a Monte Carlo simulation study to determine degree of quarantining required to achieve acceptable in-control ARL performance. We simulated MVN data, CR multivariate gamma data with shape parameter values  $\theta_0 = 4, 2, 1/2$ , multivariate transformed gamma data with  $\theta_0 = 1/2$ , and multivariate Cauchy data. For all simulations, we set the nominal in-control ARL to  $1/\alpha = 500$ . We simulated equicorrelated data vectors of dimension  $p = 2, 5, 10$  with correlation  $\rho = 0, 0.5, 0.9$  and imposed quarantine values of  $c = 0, 9$  when  $p = 2$  and  $c = 0, 15$  when  $p = 5, 10$ . For each distribution and combination of  $p, \rho$ , and  $c$ , we simulated 10,000 sequences to estimate the in-control ARL. If the estimated in-control ARL was within approximately 10% of nominal, we concluded that the in-control performance is acceptable for practical purposes. Table 2 contains the degree of quarantine sufficient to achieve acceptable in-control performance, with a few exceptions.

Under the multivariate normal distribution, no quarantine is necessary to achieve acceptable in-control ARL for  $p = 2$ . When  $p = 5$  or 10, a quarantine value of  $c = 15$  is sufficient to achieve acceptable in-control ARL.

Under the CR multivariate gamma distribution with  $p = 2$ , a quarantine value of  $c = 9$  is sufficient to achieve an acceptable in-control ARL for all values of  $\theta_0$ . A quarantine value of  $c = 15$  is sufficient to achieve acceptable in-control ARL when  $p = 5$  for all of the scenarios except when  $\rho = 0.9$  and  $\theta_0 = 1/2$ , or 2. When  $p = 10$ , a quarantine of  $c = 15$  is sufficient to achieve an acceptable in-control ARL for all values of  $\theta_0$ , unless the components are correlated with  $\rho = 0.9$ .

Under the multivariate transformed gamma distribution, a quarantine constant of  $c = 9$  is sufficient for dimension  $p = 2$ , and  $c = 15$  is sufficient for di-

TABLE 2. In-Control ARL of Quarantined Change Point Model Based on Directional Rank Test Statistic Estimated Using 10,000 Simulated Sequences for Each Scenario. Nominal in control ARL = 500.

		$p = 2$		$p = 5$		$p = 10$	
		$c = 0$	$c = 9$	$c = 0$	$c = 15$	$c = 0$	$c = 15$
MVN							
$\rho = 0$		505	503	501	504	493	489
$\rho = 0.9$		481	501	402	496	320	477
MVT							
$df = 1$ (Cauchy)	$\rho = 0$	443	492	307	478	174	450
	$\rho = 0.9$	252	476	54	452	27	394
$df = 5$	$\rho = 0$	493	500	450	500	378	486
	$\rho = 0.9$	421	489	205	492	92	456
CR MVG							
$\theta_0 = 4$	$\rho = 0$	489	503	459	503	418	488
	$\rho = 0.9$	306	474	74	455	42	413
$\theta_0 = 2$	$\rho = 0$	476	502	399	489	352	470
	$\rho = 0.9$	224	472	38	439	21	391
$\theta_0 = 1/2$	$\rho = 0$	407	481	253	465	174	450
	$\rho = 0.9$	96	450	12	410	6	353
Transformed MVG							
$\theta_0 = 1/2$	$\rho = 0$	412	483	253	470	175	447
	$\rho = 0.9$	439	497	263	479	155	450

mension  $p = 5$ . When  $p = 10$ ,  $c = 15$  is only sufficient to achieve in-control ARL within 10% of nominal when  $\rho = 0.9$ . However, when  $\rho = 0$ , the in-control ARL is nearly within 10% of nominal. Note that the correlation matrix of the transformed gamma data is approximately equal to the covariance matrix of the original data vectors when the original data vectors follow the MVN distribution with equicorrelated components ( $\rho = 0.9$ ).

Under the multivariate Cauchy distribution,  $c = 9, 15$  are sufficient to achieve acceptable in control ARL when  $p = 2, 5$ , respectively. When  $p = 10$ , a quarantine of  $c = 15$  observations is only sufficient to achieve acceptable in-control ARL when  $\rho = 0$ . However, if the tails of the marginal distribution of vector components are slightly lighter than Cauchy (multivariate  $T$  distribution with  $df = 5$ ), a quarantine period of  $c = 15$  is sufficient to achieve accept-

able in control ARL for process readings with  $p = 10$  dimensions.

A moderate degree of quarantining is sufficient to achieve acceptable in-control ARL performance for the phase-II procedure based on  $r_{k,n}$  for multivariate normal data and for several non-normal multivariate normal distributions. However, for some non-normal multivariate distributions with dependence structures between vector components that differ greatly from the multivariate normal distribution, the amount of quarantine used here is not sufficient to achieve acceptable in-control ARL.

### Evaluation of Out-of-Control Performance

Our evaluation of out-of-control (OOC) performance using Monte Carlo simulation focused on three areas: (1) effect of quarantine, (2) performance

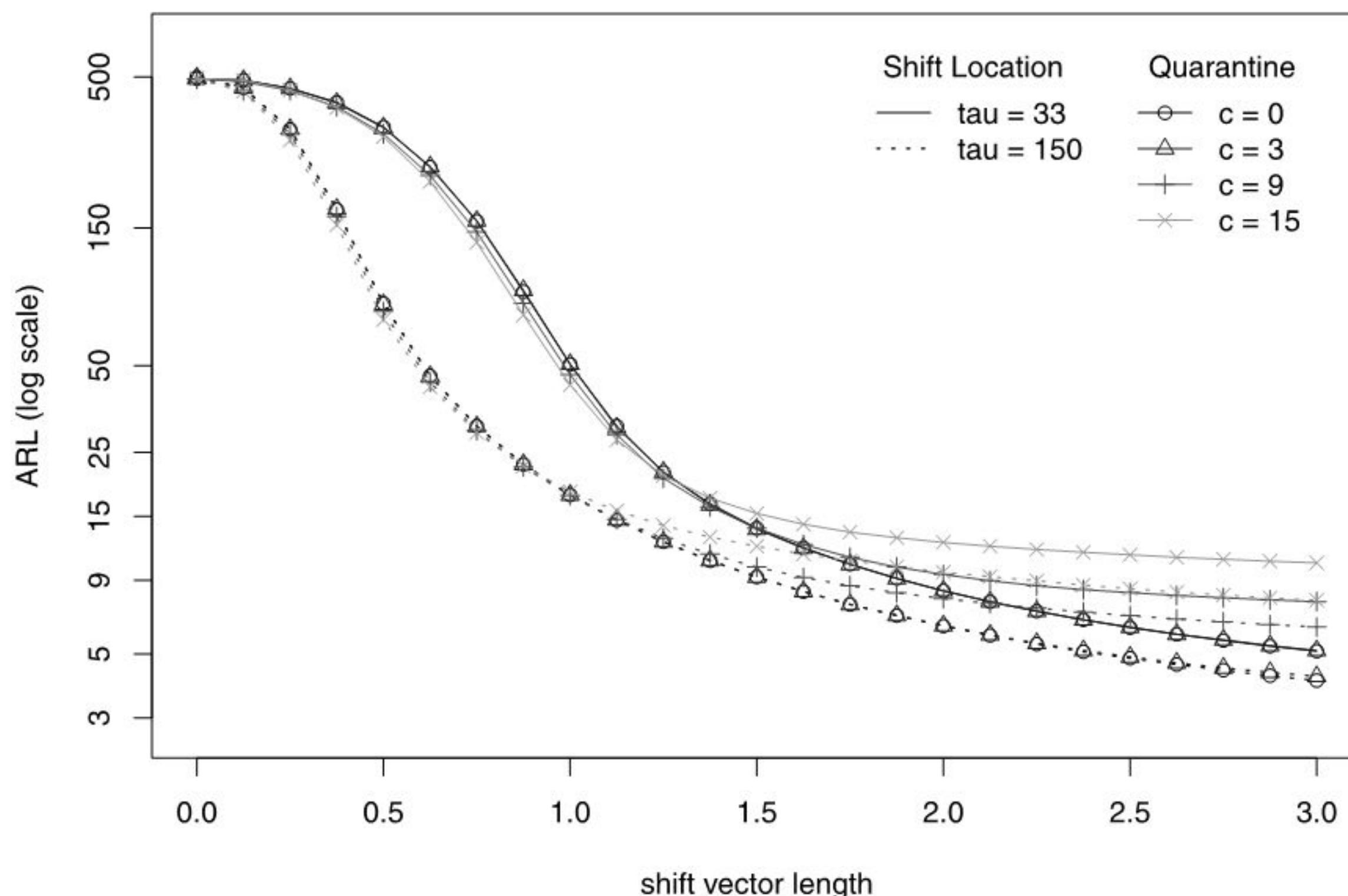


FIGURE 1. Quarantined Directional Rank Out-of-Control (OOC) ARL,  $p = 5$ . Effect of shift location and quarantine on OOC ARL for nonparametric change-point model.

of the nonparametric phase-II procedure compared to the parametric procedure (Zamba and Hawkins, 2006), (3) ability of the nonparametric change-point procedure to detect shifts in location of non-normal data.

To investigate the effect of imposing a quarantine period as well as the timing and direction of a mean vector shift on the phase-II procedure based on  $r_{k,n}$ , we simulated uncorrelated multivariate normal data ( $p = 2, 5$ , and  $10$ ), applied the procedure with  $c = 0, 3, 9, 15$ , and introduced mean vector shifts with  $\tau$  equal to the observation at which monitoring began. We simulated 100,000 data sequences for each of 25 shift vector lengths between 0 and 3, where vector length is measured using Euclidean distance. The use of a large number of simulated data sequences to estimate each OOC ARL value ensured that observed differences were not due to Monte Carlo error. We present simulation results for  $p = 5$  only, because the results for  $p = 2, 10$  are similar to those for  $p = 5$ . The out of control performance of the proposed procedure is similar if data vector components are correlated when Mahalanobis distance is used to define the length of the shift vector. The shift vector length defined using Mahalanobis distance is

$\|\delta\|_M = \sqrt{\delta^T \Sigma^{-1} \delta}$ , where  $\Sigma$  is the true covariance matrix of the simulated multivariate normal data. Each mean vector component was shifted an equal amount, but the results for multivariate normal data are identical if other shift directions are used (for example, if only 2 out of 5 mean vector components are shifted.) However, if the data vectors follow an arbitrary multivariate distribution the performance for detecting the presence of a shift in location vector may vary by the direction of the shift.

From Figure 1, we observe that imposing a quarantine does not affect the ability of the procedure to detect small to moderately sized shifts, but increasing the degree of quarantine does result in a larger ARL when a large shift occurs. Specifically, the ARL for all choices of  $c$  is nearly the same for all shift vector lengths less than 1.625. A shift vector length of 1.625 corresponds to an ARL of 12 for  $c = 0, 3, 9$  and an ARL of 14 for  $c = 15$ , when  $\tau = 33$ . In general, the quarantine period does not have an effect on out-of-control ARL for detecting the presence of a shift vector unless the magnitude of the shift is large enough that the un-quarantined procedure has ARL approximately equal to or less than the length of the quarantine period. We also observed that the

performance of the procedure improves if more in control observations are collected before a shift occurs ( $\tau = 33$  vs.  $\tau = 150$ ).

Because the ability of the proposed procedure to detect large shifts quickly is limited, the user may wish to supplement the proposed chart with a  $T^2$  chart (Hotelling, 1947) with the control limits chosen conservatively large so that the in control ARL is not reduced substantially. However, if the  $T^2$  chart is applied to data that is not elliptically symmetric, the in control performance of the  $T^2$  chart may be very poor. Because the distribution of  $T^2$  is not robust to non-normality, the in control ARL can be strongly affected by the data distribution even for control limits that are chosen to have a large in control ARL for multivariate normal data.

One may expect that the use of a quarantine would improve the detection of small shifts, because we would expect the control limit to decrease as  $c$  increases for a fixed in control ARL value. However, when the process is in control the maximum  $r_{k,n}$  value is not likely to occur near the beginning or the end of the sequence, so imposing a moderate quarantine period does not have a large impact on the distribution of  $\max r_{k,n}$  under the MVN distribution with independent random vector components. As a result, the control limits are not significantly affected for moderate quarantine values, so the detection of small shifts is not improved by introducing a moderate quarantine period.

### **Comparison With Parametric Multivariate Change Point Procedure**

We compared the performance of the nonparametric phase-II procedure to the parametric procedure (Zamba and Hawkins, 2006) using simulated multivariate normal data. We simulated equicorrelated MVN data with  $\rho = 0$  and  $\rho = 0.9$ , because  $r_{k,n}$  is not invariant to all linear transformations of the original data. We only apply the parametric procedure when  $\rho = 0$ , because this procedure is affine invariant so the performance is not affected by the covariance matrix of the original data. We measured shift vector length using Mahalanobis distance to provide a basis for performance comparison when each vector component is shifted an equal amount for uncorrelated and correlated data.

The performance of the nonparametric procedure is nearly identical for all shift sizes for uncorrelated and correlated data. Despite the fact that the data was simulated according to the multivariate normal

distribution, the nonparametric procedure detected small to moderately sized shifts *faster* on average than the parametric procedure. However, the parametric procedure detects large shifts faster on average than the nonparametric procedure. This result is consistent with the univariate case (Hawkins and Deng, 2010). The nonparametric procedure has the performance advantage for shift vectors with length less than 1.5, while the parametric procedure has the performance advantage for shift vectors with length greater than 1.5. A shift vector of length 1.5 corresponds to an ARL of approximately 15 for both procedures (see Figure 2).

Across all data dimensions ( $p = 2, 5, 10$ ), the nonparametric procedure has lower ARL for detecting shifts with magnitudes small enough that the ARL is greater than the length of the quarantine period.

### **Comparison With MEWMA**

Previous literature has claimed that the MEWMA is robust to departures from multivariate normality if the smoothing parameter is chosen to be reasonably small and the in control distribution of the data vectors is known (Stoumbos and Sullivan, 2002) and (Testik et al., 2003). Thus, it appears that the MEWMA would provide a natural comparison for the multivariate nonparametric change point model. One advantage of the proposed procedure over the MEWMA chart and other traditional charts for multivariate SPC is that Phase I analysis is not required. The proposed procedure does require the user to collect some "warm-up" cases, but this requirement is not equivalent to a small Phase I sample. The MEWMA chart based on estimated parameters from a Phase I sample only achieves the targeted in control ARL *averaged* over Phase I samples, so the performance is not guaranteed for any single data set. Substituting estimated parameters for the truth can have substantial effects on in control run length behavior for any single data set, as demonstrated in the univariate case by Hawkins et al. (2003). Unlike procedures that require estimation of in control parameters from a Phase I analysis, the change point procedure achieves desired in control run length behavior for each individual data set to which it is applied.

As demonstrated by a reviewer, when the in-control parameters of multivariate normal data are estimated from a phase I sample of  $n = 33$  observations, the MEWMA chart performs uniformly better than the proposed procedure for detecting both large

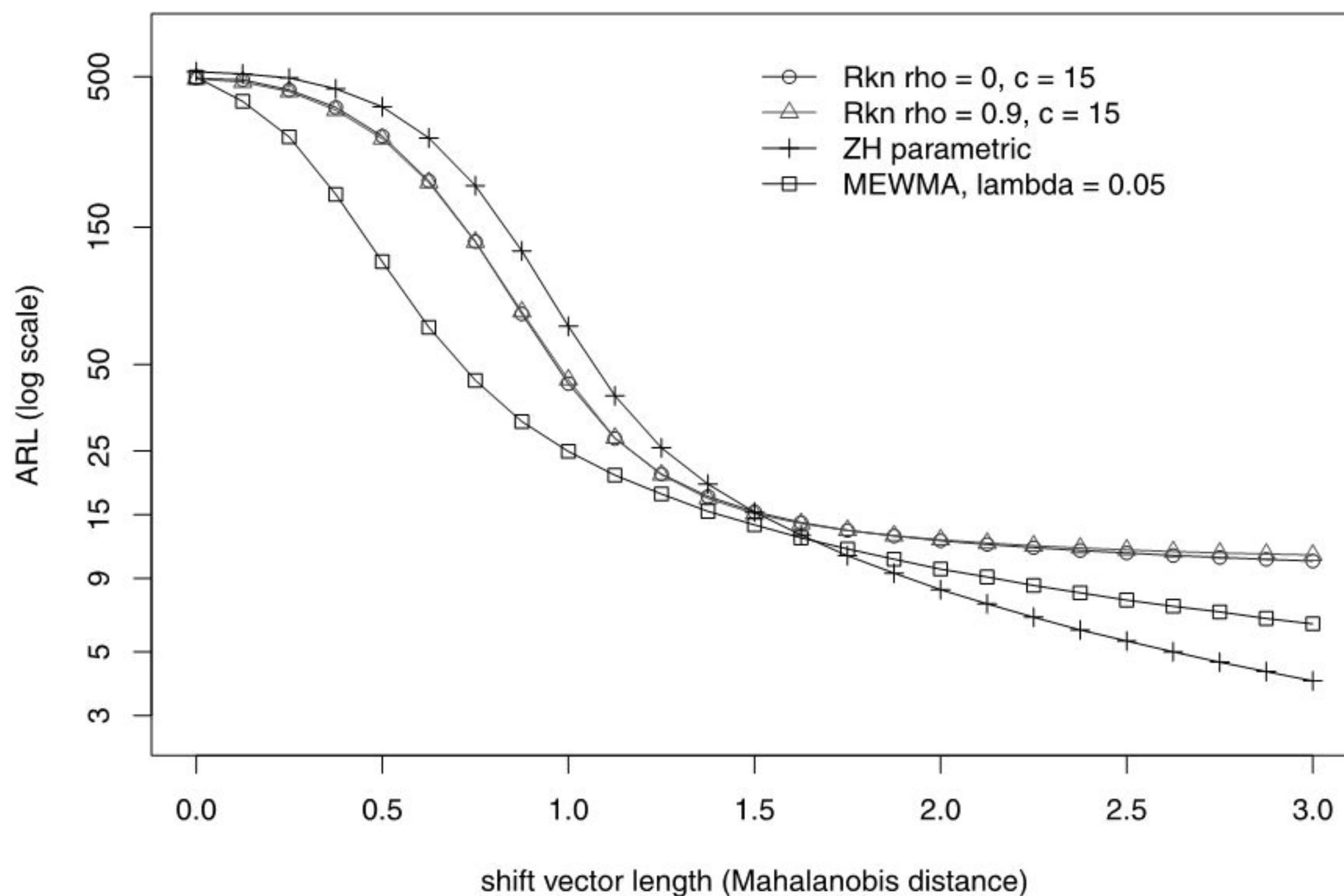


FIGURE 2. Out-of-control ARL of Quarantined Change-Point Model Based on Directional Rank Test Statistic Compared to Zamba and Hawkins (2006) Parametric Change-Point Model (ZH) and MEWMA, Shift Location  $\tau = 33$ .

and small shifts when the shift in mean is introduced immediately after monitoring begins (see Figure 2). Because the initial samples are not treated as a fixed phase I sample in the proposed procedure, the ability to rapidly detect small to moderately sized shifts in location vector improves as more in control samples are collected (see Figure 1). If in control parameters are estimated from a Phase I sample, the user must subjectively choose time points in the future to update the baseline period to enjoy a similar improvement in performance. The change point methodology also allows detection of shifts that occur during the initial warm-up period, so the user need not assume that the initial process readings are *iid* as in the case of a procedure that requires estimated in control parameters.

Table 3 contains in control average run length values for the MEWMA procedure applied to CR multivariate gamma data with  $p = 5$  dimensions. Because the in control ARL values are often well below the nominal value of 500, it is clear that the MEWMA procedure is not robust to all departures from multivariate normality. This result suggests that the MEWMA procedure may be robust when the mean and covariance matrix characterize the process well

(in the family of elliptically symmetric distributions) or if the in control distribution is known. However, the performance can be extremely poor for arbitrary multivariate distributions with dependence structure that is very different than the multivariate normal distribution, or in the common setting when the in control distribution parameters must be estimated.

TABLE 3. In-Control ARL Simulation Results for MEWMA for CR Multivariate Gamma Data with  $p = 5$  Dimensions. MEWMA parameters:  $\lambda = 0.05$ , CL = 28.049. In-control mean and covariance matrix estimated using Phase I sample of 33 observations for MEWMA procedure. Nominal in-control ARL = 500. Actual in-control ARL estimated using 10,000 simulated data sequences for each scenario.

	$\theta_0 = 4$	$\theta_0 = 2$	$\theta_0 = 1/2$
$\rho = 0$	426	400	253
$\rho = 0.5$	437	401	258
$\rho = 0.9$	245	177	86

Note that the estimated in control ARL for multivariate normal data estimated using 10,000 simulated sequences and the MEWMA parameters in Table 3 equals 499.

### OOC Performance for Non-Normal Data

To conclude our evaluation of the out-of-control performance of the nonparametric procedure, we applied the procedure to simulated multivariate transformed gamma data with shape parameter  $\theta = 1/2$ . When  $p = 2$ , the performance of the procedure is comparable for correlated and uncorrelated data for all shift sizes. When  $p = 5, 10$ , moderately sized shifts are detected faster on average for correlated data than uncorrelated data, while the OOC ARL is comparable for small or large shifts.

### Diagnostic to Select Degree of Quarantine

Because the proposed procedure is only approximately distribution free, in this section we define a diagnostic tool to aid in determining if the procedure is appropriate for a given data set if a historical sample of *iid* process readings is available. The degree of quarantine required to achieve acceptable in control run length behavior does not depend solely on the skewness or tail weight of the marginal distributions of the observation vector components or the covariance matrix. A more complete characterization of multivariate dependence is required to determine if the recommended degree of quarantine is sufficient to achieve acceptable in control ARL performance. The diagnostic tool defined below is based on the Anderson-Darling Goodness-of-Fit statistic applied to the copula function of the distribution of the process readings (Durante and Sempi, 2010). For  $p \geq 2$ , a *copula* is a  $p$ -dimensional distribution function on  $[0, 1]^p$  whose univariate marginal distributions are uniform on  $[0, 1]$ . Sklar's Theorem states that any  $p$ -dimensional distribution function is associated with a unique copula function. This copula can therefore be used to characterize the dependence between the components of a random vector. Let  $F$  be a  $p$ -dimensional distribution function with univariate margins  $F_1, F_2, \dots, F_p$ . Then there exists a copula  $C$  such that for all  $(x_1, x_2, \dots, x_p) \in \bar{\mathbf{R}}^p$ ,

$$F(x_1, x_2, \dots, x_p) = C(F_1(x_1), F_2(x_2), \dots, F_p(x_p)).$$

Furthermore,  $C$  is unique when  $F_1, F_2, \dots, F_p$  are all continuous. When  $F_1, F_2, \dots, F_p$  are continuous,  $C$  can be obtained by

$$C(u_1, u_2, \dots, u_p)$$

$$= F(F_1^{-1}(u_1), F_2^{-1}(u_2), \dots, F_p^{-1}(u_p)).$$

If the copula function of the process readings is that of a multivariate normal distribution and the marginal distributions are highly skewed gamma distributions, we have demonstrated that no additional quarantining greater than that for multivariate normal data is required (Table 2). However, for CR multivariate gamma data with the same marginal distributions but a different copula, additional quarantining is necessary. Dependent multivariate Cauchy data also requires a larger quarantine constant than multivariate normal data, but not as large as the dependent CR multivariate gamma data with  $\theta_0 = 1/2$ . The multivariate normal copula Goodness-of-Fit test adapted from Malevergne and Sornette (2003) can be used to differentiate between these three distributions, which suggests that it can be used to aid in determining if the recommended degree of quarantine is sufficient to achieve acceptable in control ARL.

Steps to compute Normal Copula Goodness-of-Fit test statistic for a sample of  $p \times 1$  iid random vectors,  $\mathbf{x}_1, \dots, \mathbf{x}_n$ :

1. Compute marginal normal scores for each component of  $\mathbf{x}_i$ . Using the formula from Royston (1982), for  $i = 1, \dots, n$  and  $j = 1, \dots, p$ ,

$$y_{ij} = \Phi^{-1} \left( \frac{r_i^{(j)} - \alpha}{n - 2\alpha + 1} \right),$$

where  $r_i^{(j)}$  is the rank of the  $j$ th component of  $\mathbf{x}_i$  among  $x_1^{(j)}, \dots, x_n^{(j)}$  and  $\alpha = 0.375$  as suggested by Royston (1982).

2. Compute sample covariance matrix of  $\mathbf{y}_i$ 's,  $\hat{\Sigma}_y$ .
3. Compute sample quadratic form deviates,

$$z_i^2 = \mathbf{y}_i' \hat{\Sigma}_y^{-1} \mathbf{y}_i.$$

If the copula of  $\mathbf{x}_1, \dots, \mathbf{x}_n$  is the multivariate normal copula, then  $z_1^2, \dots, z_n^2$  will be approximately  $\chi_p^2$  distributed.

4. Compute Anderson-Darling statistic to compare the distribution of  $z_1^2, \dots, z_n^2$  to the  $\chi_p^2$  distribution. Define

$$A_n^2 = -n - \sum_{k=1}^n \frac{2k-1}{n} \left[ \ln F_{\chi_p^2}(z_{(k)}^2) + \ln(1 - F_{\chi_p^2}(z_{(n+1-k)}^2)) \right],$$

where  $z_{(1)}^2, \dots, z_{(n)}^2$  are the order statistics of the sample  $z_1^2, \dots, z_n^2$ .

TABLE 4. Mean Goodness-of-Fit AD Test Statistic Values,  $A_n^2$ , Based on a Sample of  $n = 50$  Simulated Observations. Pairwise correlation between vector components is  $\rho = 0.9$ . Note that for the cases of uncorrelated data considered here, the random vector components are independent so the copula function is indistinguishable from the MVN copula regardless of the marginal distribution of the vector components. Nonparametric change point procedure achieves acceptable in-control ARL performance for cases above or to the left of dividing line, but does not achieve acceptable in-control ARL for cases below or to the right (see Table 2).

	$p = 2$	$p = 5$	$p = 10$
MVN	0.48	0.63	0.70
Transformed Gamma, $\theta_0 = 1/2$	0.48	0.63	0.70
MVT (df = 5)	0.56	1.42	2.71
MVT (Cauchy, df = 1)	1.81	9.48	16.79
CR Multivariate Gamma, $\theta_0 = 4$	1.23	5.89	8.00
CR Multivariate Gamma, $\theta_0 = 2$	2.16	11.39	16.42
CR Multivariate Gamma, $\theta_0 = 1/2$	4.69	28.75	46.65

From Tables 2 (in-control ARL simulation results) and 4, we can see that for a historical sample of size  $n = 50$  observations, the proposed procedure achieves acceptable in control ARL for all cases where the GOF statistic is less than 10 (except for CR multivariate gamma with  $\theta_0 = 4$ ,  $\rho = 0$ ,  $p = 10$  where the in control ARL is estimated to be 418 when the nominal value is 500), but the in control run length behavior suffers for distributions with values of the GOF statistic much larger than 10. If a historical sample size much greater than  $n = 50$  observations is available, then this diagnostic procedure would be conservative since the power of the Anderson-Darling test increases with  $n$ . For example, if a historical sample of  $n = 500$  observations following the CR MVG distribution with  $\theta_0 = 1/2$  and  $\rho = 0.9$ , the simulated mean GOF statistic is 907. Future developments may be able to improve the diagnostic procedure by reducing the dependence on  $n$ .

### Example: Analysis of Aluminum Smelter Data

The process of extracting aluminum metal begins with refining aluminum ore into alumina ( $\text{Al}_2\text{O}_3$ ).

The alumina is then reduced into metallic aluminum using an electrolysis process, which is referred to as aluminum smelting. Figure 3 shows a time series plot of a data set (kindly provided by Len Homer) measuring alumina ( $\text{Al}_2\text{O}_3$ ) content of a smelter feed along with several impurities: silica ( $\text{SiO}_2$ ), ferric oxide ( $\text{Fe}_2\text{O}_3$ ), magnesium oxide ( $\text{MgO}$ ), and calcium oxide ( $\text{CaO}$ ).

As one would expect with compositional data, the content of the compounds are negatively correlated, so a multivariate method will likely produce better results than monitoring the compound content using separate univariate charts. Zamba and Hawkins (2006) analyzed a similar data set using the change-point model for multivariate normal data. Those authors suggested transforming  $\text{SiO}_2$  to the log scale and replacing  $\text{CaO}$  with  $1/\text{CaO}$  because their procedure requires the assumption of joint normality of the compound compositions.

Although transformations exist to achieve approximate marginal normality, it is still difficult to verify that the joint distribution of the five compounds is multivariate normal. The problem of choosing a transformation to joint normality is made more difficult by the possible presence of a change-point. If a change-point is present, the desired transformation would be to a mixture of normal distributions, not a single normal distribution. We analyze the data on the original scale because our nonparametric change-point model does not require the assumption of joint normality.

We applied the change-point model based on the directional rank statistic with  $p = 5$ ,  $c = 15$  to the aluminum smelter data. Due to our choice of the quarantine constant,  $c$ , we began monitoring at observation  $n = 33$ . Figure 4 shows the change-point statistic,  $r_{\max,c,n}$ , and the control limit,  $h_{\alpha=0.002,p=5,c=15,n}$ , which corresponds to an in-control ARL of 500. At observation  $n = 44$ , the change-point statistic crosses the upper control limit. The corresponding estimate of the shift location is  $\hat{\tau} = 19$ . This illustrates that even though several “warm-up” observations are required before monitoring begins, the change point methodology can detect shifts in location vector that occur during the initial time period.

After the change-point model has signaled that a shift in location vector of the process readings has occurred, we must diagnose the signal to determine what characteristics of the process have shifted. This

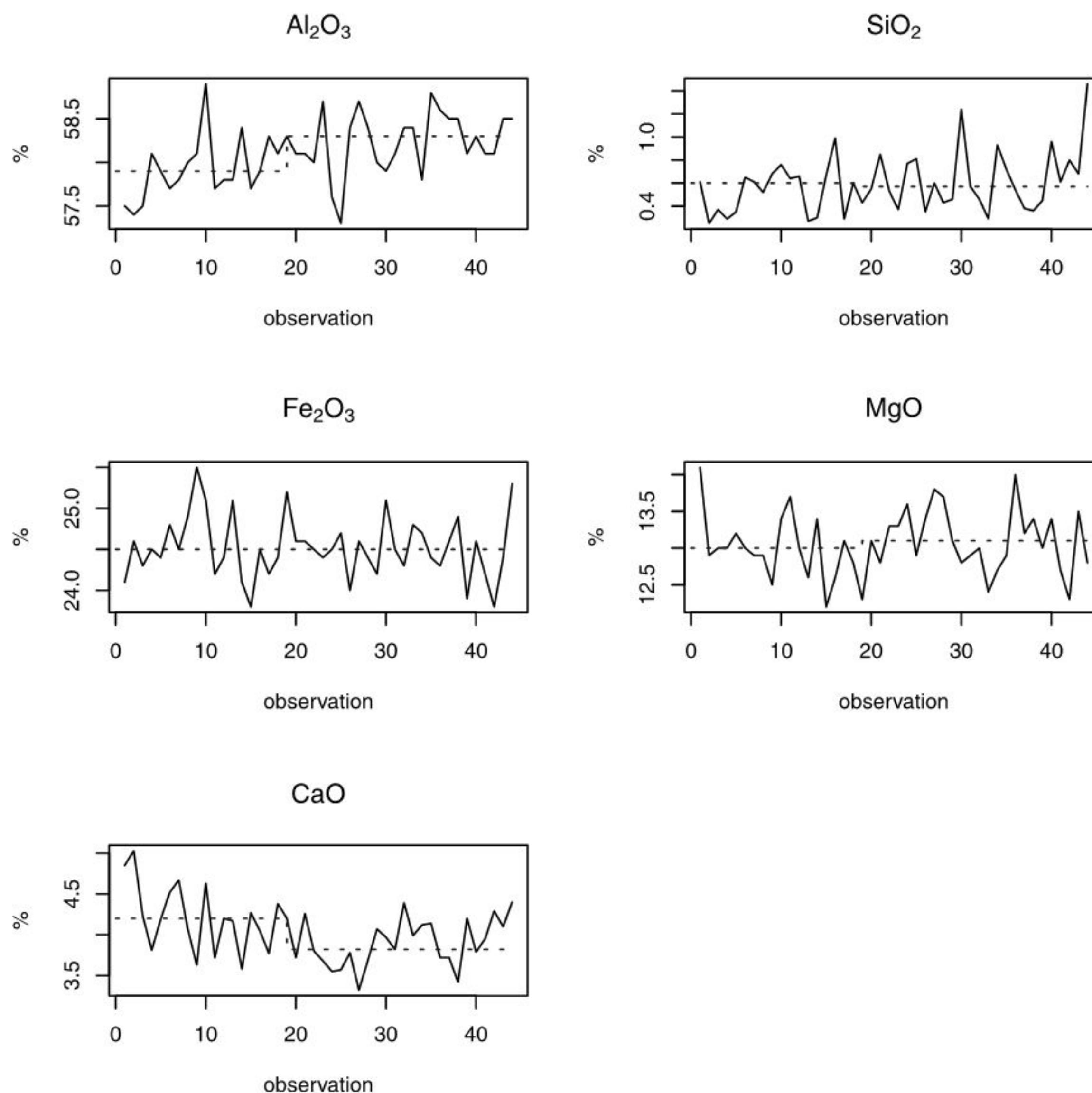


FIGURE 3. Time Series Plots of Each Component of Aluminum Smelter Feed versus Observation Number from First Observation through the First Signal. Center line represents sample mean before and after estimated shift location,  $\hat{\tau} = 19$ .

is a more challenging task for multivariate data than for univariate data. In the univariate setting, we only have to estimate the location and magnitude of the shift. However, in the multivariate setting one or more variables may exhibit a shift in location or a signal may result from a change in the relationship between the variables. Mason and Young (2002) and Hawkins (1993) discuss these issues in greater detail. We followed the post-shift diagnosis procedure of Zamba and Hawkins (2006), but replaced parametric test statistics with nonparametric analogues.

First, we conducted a set of univariate two-sample Wilcoxon-Mann-Whitney rank sum tests for differ-

ence in location for each component. Zamba and Hawkins (2006) suggest arranging the variables in ascending order of mean composition in the alumina and conducting regression adjusted step-down *t*-tests of each variable conditioned on those with lower mean concentration. We conducted a series of  $\chi^2$  tests against a restricted alternative (Ferguson, 1996) using the component-wise rank test statistic to test for a difference in each compound concentration after accounting for those with lower concentrations. For example, the restricted alternative test for CaO is  $\tilde{R}_{\hat{\tau},n}(\text{SiO}_2, \text{CaO}) - \tilde{R}_{\hat{\tau},n}(\text{SiO}_2)$ , where  $\tilde{R}_{\hat{\tau},n}(\text{SiO}_2, \text{CaO})$  denotes the bivariate component-wise rank test statistic computed using data vec-

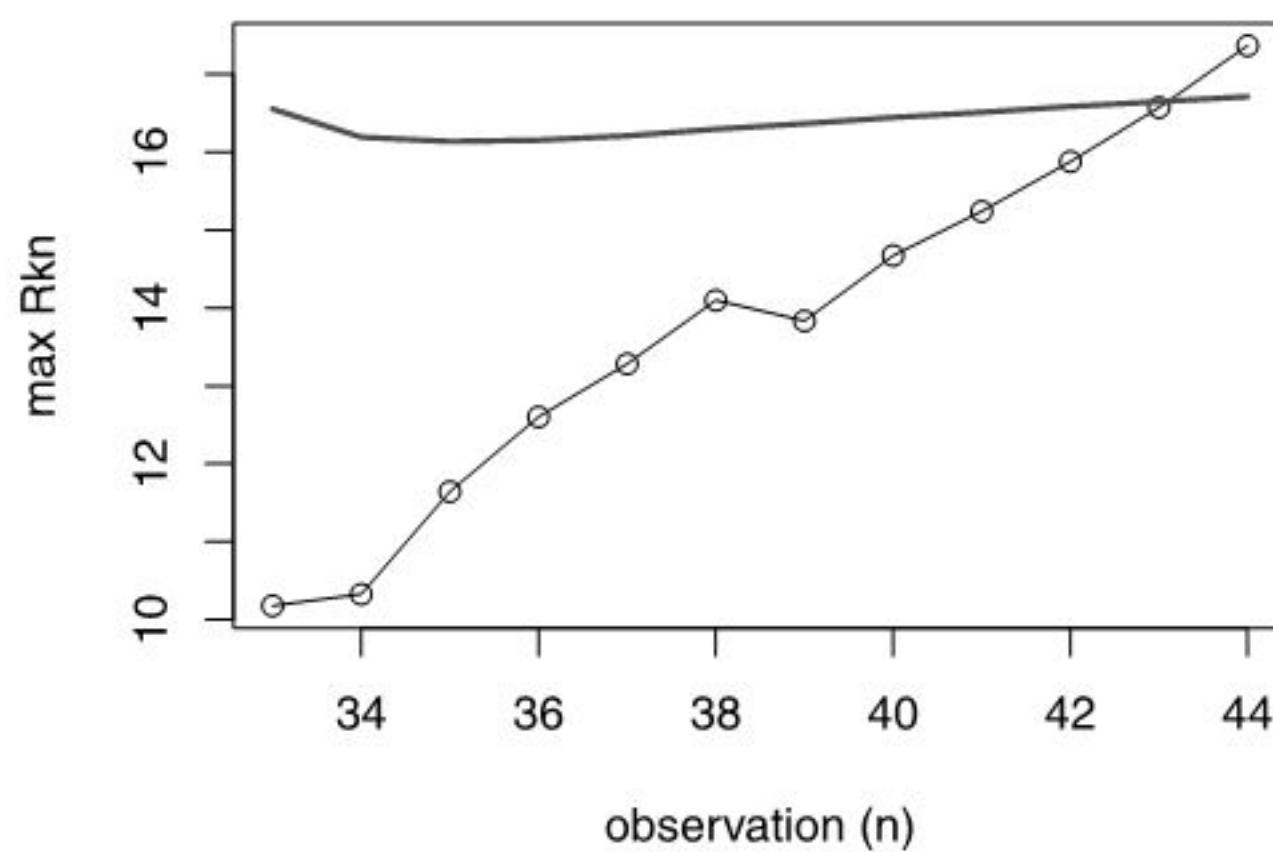


FIGURE 4. Analysis of Aluminum Smelter Data. Plot of control chart statistic  $r_{\max,c,n}$  and control limit  $h_{\alpha}=0.001, p=5, c=15, n$ .

tor components  $\text{SiO}_2$  and  $\text{CaO}$  and  $\tilde{R}_{\hat{\tau},n}(\text{SiO}_2)$  denotes the univariate component-wise rank test statistic computed using only  $\text{SiO}_2$ . The asymptotic null distribution of the difference of the test statistics is  $\chi^2$ . We used the asymptotic null distribution to obtain approximate  $p$ -values for this series of tests.

The univariate Wilcoxon-Mann-Whitney tests (Hollander and Wolfe, 1999) suggest that there was a decrease in the median  $\text{MgO}$  concentration and an increase in the median  $\text{Al}_2\text{O}_3$  concentration. The multivariate restricted alternative  $\chi^2$  tests show that  $\text{CaO}$  increased in concentration relative to the lower concentration impurities, and that the increase in median concentration of alumina ( $\text{Al}_2\text{O}_3$ ) is significant relative to the concentration of the impurities in the smelter feed.

## Conclusion

Existing multivariate phase-II SPC procedures require the user to make assumptions about the in-

control or out-of-control distributions, or to collect a large phase-I training sample before beginning to monitor the process. We proposed a multivariate phase-II SPC procedure that does not require the user to assume that process readings follow the normal distribution, to have knowledge of in- or out-of-control distribution parameters, or to collect a large phase-I training sample before monitoring the process for a shift in distribution. Because the multivariate test statistic is not distribution-free for small to moderate sample sizes, we suggest the use of a quarantine period. Imposing a quarantine only adversely affects the performance for detecting large shifts. This is an acceptable penalty to incur, because we designed the procedure for use in detecting small to moderate shifts that would not be detected when only a few post-shift observations have been collected and because nonparametric procedures are inherently less powerful than parametric for detecting large shifts in small samples. The nonparametric procedure proposed here detected small to moderately sized shifts faster on average than the parametric method for MVN data, as observed in the univariate case by Hawkins and Deng (2010).

While the MEWMA chart has been considered to be robust to non-normality, we have shown that for distributions with dependence structures that differ greatly from multivariate normal the MEWMA does not demonstrate acceptable in-control ARL performance. For multivariate normal data, the MEWMA chart detects both large and small shifts more rapidly than the proposed procedure when in-control parameters are estimated from a small phase I data set. However, the MEWMA chart based on estimated in-control parameters only achieves targeted in-control ARL averaged across data sets, while the proposed method achieves desired in-control ARL behavior for each data set to which it is applied.

TABLE 5. Diagnosis After Signal. Univariate Wilcoxon-Mann-Whitney (WMW) tests computed following Hollander and Wolfe (1999), “step-down”  $\chi^2$  test refers to multivariate restricted alternative  $\chi^2$  test, and  $\tilde{x}_{m_1, \dots, m_2} = \text{median}\{x_{m_1}, \dots, x_{m_2}\}$

Analyte	$\tilde{x}_{1, \dots, \hat{\tau}}$	$\tilde{x}_{\hat{\tau}+1, \dots, n}$	WMW	$p$	“step-down” $\chi^2$	$p$
$\text{SiO}_2$	0.60	0.57	179.5	0.17	—	—
$\text{CaO}$	4.20	3.82	341	0.01	7.16	0.007
$\text{MgO}$	13.00	13.10	193.5	0.30	0.11	0.74
$\text{Fe}_2\text{O}_3$	24.50	24.50	241.5	0.92	0.23	0.63
$\text{Al}_2\text{O}_3$	57.90	58.30	120.5	0.005	5.77	0.02

High dimensional data may pose a limitation for use of the proposed method, because a larger quarantine value may be required which could limit the utility of the method for detecting moderate sized shifts. However, common applications of multivariate SPC do not have data dimension much larger than studied here. Furthermore, large data sets would be required to perform multivariate SPC in general in very high dimensions and the proposed method is best suited for applications when a limited amount of historical data is available.

## Appendix: Computationally Efficient Algorithm

We enumerate the steps to compute the test statistic  $r_{k,n+1}$  for all  $k = 1, \dots, n$  when observation  $\mathbf{x}_{n+1}$  is added to the data set using a method for which computation time scales linearly in  $n$ . For this method, we store and update the centered rank vectors,  $\mathbf{R}_n(\mathbf{x}_i)$  each time a new observation is acquired. This requires storage of a  $p \times n$  matrix, which is not a restrictive requirement for realistic values of  $n$  and  $p$ .

### Update procedure

1. For  $i = 1, \dots, n$ , update

$$\begin{aligned}\mathbf{R}_{n+1}(\mathbf{x}_i) &= \sum_{j=1}^{n+1} \mathbf{h}(\mathbf{x}_i, \mathbf{x}_j) \\ &= \left( \sum_{j=1}^n \mathbf{h}(\mathbf{x}_i, \mathbf{x}_j) \right) + \mathbf{h}(\mathbf{x}_i, \mathbf{x}_{n+1}) \\ &= \mathbf{R}_n(\mathbf{x}_i) + \mathbf{h}(\mathbf{x}_i, \mathbf{x}_{n+1}).\end{aligned}$$

2. Compute

$$\mathbf{R}_{n+1}(\mathbf{x}_{n+1}) = \sum_{j=1}^{n+1} \mathbf{h}(\mathbf{x}_{n+1}, \mathbf{x}_j).$$

3. Compute

$$\hat{\Sigma}_{n+1} = \frac{1}{n} \sum_{i=1}^{n+1} \mathbf{R}_{n+1}(\mathbf{x}_i) \mathbf{R}_{n+1}(\mathbf{x}_i)' \text{ and } \hat{\Sigma}_{n+1}^{-1}.$$

4. For  $k = 1$ , compute

$$\bar{\mathbf{r}}_{n+1}^{(1)} = \mathbf{R}_{n+1}(\mathbf{x}_1).$$

For  $2 \leq k \leq n$ , compute

$$\bar{\mathbf{r}}_{n+1}^{(k)} = \frac{1}{k} \sum_{i=1}^k \mathbf{R}_{n+1}(\mathbf{x}_i)$$

$$\begin{aligned}&= \frac{1}{k} \left[ \left( \sum_{i=1}^{k-1} \mathbf{R}_{n+1}(\mathbf{x}_i) \right) + \mathbf{R}_{n+1}(\mathbf{x}_k) \right] \\ &= \frac{(k-1)}{k} \bar{\mathbf{r}}_{n+1}^{(k-1)} + \mathbf{R}_{n+1}(\mathbf{x}_k).\end{aligned}$$

5. Compute

$$\hat{\Sigma}_{k,n+1}^{-1} = \frac{(n+1)k}{(n+1-k)} \hat{\Sigma}_{n+1}^{-1}.$$

6. Compute

$$r_{k,n+1} = \bar{\mathbf{r}}_{n+1}^{(k)'} \hat{\Sigma}_{k,n+1}^{-1} \bar{\mathbf{r}}_{n+1}^{(k)}.$$

Steps 1–3 only need to be executed once each time a new observation is added to the data set. All of the operations in these steps scale linearly in  $n$ . We use the unpooled covariance estimator because the inverse of the covariance matrix only needs to be computed once to obtain test statistics for each possible split point  $k = 1, \dots, n-1$ , and within-group rank vectors do not need to be computed. Steps 4–6 need to be executed for each  $k = 1, \dots, n$  each time a new observation,  $\mathbf{x}_{n+1}$ , is collected. For each  $k$ , the number of operations necessary to execute steps 4–6 once does not depend on  $n$  or  $k$ . Thus, the entire update procedure scales linearly in  $n$ . Because the computation time scales linearly in  $n$ , we were able to employ large scale Monte Carlo simulation studies to obtain control limits and evaluate the performance of the proposed procedure.

## References

- CHOI, K. and MARDEN, J. (1997). "An Approach to Multivariate Rank Tests in Multivariate Analysis of Variance". *Journal of the American Statistical Association* 92(440), pp. 1581–1590.
- CROSIER, R. B. (1988). "Multivariate Generalizations of Cumulative Sum Quality Control Schemes". *Technometrics* 30(3), pp. 291–303.
- DURANTE, F. and SEMPI, C. (2010). "Copula Theory: An Introduction". In *Copula Theory and Its Applications*, Jaworski, P.; Durante, F.; Hardle, W.; and Rychlik, T., eds., vol. 198 of *Lecture Notes in Statistics*, pp. 3–31. Berlin, Heidelberg: Springer-Verlag.
- FERGUSON, T. S. (1996). *A Course in Large Sample Theory*. Boca Raton, FL: Chapman and Hall.
- HAWKINS, D. M. (1993). "Regression Adjustment for Variables in Multivariate Quality Control". *Journal of Quality Technology* 25, pp. 170–182.
- HAWKINS, D. M. and DENG, Q. (2010). "A Nonparametric Change-Point Control Chart". *Journal of Quality Technology* 42(2), pp. 165–173.
- HAWKINS, D. M.; QIU, P.; and KANG, C. W. (2003). "The Changepoint Model for Statistical Process Control". *Journal of Quality Technology* 35(4), pp. 355–366.
- HOLLAND, M. D. (2013). *NPMVCP: Nonparametric Multivariate Change Point Model*. R package version 1.0.

- HOLLANDER, M. and WOLFE, D. A. (1999). *Nonparametric Statistical Methods*, 2nd edition. New York: John Wiley & Sons.
- HOTELLING, H. (1947). "Multivariate Quality Control—Illustrated by the Air Testing of Bombsights". In *Techniques of Statistical Analysis*, Eisenhart, M. W.; Hastay, M. W.; and Wallis, W. A., eds., pp. 111–184. New York: McGraw Hill.
- KOTZ, S.; BALAKRISHNAN, N.; and JOHNSON, N. L. (2000). *Continuous Multivariate Distributions*, 2nd edition. New York: John Wiley & Sons.
- LOWRY, C. A.; WOODALL, W. H.; CHAMP, C. W.; and RIGDON, S. E. (1992). "A Multivariate Exponentially Weighted Moving Average Control Chart". *Technometrics* 34(1), pp. 46–53.
- MALEVERGNE, Y. and SORNETTE, D. (2003). "Testing the Gaussian Copula Hypothesis for Financial Assets Dependences". *Quantitative Finance* 3, pp. 231–250.
- MASON, R. L. and YOUNG, J. C. (2002). *Multivariate Statistical Process Control With Industrial Application*. Philadelphia, PA: SIAM.
- QIU, P. (2008). "Distribution-Free Multivariate Process Control Based on Log-Linear Modeling". *IEEE Transactions* 40(7), pp. 664–677.
- QIU, P. and HAWKINS, D. M. (2001). "A Rank-Based Multivariate CUSUM Procedure". *Technometrics* 43(2), pp. 120–132.
- QIU, P. and HAWKINS, D. M. (2003). "A Nonparametric Multivariate Cumulative Sum Procedure for Detecting Shifts in All Directions". *The Statistician* 52(2), pp. 151–164.
- ROYSTON, J. P. (1982). "Algorithm AS 177: Expected Normal Order Statistics (Exact and Approximate)". *Journal of the Royal Statistical Society, Series C: Applied Statistics* 31(2), pp. 161–165.
- STOUMBOS, Z. G. and SULLIVAN, J. H. (2002). "Robustness to Non-Normality of the Multivariate EWMA Control Chart". *Journal of Quality Technology* 34(3), pp. 260–276.
- TESTIK, M. C.; RUNGER, G. C.; and BORROR, C. M. (2003). "Robustness Properties of Multivariate EWMA Control Charts". *Quality and Reliability Engineering International* 19, pp. 31–38.
- ZAMBA, K. D. and HAWKINS, D. M. (2006). "A Multivariate Change-Point Model for Statistical Process Control". *Technometrics* 48(4), pp. 539–549.

