



HAL
open science

Performance Evaluation Based on the Robust Mahalanobis Distance and Multilevel Modelling Using Two New Strategies

Ghazi Shukur, Shakir Hussain, Mohammed Mohammed, Roger Holder,
Abdullah Almasri

► **To cite this version:**

Ghazi Shukur, Shakir Hussain, Mohammed Mohammed, Roger Holder, Abdullah Almasri. Performance Evaluation Based on the Robust Mahalanobis Distance and Multilevel Modelling Using Two New Strategies. *Communications in Statistics - Simulation and Computation*, 2008, 37 (10), pp.1966-1980. 10.1080/03610910802311692 . hal-00514331

HAL Id: hal-00514331

<https://hal.science/hal-00514331>

Submitted on 2 Sep 2010

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Performance Evaluation Based on the Robust Mahalanobis Distance and Multilevel Modelling Using Two New Strategies

Journal:	<i>Communications in Statistics - Simulation and Computation</i>
Manuscript ID:	LSSP-2008-0045.R1
Manuscript Type:	Original Paper
Date Submitted by the Author:	30-Jun-2008
Complete List of Authors:	shukur, ghazi; Economics and Statistics, School of Management and Economics; Jönköping (Jonkoping) International Business School, Economics and Statistics Hussain, Shakir; School of Medicine, University of Birmingham, Division of Primary Care and General Practice Mohammed, Mohammed; University of Birmingham, Department of Public Health Holder, Roger; School of Medicine, University of Birmingham, Division of Primary Care and General Practice Almasri, Abdullah; Karlstad University, Department of Economics and Statistics
Keywords:	Ranking indicators, performance evaluation, robust statistics, multilevel estimation
Abstract:	In this paper we propose a general framework for performance evaluation of organisations and individuals over time using routinely collected performance variables or indicators. Two new double robust and model-free strategies are used for evaluation (ranking) of sampling units. Strategy 1 can handle missing data using (RML) at stage two, while strategy two handle missing data at stage one. Strategy 2 has the advantage that overcomes multicollinearity problem. Strategy one requires independent indicators for the construction of the distances, where strategy two does not. Two different domain examples are used to illustrate the application of

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60



For Peer Review Only

Performance Evaluation Based on the Robust Mahalanobis Distance and Multilevel Modelling Using Two New Strategies

By

S. Hussain¹; M. A. Mohamed²; R. Holder³; A. Almasri⁴; and G. Shukur⁵.

^{1&3} Division of Primary Care and General Practice, School of Medicine, University of Birmingham,

² Department of Public Health, University of Birmingham, UK

⁴ Department of Economics and Statistics, Karlstad University, Sweden

⁵ Department of Economics and Statistics, Jönköping University, Sweden, and Centre for Labour
Market Policy (CAFO), Department of Economics and Statistics, Växjö University, Sweden

Abstract.

In this paper we propose a general framework for performance evaluation of organisations and individuals over time using routinely collected performance variables or indicators. Such variables or indicators are often correlated over time, with missing observations, and often come from heavy tailed distributions shaped by outliers. Two new double robust and model-free strategies are used for evaluation (ranking) of sampling units. Strategy 1 can handle missing data using residual maximum likelihood (RML) at stage two, while strategy two handle missing data at stage one. Strategy 2 has the advantage that overcomes the problem of multicollinearity. Strategy one requires independent indicators for the construction of the distances, where strategy two does not. Two different domain examples are used to illustrate the application of the two strategies. Example one considers performance monitoring of gynaecologists and example two considers the performance of industrial firms.

Key Words: Ranking indicators, performance, robust statistics, multilevel estimation, Mahalanobis distance.

1. Introduction

The development of performance or ranking indicators has a wide range of application in several disciplines. It constitutes a central topic in (e.g., economics, health, sociology, education) and modern formalized modeling has dealt with it for some considerable time. The formalized modeling of this topic has led the way towards statistical application and evaluation, and it is such application and evaluation that form the main theme of this paper. At this stage it is important to stress that each and every area of application has its own characteristics regarding the nature and definition of the problem and the observations and data collected thereby.

1
2
3 In the near past, numerous studies in many areas, starting with education; economics; health;
4 and other social sciences, have been conducted in the developing of performance indicators
5 where quantitative comparisons between institutions have been developed to introduce
6 efficiency among the activities in those areas. Many of these studies, however, did not pay
7 attention to possible complexities associated with real data like, for example missing values
8 and the nature of this missingness; heavy tailed data as a result of outliers; among the
9 *independent* variables (note that in many cases there is no dependent variable involved in the
10 analysis); possible linear trends. However, Goldstein and Spiegelhalter (1996) addressed in
11 some detail the need to take account of model based uncertainty in making comparisons, to
12 establish appropriate measures of institutional outcomes and base-line measures, and to
13 exercise care and sensitivity when interpreting apparent differences.
14
15
16
17
18
19
20
21
22
23

24 Among the most recent studies on performance are those of Harley et al (2005) in the area of
25 health and Behrenz and Althin (2005) in the area of labour market and unemployment. In the
26 first paper, the authors have used a two stages statistical method for evaluating health
27 indicators using routine hospital data to identify gynecologists' performance. The robust
28 Mahalanobis distance is first used to reduce the dimension of the indicators, after which the
29 robust residuals from the level-two estimation were investigated to indicate the outliers. In the
30 second paper the authors studied the efficiency and productivity of employment offices in
31 Sweden. They applied the so-called Malmquist productivity index, which measures the
32 distance from every office to the production frontier and the shifting of the frontier over time.
33 Note that, in the second paper, the analysis has been done in the micro level and no
34 consideration have been taken for the possible existence of any hierarchy structure in the data,
35 while in the first paper possible hierarchy is taken into account.
36
37
38
39
40
41
42
43
44
45
46

47 However, increasing attention is nowadays paid to multilevel modelling especially in
48 modelling dynamic relationships of hierarchical structures. The purpose of this study is to
49 provide two strategies, for performance evaluation of individuals and organisations over time,
50 that take into consideration the problems mentioned associated with real data which other
51 previously studies did not. Here we combine robust multivariate methods together with
52 multilevel estimation methods that allow for missing observations among the data. By
53 following such strategies one can avoid inadequate methods that might lead to extremely
54 misleading results and inferences. To demonstrate these strategies, we use two examples from
55 two different areas. The first example is regarding evaluating health indicators to assist health
56
57
58
59
60

1
2
3 authorities in decision-making, while the second example evaluates the performances of
4 industrial firms. These examples illustrate how our approach can disentangle issues regarding
5 producing performance indicators when the data under consideration suffer from missing
6 values, outliers, are fat tailed, and different levels of multicollinearity in the main structure.
7 Moreover, it is important to mention that one of the crucial advantages of our approach is that
8 our ranking strategies are totally model-free in the sense that we do not specify a model of
9 dependent and independent variables. What our technique needs is only available necessary
10 performance indicators or variables. Processing in this manner, we avoid problems of
11 specification and misspecifications of models. This makes our approach non-comparable to
12 other performance or ranking methods such as those mentioned above or the GLM or the
13 Mixed procedures in SAS. According to our knowledge, no such methods are yet available in
14 the literatures. In short, our approach add the following important contributions to the readily
15 available methods in the sense that it does not require any model specification, it is double
16 robust, it captures dynamics of performances over time and it is applicable on highly collinear
17 or independent indicators or variables.
18
19
20
21
22
23
24
25
26
27
28
29
30

31 The rest of the paper is organised as follows: Section 2 presents a description of the data and
32 the theoretical assumptions; Section 3 describes the methodology and estimation procedures;
33 while the results are presented in Section 4. Section 5, finally, gives conclusions and a
34 summary of findings.
35
36
37
38
39

40 2. Data Description

41 In this section we describe the nature of the two data sets that we use to illustrate our
42 modelling strategy.
43
44
45
46
47

48 Example one:

49 The data set here consists of routine hospital data of seven clinically relevant indicator
50 variables from hospital episode statistics for 143 gynaecologists. These indicators are: % of
51 finished gynaecologist episodes with complications; mean length of spell (in days); % of
52 finished gynaecologist episodes with more than two operations; % of finished gynaecologist
53 episodes where spell is longer than episode; % of finished gynaecologist episodes for
54 dilatation and curettage on women aged less than 40 years; % of finished gynaecologist
55 episodes for sterilisation on women aged less than 25 years; and % of finished gynaecologist
56
57
58
59
60

1
2
3 episodes for hysterectomy on women aged less than 30 years. The data are collected over five
4 years period (from 1991:2 to 1995:6, but only 68 gynaecologists appear in all five years. High
5 proportions may potentially indicate sub-optimal performance, although this requires further
6 study.¹
7
8
9

10 11 **Example two:**

12 The dataset contains all Swedish industrial firms that are registered on the stock market and
13 with available accounts during the recent five years (1999 to 2003). Data about 2004 are
14 unfortunately not completely recorded yet. With help of the database OSIRIS, 57 firms were
15 identified.² The following seven variables were formulated: cost of goods sold (COGS) to
16 sales; free cash flow to current liabilities; growth in gross investment; growth in net turnover;
17 growth in total debt; R&D to sales and cost of employees.
18
19
20
21
22
23
24
25

26 *COGS to sales* is a profitability ratio where COGS can be seen as the cost of doing business,
27 for example the cost of raw materials. The variable is weighted by sales and measured in
28 percentage. *Free cash flow to current liabilities* is a liquidity ratio. Free cash flow is how
29 much cash a firm has after paying its bills for ongoing activities and growth. *Growth in gross*
30 *investment* measures the growth in capital goods; the variable is measured in percentage
31 terms. *Growth in net turnover* measures the growth in net income or revenue from the sale of
32 goods and services, in percentage. *Growth in total debt* is a variable that measure the growth
33 in percentage of the total debt, i.e. the sum of short-term and long-term debt. *R&D to sales* is
34 a variable in percentage terms that quantify the relative importance of R&D investments in
35 the firm.
36
37
38
39
40
41
42
43
44
45

46 **3. Methodology and estimation procedures**

47 This paper is designed to provide frameworks that encompass two stages statistical methods
48 that evaluate health and firm indicators thereby assisting health- and other authorities in
49 decision-making. Two stages of data reductions are proposed using two strategies. These
50 strategies consist of two major components, namely the multivariate Mahalanobis distance
51 (MD) and the level two standardised residuals estimates from the multilevel modelling. In the
52 following we present a brief description of these two components.
53
54
55
56
57
58
59

60 ¹ The data however is described in more details in Harley et al (2005).

² OSIRIS is a comprehensive database of listed companies, banks and insurance companies around the world.

Suppose that $x = (x_1, x_2, \dots, x_n)^T$ is a set of n observations on p random variables. The classical MD is defined as:

$$MD(x_i) = \sqrt{(x_i - \mu)V^{-1}(x_i - \mu)^t} \quad (1)$$

Where μ is the arithmetic mean vector and V is the covariance matrix. The classical squared Mahalanobis distance $(MD)^2$ is not ideally suited to multivariate outlier detection because it is not resistant to outliers. Rousseeuw and Leroy (1987) recommend using distance based on robust estimators of multivariate location and scatter (μ_R, V_R) to avoid masking effect. A cutoff point of $\sqrt{\chi_{p,0.975}^2}$, (p is the number of the variables), used to determine points above as outliers. The minimum covariance determinant (MCD) method of Rousseeuw (1985) aims to find h observations out of (n) whose covariance matrix (V_R) has the lowest determinant. The development of MCD is mentioned in Rousseeuw and Van Driessen (1999) under the name Fast MCD algorithm, where lower determinant of MCD can be approximated from the initial MCD.

Assuming that the fraction of outliers is at most α , ($0 < \alpha \leq 1/2$), e.g. 50%, then α can be chosen to equal $\chi_{p,0.50}^2$ where, except for the extreme values cases, we expect the majority of the data to come from a normal distribution. Let the halve contain $h = (n + p + 1) / 2$ observations with (n) total sample size and (p) number of variables, however the determinants of the covariance matrix (V_R) will be minimized subject to the inequality

$$\{i, (x_i - \mu_R)V_R^{-1}(x_i - \mu_R)^t \leq \alpha^2\} \geq h \quad (2)$$

Finally the robust MCD distance can be written as

$$RMD(x_i) = \sqrt{(x_i - \mu_R)V_R^{-1}(x_i - \mu_R)^t} \quad (3)$$

Where μ_R is now our first moment vector and V_R is the robust covariance matrix. Rocke and Woodruff (1996) proposed the robust M estimator that uses the fast minimum covariance determinant estimator as an initial robust estimate then the estimate refines with M iterations using the translated bi-weight function that is described in Rocke (1996). In computing the Robust Distance we use either the Fast MCD or the M estimator that is available in the robust library of S+ 8.

The data we study here consists of repeated performance measures of random samples of item (gynaecologists / firms). We start by writing a simple empty model

$$y_{ij} = \bar{\beta} + \beta_{0j} + e_{ij} \quad (4)$$

where y_{ij} denotes the MD of the i^{th} year of the j^{th} gynaecologist or firm, $\bar{\beta}$ represent the mean MD distance, β_{0j} is a random variable representing “between-items” variability, and e_{ij} is a random variable representing “within- items” variability. The distributions of the random variables are

$$\beta_{0j} \sim N(0, \sigma_b^2) \quad e_{ij} \sim N(0, \sigma_e^2), \quad (5)$$

where σ_b^2 and σ_e^2 are the residuals variances of the between items (level two) and within items (level one) effects, respectively. The so-called Huber/White or sandwich estimator is then used to obtain robust tests and confidence intervals by correcting the asymptotic standard errors.

Langford and Lewis (1998) propose downward residual checking starting from the heights level and continuing with each next lower level for the purpose of outlier inspection and ranking. Here, we consider both level one residuals, e_{ij} and level two residuals β_{0j} for this purpose.

The underlying assumptions of the model in formulas 4 and 5 suffer from lack of robustness against outlying observations since both residuals are based on the Gaussian distribution, Seltzer and Choi (2003). According to Pinheiro et al (2001), an interesting feature of mixed-effects models is that outlier may occur either at the level of the within subject error e_{ij} called e-outliers, or at the level at random effect β_{0j} called β -outliers. In this paper our concern is the level two subjects β -outliers. Empirical Bayesian (EB) method is used to predict level two residuals. One disadvantage is its strong dependence on the model assumptions and the other disadvantage is when the number of sampling units in the level two is small, the EB approach in this case can result in underestimating of uncertainty.

The Fully Bayesian (FB) approach bases inferences on the marginal posterior distribution of all the parameters in the model and the 0.025 and 0.975 quantiles of this distribution would provide the Bayesian analogue of a 95% confidence intervals, see Seltzer and Choi (2003).

Using the scale mixture of normal representation presented by Seltzer and Choi (2003), the normal representation of the t distribution with mean 0, scale parameter equals to 1 and γ degrees of freedom, $t(0,1,\gamma)$ can be expressed by $\frac{Z}{\omega^{1/2}}$. Z here is standard normal with mean 0 and variance 1 and (ω) is chi-squared distributed divided by its degrees of freedom and the distribution of $\chi^2_\gamma / \gamma = \text{Gamma}(a,b)$ where $a = \frac{\gamma}{2}$ and $b = \frac{\gamma}{2}$, see Gelman et al (1995). Using the above argument and if we assume that the random effects has $t(0, \sigma_b^2, \gamma)$ then β_{0j} in (2) has the form:

$$\beta_{0j} \sim N\left(0, \frac{\sigma_b^2}{s_i}\right), \text{ where } s_i \sim \text{Gamma}(\gamma/2, \gamma/2). \quad (6)$$

For the level one residuals, e_{ij} , if we assume $t(0, \sigma_e^2, \gamma_1)$ then the definition under normality assumption is

$$y_{ij} \sim N\left(\bar{\beta}, \frac{\sigma_e^2}{r_{ij}}\right), \text{ where } r_{ij} \sim \text{Gamma}(\gamma_1/2, \gamma_1/2). \quad (7)$$

Note that the mean $\bar{\beta}$ is a constant that may take any value based on the data we analyse.

Here we will conduct a simple MCMC simulation study to show the difference between normal and the scale mixture (see Peel and McLachlan, 2000).

The construction of the proper normal prior for the above simulation study will consider minimally informative prior for $\bar{\beta}$ which we choose to be

$$\bar{\beta} \sim N(0, 1*10^6). \quad (8)$$

This type of normal prior adds no information to the data since the data have a range of {0.88, 12.72}.

Now, in the first stage of strategy one, a robust multivariate MD (RMD) is computed from several independent indicators and assigned to each subject over the years. These distances are used in the second stage as the outcome in a multilevel model (level 1: design variable for distances over time, level 2: design variable for sampling units) to obtain ranks of the units. In other words, the rank from the second level robust standardized residual of a multilevel model

1
2
3 of the repeated RMD is computed for each individual or firm. This can be done according to
4 the following:
5

- 6
- 7 1- Compute the RMD for all the variables over the years.
- 8
- 9 2- Use these RMD as the outcome in the robust multilevel model (level one is the repeated
10 RMD while level two is the gynaecologist).
- 11
- 12 3- Obtain the rank from the robust level two standardised residuals.
- 13
- 14 4- Compute the uncertainty of each rank by sampling from the posterior distribution of the
15 MCMC.
16
17
18

19
20
21
22 For strategy two, in the first stage and to avoid possible multicollinearity, we consider one
23 variable at a time. For each variable, we apply the multilevel modelling to achieve a reduction
24 over time to one single point. In order to obtain robust standardised residuals for a specific
25 gynaecologist or firm, the level two subjects β -outliers will be the outcome of this
26 modelling. We then repeat this procedure for the rest of the variables and obtain as many
27 points as the number of variables. In the second stage, we compute a RMD to obtain the rank
28 of those units. Processing in this manner we avoid the problem of dealing with dependency
29 among the variables. Strategy two can be summarised according to the following:
30
31
32
33
34
35

- 36
- 37
- 38
- 39 1- Compute the robust level two residuals for each indicator or variable separately over the
40 years.
- 41
- 42 2- Use these level two residuals in the second stage to compute the RMD.
- 43
- 44 3- Obtain the rank from the RMD.
- 45
- 46 4- Compute the uncertainty of each rank by sampling from the posterior distribution of the
47 MCMC.
48
49
50
51
52
53

54 Note that, one can also apply the idea of strategy one to dependent data but MD requires
55 independent indicators for the construction of the distances. However, robust principle
56 components can help us to decide which strategy to use by detecting independence in the data
57 structure. In what follows, we confine ourselves to brief descriptions of the methodological
58 issues used in this study, and further details are found in cited references.
59
60

3.1. Robust Multivariate Outliers Detection

The classical squared Mahalanobis Distance (MD)² is one approach to multivariate outlier detection based on arithmetic means and covariance matrix, MD measures how far a random vector is from the middle of its distribution. This will provide a reasonable summary measure of the distance of each item (individual or firm) from the mean. Points in multivariate data with large MD and greater than $\sqrt{\chi^2_{p,0.95}}$ are approximately considered outliers, where p here denotes the number of variables with 0.95 χ^2 quantile (see Rousseeuw and Leroy, 1987).

Note that occasionally one might be faced with cases where one outlier is too extreme that it may mask other outliers as a result. To overcome this problem in MD, Rousseeuw and Leroy (1987) proposed a robust estimation of covariance (M-type robust of covariance). If the covariance matrix for example is not estimated robustly then the underlying structured parameters are not robust. Using a method called “C-step”, Rousseeuw and Van Driessen (1999) developed a fast algorithm for Minimum Covariance Determinant (MCD) to approximate the minimum covariance determinant estimator of the MD (see the S+ 8 Robust Library (2007) for use of MCD algorithm and other algorithms).

3.2. MANOVA and Multilevel Model

The data collected in this study are unbalanced repeated measures over time. For strategy one, for example, when modelling the outcome measure MD, this repeated measurement data could be viewed as multilevel data with repeated MD nested within subjects. This leads to a two level model with the series of repeated MD as the lowest level and the individual subjects as the highest level. Multivariate ANalysis Of VAriance (MANOVA) may be used to model the MD data. The advantages of MANOVA are that no assumptions about the covariance matrix need to be made and that under normality assumption it yields exact tests. The assumptions are related to independent and identical distributions within treatment groups, homoscedasticity between groups, multivariate normality and complete data. When group sizes are different MANOVA will suffer from lack of uniqueness, Searle, Casella and McCulloch (1992). The assumption of homoscedasticity is unlikely to occur in the structure of the data that we are dealing with. However, since a group of subjects may show more variation over time than other subjects, this variation may cause the observation to be an outlier. This example shows the violation of the homoscedasticity assumption.

1
2
3 The multilevel model in this study is composed of three components; first, the fixed part that
4 represents the fixed effect of the intercept and the trend, second, the random effect at level
5 one, and finally we have the random effect at level two. The measurement occasions are
6 nested within subjects, level one units are occasions and level two units are subjects. The
7 random coefficient approach to repeated measures is usually based on polynomial trend to a
8 model that is a polynomial function of time. Depending on the research topic one may use any
9 other linearly independent set of functions. At this stage it is important to mention that there
10 are two estimation procedures for these purposes, namely, the Maximum Likelihood (ML)
11 and the Residual Maximum RML. RML estimation takes into account the loss of degree of
12 freedom resulting from the estimation of the parameters of the fixed part and also allows for
13 missing values. Hence, in this paper, we are considering the RML as our estimation method.
14
15
16
17
18
19
20
21
22
23

24 **4. Results**

25
26
27 In this section we present our main results for our methodology for performance evaluation
28 and giving our two empirical examples. For strategy one, when we use the assumption of
29 normality and in the case of almost independent or weak dependent data, we start the analysis
30 by obtaining the MD that we then use as the outcome in a multilevel model to obtain ranks of
31 sampling units, (the outcomes from different time point are unbalanced). On the other hand,
32 and when some dependency is existing among the variables, we start by calculating the
33 principle components (to detect independence) and then we obtain the RMD that we then use
34 as the outcome in a multilevel model using RML to obtain ranks of sampling units.
35
36
37
38
39
40
41

42 Strategy two starts with the multilevel modelling using RML first and then the RMD in the
43 case when normality is assumed, while we use the robust counterpart in the situation of fat
44 tailed distribution. The only difference between this strategy and strategy one is that here we
45 ignore the dependency between the variables since we consider only one variable at a time
46 (missing time points is possible) by applying multilevel modelling to reduce the points over
47 the time for the respective variable to one single point. We then repeat this procedure for the
48 rest of variables and obtain as many standardized points as the number of variables. Finally
49 we obtain the RMD, or the robust version of it in the presence of fat tailed, to obtain the rank
50 of those units.
51
52
53
54
55
56
57
58
59
60

Example one

For this example, and since the data is not highly collinear, we apply both strategies to maintain the ranking performance. We first follow strategy two to achieve the results of ranking indicators of the gynaecologists and compare the last 10% of the ranking distribution with those obtained by means of strategy one. If the two strategies lead to similar results, we then conclude that they converge to each other with respect to ranking the best/worst individual in the sample. The results of the two strategies are presented in Figures 1 and 2. These figures are constructed using two different functions the matter that makes them looking different in construction. In these figures, we show all the results from our analyses but we only rank a subset of them that are quite different in performance than the others. For example using strategy one, we in Figure 1 show the results from about the upper 10th percentile ranking of the indicators (16 gynaecologists) that might give us a picture of bad performance for the ranked individuals.

Figure 1. Upper 10th percentile ranking of the indicators using strategy one

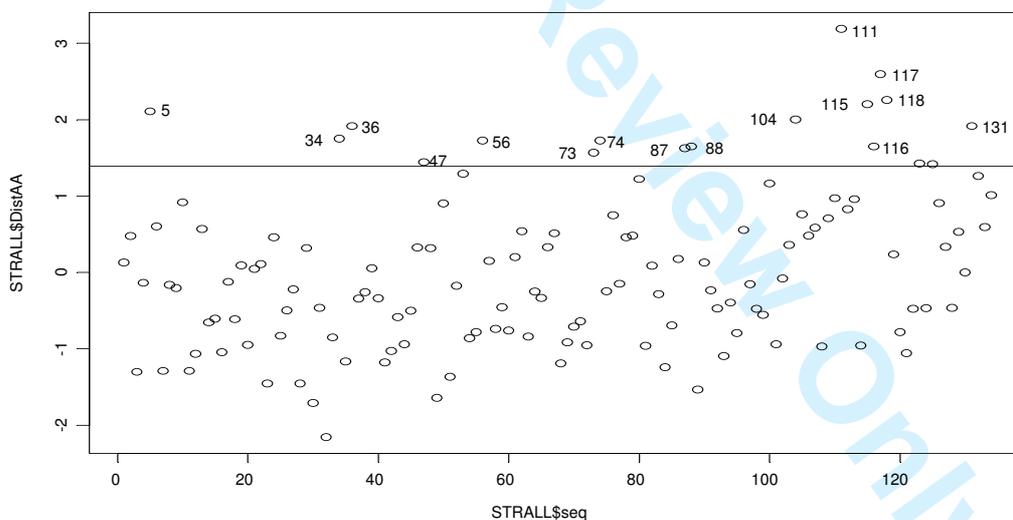
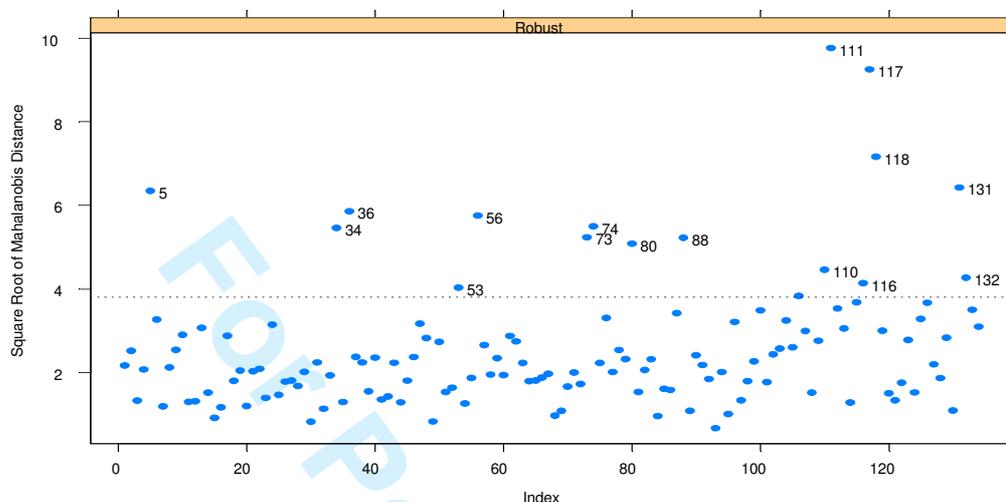


Figure 2 also shows results about the upper 10th percentile ranking of the indicators (16 gynaecologists), but we use strategy two to produce these ranking indicators. The results from the two figures reveal that both strategies almost agree in this case regarding the number of individual that belong to the upper 10th percentile of the ranking distribution of the gynaecologists (an agreement of about 80% in ranking and a correlation of about 0.75%). However, both strategies totally agree with respect to the last three gynaecologists with the worst performances in the sample, namely gynaecologists with ranks, 111, 117, and 118.

Figure 2. Upper 10th percentile ranking of the indicators using strategy two

A Markov Chain Monte Carlo (MCMC) simulation study has been conducted using WinBugs, version 1.4 to compute the uncertainty of the ranks. We compare between two types of distributions for the level 1 and level 2 residuals, namely the Normal and the Scale Mixture of the normal representation of the t-distributions. The results of this simulation study are presented in the following Table 1. These results are for the cases when there exist some kind of disagreement, when using these two distributions, regarding indicating the outlierness of the individuals. In this table, we only include the cases of individuals that considered as outliers by at least one of used distributions (i.e., when the values of the level two residuals are greater than 2). The results, however, indicate the following: The two methods show negative slope over the studied period, i.e. the fixed linear trend = $\{-0.1136, -0.0735\}$ indicating progressive shrink of the MD distance over time. The Inter Class Correlation (ICC) shows stronger clustering with the case of the Normal distribution (0.614) and weaker in the case scale mixture (0.429). The Tau in the table stands for the within level 1 residuals variance over the number of years, while the Tau.u2 is the between level 2 residuals variance. These variances are significant in both cases.

Using Scale Mixture, the level 2 residual value (of sample size 5 years) for the individual number [75] strongly classifies him/her as an outlier, indicating consistent low performance for this individual. On the other hand, when using the Normal distribution this individual is not considered as an outlier. This indicates that scale mixture show robust outlier detection.

Table 1.

Results that disagree regarding indicating the outlieriness of the individuals using the two distributions (values, standard deviations and confidence intervals are reported)

	Normal	Scale Mixture	Sample size
Grand mean	3.584 0.1762 (3.24, 3.934)	3.039 0.146 (2.751, 3.330)	-
Fixed linear trend	-0.1136 0.038 (-0.187, -0.038)	-0.0735 0.031 (-0.135, -0.013)	-
Tau	0.825 0.0642 (0.704, 0.954)	1.739 0.173 (1.431, 2.099)	-
Tau.u2	0.524 0.086 (0.374, 0.718)	1.324 0.276 (0.853, 1.921)	-
Individual [75]	1.871 0.483 (0.914, 2.835)	2.417 0.685 (1.00, 3.613)	5
Individual [90]	2.167 0.592 (1.001, 3.35)	1.721 0.847 (0.192, 3.504)	3
Individual [106]	2.594 0.888 (0.843, 4.304)	2.613 1.747 (-0.661, 5.405)	1
Individual [111]	5.952 0.982 (3.99, 7.86)	6.194 4.396 (-1.006, 11.19)	1
Individual [115]	2.765 0.887 (1.050, 4.531)	3.021 1.872 (-0.560, 5.81)	1
Individual [116]	2.206 0.599 (1.02, 3.36)	1.101 0.852 (-0.348, 2.922)	3
Individual [117]	3.78 0.91 (2.02, 5.599)	3.813 2.70 (-0.835, 7.63)	1
Individual [118]	2.876 0.895 (1.186, 4.668)	3.062 1.99 (-0.578, 6.06)	1
Individual [131]	2.199 0.876 (0.549, 3.977)	2.342 1.515 (-0.558, 4.804)	1
DIC	1593.74	1400.50	
ICC	0.614	0.429	

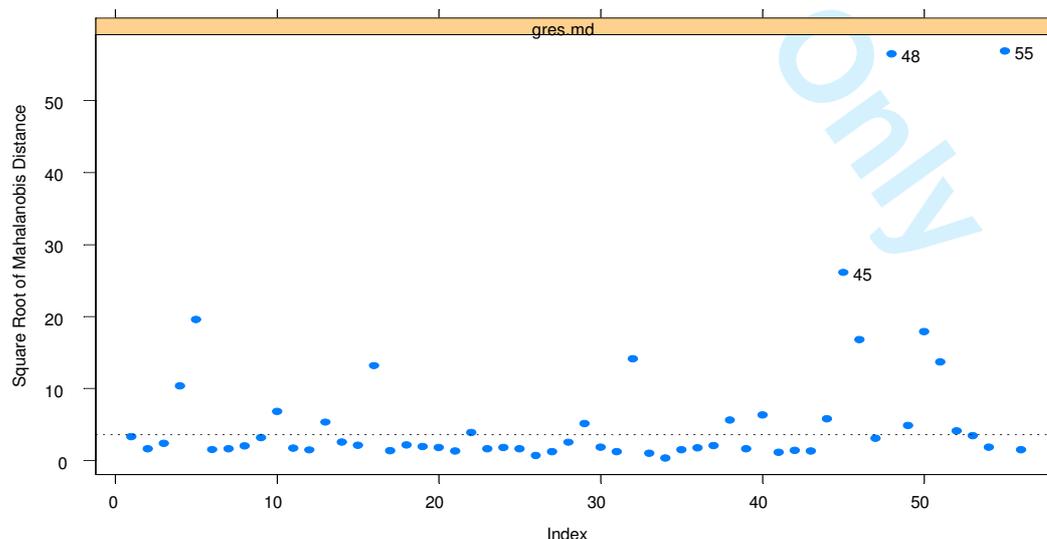
In the table, we have 6 individuals with only one year sample size available, namely [106], [111], [115], [117], [118] and [131]. Using the Normal distribution, these individuals have been significantly considered as outliers. On the other hand, when applying the Scale Mixture,

these individuals are still considered as outliers, but with higher uncertainty rendering the results statistically insignificant. The confidence intervals associated with the estimated residuals for these individuals cover the value of zero. This means that making analysis on these individuals by means of only one observation of time might lead to misleading conclusions. This has been discovered when using the more robust Scale Mixture but not the Normal distribution. The same is true for individuals [90] and [116] for whom the level 2 residual values (of sample size 3 years) significantly indicate them as outliers using the Normal distribution. On the other hand, the Scale Mixture significantly do not indicate the first as an outlier, the second however is not statistically significant.

Example two

In this example, and since the variables are highly correlated (the correlations between the variables are between 0.60 - 0.90), we only apply strategy two. Applying the first strategy, by first computing the PC to detect independence, will significantly reduce the information used in the analysis which might lead to the problem of omitting relevant information. The results shown are obtained using the variables *COGS to sales*, *Free cash flow to current liabilities*, *Growth in gross investment*, *Growth in net turnover*, *Growth in total debt*, *R&D to sales* and *cost of employees*. The results of the ranking indicators are shown in Figure 3 below.

Figure 3. Ranking of the companies indicators regarding best performances using strategy two



1
2
3 Our analysis indicated the firms with the best performances among the others, these are the
4 observations 45, 48 and 55 that stand for the firms Sandvik, Securitas and Volvo, respectively
5 (but that Sandvik and Securitas clearly have the best performance). A closer examination of
6 the figure, reveal that some other firms are also indicated to have fairly good performances,
7 namely those that are associated with observations 5, 16, 32, 46, 50 and 51 (not specified in
8 the figure) that stand for the firms Atlas Copo, Fingerprint Card, NCC, Scania, Skanska and
9 SKF. Note that all these indicated firms are heavy industry, building and computer firms in
10 Sweden.
11
12
13
14
15
16
17
18
19
20

21 **5. Summary and Conclusions**

22
23
24
25 The main purpose of this paper is to propose a general framework for performance evaluation
26 of organisations and individuals over time using routinely collected performance variables or
27 indicators. Two model-free approaches or strategies are recommended depending on the
28 nature of the data under consideration. It is not unusual that the data are often interdependent,
29 correlated over time, with missing observations, or used to come from heavy tailed
30 distribution shaped by outliers. Based on this fact, and the strength of possible
31 interdependence between variables, we introduce two strategies for evaluation (ranking) of
32 sampling units. In cases when the dependency structure between the variables is weak or the
33 variables are almost orthogonal to each other, the first strategy starts by computing the
34 Mahalanobis distance (MD) for each sampling unit (if these units are normally distributed), or
35 the RMD (in the case of heavy tails distributions) for all indicators over time. These distances
36 are then used in the second stage as the outcome in a multilevel model using RML to obtain
37 ranks of sampling units. The second strategy should be used when the dependency structure
38 between the variables is high (but it also works in cases with weak dependency or if the
39 variables are almost orthogonal). It starts by first applying the multilevel model with
40 indicators as the outcome to derive robust standardised residual for each indicator variable in
41 the first stage, and then compute an MD or RMD of all indicators for each sampling unit in
42 the second stage.
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57

58 To summarise, strategy one can handle missing data using robust residual maximum
59 likelihood (RML) at stage two, while strategy two handle missing data at stage one. Running
60

1
2
3 one indicator at a time in strategy two solve the multicollinearity problem. Strategy one
4 requires independent indicators for the construction of the distances (this imply a great loss of
5 information when the variables are not independent, however when calculating principle
6 components and excluding dependent data to achieve independency), where as strategy two
7 does not. Two different domain examples are used to illustrate the application of the two
8 strategies. Example one considers performance monitoring of surgeons and example two
9 considers the performance of industrial firms. The code for the program is included at the end
10 of this paper and the data we used together with the initial values are available upon request
11 from the authors.
12
13
14
15
16
17
18
19
20

21 The results from the first example reveal that both strategies almost agree regarding the
22 number of individual that belong to the upper 10th percentile of the ranking distribution of the
23 gynaecologists (an agreement of about 80% in ranking and a correlation of about 0.75).
24 However, both strategies totally agree with respect to the last three gynaecologists with the
25 worst performances in the sample. In the second example, however, we are intending to rank
26 the firms with the best performances. Since the analysed variables are very collinear, we
27 applied strategy two in our analysis. The results show that three firms are ranked as having the
28 best performance with respect to the sample period. These are Sandvik, Securitas and Volvo.
29 Other firms have also shown to have good performances, namely; Atlas Copo, Fingerprint
30 Card, NCC, Scania, Skanska and SKF. These firms are known to be heavy industry, building
31 and computer firms in Sweden.
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

References

- Behrenz, L. and Althin, R. (2005), "Efficiency and productivity of employment offices: Evidence from Sweden", *International Journal of Manpower*, Vol. **26**, No. 2.
- Gelman, A., Carlin, J., Stern, H., and Rubin, D (1995), "Bayesian data analysis". London: Chapman & Hall
- Goldstein, H., and Spiegelhalter, D. J (1996), "League Tables and Their Limitations: Statistical Issues in Comparisons of Institutional Performance", *J. R. Statist. Soc. A*, **159** (1996), pp. 385-443.
- Harley, M., M. A. Mohammed, S. Hussain, J. Yates, and A. Almasri (2005), "Was Rodney Ledward a Statistical outlier? Retrospective Analysis Using Routine Hospital Data to Identify Gynaecologist's Performance", *BMJ*, doi: 10.1136 / bmj.38377.675440.8F (Published 15 April 2005).
- Langford, I. H and Lewis, T. (1998), "Outliers in multilevel data", *J. R. Statist. Soc. A*, **161**: pp 121-160, 1998.
- Peel, D. and McLachlan, G. J. (2000), "Robust mixture modelling using the t distribution", *Statistics and Computing*, **10**, No.4/October, 2000 pp 339-348.
- Pinheiro, Jose C; Chuanhai Liu; Ying Nian Wu (2001), "Efficient Algorithms for Robust Estimation in Linear Mixed-Effects Models Using the Multivariate t-Distribution", *Journal of Computational and Graphical Statistics*, **10**, No. 2, pp. 249-276.
- Rocke, D. M. (1996), "Robustness Properties of S-Estimators of Multivariate Location and Shape in high Dimension", *Annals of Statistics*, Vol. **24**, No. 3, pp. 1327-45.
- Rocke, D. M. and Woodruff, D. L. (1996), "Identification of Outliers in Multivariate Data", *JASA*, **91**, pp. 1047-61.
- Rousseeuw, P. J. (1985), "Multivariate Estimation with High Breakdown Point", *Mathematical Statistics and Applications B* (W. Grossmann, G. Pflug, I. Vineze and W. Werz, eds.) pp. 283-297
- Rousseeuw, P. J. and Van Driessen (1999), "A Fast Algorithm for the Minimum Covariance Determinant Estimator", *Technometrics*, **41**, No. pp. 212-223.
- Rousseeuw, P. J. and A. M. Leroy (1987), "*Robust Regression and Outlier Detection*", New York, John Wiley.
- Searle, R. S., Casella G., and McCulloch. C. E. (1992), "*Variance Components*", John Wiley and Sons, INC.
- Seltzer, M. and K. Choi (2003), "Sensitivity Analysis for Hierarchical Models: Down weighting and Identifying Extreme Cases Using the t distribution", edited by Steven P. Reise and Naihua Duan, "Multilevel Modeling Methodological Advances, Issues, and Application" LONDON, LAWRENCE ERLBAUM ASSOCIATES, PUBLISHERS 2003.

```

1
2
3
4      NORMAL
5
6      model
7      {
8      # Level 1 definition
9      for(i in 1:N) {
10     md[i] ~ dnorm(mu[i],tau)
11     mu[i]<- beta[1]*cons[i]+beta[2]*rep[i]+ u2[ID[i]]*cons[i]
12     }
13     # Higher level definitions
14     for (j in 1:n2) {
15     u2[j] ~ dnorm(0,tau.u2)
16     }
17     # Priors for fixed effects
18     for (k in 1:2) { beta[k] ~ dflat() }
19     # Priors for random terms
20     tau ~ dgamma(1.5,3.012)
21     sigma2 <- 1/tau
22     tau.u2 ~ dgamma(1.5,4.656)
23     sigma2.u2 <- 1/tau.u2
24     ICC<-sigma2.u2/(sigma2+sigma2.u2)
25     }
26     -----
27
28     Scale Mixture
29
30     model
31     {
32     # Level 1 definition
33     for(i in 1:N) {
34     md[i] ~ dnorm(mu[i],winvsig2[i])
35     mu[i]<- beta[1]*cons[i]+beta[2]*rep[i]+ u2[ID[i]]*cons[i]
36     w[i] ~ dgamma(2.0,2.0)
37     winvsig2[i]<-w[i]*sig2inv
38     }
39     # Higher level definitions
40     for (j in 1:n2) {
41     u2[j] ~ dnorm(0,qinvtau[j])
42     q[j] ~ dgamma(2.0,2.0)
43     qinvtau[j]<-q[j]*taullinv
44     }
45     # Priors for fixed effects
46     for (k in 1:2) { beta[k] ~ dflat() }
47     # Priors for random terms
48     sig2inv ~ dgamma(1.5,2.259)
49     sig2 <- 1/sig2inv
50     taullinv ~ dgamma(1.5,3.492)
51     taull<- 1/taullinv
52     ICC<-taullinv/(taullinv+sig2inv)
53     }
54     -----
55
56
57
58
59
60

```