

Document downloaded from:

<http://hdl.handle.net/10251/182472>

This paper must be cited as:

González-Cebrián, A.; Arteaga, F.; Folch-Fortuny, A.; Ferrer, A. (2021). How to Simulate Outliers with the Desired Properties. *Chemometrics and Intelligent Laboratory Systems*. 212:1-16. <https://doi.org/10.1016/j.chemolab.2021.104301>



The final publication is available at

<https://doi.org/10.1016/j.chemolab.2021.104301>

Copyright Elsevier

Additional Information

1 How to Simulate Outliers with the Desired Properties

2 Alba González-Cebrián^a, Francisco Arteaga^b, Abel Folch-Fortuny^c, Alberto
3 Ferrer^a

4 ^a*Multivariate Statistics Engineering Research Group, Universitat Politècnica de*
5 *València, 46022 València, Spain*

6 ^b*Universidad Católica de València, 46001 València, Spain*

7 ^c*DSM Biotechnology Center, 2613 AX Delft, The Netherlands*

8 Abstract

Deviating multivariate observations are used typically to test the performance of outlier detection methods. Yet, the generation of outlying cases itself usually appears as a secondary methodological step in methods comparison. In the literature, outliers are defined using certain distribution parameters which differ from those of the clean or reference data. However, these parameters change among authors, leading to a lack of a standard and measurable definition of the characteristics simulated outliers. This makes the comparison between methods hard and its results dependent on the procedure followed to simulate the data. In order to set a standard procedure, a framework to simulate outliers is defined here. Since it is based on certain specifications for both the Squared Prediction Error (*SPE*) and Hotelling's T^2 statistics from a Principal Component Analysis (PCA) model, tuning them becomes a simple and efficient task. This procedure has been implemented in a set of Matlab functions.

9 *Keywords:* PCA, Outliers, Squared Prediction Error, Hotelling's T^2 ,
10 Simulation, Matlab

11 1. Introduction

12 Principal component analysis (PCA) models are specially useful in the
13 context of highly correlated data sets, given its dimensionality and noise
14 reduction power. This is accomplished by obtaining the A latent variables or
15 principal components (PCs) that are linear combinations of the K original
16 variables (usually with $A \ll K$). These A components explain most of the

17 variance of the K original variables. Beyond the use of PCA as a model
18 itself, it is also widely used for Exploratory Data Analysis (EDA), given
19 the effectiveness offered by the compression of a high-dimensional space to
20 a lower dimension representation that retains most of its variability. One of
21 the reasons why EDA is a good practice before any further use of a data
22 set, is that during this prior steps one can deal with events such as missing
23 data or the potential existence of rare events, also named outliers [1, 2, 3].
24 When PCA is used in an EDA framework, a model is built, which is known
25 as the PCA Model Building (PCA-MB) stage. In its basic definition, PCA
26 uses least squares parameters which can be very distorted by the influence
27 of outliers. In order to deal with this issue, several approaches that avoid
28 this negative effect have been proposed in literature, assembled in what are
29 known as robust PCA methods. There are plenty of strategies to conduct
30 PCA in a robust way. However, beyond the particularities of each proposal,
31 what basically defines these algorithms is their ability to neglect the influence
32 of potential outliers during the PCA-MB stage. To develop methodological
33 work on how to detect and how to treat outliers, it is often useful to simulate
34 this type of data.

35 Most approaches used to simulate outliers assume the paradigm of row-
36 wise outliers. This paradigm defines an outlier as a whole observation or
37 row in a matrix. Probably, the most famous model in order to define this
38 situation is the classical Tukey-Huber Contamination Model (THCM) [4]. In
39 these scenarios the observed data \mathbf{X} is thus a mix of unobserved distribu-
40 tions defining two different submatrices \mathbf{Y} and \mathbf{Z} , representing data from two
41 different populations. As one could expect, the election of the distributions
42 to simulate both the contaminated and clean parts of the data is a critical
43 procedure step, as it creates the conditions under which the performance of
44 different methods are evaluated and compared. Examining literature, one
45 can notice that the task of simulating the data sets and outliers in the frame-
46 work of PCA-MB has been addressed differently [5, 6, 7, 3]. In general terms,
47 what remains in common among most proposals is that outliers are defined
48 by setting the parameters of the population to which they belong. Thus,
49 observations are classified as outliers because they are drawn from a distri-
50 bution that is different from the one which describes the clean data. However,
51 it is not straightforward to stablish the relationship between the chosen pa-
52 rameters for the distribution of the outliers and the resulting properties of
53 the simulated observations. As a result, simulating observations with the
54 desired distance from the reference data set by setting different parameters

55 of the data distribution, becomes practically unfeasible. Moreover, working
56 with this simulation paradigm means to make assumptions about the distri-
57 butions that describe both the reference and outlying data set. Usually, a
58 multivariate normal distribution is assumed and the mean vector or the co-
59 variance matrix are altered in order to generate outlying observations. Yet,
60 assuming a particular probability distribution might not be that simple in
61 case that one wants to simulate outliers for a real reference data set. For
62 these reasons, though the traditional paradigm is technically correct, our be-
63 lief is that one could further exploit the information offered by a PCA model
64 in order to generate outliers with more control of their properties based on
65 two statistics: the Squared Prediction Error (*SPE*) and the Hotelling T^2
66 (T^2). In this work we propose a standard framework for outliers definition
67 and simulation based on its characterization in terms of these statistics.

68 Firstly, the conceptual framework is introduced, defining the PCA model
69 and the aforementioned pair of statistics. Afterwords, the methodology to
70 generate moderate and severe perturbations, based on shift directions of the
71 *SPE* and the T^2 , is explained. Later on, the proposed variants of the algo-
72 rithm to simulate outliers are introduced, and some examples of how to simu-
73 late controlled outliers are shown. Moreover, some practical applications will
74 be provided to illustrate the potential of the proposed method as standard
75 framework to simulate outliers. In these examples, our procedure to simulate
76 controlled outliers will be configured to emulate other strategies of outliers
77 generation from literature on PCA models. Additionally, the consistency of
78 the outlying properties will be assessed by projecting our simulated outliers
79 onto a robust PCA model. Finally, a summary of the main conclusions is pro-
80 vided. The Matlab code and documentation for outliers generation are avail-
81 able in the GitHub repository <https://github.com/albagc/SCOUTer.git>.
82 Detailed code lines to reproduce the results from Section 3 are available in the
83 *howto.pdf* document on the repository and further details about references
84 for the outliers simulation are provided in Appendix A.

85 **2. Materials and methods**

86 *2.1. The PCA model framework*

87 Let \mathbf{X} be a matrix with N observations on K variables. After some pre-
88 processing such as mean-centering and/or unit variance scaling, a PCA model
89 is estimated. This is done by compressing the high-dimensional \mathbf{X} matrix
90 into a low-dimensional subspace of dimension A (with $A \leq \text{rank}(\mathbf{X})$). PCA

91 is based on the bilinear decomposition of \mathbf{X} in $\mathbf{X} = \mathbf{TP}^\top + \mathbf{E}$, where \mathbf{T} is an
 92 $N \times A$ matrix of *scores* and \mathbf{P} is a $K \times A$ matrix of *loadings* (Figure 1).

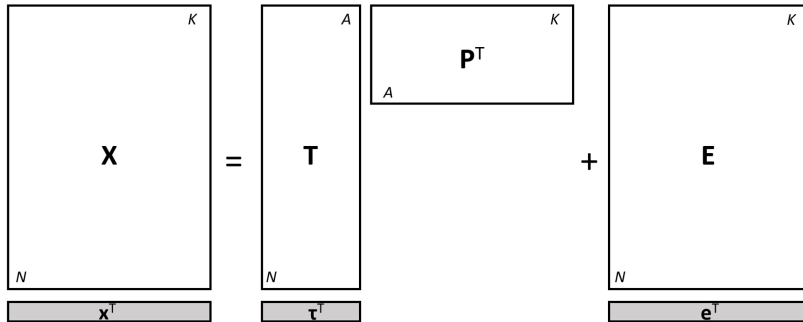


Figure 1: Visual representation for the PCA model.

93 The A columns of the loading matrix \mathbf{P} are the *loading* vectors \mathbf{p}_a , with
 94 $a = 1, 2, \dots, A$. The *score* matrix \mathbf{T} can be considered as a collection of row
 95 vectors $\boldsymbol{\tau}^\top$ (scores of an observation) or column vectors \mathbf{t}_a (latent variables,
 96 with $\mathbf{t}_a = \mathbf{X}\mathbf{p}_a$ and $a = 1, 2, \dots, A$). The score matrix can be obtained as
 97 $\mathbf{T} = \mathbf{XP}$, that is, as the projection of the \mathbf{X} matrix on the A -dimensional
 98 space of the PCA model (i.e., columns of \mathbf{P} matrix). Analogously, given an
 99 observation \mathbf{x} of the original K -dimensional space, its projection $\boldsymbol{\tau}$ onto the
 100 subspace of the model can be obtained using the projection matrix \mathbf{P} as well
 101 by $\boldsymbol{\tau} = \mathbf{P}^\top \mathbf{x}$.

102 From the scores matrix one can recall the explained part of \mathbf{X} in the
 103 PCA model as $\hat{\mathbf{X}} = \mathbf{TP}^\top$. This notation can be used as well for individual
 104 observations, where $\hat{\mathbf{x}} = \mathbf{P}\boldsymbol{\tau}$. The original observation can be decomposed by
 105 the part explained (i.e., predicted) by the model (signal or $\hat{\mathbf{x}}$) and the error
 106 not considered in any of the A latent variables (noise or \mathbf{e}). Thus, for a given
 107 observation we have $\mathbf{x} = \mathbf{P}\boldsymbol{\tau} + \mathbf{e} = \hat{\mathbf{x}} + \mathbf{e}$. From the previous expressions it
 108 can be seen that $\mathbf{E} = \mathbf{X}(\mathbf{I} - \mathbf{PP}^\top)$ and then $\mathbf{e} = (\mathbf{I} - \mathbf{PP}^\top) \mathbf{x}$.

109 2.2. Outliers in the PCA model

110 An observation can be considered an outlier in terms of a PCA model,
 111 according to its values for the Squared Prediction Error (*SPE*) and the
 112 Hotelling's T^2 (T^2 , or more specifically, T_A^2 for a PCA model with A compo-
 113 nents). These statistics, obtained from the residuals and the scores respec-
 114 tively, offer complementary information about the distance of an observation
 115 to the PCA model and the majority of data. In [8], there is a comprehensive

116 explanation about the properties of these statistics and their use to detect
 117 outlying observations. Following their work, some mathematical aspects of
 118 SPE and the T_A^2 are given in this section.

119 The SPE is the squared Euclidean (perpendicular) distance from the
 120 observation \mathbf{x} to the A -dimensional subspace of the model, that is $SPE =$
 121 $\mathbf{e}^\top \mathbf{e}$, where \mathbf{e} is the error vector of the observation \mathbf{x} . From the previous
 122 expression, the SPE can be rewritten as $SPE = \mathbf{x}^\top (\mathbf{I} - \mathbf{P}\mathbf{P}^\top)^\top (\mathbf{I} - \mathbf{P}\mathbf{P}^\top) \mathbf{x}$.
 123 Since $(\mathbf{I} - \mathbf{P}\mathbf{P}^\top)$ is symmetric and idempotent matrix:

$$SPE = \mathbf{x}^\top (\mathbf{I} - \mathbf{P}\mathbf{P}^\top) \mathbf{x} \quad (1)$$

124 Assuming that residuals follow a multivariate normal distribution, [9],
 125 [10] and [11], derived approximate distributions for such quadratic forms.

126 On the other hand, the Hotelling- T_A^2 statistic for an observation is defined
 127 as

$$T_A^2 = \boldsymbol{\tau}^\top \boldsymbol{\Theta}^{-1} \boldsymbol{\tau} = \sum_{a=1}^A (\tau_a^2 / \lambda_a) \quad (2)$$

128 where $\boldsymbol{\Theta}(A \times A)$ is the covariance matrix of \mathbf{T} (diagonal matrix of the high-
 129 est A eigenvalues $\{\lambda_1, \dots, \lambda_A\}$). It represents the estimated squared Maha-
 130 lanobis distance from the center of the latent subspace to the projection of
 131 an observation onto this subspace.

132 When diagnosing which variables yield the obtained values for the SPE
 133 and the T^2 it can be useful to check the contributions of each variable to
 134 each statistic [8].

135 From these two statistics (the SPE and the T^2), two complementary
 136 control metrics are obtained. Firstly, with an appropriate reference set of
 137 data, the in-control PCA model is built. The control limits are defined as
 138 well using the reference distributions for each statistic.

139 Regarding the Upper Control Limit (UCL) for the SPE , several proce-
 140 dures can be used. In [10] it is shown that an approximate SPE critical
 141 value at significance level α is given by

$$UCL(SPE)_\alpha = \theta_1 \left[z_\alpha \sqrt{2\theta_2 h_0^2 / \theta_1 + 1 + \theta_2 h_0 (h_0 - 1) / \theta_1^2} \right]^{1/h_0} \quad (3)$$

142 where $\theta_k = \sum_{j=A+1}^{rank(\mathbf{X})} (\lambda_j)^k$, $h_0 = 1 - 2\theta_1 \theta_3 / 3\theta_2^2$, λ_j are the eigenvalues of the
 143 PCA residual covariance matrix $\mathbf{E}^\top \mathbf{E} / (N - 1)$, and z_α is the 100(1 - α)%
 144 percentile of a standard normal variable.

145 Alternatively, one can use an approximation based on the weighted chi-
 146 squared distribution ($g\chi_h^2$) proposed by [9]. In [12] authors suggested a simple
 147 and fast way to estimate parameters g and h which is based on matching
 148 moments between a $g\chi_h^2$ distribution and the sample distribution of SPE .
 149 The mean ($\mu = gh$) and variance ($\sigma^2 = 2g^2h$) of the $g\chi_h^2$ distribution are
 150 equated with the sample mean (b) and variance (v) of the SPE sample.
 151 Hence, the Upper SPE Control Limit at significance level α is given by

$$UCL(SPE)_\alpha = v\chi_{(2b^2/v),\alpha}^2/(2b) \quad (4)$$

152 where $\chi_{(2b^2/v),\alpha}^2$ is the $100(1-\alpha)\%$ percentile of the corresponding chi-squared
 153 distribution with $2b^2/v$ degrees of freedom.

154 Upper Control Limits (UCL) for the T_A^2 at a significance level (type I)
 155 risk α can be obtained assuming that the statistic follows an F distribution

$$T_A^2 \sim A(N^2 - 1) F_{A,(N-A)}/(N(N - A)) \quad (5)$$

156 Thus, the corresponding UCL from Equation 5 is given by

$$UCL(T_A^2)_\alpha = A(N^2 - 1) F_{(A,(N-A)),\alpha}/(N(N - A)) \quad (6)$$

157 According to the aforementioned conceptual meaning of these multivari-
 158 ate statistics (SPE and T_A^2), observations above their associated UCL will
 159 be representing different types of outliers.

160 The first type of outliers, with high SPE , occurs when the correlation
 161 structure between variables is different from the observed one during the
 162 model fitting with the clean data set. Using these observations to fit the
 163 model can lead to dramatic distortions on the correlation structure captured
 164 by the PCs. These perturbations are named “moderate outliers” or “anoma-
 165 lous observations” and are caused by unusual variations outside the model.

166 The second type of outliers, with high values of the T_A^2 , appears when
 167 the correlation structure between measured variables remains constant but
 168 their absolute values differ from the expected ones. These perturbations are
 169 named “severe outliers” or “extreme observations” and they are usually rep-
 170 resenting unusual shifts in the model (i.e. shifts that respect the correlation
 171 structure of the model). This leads to extreme values in the projection of
 172 these observations with respect to the ones obtained for the clean data set.

173 These links between distances and types of outliers, or outlying properties,
 174 can be described using the Squared Prediction Error and Hotelling’s T_A^2 . On

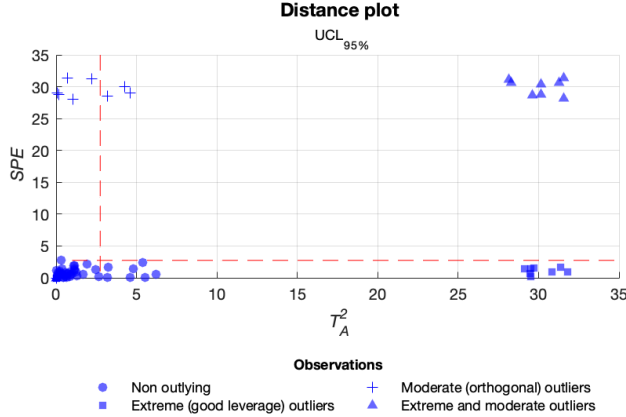


Figure 2: Types of outliers according to the PCA model built with a reference data set (blue dots). Red discontinuous lines are the 95% control limits for the SPE and Hotelling's T^2 .

175 one hand, moderate outliers will present high values for the SPE statistic,
 176 which is the reason why they are also known as orthogonal outliers. On
 177 the other hand, extreme outliers are observations with high values for the
 178 T_A^2 . They are named as well good leverage observations, since their presence
 179 does not distort the correlation structure of the model. There can be as well
 180 observations that are both moderate and extreme outliers (Figure 2).

181 In conclusion, outliers in the context of a PCA model will be associated
 182 with large values of the SPE , the T_A^2 or both distances. Using this pair of
 183 statistics to describe observations provides more meaningful criteria to define
 184 outliers than setting their distribution parameters in the original or latent
 185 space. Setting not only the type, but also how far outliers will be from the
 186 reference data set, is something plausible when using the SPE and T_A^2 as
 187 targets to simulate outliers. This idea of representing all types of outliers as
 188 combinations of SPE and T_A^2 values, is the basis of the simulation approach
 189 presented in this work.

190 2.3. Framework to generate outliers

191 Our proposal for the generation of outliers is to transform an observation
 192 \mathbf{x} , with given SPE and T^2 values ($SPE_{\mathbf{x}}$ and $T_{\mathbf{x}}^2$, respectively), into a new
 193 observation with an SPE and/or T^2 values specified by the user ($SPE_{\mathbf{y}}$ and

194 $T_{\mathbf{y}}^2$, respectively). The transformation will consist in a shift of the observation
 195 following certain direction in the space of the original variables.

196 Moving the observation \mathbf{x} in the direction \mathbf{v} to obtain a new observation
 197 $\mathbf{y} = \mathbf{x} + \mathbf{v}$, we can calculate the new value of the SPE and the T^2 statistics,
 198 based on the original values:

$$SPE_{\mathbf{x}+\mathbf{v}} = (\mathbf{x} + \mathbf{v})^\top (\mathbf{I} - \mathbf{P}\mathbf{P}^\top) (\mathbf{x} + \mathbf{v}) = SPE_{\mathbf{x}} + \mathbf{v}^\top (\mathbf{I} - \mathbf{P}\mathbf{P}^\top) (2\mathbf{x} + \mathbf{v}) \quad (7)$$

$$T_{\mathbf{x}+\mathbf{v}}^2 = (\mathbf{x} + \mathbf{v})^\top \mathbf{P}\mathbf{\Theta}^{-1}\mathbf{P}^\top (\mathbf{x} + \mathbf{v}) = T_{\mathbf{x}}^2 + \mathbf{v}^\top \mathbf{P}\mathbf{\Theta}^{-1}\mathbf{P}^\top (2\mathbf{x} + \mathbf{v}) \quad (8)$$

200 The next issue is how to choose the direction \mathbf{v} . An obvious choice is to
 201 shift the observation in the direction that joins it with the origin of coordi-
 202 nates in the original data space, taking $\mathbf{v} = c\mathbf{x}$. In this case, it is easy to
 203 calculate the change in both statistics:

$$SPE_{\mathbf{x}+c\mathbf{x}} = (1 + c)^2 SPE_{\mathbf{x}} \quad (9)$$

$$T_{\mathbf{x}+c\mathbf{x}}^2 = (1 + c)^2 T_{\mathbf{x}}^2 \quad (10)$$

205 However, we are interested in finding directions in which we can control
 206 the change that occurs in each statistic. For example, there are specific
 207 directions that allow the change in one of both statistics, without affecting
 208 the other. In particular, we can move the observation in the direction of its
 209 residual vector in the PCA model: $\mathbf{e} = (\mathbf{I} - \mathbf{P}\mathbf{P}^\top) \mathbf{x}$, so that a change in the
 210 SPE will occur, without modifying the T^2 . Similarly, we can move it in the
 211 direction that joins the projection of the observation on the model with the
 212 origin (i.e. the direction of the predicted observation $\hat{\mathbf{x}}$): $\mathbf{P}\mathbf{P}^\top \mathbf{x}$, so that there
 213 will be a change in T^2 , without modifying the SPE . As both directions are
 214 orthogonal, we can compose both displacements in one operator, with control
 215 over the amount by which each of them increases. This will be illustrated in
 216 following sections.

217 2.3.1. Shift of the SPE statistic

218 If we move the observation \mathbf{x} in the direction given by its residual vector
 219 (according to the PCA model): $\mathbf{e} = (\mathbf{I} - \mathbf{P}\mathbf{P}^\top) \mathbf{x}$, multiplied by a scalar a ,
 220 we get, from Equation 7 and Equation 8:

$$SPE_{\mathbf{x}+a(\mathbf{I}-\mathbf{P}\mathbf{P}^\top)\mathbf{x}} = SPE_{\mathbf{x}+a\mathbf{x}^\top (\mathbf{I} - \mathbf{P}\mathbf{P}^\top) (2\mathbf{x} + a (\mathbf{I} - \mathbf{P}\mathbf{P}^\top) \mathbf{x})} = (1+a)^2 SPE_{\mathbf{x}} \quad (11)$$

221

$$T_{\mathbf{x}+a(\mathbf{I}-\mathbf{P}\mathbf{P}^\top)\mathbf{x}}^2 = T_{\mathbf{x}}^2 + a\mathbf{x}^\top (\mathbf{I} - \mathbf{P}\mathbf{P}^\top) \mathbf{P}\boldsymbol{\Theta}^{-1}\mathbf{P}^\top (2\mathbf{x} + a(\mathbf{I} - \mathbf{P}\mathbf{P}^\top)\mathbf{x}) = T_{\mathbf{x}}^2 \quad (12)$$

222 We can choose the value a to achieve a target value for the SPE statistic,
223 say $SPE_{\mathbf{y}}$:

$$(1 + a)^2 SPE_{\mathbf{x}} = SPE_{\mathbf{y}} \rightarrow a = \sqrt{SPE_{\mathbf{y}}/SPE_{\mathbf{x}}} - 1 \quad (13)$$

224 Note that the selected direction is the one that maximizes the change in
225 the SPE , because the gradient of this statistic is: $\nabla(SPE)(\mathbf{x}) = 2(\mathbf{I} - \mathbf{P}\mathbf{P}^\top)\mathbf{x}$.

226 2.3.2. Shift of the T^2 statistic

227 If we move the observation \mathbf{x} in the direction $\mathbf{P}\mathbf{P}^\top\mathbf{x}$, multiplied by a
228 scalar b , we get, from Equation 7 and Equation 8:

$$SPE_{\mathbf{x}+b\mathbf{P}\mathbf{P}^\top\mathbf{x}} = SPE_{\mathbf{x}} + b\mathbf{x}^\top \mathbf{P}\mathbf{P}^\top (\mathbf{I} - \mathbf{P}\mathbf{P}^\top) (2\mathbf{x} + b\mathbf{P}\mathbf{P}^\top\mathbf{x}) = SPE_{\mathbf{x}} \quad (14)$$

229

$$T_{\mathbf{x}+b\mathbf{P}\mathbf{P}^\top\mathbf{x}}^2 = T_{\mathbf{x}}^2 + b\mathbf{x}^\top \mathbf{P}\boldsymbol{\Theta}^{-1}\mathbf{P}^\top (2\mathbf{x} + b\mathbf{P}\mathbf{P}^\top\mathbf{x}) = (1 + b)^2 T_{\mathbf{x}}^2 \quad (15)$$

230 We can choose the value b to achieve a target value for the T^2 statistic,
231 say $T_{\mathbf{y}}^2$:

$$(1 + b)^2 T_{\mathbf{x}}^2 = T_{\mathbf{y}}^2 \rightarrow b = \sqrt{T_{\mathbf{y}}^2/T_{\mathbf{x}}^2} - 1 \quad (16)$$

232 We can also select the direction that maximizes the change in the T^2
233 statistic, without changing the SPE statistic, choosing the gradient of the
234 T^2 statistic: $\nabla(T^2) = 2\mathbf{P}\boldsymbol{\Theta}^{-1}\mathbf{P}^\top\mathbf{x}$. We do not use this direction because it
235 is difficult to parametrise the amount of change in the T^2 statistic.

236 2.3.3. Shift both statistics simultaneously

237 If we have an observation \mathbf{x} with statistics $SPE_{\mathbf{x}}$ and $T_{\mathbf{x}}^2$, we can trans-
238 form it into a new observation with statistics $SPE_{\mathbf{y}}$ and $T_{\mathbf{y}}^2$ combining the
239 aforementioned transformations:

$$\mathbf{y} = \mathbf{x} + a(\mathbf{I} - \mathbf{P}\mathbf{P}^\top)\mathbf{x} + b\mathbf{P}\mathbf{P}^\top\mathbf{x} \quad (17)$$

240 With $a = \sqrt{SPE_{\mathbf{y}}/SPE_{\mathbf{x}}} - 1$ and $b = \sqrt{T_{\mathbf{y}}^2/T_{\mathbf{x}}^2} - 1$, as seen in Equa-
241 tion 13 and Equation 16. The procedure to build a new observation with
242 desired SPE and T^2 statistics, based on an arbitrary prior observation \mathbf{x} ,

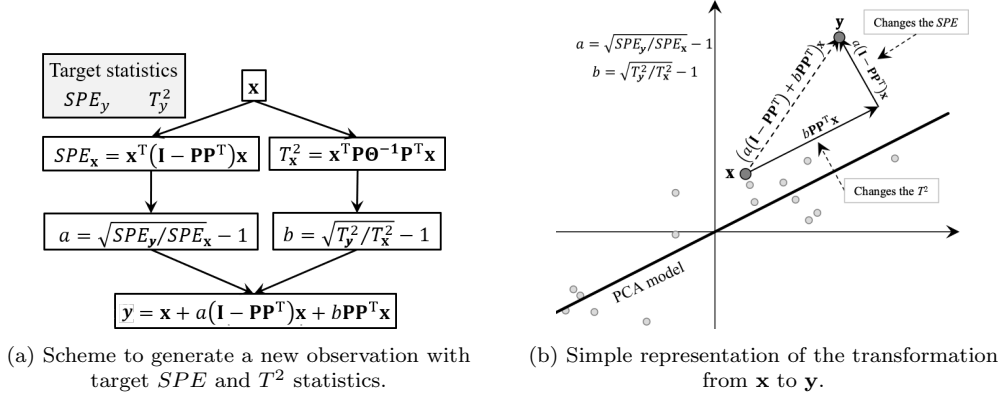


Figure 3: Representation of the algorithm to generate \mathbf{y} .

243 is illustrated in Figure 3a. The visual representation of the algorithm with
 244 a model of only one PC, for an original space with only two variables is
 245 represented in Figure 3b.

246 Furthermore, there is another aspect that can be used to control the outly-
 247 ing behaviour of the new observations. Given the reference and target values
 248 of a statistic, one can generate a series of $M - 1$ intermediate observations
 249 between the reference and the target one: $\{\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_{M-1}\}$. Mathemati-
 250 cally, the expected value of a statistic H_m as a result of a transition from the
 251 reference H_0 to the target value H_M :

$$H_m = H_0 + (m/M)^\gamma (H_M - H_0) \quad m = 1, 2, \dots, M - 1 \quad (18)$$

252 Thus, SPE_m and T_m^2 will follow a pattern of gradual change according
 253 not only to the number of steps, but also to the spacing between them.
 254 This spacing is regulated in Equation 18 by the γ parameter. As it can be
 255 appreciated in Figure 4, when this parameter is set to 1, the spacing between
 256 steps is linear, shifting towards a non-linear dynamic as it drifts from 1.

257 Given that both parameters (γ_{SPE} and γ_{T^2}) can be shifted simultaneously,
 258 this gives to the user the flexibility to simulate a wider variety of trajectories
 259 for each possible combination of values along the spacing of the two param-
 260 eters. Performing simultaneous shifts with some values for the parameters,
 261 results in the curves of Figure 5.

262 This framework, including the possibility of controlling the distance be-
 263 tween intermediate observations in series of outliers, can be useful in order
 264 to study and compare the sensitivity of different robust PCA approaches

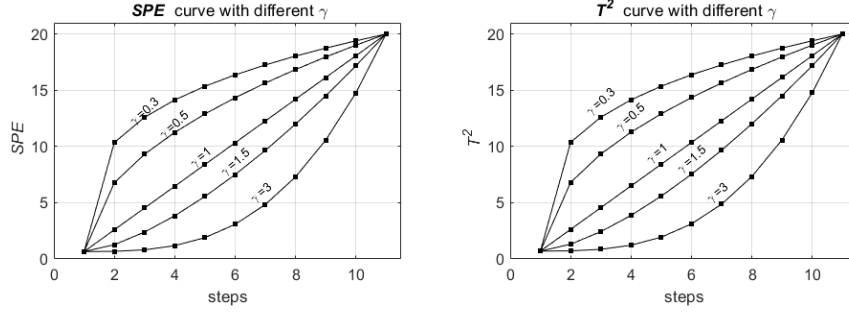


Figure 4: Curves for the SPE (left) and T^2 (right) statistics along the shift in 20 steps for different values of their spacing parameters γ .

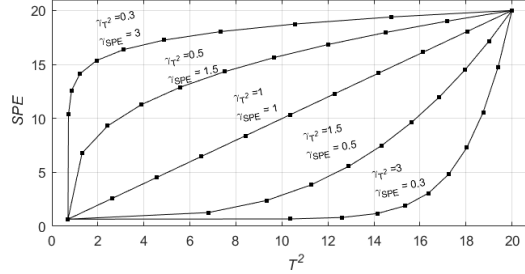


Figure 5: Curves for the SPE and T^2 statistics along the shift in 20 steps for different combinations of their γ parameters.

265 or methods for outlying detection. Thus, one could know not only for what
 266 type of outliers, but also at which step, one method performs differently from
 267 others. Finally, considering all these parameters one has the complete flux
 268 diagram of the procedure in Figure 6.

269 If a given observation \mathbf{x} is moved in different directions, it will be appre-
 270 ciated both in the SPE and T^2 statistics, and also in the scores. Figure 7
 271 illustrates different shifts on a five dimensional observation \mathbf{x} according to a
 272 reference PCA model.

273 In Figure 7a, red dashed lines represent the UCL for the T^2 and SPE
 274 statistics. The ellipse represented in the score plot from Figure 7b, is the
 275 contour curve of the confidence ellipsoid for the T^2 statistic, calculated for a
 276 confidence level of $(1-\alpha) \times 100\%$. From Equation 6, it is obtained an ellipsoid
 277 delimited in each dimension (i.e. PC) of the latent subspace. The contour of
 278 that ellipsoid represents a region of the space which holds $T^2 = T^2_{100(1-\alpha)\%CL}$
 279 for each observation lying on that contour. Since the score plot is a bi-

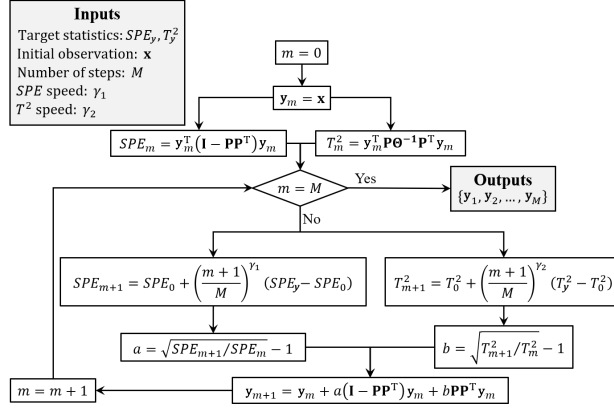
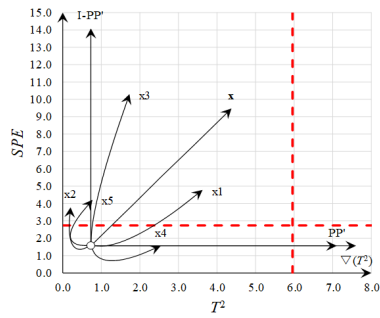
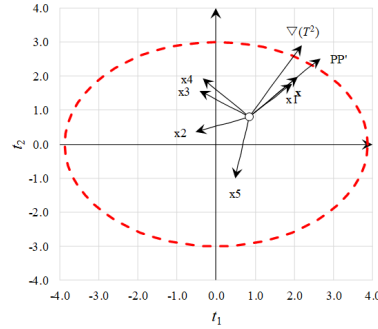


Figure 6: Flux diagram of simulation algorithm including all the parameters.



(a) Distance plot with five different shift directions for an observation.



(b) Score plot with five different shift directions for an observation.

Figure 7: Illustration of how moving in different directions would affect the SPE and T^2 (a) and the scores (b).

280 dimensional plot, the bi-dimensional representation of the confidence ellipsoid
 281 turns into a confidence ellipse. Therefore, observations lying outside the
 282 ellipse will be over passing the UCL for the T^2 statistic.

283 The first set of directions are those that correspond to the five variables
 284 (labelled as x_1, \dots, x_5). The trivial direction ($\mathbf{v} = \mathbf{x}$) is also considered.
 285 The direction corresponding to the residual vector ($\mathbf{v} = (\mathbf{I} - \mathbf{P}\mathbf{P}^T)\mathbf{x}$) is easy
 286 to recognize, since it causes an increase in the SPE without affecting the
 287 T^2 statistic. In the distance plot (Figure 7a) it is represented as a vertical
 288 arrow, whereas it does not appear in the score plot (Figure 7b), given that
 289 the projection of $\mathbf{x} + a(\mathbf{I} - \mathbf{P}\mathbf{P}^T)\mathbf{x}$ in the model space is the same as that of
 290 \mathbf{x} , for all a values.

291 The last two directions are $\mathbf{P}\mathbf{P}^\top \mathbf{x}$ and $\mathbf{P}\mathbf{\Theta}^{-1}\mathbf{P}^\top \mathbf{x}$ (labelled as $\nabla(T^2)$ in
 292 Figure 7). These two directions are in the model plane and this means that
 293 the SPE will not be affected, which can be appreciated by the horizontal
 294 arrows in Figure 7a. The magnitude of the shift in the T^2 value is bigger
 295 for the $\mathbf{P}\mathbf{\Theta}^{-1}\mathbf{P}^\top \mathbf{x}$ direction, since it corresponds to the gradient of the T^2
 296 statistic. The trajectory described by the scores when the direction $\mathbf{P}\mathbf{P}^\top \mathbf{x}$
 297 is chosen, is an extension of the segment that joins the origin (0,0) with the
 298 scores of \mathbf{x} (i.e. the direction of the predicted observation $\hat{\mathbf{x}}$). The trajectory
 299 followed when the shift is performed in the direction $\mathbf{P}\mathbf{\Theta}^{-1}\mathbf{P}^\top \mathbf{x}$ ($\nabla(T^2)$) is
 300 perpendicular to the $(1 - \alpha) \times 100$ confidence level Hotelling's T^2 ellipse,
 301 which is defined as the level curve for the T^2 statistic.

302 3. Results

303 In this section, some examples of how to simulate outliers with the desired
 304 properties are shown. This section is divided in two main parts. The first
 305 part will present results for three different scenarios of outliers simulation.
 306 Afterwards, four examples of outliers simulation extracted from literature are
 307 emulated using the framework proposed in this work. The aim of this exercise
 308 is to show how the technique described in this work can comprise other
 309 particular simulation settings. Finally, an assessment about the properties
 310 of the simulated outliers in terms of a robust PCA model is provided as well.

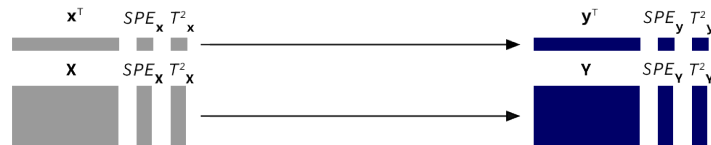
311 3.1. Cases of use of the proposed method.

312 These results illustrate three generic simulation scenarios: generating out-
 313 liers in one step, generating a sequence of outliers, and generating a grid of
 314 outliers. For this purpose, a reference matrix \mathbf{X} of $n = 50$ observations and
 315 $k = 5$ normally distributed variables is simulated. The PCA model based on
 316 \mathbf{X} is built with two PCs, assuming a type I risk α of 5% and performing a
 317 mean centering. All functions along their documentation and the script to
 318 reproduce the following scenarios can be downloaded from the github reposi-
 319 tory <https://github.com/albagc/SCOUTer.git>. A detailed explanation
 320 about the obtention of the following results can be found in the *howto.pdf*
 321 file.

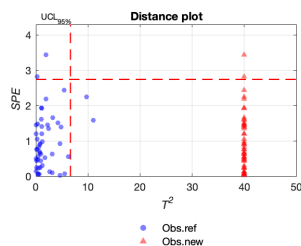
322 3.1.1. Case I: One-step simulation of outliers.

323 This is the simplest case, in which from an initial observation \mathbf{x} with
 324 reference values $SPE_{\mathbf{x}}$ and $T_{\mathbf{x}}^2$, a new observation \mathbf{y} is obtained, with the

325 desired $SPE_{\mathbf{y}}$ and $T_{\mathbf{y}}^2$ values (8a). The aforementioned scheme can be easily
 326 generalised for a set of observations. In the following example, the original
 327 \mathbf{X} matrix will be drifted from its initial coordinates. In this scenario a set of
 328 one-step outliers is generated by increasing only the T^2 value (i.e. extreme
 outliers). The SPE remains at its reference value.



(a)



(b)

Figure 8: (a) Illustration of a one-step simulation of controlled outliers. (b) Distance plot with the reference (blue circles) and the shifted (red triangles) data sets, performing a single step keeping the initial $SPE_{\mathbf{X}}$ value, but setting a target value $T_{\mathbf{Y}}^2 = 40$ for all the observations.

329 As it can be seen in Figure 8b, all observations have been shifted in their
 330 distance to the center on the model plane, drawing a contour on the score plot
 331 for the value $T_A^2 = 40$, whereas they have kept their values on the SPE
 332 statistic. In other words, this is an example of a how to simulate a set of
 333 extreme observations.
 334

335 3.1.2. Case II: Step-wise simulation of outliers

336 In this scenario, the transition between the reference and the target values
 337 for the statistics is performed with a spacing of n steps between them. From
 338 a reference observation \mathbf{x} (or set of observations \mathbf{X}) with reference values
 339 $SPE_{\mathbf{x}}$ and $T_{\mathbf{x}}^2$ (or $SPE_{\mathbf{X}}$ and $T_{\mathbf{X}}^2$), a series of $M - 1$ new sets of observations
 340 up to \mathbf{y} (or \mathbf{Y}) with the desired $SPE_{\mathbf{y}}$ and $T_{\mathbf{y}}^2$ (or $SPE_{\mathbf{Y}}$ and $T_{\mathbf{Y}}^2$) values is
 341 generated (9a).

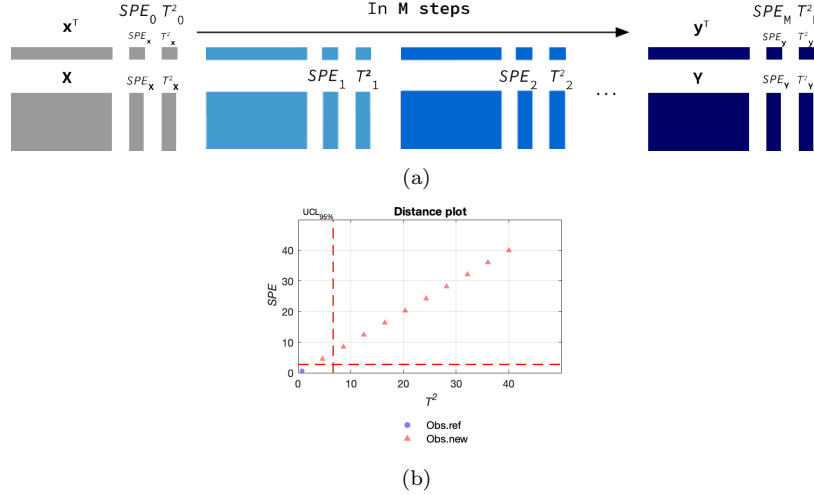


Figure 9: (a) Illustration of a M -step simulation of controlled outliers. (b) Distance plot after performing a 10-step shift both in the SPE_x and the T^2 values from one initial observation \mathbf{x} (blue circle).

342 In the example from above (Figure 9b), there is a linear spacing between
 343 steps for the SPE and the T^2 . However, the spacing between steps can be
 344 tuned, as seen in Figure 4 and Figure 5 from Section 2.3.

345 3.1.3. Case III: Grid-wise simulation of outliers

346 With the step-wise approach the same number of steps is performed for
 347 both statistics. Finally, the grid-wise case enables a different number of
 348 steps for each statistic. Starting from an initial data set \mathbf{x} (or \mathbf{X}) with
 349 reference values SPE_x and T_x^2 (or $SPE_{\mathbf{X}}$ and $T_{\mathbf{X}}^2$), a grid of new observations
 350 combining each step of the statistics is obtained (Figure 10a). As a result,
 351 there are as many data sets simulated as combinations between the steps of
 352 the statistics.

353 In this last case, a grid with 2 steps for the SPE and 3 steps for the T^2
 354 has been produced, setting different spacing parameters for each parameter
 355 as well (Figure 10b).

356 3.2. Comparison to other simulation methods and PCA frameworks

357 The aim of this section is to address two important questions about the
 358 simulation method proposed in this work: i) the proposed simulation frame-
 359 work can encompass other existing simulation strategies, and ii) if the prop-

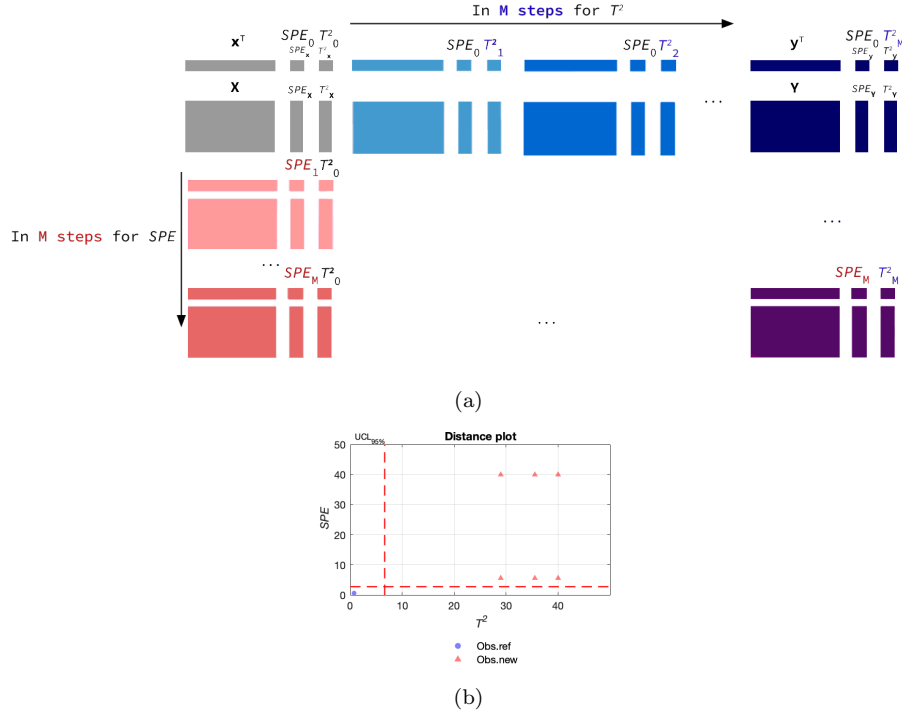


Figure 10: (a) Illustration of a grid-case simulation with M -step shifts for the SPE and the T^2 . (b) Distance plot after performing two steps for the SPE with $\gamma_{SPE} = 3$ and three steps for the T^2 with $\gamma_{T^2} = 0.3$ from one reference observation \mathbf{x} (blue circle).

360 eries of the simulated outliers will be maintained when they are projected
 361 onto PCA models fitted with other algorithms rather than the classical least
 362 squares version.

363 3.2.1. Simulation of other outlier generation strategies

364 With the aim of assessing if the proposed method can be seen as a general
 365 simulation framework, four strategies to simulate outliers extracted from
 366 literature [5, 6, 13, 3] will be redefined in terms of the proposed simulation
 367 framework. Figure 11 provides a graphical comparison between the simu-
 368 lated outliers following the original strategy from the aforementioned works
 369 and using the algorithm proposed in this article. Each simulation procedure
 370 and all the details to get the results presented in this section, are further
 371 explained in the Appendix A.

372 At first glance, one can notice in Figure 11 that despite sharing the pur-
 373 pose of simulating outliers, each strategy leads to very different outliers in

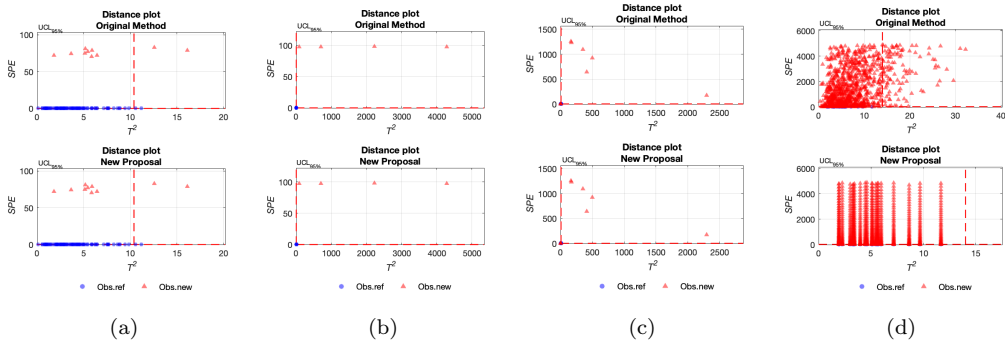


Figure 11: Distance plots of the observations simulated as in [5] (a), as in [6] (b), as in [13] (c) and as in [3] (d) using the approach from the original work (first row) and the proposed algorithm controlling the outlier properties (second row).

374 qualitative and quantitative terms. In Figures 11a and 11d outliers are far in
 375 terms of their orthogonal distance but their projection onto the model plane
 376 seems still under control limits. These plots differ from the ones reported in
 377 Figures 11b and 11c, where outliers are distant both in terms of the T^2 and
 378 in the SPE .

379 Furthermore, the simulation procedure from Figure 11c differs strategi-
 380 cally from the others, since the same set of observations is shifted apart in
 381 50 steps from their reference values. In Figure 11d, it can be appreciated the
 382 gradual shift of the same set of observations increasing their SPE and ran-
 383 domly shifting the T^2 . This can be seen as well in Figure A.13. It also stands
 384 out the difference between the upper and lower distance plots in Figure 11d.
 385 This is because we considered that variations of the T^2 in their simulated
 386 outliers were not a strategic feature of the simulation. This is explained in
 387 detail in the Appendix A.

388 Comparing the original methods to simulate outliers (upper row of plots
 389 in Figure 11), it can be seen that all of them increase the SPE of the out-
 390 liers. This is because in the end, despite following different strategies, all
 391 procedures to simulate outliers rely on breaking the correlation structure de-
 392 scribed by the reference data set. This is done differently by each author. In
 393 [5] the simulation strategy relies on adding noise to the outlying observations,
 394 whereas in [3], the noise is introduced as the new mean vector of the outly-
 395 ing distribution. This results in outliers with an increased SPE but with a
 396 moderate T^2 , as it can be seen in Figures 11a and 11d. In [6], outliers are
 397 generated by altering the variance of variables, which leads to an increase in

398 the T^2 (Figure 11b). The mean vector of the outliers distribution is changed
399 as well, in such a way that the correlation pattern is not respected anymore,
400 which leads to the increase of the SPE . Finally, in [13], authors shift the
401 sign of randomly selected cells. As a consequence, they are clearly breaking
402 the correlation structure and this can lead as well to an increase in the T^2 of
403 the outlying observations (Figure 11c).

404 The comparison between plots from the upper and lower row in Fig-
405 ure 11, shows that results obtained by the proposed algorithm to simulate
406 outliers with the desired properties, are fairly similar to the ones obtained by
407 other simulation settings. Furthermore, some limitations of the traditional
408 paradigm to simulate outliers can be seen as well. This traditional frame-
409 work relies on changing the parameters of the distribution that describes the
410 outlying population, but there is not a direct and clear relationship between
411 the new parameters of the outlying distribution and their effect on the SPE
412 or the T^2 . Consequently, it is difficult to control how this new distribution
413 will affect to the outlying properties of the outliers when they are projected
414 onto the reference PCA model. This can be appreciated in the fact that most
415 simulation strategies easily increase the SPE of their observations, but with-
416 out controlling its value and without having the same control over the T^2 of
417 the outliers. In fact, the T^2 seems to be a more uncontrolled parameter and
418 any of the proposals includes specific outliers for the T^2 . This is probably
419 because it is not trivial how to find a new mean vector for the distribution of
420 the outliers that still respects the correlation structure of the reference data
421 set.

422 The change from the traditional simulation paradigm, to the new one pro-
423 posed in this work, simplifies the relationship between the simulation setup
424 and the properties of the resulting outliers. The algorithm proposed in this
425 work does not rely on the distribution of the reference and the outlying ob-
426 servations and it has an independent control over the SPE and the T^2 . This
427 results in a new simulation approach that is versatile enough to encompass
428 other particular simulation strategies (Figure 11). Besides, differences be-
429 tween simulation settings can be directly measured in terms of the target
430 SPE and T^2 of the outliers.

431 3.2.2. *Properties of the simulated outliers in a robust PCA model*

432 The second aspect to assess in this comparison is to what extent (just
433 quantitative or also qualitative) outliers simulated by the proposed algorithm
434 behave as outliers in terms of other detection techniques. In this sense, it is

435 also interesting to assess if the properties of simulated outliers change when
436 they are expressed in terms of different distance metrics. For instance, some
437 robust PCA techniques differ not only in the core algorithm to calculate the
438 principal components, but also in terms of the statistics that measure the
439 distance of an observation to the model. Hence, the whole basis used by our
440 proposed framework to define the outliers, is different in these cases. This
441 may affect the properties of simulated observations when they are defined in
442 these new terms.

443 For this purpose, simulation scenarios from sections 3.1.1, 3.1.2 and 3.1.3,
444 will be projected onto a robust PCA model calculated with MacroPCA[3].
445 This technique can be considered as an ensemble of several outlier detection
446 methods. It includes the *Detect Deviating Cells* (DDC) [14] algorithm as first
447 step in order to detect outlying cells, which itself, can be regarded as an out-
448 lier detection technique. Later on, MacroPCA algorithm fits a robust PCA
449 model using a version of the ROBPCA algorithm [15]. In this ROBPCA step,
450 they include the detection of outlying observations in several steps. Firstly,
451 in the Projection Pursuit step, to rank rows according to their outlyingness.
452 Secondly, after the iterative subspace estimation, they apply a filter on ob-
453 servations based on their orthogonal distance to the model. Thirdly, they
454 apply the DetMCD method [16], for the covariance matrix estimation, which
455 also includes intermediate distance calculations to use the least distant ob-
456 servations for the covariance matrix computation. Finally, when the PCA
457 model has been estimated, they perform a last outlier detection based on two
458 robust distance metrics: the orthogonal distance and the score distance.

459 It is worth to highlight that although the distance metrics used in [3]
460 do not coincide with the SPE and T^2 , their conceptual meaning is equiva-
461 lent, since they represent the orthogonal distance and the Mahalanobis dis-
462 tance on the model plane, respectively. Thus, we considered that MacroPCA
463 was clearly representative as a state-of-the-art outlier detection method and
464 as a robust PCA model building algorithm. Moreover, its good perfor-
465 mance in outliers detection and the comparable meaning of its distance
466 metrics (orthogonal and score distances) to ours (the SPE and the T^2),
467 were considered as interesting factors for the comparison. Results shown
468 in Figure 12 were obtained using the *cellWise* package in R (available in
469 <https://CRAN.R-project.org/package=cellWise>).

470 As it can be seen in Figure 12, qualitative properties of the simulated out-
471 liers are still met in terms of alternative PCA models and distance metrics.
472 However, there are some differences in the distance values and their Upper

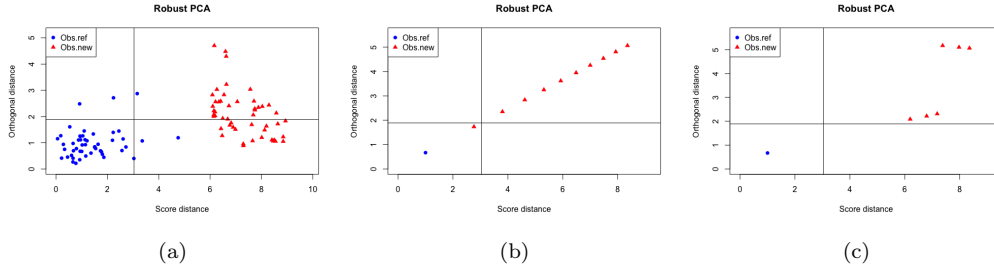


Figure 12: Distance plots of the observations simulated in Figure 8b (a), in Figure 9b (b) and in Figure 10b (c), when they are projected onto the PCA model fitted using MacroPCA with the reference data set. Blue circles represent reference observations, whereas red triangles represent the simulated outliers. Black lines represent the Upper Control Limits for the Orthogonal Distance (ordinate) and the Score Distance (abscissa).

473 Control Limits, which is reasonable given that the Orthogonal Distance and
 474 Score Distance are not exactly the SPE nor the T^2 . Results in Figure 12a
 475 show also an increase of the simulated outliers in terms of the orthogonal
 476 distance. Given the robust estimation of the covariance determinant in the
 477 detMCD step of MacroPCA, extreme observations in the T^2 were detected
 478 as outliers, and excluded for the computation of the final PCA model param-
 479 eters. As a result, since these observations were excluded at some point for
 480 the PCA model building, we find reasonable that they also increased their
 481 distance to the model. Nonetheless, in all cases simulated outliers keep the
 482 outlying character that they were asked to represent in first instance. This
 483 can be appreciated by their position above the cut-off values for the distances
 484 in all distance plots, indicating the persistence of their outlying properties.

485 4. Conclusion

486 In this work, a new framework to simulate outliers directly controlling
 487 their outlying properties has been proposed. This approach is based on the
 488 use of a well-known pair of statistics, the SPE and the Hotelling- T^2 from a
 489 PCA model, which evaluate in a complementary way how far an observation
 490 is from the majority of the data set (i.e. the outlyingness degree).

491 Given an observation with initial values for the statistics, a PCA model
 492 and target values for the pair of statistics, our simulation method drifts
 493 the aforementioned observation in a direction that shifts the initial SPE
 494 and Hotelling- T^2 until reaching their target values. This shift direction is a

495 combination of two orthogonal directions, each one independently controlling
496 the shift on the SPE and the Hotelling- T^2 .

497 This feature is a key factor, since it enables a specific control over the
498 two properties that define multivariate outliers in terms of a PCA model.
499 This becomes critical specially when simulating anomalous data, which is
500 a extremely common procedure when testing the performance of different
501 statistical methods handling datasets with outlying observations. However,
502 the outliers generation is usually an ad hoc procedure, with a lack of standard
503 protocols and being based most of the times, even when working with PCA
504 models, on distributions and parameters that do not tune neither how nor
505 how much an observation is outlying. This makes the supposed benefits
506 of the different statistical methods depend on the nature of the simulated
507 outliers and consequently, the comparison of the different methods reported
508 in the literature becomes difficult or impossible. Moreover, most simulation
509 methods require an assumption about the distribution of the reference data
510 set, and simulate outliers by changing one of its parameters, such as the
511 mean or the covariance matrix. This simulation paradigm might not be
512 feasible to implement with real data sets, when the distribution is unknown.
513 Furthermore, the relationship between the new parameters of the distribution
514 and the outlying properties of the simulated observations is not simple and
515 direct.

516 In Section 3.2.1 we showed how the methodology proposed in this arti-
517 cle, successfully encompasses particular simulation strategies proposed in the
518 literature in a common framework. Consequently, the comparison between
519 approaches can be easily measured in terms of target specifications or in
520 terms of the strategy followed to shift the outliers, i.e.: one step, step-wise or
521 grid-wise. Besides, we also illustrated the shortage of extreme (good lever-
522 age) outliers simulated in the literature given the difficulty of modifying the
523 reference distribution while respecting its covariance structure, which is eas-
524 ily achieved by the simulation framework proposed in this paper (Figure 8b).
525 Moreover, in Section 3.2.2 we also showed how the outlying properties are, at
526 least, qualitatively consistent when the simulated outliers are projected on a
527 robust PCA model.

528 However, the proposed method has some limitations, which are further
529 addressed in Appendix B. The simulation procedure does not set any restric-
530 tion in case that binary or categorical variables are present in the matrix.
531 Naturally, this framework is also restricted by the same limitations as the
532 PCA model is, such as the inability to model non-linear relations between

533 variables (see Appendix B.2).

534 In summary, the framework proposed in this paper offers the possibility
535 of generating outlying observations with a wide range of desired properties,
536 given that the user can control the pair of statistics that essentially define
537 the outlyingness degree: the *SPE* and the Hotelling- T^2 . This procedure has
538 been implemented in Matlab, providing a set of functions to perform the PCA
539 Model Building and the simulation of controlled outliers. Further details
540 about the Matlab code can be found in the documentation file available in
541 the GitHub repository.

542 Computational details

543 The results have been obtained executing the functions from [https://](https://github.com/albagc/SCOUTer.git)
544 github.com/albagc/SCOUTer.git in Matlab version R2020a 9.8.0.1323502.
545 Further information about the functions can be found in the *documentation*
546 and *howto* documents on the repository.

547 Acknowledgements

548 Financial support was granted by the UPV as a part of the FPI Grant
549 (Subprogramme I) under the UPV Research and Development support pro-
550 gramme and also by the Spanish Ministry of Economy and Competitiveness
551 under the project DPI2017-82896-C2-1-R.

552 References

- 553 [1] F. Arteaga, A. Ferrer, Dealing with missing data in MSPC: Several
554 methods, different interpretations, some examples, *Journal of Chemo-*
555 *metrics* 16 (2002) 408–418. doi:10.1002/cem.750.
- 556 [2] A. Smoliński, B. Walczak, J. W. Einax, Exploratory analysis of data
557 sets with missing elements and outliers, *Chemosphere* 49 (2002) 233–
558 245. doi:10.1016/S0045-6535(02)00326-0.
- 559 [3] M. Hubert, P. J. Rousseeuw, W. Van den Bossche, W. V. den
560 Bossche, MacroPCA: An all-in-one PCA method allowing for miss-
561 ing values as well as cellwise and rowwise outliers, *Technomet-*
562 *rics* (2018) 1–18. URL: [https://doi.org/10.1080/00401706.2018.](https://doi.org/10.1080/00401706.2018.1562989)
563 [1562989](https://doi.org/10.1080/00401706.2018.1562989). doi:10.1080/00401706.2018.1562989. arXiv:1806.00954.

- 564 [4] P. J. Huber, Robust Estimation of a Location Parameter, *The An-*
565 *nals of Mathematical Statistics* 35 (1964) 73–101. doi:10.1214/aoms/
566 1177703732.
- 567 [5] I. Stanimirova, M. Daszykowski, B. Walczak, Dealing with missing val-
568 ues and outliers in principal component analysis, *Talanta* 72 (2007)
569 172–178. doi:10.1016/j.talanta.2006.10.011.
- 570 [6] S. Serneels, T. Verdonck, Principal component analysis for data contain-
571 ing outliers and missing elements, *Computational Statistics and Data*
572 *Analysis* 52 (2008) 1712–1727. doi:10.1016/j.csda.2007.05.024.
- 573 [7] C. Agostinelli, A. Leung, V. J. Yohai, R. H. Zamar, Robust estima-
574 tion of multivariate location and scatter in the presence of cellwise
575 and casewise contamination, *Test* 24 (2015) 441–461. doi:10.1007/
576 s11749-015-0450-6. arXiv:1406.6031.
- 577 [8] A. Ferrer, Multivariate Statistical Process Control Based on Principal
578 Component Analysis (MSPC-PCA): Some Reflections and a Case Study
579 in an Autobody Assembly Process, *Quality Engineering* 19 (2007) 311–
580 325. doi:10.1080/08982110701621304.
- 581 [9] G. E. P. Box, Some Theorems on Quadratic Forms Applied in
582 the Study of Analysis of Variance Problems, I. Effect of Inequal-
583 ity of Variance in the One-Way Classification, *Ann. Math. Statist.*
584 25 (1954) 290–302. URL: [https://projecteuclid.org:443/euclid.](https://projecteuclid.org:443/euclid.aoms/1177728786)
585 [aoms/1177728786](https://projecteuclid.org:443/euclid.aoms/1177728786). doi:10.1214/aoms/1177728786.
- 586 [10] J. E. Jackson, G. S. Mudholkar, Control Procedures for Residu-
587 als Associated with Principal Component Analysis, *Technometrics*
588 21 (1979) 341–349. URL: <http://www.jstor.org/stable/1267757>.
589 doi:10.2307/1267757.
- 590 [11] L. Eriksson, T. Byrne, E. Johansson, J. Trygg, C. Vikström, Multi-
591 and Megavariate Data Analysis: Principles and Applications., 3rd ed.,
592 Umetrics Academy, 2001.
- 593 [12] P. Nomikos, J. F. MacGregor, Multivariate SPC Charts for
594 Monitoring Batch Processes, *Technometrics* 37 (1995) 41–59.
595 URL: [https://www.tandfonline.com/doi/abs/10.1080/00401706.](https://www.tandfonline.com/doi/abs/10.1080/00401706.1995.10485888)
596 [1995.10485888](https://www.tandfonline.com/doi/abs/10.1080/00401706.1995.10485888). doi:10.1080/00401706.1995.10485888.

- 597 [13] A. Folch-Fortuny, A. F. Villaverde, A. Ferrer, J. R. Banga, Enabling
598 network inference methods to handle missing data and outliers, *BMC*
599 *Bioinformatics* 16 (2015) 1–12. URL: <http://dx.doi.org/10.1186/s12859-015-0717-7>. doi:10.1186/s12859-015-0717-7.
600
- 601 [14] P. J. Rousseeuw, W. V. D. Bossche, Detecting Deviating Data
602 Cells, *Technometrics* 60 (2018) 135–145. URL: <https://www.tandfonline.com/doi/full/10.1080/00401706.2017.1340909>.
603 doi:10.1080/00401706.2017.1340909.
604
- 605 [15] M. Hubert, P. J. Rousseeuw, K. Vanden Branden, ROBPCA:
606 A New Approach to Robust Principal Component Analysis,
607 *Technometrics* (2005). URL: <https://pdfs.semanticscholar.org/250b/4f05982b491ad80ba8b986d958eedb69a6be.pdf>.
608 doi:10.1198/004017004000000563.
609
- 610 [16] M. Hubert, P. J. Rousseeuw, T. Verdonck, A deterministic algorithm for
611 robust location and scatter, *Journal of Computational and Graphical*
612 *Statistics* 21 (2012) 618–637. doi:10.1080/10618600.2012.672100.
- 613 [17] N. Locantore, J. S. Marron, D. G. Simpson, N. Tripoli, J. T. Zhang,
614 K. L. Cohen, Robust principal component analysis for functional data,
615 *Test* 8 (1999) 1–73.
- 616 [18] F. Arteaga, A. Ferrer, How to simulate normal data sets with the
617 desired correlation structure, *Chemometrics and Intelligent Labora-*
618 *tory Systems* 101 (2010) 38–42. URL: <http://dx.doi.org/10.1016/j.chemolab.2009.12.003>.
619 doi:10.1016/j.chemolab.2009.12.003.

620 **Appendix A. Comparison of the simulation method to other sim-** 621 **ulation strategies**

622 This section contains information about how to replicate the strategies
623 to simulate outliers present in different articles of PCA-MB dealing with
624 outliers. In each case there are two main items to simulate: the reference
625 data set used to fit the PCA model and the outlying observations.

626 The aim of this section is to provide a brief summary about the simulation
627 strategies followed in each reference and to give the details about the set up
628 of our proposed algorithm to imitate them, getting the results of Figure 11.
629 The following table provides information about the method used in each

630 referred work to simulate the reference data set and the outlying observations.
631 Some notation has been changed from the original works to avoid potential confusions with other terms used in this paper.

Ref.	Simulation of reference data set	Simulation of outliers
1 - [5]	$\mathbf{X}_0 \sim N_n(\mathbf{0}_n, \mathbf{I}_n) \rightarrow \mathbf{X}_0 = \mathbf{T}_A \mathbf{P}_A^\top + \mathbf{E}_0$ $\mathbf{E}_1 \sim N(\mathbf{0}, \mathbf{1}) \cdot 0.1$ $\mathbf{X}_1 = \mathbf{T}_{1,A} \mathbf{P}_A^\top + \mathbf{E}_1$ $n = 98; k = 20; A = 4$	$\mathbf{X}_2 = \mathbf{T}_{2,A} \mathbf{P}_A^\top + \mathbf{E}_2$ $\mathbf{E}_2 \sim N(\mathbf{10}, \mathbf{1})$
2 - [6]	$\mathbf{T}_A \sim N_A(\mathbf{0}_A, \mathbf{I}_A)$ $\mathbf{P}_A : \perp k \times A$ uniformly distributed pseudorandom numbers $\mathbf{E}_k \sim N_k(\mathbf{0}_k, \mathbf{1}_k)/100$ $\mathbf{X}_1 = \mathbf{T}_A \mathbf{P}_A^\top + \mathbf{E}$	$\mathbf{X}_2 \sim N_A(\mathbf{15}_A, 8 * \mathbf{I}_A)$
3 - [13]	\mathbf{X}_1 : Data reconstructing the metabolic network of the benchmark problem 4 from the original work	\mathbf{X}_2 : outliers $\forall i \in 1, \dots, n_2$ $\forall j \in 1, \dots, k$ $m_j = \text{mean}(\mathbf{x}_j)$ $s_j = \text{std}(\mathbf{x}_j)$ if $ x_{ij,1} \leq m_j + 1.5s_j$: $x_{ij,2} = -x_{ij,1}$
4 - [3]	$\mathbf{X}_1 \sim N(\mathbf{0}, \boldsymbol{\Sigma}_{A09})$ $A = 6PCs, N = 100, K = 200$ $\boldsymbol{\Sigma}_{A09} = \mathbf{V}_{A09} \mathbf{D}_{A09} \mathbf{V}_{A09}^\top$ $\mathbf{D}_{A09} = \text{diag}(30, 25, \dots, 5, 0.098, 0.0975, \dots, 0.005)$	$\mathbf{X}_2 \sim N(m\boldsymbol{\nu}_{A+1}, \boldsymbol{\Sigma}_{A09})$ $m \in 1, \dots, 50$ $\boldsymbol{\nu}_{A+1} = \mathbf{V}_{A09}[:, A+1]$

Table A.1: Strategies followed by different authors to simulate the reference data sets and the outlying observations.

632

633 Appendix A.1. Reference 1

634 In [5], an adaptation of the classical Expectation Maximization PCA
635 (EM-PCA) is provided. The core least squares PCA is substituted by the
636 robust spherical PCA. [17] Their reference data set is a matrix of dimen-
637 sions $n = 98$ and $k = 20$, with its variables following a multivariate normal
638 distribution. After its reconstruction with A principal components, an error
639 term following a $N(0, 1)$ distribution is added to regular observations. A
640 certain percentage of observations is randomly sampled and instead, their

641 added error term follows a $N(10, 1)$ distribution. Before its addition to the
 642 reconstructed matrix, the error term is multiplied by a 0.2 factor.

643 Following one of the simulation settings from the original work, we set
 644 a number of $A = 4$ PCs and selected 10% of the observations to transform
 645 them into outliers. After building a PCA on the clean data, the outliers
 646 simulated as in the original work, were projected onto the model, obtaining
 647 their SPE and T^2 . These metrics were used as an input to our simulation
 648 function, setting them as target values for the algorithm.

649 The code used to generate the outliers is the following one:

```
650 pcamodel_ref1 = pcamb_classic(Xref1, 4, 0.05, 'cent');
651 pcaxout = pcame(Xoutref1, pcamodel_ref);
652 Xout1 = scout(Xoutref1_0, pcamodel_ref, 'simple', 'spey',
653 pcaxout.SPE, 't2y', pcaxout.T2);
```

654 The elements `Xref1`, `Xoutref1` and `Xoutref1_0` are matrices containing the
 655 reference data set used to build the PCA model, the outliers simulated as
 656 in the original work and the outlying observations before being transformed
 657 into outliers, respectively.

658 *Appendix A.2. Reference 2*

659 In [6], the simulation procedure begins by simulating the latent subspace.
 660 On one hand, scores are simulated as independent normally distributed vari-
 661 ables with zero mean and unitary variance. On the other hand, loadings are
 662 simulated as orthogonal vectors with pseudo-random uniformly distributed
 663 pseudo-random numbers. In third place, the error matrix is simulated as
 664 normally distributed random noise divided by 100. Using this terms, the
 665 reference matrix (\mathbf{X}_1) can be reconstructed.

666 Afterwards, the outlying data set is simulated as a matrix \mathbf{X}_2 , with $n_2 =$
 667 $0.1 \cdot n_1$ observations, and which follows a normal distribution $N_A(\mathbf{15}_A, 8\mathbf{I}_A)$,
 668 where $\mathbf{15}_A$ is a vector containing A elements equal to 15.

669 Following this simulation procedure to generate outliers, authors create
 670 three different setups (C_1, C_2 and C_3) varying the number of observations
 671 ($n_{C_1} = 100, n_{C_2} = 40, n_{C_3} = 40$), the number of variables ($k_{C_1} = 5, k_{C_2} =$
 672 $10, k_{C_3} = 200$) and maintaining the number of PCs ($A_{C_1} = A_{C_2} = A_{C_3} = 2$).
 673 Results shown in Figure 11b are the ones obtained with the configuration B .

674 The following lines of code were used to replicate the PCA model of the
 675 reference data set and the outliers simulated in this work:

```

676 pcamodel_ref2 = pcamb_classic(Xref2, 2, 0.05, 'cent');
677 pcaxout= pcame(Xoutref2, pcamodel_ref2);
678 Xout2 = scout(Xoutref2_0, pcamodel_refB, 'simple', 't2y',
679   pcaxout.T2,'spey',pcaxout.T2);

```

680 The elements `Xref2`, `Xoutref2` and `Xoutref2_0` are matrices containing the
681 reference data set used to build the PCA model, the outliers simulated as
682 in the original work and the outlying observations before being transformed
683 into outliers, respectively. In this case, observations from `Xoutref2_0` were
684 observations following the same distribution as the reference data set.

685 *Appendix A.3. Reference 3*

686 In [13], authors provide a solution based on Trimmed Scores Regression
687 (TSR) to enable network inference methods work in presence of missing data
688 and outliers. Outlying observations are simulated by shifting the sign of cells
689 above the variable average plus 1.5 times the standard deviation, or below
690 the mean minus 1.5 times the standard deviation. Thus, the correlation
691 pattern between variables is broken for these outliers. In order to illustrate
692 the results, we show the outliers generated for the benchmark problem 4, one
693 of the five benchmark problems addressed in the original paper.

694 The original work provides the information about the data used in the
695 article. After downloading it, we built a PCA model based on the `Xref3` ma-
696 trix, and projected. Afterwards, a random selection of rows determined the
697 observations that were transformed to outliers using the original procedure
698 in [13] and the algorithm proposed in this article.

699 The process to achieve this simulation framework is very similar to the
700 ones from previous references, where the *SPE* and the T^2 of the outliers
701 generated following the procedure from the original work are used as target
702 values in our simulation function.

```

703 pcamodel_ref3 = pcamb_classic(Xref3, 3, 0.05, 'cent');
704 pcaxout = pcame(Xoutref3, pcamodel_ref3);
705 Xout3 = scout(Xoutref3_0, pcamodel_ref3, 'simple', 'spey',
706   pcaxout.SPE, 't2y',pcaxout.T2);

```

707 The elements `Xref3`, `Xoutref3` and `Xoutref3_0` are matrices containing
708 the reference data set used to build the PCA model, the outliers simulated as
709 in the original work and the outlying observations before being transformed
710 into outliers, respectively.

711 *Appendix A.4. Reference 4*

712 In [3] authors propose an adaptation of their previous robust PCA algo-
 713 rithm to deal with missing data and cellwise outliers. However, in this work
 714 we are focusing exclusively in the comparison between rowwise outliers, i.e.
 715 anomalous observations. In this case the reference data set is generated as
 716 a matrix whose columns follow a multivariate normal distribution $N(\mathbf{0}, \mathbf{\Sigma})$.
 717 Two different covariance matrices (A09 and ALYZ) are used and their singu-
 718 lar values are adapted in order to reach over the 80% of explained variance
 719 with the first 6 principal components.

720 Later on, a certain percentage of rows is randomly sampled and changed
 721 by new observations that follow the distribution $N(m\boldsymbol{\nu}_{A+1}, \mathbf{\Sigma})$. In the pre-
 722 vious expression, A is the number of principal components and the term $\boldsymbol{\nu}_j$
 723 refers to the j th eigenvector of the covariance matrix. The factor m that
 724 multiplies the new mean vector ranges from 1 to 50, leading an increasing
 725 noise introduced in the outliers along with the increase in the m parameter.
 726 This is equivalent to make outliers more distant to the model hyperplane.

727 The following code lines show the procedure used to imitate the simulation
 728 with the A09 covariance matrix. The matrix dimensions are $n = 100$ observa-
 729 tions and $k = 200$ variables, with $A = 6$ principal components, $m = 1, \dots, 50$
 730 and 20 rows randomly selected to transform them into outliers.

```
731 pcamodel_ref4 = pcamb_classic(Xref4, 6, 0.05, 'cent');
732 pcaxout4 = pcame(Xout4, pcamodel_ref4);
```

733 The matrix X_{out3} contains the outlying rows for all the values of the
 734 step parameter m , i.e. it is a matrix of 1020 rows (20×51 , the 20 original
 735 observations and their progressive 50 shifts). Vectors \mathbf{SPE}_y^2 and \mathbf{T}_y^2 contain
 736 the SPE and T^2 values of the 20 outliers along the 50 steps.

737 A characteristic aspect of this simulation is the gradual shift described by
 738 the outliers. In terms of our proposed procedure, this is equivalent to use the
 739 step-wise generation of outliers as in Section 3.1.2. For this purpose, we need
 740 the final values of the statistics at the $m = 50$ step for all the observations,
 741 but also the step parameter γ . In order to study the progression for the SPE
 742 and the T^2 along the m steps, we plotted their evolution. in Figure A.13.

743 After visualising Figure A.13, it stands out a clear difference between the
 744 growing patterns of the SPE and the T^2 along m . Whereas the SPE trajec-
 745 tory for the outliers draws a clear ascending pattern for all the observations,
 746 the T^2 does not seem to do so.

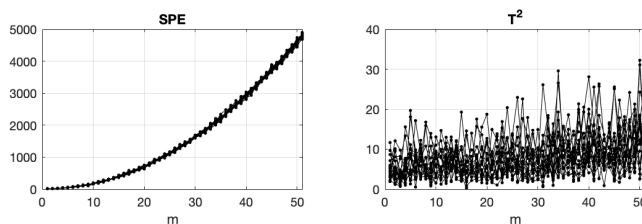


Figure A.13: SPE (left) and T^2 evolution of the outliers simulated in [3] along the m steps.

747 Moreover, when the outliers are visualised in the distance plot from Fig-
 748 ure 11c, barely any of them is above the UCL for the T^2 . This lead us to
 749 consider that changes in the T^2 among the outliers were more an artifact than
 750 a desired outcome of the simulation. Hence, we focused on calculating the
 751 γ_{SPE} parameter only. This parameter appears in Figure 3a and Equation 18
 752 tuning the spacing between the SPE steps in the following expression.

753 In order to fit the γ_{SPE} parameter, we used the MatLab function `lsqnonlin.m`:

```

754 sperw = reshape(pcaout4.SPE,20,51)';
755 xg_spe = nan(20,1);
756 for i = 1:20
757     HM = sperw(end,i);
758     H0 = sperw(1,i);
759     Hm = sperw(:,i);
760     M = 50;
761     m = 0:50;
762     gfun = @(gamma)H0 + (m./M).^gamma*(HM - H0) - Hm;
763     xg_spe(i) = lsqnonlin(gfun,3);
764 end
765 g_spe_mean = mean(xg_spe);

```

766 After calculating the γ_{SPE} for each observation, the mean value is calcu-
 767 lated and stored in the `g_spe_mean` variable. Figure A.14 shows the estimated
 768 trajectory using the average $\gamma_{SPE} = 2.6348$ value. This parameter used later
 769 as an input to the `scout.m` function:

```

770 Xout4 = scout(Xoutref4_0, pcamodel_ref4, 'steps', 'spey',
771 sperw(end,:), 'nsteps', 50, 'gspe', g_spe_mean);

```

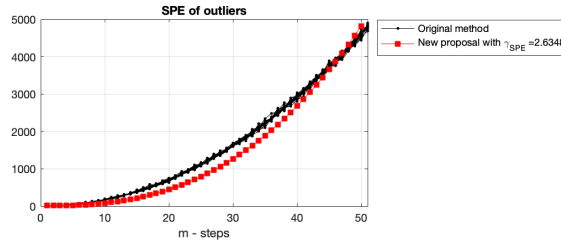


Figure A.14: SPE of outliers simulated by the strategy from the original work [3] (black) and SPE of the outliers simulated by the proposed algorithm (red) along the m steps.

772 Appendix B. Limitations of the proposed algorithm

773 This section addresses in further detail the results obtained with the
 774 method to simulate outliers with desired properties, when it is used on a
 775 matrix with non-linearities or with binary data.

776 The reference matrix \mathbf{X}_0 is simulated using the functions from [18]. The
 777 following code lines are the ones used to generate the reference matrix:

```
778 [X,S,srnd] = simdataset(100,10,[6,3],ones(1,10));
779 [X_0,srndn]=randnm(S,100,srnd);
```

780 The resulting matrix has 100 observations, 10 variables normally dis-
 781 tributed and two principal components which explain above the 80% of the
 782 variance.

783 *Appendix B.1. Non linearities*

784 In this case, the matrix will present relations between variables that the
 785 classical PCA model will not be able to capture. In order to study to what
 786 extent this limitation of the PCA model would affect the simulations, we
 787 carried out a generation of outliers with a reference matrix that contained
 788 non-linearities and increasing only the T^2 of the outliers. This means that the
 789 generated observations should nor break the correlation pattern described by
 790 variables.

791 The new matrix \mathbf{Y} is the result of concatenating the original matrix \mathbf{X}_0 ,
 792 and a set of non-linear variables generated from the original ones in \mathbf{X}_0 . The
 793 non-linear relations included in each variable are:

```
794 rng(1101)
795 varind = randperm(10,8);
```

```

796 Y_11 = X_0(:,varind(1)).^2;
797 Y_12 = X_0(:,varind(2)).^3;
798 Y_13= exp(X_0(:,varind(3)));
799 Y_15 = rand(1,1) + X_0(:,varind(5)) + X_0(:,varind(5)).^2;
800 Y_16 = X_0(:,varind(2)).*X_0(:,varind(4));
801 Y_17 = X_0(:,varind(6)).*X_0(:,varind(7)).^2;
802 Y_18 = exp(X_0(:,varind(3))).^(X_0(:,varind(7)) + X_0(:,varind(8)));
803 Y_19 = X_0(:,varind(3))*2;
804
805 Y = [X_0,Y_11,Y_12,Y_13,Y_14,Y_15,Y_16,Y_17,Y_18,Y_19];
806

```

807 As one can notice, the selection of the variables that were non-linearly
808 combined was perform randomly. Also, a linearly generated variable (y_{19})
809 was included in the set, to compare if the outliers on this variable still followed
810 their analytic relation with the column. used to generate them.

811 As we aforementioned, some outliers on the T^2 were generated to keep
812 the original correlation structure between variables. In order to do so, the
813 PCA reference model based on \mathbf{Y} had to be calculated. By setting “0” as
814 the second input argument in the PCA-MB function, it returns a suggestion
815 about the number of PCs to consider:

```

816 pcamodel_ref = pcamb_classic(Y, 0, 0.05, 'cent');
817
818 Sugested number of PCs:
819 - Singular values of covariance matrix > 1 = 6
820 - Minimum PCs to reach cummulative variance > 80 % = 3
821 Select the number of PCs: 3
822

```

823 A number of 3 PCs was selected. Then, outliers on the T^2 were generated
824 setting the same target value for all of them in the *scout.m* function:

```

825 T2target = 60*ones(size(Y, 1), 1);
826 Yextreme = scout(Y, pcamodel_ref, 'simple', 't2y', T2target);
827 Yall = [Y; Yextreme.X];
828

```

829 The resulting outliers are represented in Figure, B.15 where it can be seen
830 that the new observations accomplish the specified target values for the T^2 .

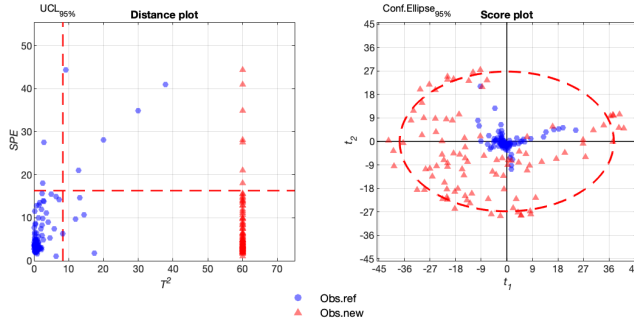


Figure B.15: Distance (left) and score (right) plot for the reference (blue circles) and the outliers (red triangles) generated.

831 However, the relations between the non-linear variables and the original
 832 columns use to generated have been distorted. In Figure B.16 there is a clear
 difference between blue and red observations.

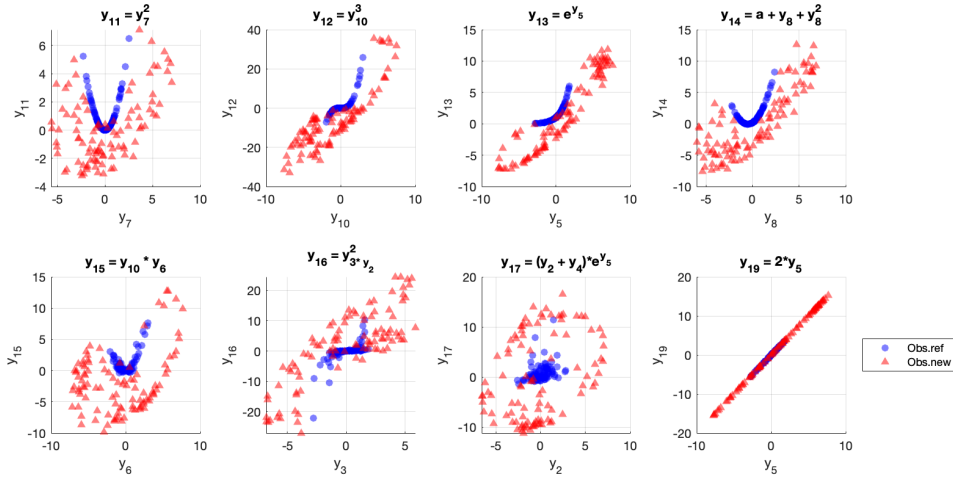


Figure B.16: Scatter plots with the reference (blue circles) and new (red triangles) observations for all the new variables in \mathbf{Y} generated as combinations of the variables in \mathbf{X}_0 .

833
 834 Whereas the blue circles perfectly describe the analytical relation used
 835 to generate them, that is not the case for red triangles, since they clearly
 836 break the relative pattern between variables. This is not the case for the
 837 last variable (\mathbf{x}_{19}), which was generated as a linear combination. This result

838 reinforces the limitation that is produced when the method has to take into
839 account non-linear relations between the variables.

840 *Appendix B.2. Binary variables*

841 In this second example the purpose is to show the changes produced on
842 categorical variables when the algorithm is used on a mixed matrix with
843 continuous and categorical data.

844 In this case, four binary variables with different percentage of 0s and 1s
845 are simulated. The resulting matrix \mathbf{Y} has the original variables from \mathbf{X}_0
846 and the four additional binary columns.

```
847 rng(1101)
848 Y = [X_0,zeros(size(X_0,1),4)];
849 Y(randperm(size(X_0,1),round(0.2*size(X_0,1))),11) = 1;
850 Y(randperm(size(X_0,1),round(0.4*size(X_0,1))),12) = 1;
851 Y(randperm(size(X_0,1),round(0.6*size(X_0,1))),13) = 1;
852 Y(randperm(size(X_0,1),round(0.8*size(X_0,1))),14) = 1;
```

853 Similarly as in Appendix B.1, a PCA model is fitted with \mathbf{Y} , but in this
854 case, two PCs were selected.

```
855 pcamodel_ref = pcamb_classic(Y, 0, 0.05, 'cent');
856 Suggested number of PCs:
857 - Singular values of covariance matrix > 1 = 2
858 - Minimum PCs to reach cummulative variance > 80 % = 2
```

859 In this case we generated outliers increasing the SPE and the T^2 , im-
860 posing a target value of 50 for both of them and for all the data points. As
861 it can be seen in Figure B.17, the set of new observations has the specified
862 values for both statistics.

```
863 T2target = 50*ones(size(Ybin, 1), 1);
864 SPETarget = 50*ones(size(Ybin, 1), 1);
865 Yout = scout(Ybin, pcamodel_ref, 'simple', 't2y', T2target,'spey',SPETarget);
866 Yall = [Ybin; Yout.X];
```

867 Nonetheless, it is easy to see in Figure B.18 that new observations are
868 outside the range of accepted values for binary variables. This artefact is
869 produced because the simulation algorithm assumes to work with continu-
870 ous variables. Consequently, it does not include any constraint in the data
871 generation to respect the binary or qualitative nature of variables.

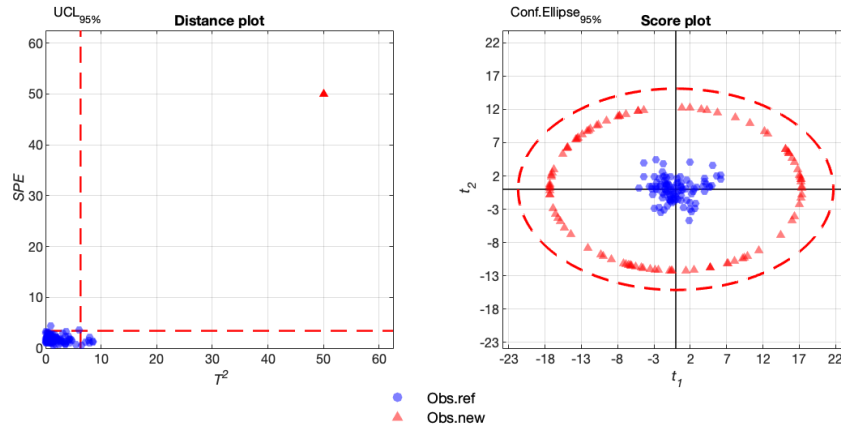


Figure B.17: Distance (left) and score (right) plot for the reference (blue circles) and the outliers (red triangles) generated.

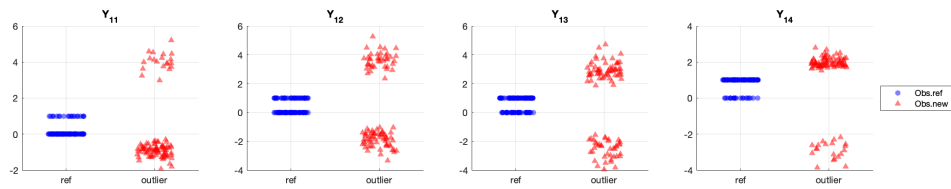


Figure B.18: Distance (left) and score (right) plot for the reference (blue circles) and the outliers (red triangles) generated.