



**HAL**  
open science

# Aberrant measurements: Detection, localization, suppression, acceptance and robustness

José Ragot

► **To cite this version:**

José Ragot. Aberrant measurements: Detection, localization, suppression, acceptance and robustness. Measurement, 2021, 172, pp.108872. 10.1016/j.measurement.2020.108872 . hal-03168381

**HAL Id: hal-03168381**

**<https://hal.science/hal-03168381>**

Submitted on 2 Jan 2023

**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.



Distributed under a Creative Commons Attribution - NonCommercial 4.0 International License

# Aberrant measurements: detection, localization, suppression, acceptance and robustness

José Ragot<sup>a</sup>

<sup>a</sup>Centre de Recherche en Automatique de Nancy,  
UMR 7039 - Nancy-Université, 2, Avenue de la Forêt de Haye, 54516 Vandœuvre, France.

5

---

## Abstract

The detection of outliers in a series of measurements, but even more so their location, is a necessity when these measurements are to be used in a monitoring system. This detection/localisation can only be done if redundant information is available, which may be based on the model of the system on which the measurements were collected.

10

In some cases, however, it is not necessary to detect and locate outliers. Instead, a robust approach to their use may be preferred, one that minimizes the influence of these outliers, such as using a median rather than a mean.

In this paper, the focus will be on the notion of robustness through a few examples and notably by proposing extensions to two well-known data processing techniques (data reconciliation and principal component analysis). The numerical examples proposed clearly show how to implement these two techniques and how to use them in a system monitoring procedure.

15

*Keywords:* Outliers, detection, deletion, acceptance, redundancy, robustness

---

## 1. Introduction

Some general considerations on the problems related to the presence of outliers will be followed by a short presentation of work related to the treatment of outliers, and then the plan adopted for this paper.

20

### 1.1. Positioning of the problem to be solved

Due to the intensive use of data (especially from sensors), but also to the increase in their volume and their use in monitoring and control tools, the problems resulting from the presence of outliers have taken on considerable importance in recent decades [54, 69].

25

How to recognize outliers? They are commonly defined as observations that appear to be inconsistent with the main part of the data set, or as observations that deviate significantly from the model postulated. Some historical definitions of outliers, which are indeed still relevant, are also worth noting. Barnett and Lewis [5] indicate that an outlier is an observation that appears to deviate significantly from other members of the sample in which it occurs. Hawkins [31] defines an outlier as an observation that deviates so far from other observations that it raises suspicions that it was generated by a different mechanism.

30

**Definition 1** (Residues as indicators of outliers). *The fundamental principle of model-based fault detection is based on the estimation of the state of the system from the available measurements. The resulting estimation error constitutes the residual vector. Subsequently, a decision on whether or not a measurement inconsistency is present is made by comparing this residue to a given threshold. For a signal generated from measurements collected on a system to be a true residue, it must be sensitive to measurement inconsistencies. In a general way, the generation of a residual is based on a data transformation to generate a residual capable of revealing at best the deviation from a reference situation. In the following, the generation of the residue will be at the heart of the approaches presented.*

35

40

---

Email address: [jose.ragot@univ-lorraine.fr](mailto:jose.ragot@univ-lorraine.fr) (José Ragot)

In any case, the detection of outliers can only be done if redundant information is available, which may be of hardware or software origin. Hardware redundancy comes from the simultaneous use of several sensors to measure the same quantity. The detection of an aberrant measurement, thus coming from a faulty sensor, is then based on the comparison of the measurements between them and the use of a majority vote. Software redundancy exploits the properties of the model of the system on which the data are collected. Subject to structural conditions to be met, the adequacy of the measurements with respect to the model can be tested to determine whether or not there are any measurement errors.

In the literature, there are currently three research communities dealing with the problem of fault diagnosis: the FDI (Fault Detection and Isolation) community, whose methodological tools are largely based on the synthesis of dynamic diagnostic filters, the DX (Diagnosis) community, whose foundations come from the fields of computer science and artificial intelligence, and the signal processing community, whose tools are based on statistical signal processing and also on pattern recognition. Even if there are common principles between the three communities, such as the use of models or the generation of alarm signals, each has focused on the development of its own terminologies, calculation tools and methodological approaches, guided by different constraints and objectives. Similarly, the modeling formalisms are very different for each domain; for example, the models of the FDI community are often based on differential algebra, while those of the DX community are mostly symbolic and qualitative. Although links exist between these three communities, which obviously have very similar objectives in the field of diagnosis, it is nevertheless true that specific vocabularies remain, (to this day, there is still no consensus on the terms used) such as for example: outliers, faults, defects, anomalies, discordant observations, peculiarities or contaminants, which nevertheless designate very similar facts.

### 1.2. Related works

The wide variety of available methods based on well-established statistical tools [53, 67], has made it possible, in many practical applications, to use techniques to handle measurement inconsistencies. This article is not intended to be an inventory of usable tools, but simply to raise the reader's awareness of the problem of outliers through a few simple examples and situations, in particular by focusing on two model-based methods, namely measurement reconciliation and principal component analysis. Of course, many other model-based techniques would need to be examined, including those using linear regression models [57] and their variants [15], nonlinear regression techniques [41], machine learning [43], Bayesian models [54]. Among the methods that have been developed to deal with outliers is that of Rousseeuw and Hubert [59]. The latter consists first of all in constructing an adjustment that is robust to them, generating residuals with respect to this adjustment, and then analyzing these residuals to identify outliers. A very large number of works have been published on this subject, and among the most recent are the following : [68, 49, 17, 60, 26, 34].

Moving beyond outlier detection, the reader may be interested in the more general problem of detecting anomalous series in relation to a set of series. In [10] the author has been interested in the detection of unsupervised anomalies in uni- and multi-variate time series with a particular application for masses of data in the field of tyres. In [30] the author deals with the detection of multiple breaks in a signal with an extension to the case of multiple breaks in several synchronized time series. Some authors [3, 21] supplement this classification of methods with techniques based on proximity concepts (data classification techniques for example). In [11] a taxonomy is presented based on the main aspects that characterize an outlier detection technique.

For industrial applications and the processing of large volumes of data, the emphasis is often placed on the online detection of outliers directly related to security, safety and production quality monitoring issues. As an example, in a wide variety of fields, we can refer to the following work [6] for telecommunications, [70] for tyre quality assessment, [62] in the field of aeronautics, [52] for sensor networks, [55] in the field of mineral processing, [1] in the field of chemical engineering. Much more broadly, it can be argued that no area is unaffected by the problem of dealing with outliers, particularly in the context

of societal applications : fraud detection, intrusion detection, face detection, video surveillance, social media analysis. 90

### 1.3. Positioning of the present work

In the following, the treatment of outliers will be approached through three complementary functions : fault detection (the determination of the presence of faults in a system and the time of occurrence of these faults), fault isolation (determination of the exact location of a fault), fault identification (determination of the size of a fault), which are often performed sequentially, but which, depending on the technique used, can also be performed concomitantly. 95

It is also necessary to specify, on the one hand, the nature of the systems to which this presentation is addressed and, on the other hand, the assumptions made about outliers. The systems will be considered stationary as well as the noises that affect their measurements. Therefore, situations where the internal structure of the systems may vary or variations in the internal parameters of the systems are excluded from the proposed approaches. As for outliers, they concern amplitude biases on their measurements. We therefore exclude outliers that may come from other sources of disturbance such as changes in variance or frequency. 100

As the objective here is to detect/locate/identify outliers in interdependent multivariate systems, these measures need to be treated globally, which obviously excludes single signal approaches. The strategies generally used are all based on the synthesis of residuals indicating the presence of outliers ; the precise localisation of the outliers needs a particular structuring of these residuals. As shown in the two examples in sections (3.3) and (4.5), the analysis of these residuals is done by fairly classical approaches, which may use jump tests or classification techniques. 105

Section 3 presents this possibility through a procedure known as data reconciliation. The presentation aims, on the one hand, to give some specific references to the users of the measurements and, on the other hand, to insist on the robustness of certain techniques for processing the measurements [44, 24, 2, 28, 37]. The robust reconciliation procedure that we propose here has the advantage of being carried out in conjunction with the detection/localization of outliers. 110

Section 4 is based on a well-known tool, Principal Component Analysis, but is revisited by proposing a relatively unknown technique for reconstructing variables from a selection of available variables and projecting them into the so-called residual space in order to obtain structured indicators for detecting and isolating outliers. 115

## 2. Some classical approaches for abnormal values detection

This section illustrates, using simple examples, some approaches to highlighting abnormal data. After illustrating some types of outliers (section 2.1), the academic example of a two-equation system introduces the notion of residual and structured residual (section 2.2), which then allows the general case of a multivariate system in a static regime to be presented (section 2.3). The next two sections 2.4 and 2.5 specify the approaches for replacing outliers and then for their acceptance. 120

### 2.1. Different types of outliers 125

Figure 1 shows two types of error [5] in the case of a two-dimensional variable. Bias affecting both directions of measurement and random errors are shown. The sensors that delivered the measurements can be characterized by their accuracy and fidelity (these two qualities being possibly associated) as shown in Table 1. Depending on the use made of them, these measurements can be processed in such a way as to reduce the influence of the two sources of error. Obviously, the situation is more complex when we are interested in a network of sensors equipping a physical process, as the faults that affect them may be related to each other. 130

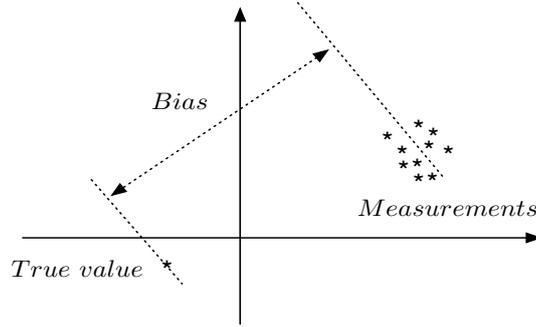


Figure 1: Various types of errors.

Sensor	Systematic errors	Random errors
true and unfair	minor	important
wrong and fair	important	minor
true and fair (precise)	minor	minor
wrong and unfair	important	important

Table 1: Quality of a sensor

## 2.2. Active approach: detection of outliers in multiple time series

This term refers to a set of techniques that can detect and locate outliers in a series of observations.

135 Once localized, these can either be removed, which may subsequently cause some processing difficulties, or replaced by so-called substitution values obtained for example by interpolation using healthy measurements close to those that have been removed. Numerous statistical tests have been developed for outlier detection (Dixon, Grubbs, Cochran, Tukey, Chauvenet, Tietjen-Moore, Student, Thompson ...)

140 techniques and the reader is invited to refer to the historical references [19, 27, 7, 66, 25], and to more recent ones [58, 14, 45, 35, 48, 63], but also works of synthesis [31, 5, 4]. The techniques mentioned above are generalized to the multivariate case where the detection involves the analysis of measurements of several variables coupled by a model.

As an example let us consider the system characterized by four variables and described by the model :

$$\begin{aligned} x_1^2 - x_2 + 2 \log x_3 &= 0 \\ x_2 + x_3 - x_4 &= 0 \end{aligned} \quad (1)$$

145 The twenty available measures of  $x_i, i = 1, \dots, 4$  are grouped in Table (2) and we wish to establish a diagnosis of the consistency of these data. As this is a simulation, one fault affects  $x_2$  for observation 12 and another affects  $x_3$  for observation 6. To make a diagnosis, the model residuals were calculated from the measurements, i.e. :

$$\begin{aligned} r_1 &= x_1^2 - x_2 + 2 \log x_3 \\ r_2 &= x_2 + x_3 - x_4 \end{aligned} \quad (2)$$

In order to improve diagnostic efficiency, the equations (2) can be combined in an additive manner, which removes the variable  $x_2$  (a similar approach is used for the variable  $x_3$ , which could also be removed):

$$x_1^2 + 2 \log x_3 + x_3 - x_4 = 0 \quad (3)$$

and makes it possible to evaluate the residue  $r_3 = x_1^2 + 2 \log x_3 + x_3 - x_4$  but without using that of  $x_2$ . The validity of the measurements is highlighted by the analysis of the magnitude of the absolute values of the three residuals ( $r_i, i = 1, 2, 3$ ) whose values are shown in Table (2) and graphically illustrated in

	$x_1$	$x_2$	$x_3$	$x_4$	$ r_1 $	$ r_2 $	$ r_3 $
1	0.86	0.61	0.95	1.54	0.02	0.02	0.03
2	1.24	1.49	0.98	2.45	0.01	0.01	0.02
3	2.15	2.01	0.29	2.30	0.12	0.00	0.12
4	1.62	2.17	0.80	2.96	0.01	0.01	0.02
5	1.50	2.05	0.90	2.95	0.00	0.00	0.00
6	1.69	1.81	0.95	2.39	0.93	0.37	1.30
7	1.31	1.48	0.89	2.37	0.02	0.00	0.02
8	1.13	1.17	0.94	2.13	0.01	0.02	0.02
9	1.45	0.91	0.55	1.45	0.01	0.00	0.00
10	1.15	0.68	0.73	1.40	0.02	0.01	0.03
11	1.27	0.50	0.58	1.07	0.04	0.00	0.04
12	2.77	0.86	0.03	0.38	0.49	0.50	0.00
13	1.37	0.25	0.44	0.69	0.00	0.00	0.01
14	1.02	0.18	0.65	0.82	0.00	0.01	0.01
15	1.20	0.12	0.52	0.64	0.01	0.00	0.01
16	1.45	0.09	0.37	0.46	0.04	0.00	0.04
17	0.43	0.06	0.94	0.99	0.00	0.00	0.00
18	0.64	0.04	0.83	0.87	0.01	0.00	0.01
19	0.59	0.03	0.85	0.88	0.00	0.00	0.01
20	1.42	0.02	0.37	0.39	0.01	0.00	0.01

Table 2: Measures available over time and model residuals

	$\delta x_1$	$\delta x_2$	$\delta x_3$	$\delta x_4$
$r_1$	×	×	×	.
$r_2$	.	×	×	×
$r_3$	×	.	×	×

Table 3: structured faults signature

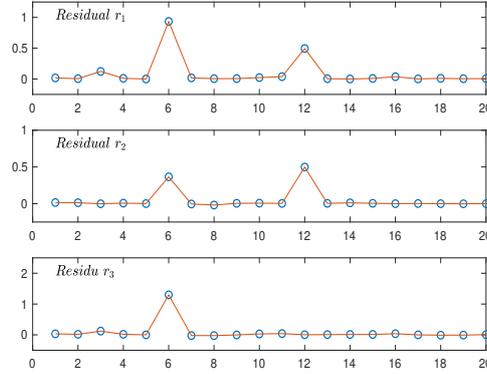


Figure 2: Model Residuals

the figure (2). The three model residuals are essentially zero except for the two observations 6 and 12. **155**  
The 6 observation triggers all three residuals significantly, while the 12 observation triggers only the  $r_1$   
and  $r_2$  residuals. This is explained by the table (3) of signatures of possible faults  $\delta x_i$  whose role is  
fundamental to characterize the detectability and the isolability of the measurement faults. The crosses  
show the sensitivity of the residues to the faults, the absence of sensitivity being marked by a dot. The  
sensitivities of the residuals with respect to the faults  $\delta x_i$  are all different which shows the isolability of **160**  
these measurement faults, which would not be possible if only the first two residuals  $r_1$  and  $r_2$  had been

used because in this case  $\delta x_2$  and  $\delta x_3$  would have had the same signature. Concerning the detection phase, the visual examination of the figure (2) can be advantageously replaced by a technique of jump detection or extreme value detection (Dixon's test for example).

**165** *2.3. Active approach: Outlier detection in a multisensor system*

The previous example can be easily generalized to any dimension system. In what follows, we summarize the technique of the parity space [56] which is based on the construction of the parity vector whose structure is established from the equations of the system whose consistency we want to monitor. Let us consider, at a particular instant, the measurement linear system :

$$\begin{aligned} x_m &= C x + \varepsilon + F d \\ x \in \mathcal{R}^n, x_m &\in \mathcal{R}^m, d \in \mathcal{R}^p, \varepsilon \in \mathcal{R}^n, m > n \end{aligned} \quad (4)$$

**170** where  $x_m$  is the known measurement vector,  $x$  the vector of the variables to be measured,  $d$  the vector of the unknown faults and  $\varepsilon$  the vector of the measurement noise.  $C \in \mathcal{R}^{m \times n}$  is the assumed full row matrix characterizing the measurement system and  $F$  is the matrix of the fault directions. The constraint  $m > n$  reflects redundancy of information and comes from the fact that there are more measurements than variables. To detect the presence of faults, we seek to establish analytical redundancy relations between  
**175** the measurements which are independent of the unknown quantities  $x$  but which remain sensitive to the faults  $d$ . For this, we define the parity vector :

$$p = W x_m \quad (5)$$

where  $W \in \mathcal{R}^{(m-n) \times n}$  is the projection matrix orthogonal to  $C$  resulting from (4) by simple multiplication by  $W$  :

$$p = W \varepsilon + W F d \quad (6)$$

The expression (5) is the so-called "computation" form of the parity vector from the  $x_m$  measurements  
**180** while the expression (6) explains the influence of the  $d$  faults on the parity vector through the  $WF$  matrix. In the absence of measurement noise  $\varepsilon$  and failure  $d$  the parity vector  $p$  is null. In this particular situation, the equation (5) then translates the set of redundancies that link the measurements  $x_m$  :

$$W x_m = 0 \quad (7)$$

Given the expression (6) the capability to isolate  $d$  faults affecting the measurements is directly related to the structure of the  $WF$  matrix and in particular to its rank. Let us consider, for example, the  
**185** system of measurements subject to two faults affecting some of them :

$$x_m = \begin{bmatrix} 1 & 2 & 1 \\ 1 & 0 & 2 \\ 1 & 1 & 1 \\ 1 & 0 & 1 \\ 2 & 0 & 2 \end{bmatrix} x + \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \\ 1 \end{bmatrix} \varepsilon + \begin{bmatrix} 1 & 0 \\ 0 & 0 \\ 0 & 0 \\ 1 & 1 \\ 0 & 0 \end{bmatrix} d \quad (8)$$

Solving  $W C = 0$  leads to:

$$W = \begin{bmatrix} -1 & 0 & 2 & -1 & 0 \\ -2 & 0 & 4 & 0 & -1 \end{bmatrix}$$

which makes it possible to explain the parity vector in the following two forms:

$$p = \begin{bmatrix} -x_{m,1} + 2x_{m,3} - x_{m,4} \\ -2x_{m,1} + 4x_{m,3} - x_{m,5} \end{bmatrix} \quad (9a)$$

$$p = \begin{bmatrix} 0 \\ 1 \end{bmatrix} \varepsilon + \begin{bmatrix} -2 & -1 \\ -2 & 0 \end{bmatrix} d \quad (9b)$$

where  $x_{m,i}$  are the components of the  $x_m$  measurement vector. The form (9a) allows the calculation of the parity vector from the available measurements; since the  $\varepsilon$  errors are usually of zero mean values or low magnitude, the form (9b) can be used to detect and estimate possible faults  $d$ . More precisely, if the influence of  $\varepsilon$  is neglected and as the  $WF$  matrix is regular, one can easily estimate the faults  $d$  from the definition of  $p$  (9b) itself evaluated from the measurements (9a):

190

$$\begin{aligned} d &= \begin{bmatrix} -0.5p_1 \\ 0.5(p_2 - p_1) \end{bmatrix} \\ &= \begin{bmatrix} x_{m,1} - 2x_{m,3} + 0.5x_{m,5} \\ 2x_{m,3} - x_{m,1} + x_{m,4} - x_{m,5} \end{bmatrix} \end{aligned} \quad (10)$$

#### 2.4. Active approach: Outlier replacement

The aim is to eliminate and replace outliers affecting a temporal signal with a "minimal" distortion of the useful signal. A basic idea is the median of a sample which is much less sensitive to extreme values than the mean. The observations furthest from the median can then be discarded and this discarding is known as trimming in the English literature and winsorizing when the discarded values are reconstructed from the remaining values.

Specifically, a winsorized or "trimmed"  $\{r, s\}$  mean is the replacement of the smallest  $r$  observations and the largest  $s$  observations, where  $r$  and  $s$  are integers. Let us consider the values  $x_{i+j}, j = -m, \dots, m$  of a signal to be filtered where the current index  $i$  corresponds to the center of a moving window of size  $n = 2m + 1$ . The filtered value is defined by :

200

$$\hat{x}_i = \frac{1}{n} \left( r x_{i-m+r} + \sum_{j=-m+r}^{m-s} x_{i+j} + s x_{i+m-s} \right) \quad (11)$$

Rejecting extreme points simply requires setting the parameters  $r$  and  $s$  and this can be done adaptively, for example by rejecting points that deviate from the mean by more than  $k$  times the standard deviation calculated on the window considered.

**Remark 1.** In the case where the number of rejected points is not an integer, one can define a winsorized average at  $2a\%$  which implies the replacement of a given  $2a$  percentage of values at both ends of the data. In the case of a symmetric filter  $r = s$ , we have :

205

$$\hat{x}_i = \frac{1}{(1-2a)n} \left( (1-f)x_{i-m+r} + \sum_{j=-m+r+1}^{m-r-1} x_{i+j} + (1-f)x_{i+m-r} \right), \quad f = an - r \quad (12)$$

As an example, with  $m = 4, a = 0.3, r = 2$ , we get the filter:

$$\hat{x}_i = \frac{1}{3.6} (0.3x_{i-2} + x_{i-1} + x_i + x_{i+1} + 0.3x_{i+2}), \quad i > 3 \quad (13)$$

For example, the median is the most fitted statistic (nominally 50 %) because it rejects all but the most central data. In [2] an application of trimming to robust classification is provided. ■

Figure (3) is an illustration of outlier replacement (in signal *sav*) using trimming (signal *stf*) and median filtering (signal *smf*) techniques. Of course the filter window width is the key parameter of these techniques. A close look at this figure shows the removal of outliers but with a slight distortion of the rest of the signal. The right part of the figure compares the 3 signals by their differences two by two.

210

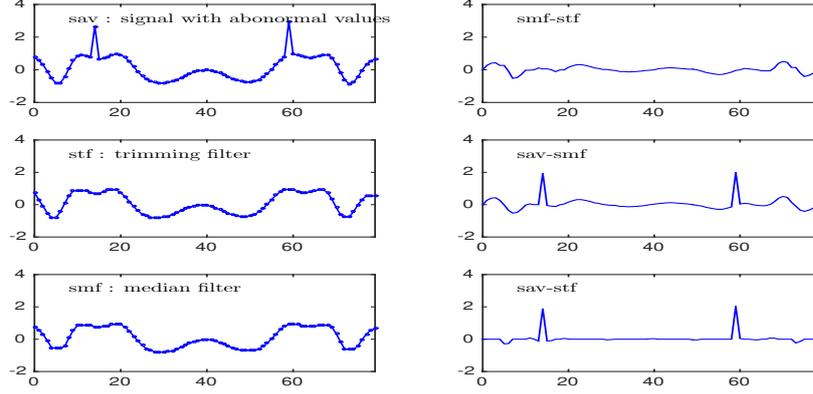


Figure 3: Outlier filtering

215 *2.5. Passive approach: acceptance of outliers and robustness*

In contrast to the previous approach, the aim here is not to eliminate outliers, but to reduce their undesirable effects during their use. For example, to identify the parameters of a system, one can try to construct an estimation algorithm that directly minimizes the influence of the outliers on these parameters. As a well-known example, let us recall the case of the robust mean using the median  
 220 filter. Two tools will be recalled to reach a certain level of robustness in the treatments: contaminated distributions and M-estimators.

- Contaminated distributions [22, 18]

A so-called contaminated model assumes that a large  $\mu$  portion of the data is generated from a classical normal error model of small magnitude. The remaining data, corresponding to the  $(1 - \mu)$  fraction of the  
 225  $N$  data set, may be affected by abnormal noise generated by a distribution of different characteristics. As an example, the distribution taking into account the **two different types of errors** can be :

$$p(\epsilon) = \mu \mathcal{N}(0, \sigma_1^2) + (1 - \mu) \mathcal{N}(0, \sigma_2^2) \quad (14)$$

This type of distribution model, after an adequate setting of  $\mu, \sigma_1$  and  $\sigma_2$ , proves to be efficient in identification in the presence of outliers. To illustrate its application, let us consider the simple case of estimating the mean of a sample size  $N$  in the presence of outliers, the aim being of course that this  
 230 estimate is not very sensitive to outliers. The likelihood function for this sample is explicit:

$$\begin{aligned} \mathcal{V} &= \prod_{i=1}^N (\mu p_1(x_i) + (1 - \mu) p_2(x_i)) \\ p_1(x) &= \frac{1}{\sqrt{2\pi}\sigma_1} \exp\left(-\frac{(x - m)^2}{2\sigma_1^2}\right) \\ p_2(x) &= \frac{1}{\sqrt{2\pi}\sigma_2} \exp\left(-\frac{(x - m)^2}{2\sigma_2^2}\right) \end{aligned} \quad (15)$$

where  $m$  is the mean to be estimated and  $\sigma_1, \sigma_2$  the standard deviations of the contaminated distribution. A few comments are necessary to justify the interest of this type of function. To do so, we can analyze its sensitivity  $g(x) = \mathcal{V}/\partial x$  compared to the data  $x$  :

$$g(x) = \frac{\frac{wp_1(x)}{\sigma_1} + \frac{(1-w)p_2(x)}{\sigma_2}}{wp_1(x) + (1-w)p_2(x)}$$

The figure (4) shows the role of the parameters  $\sigma_1, \sigma_2, w$  on the ability to take into account the values of  $x$  which will subsequently represent the outliers. For an easier interpretation of the graphs, they have been normalized, i.e. they represent the normalized functions  $\bar{g}(x) = g(x)/g(0)$ .

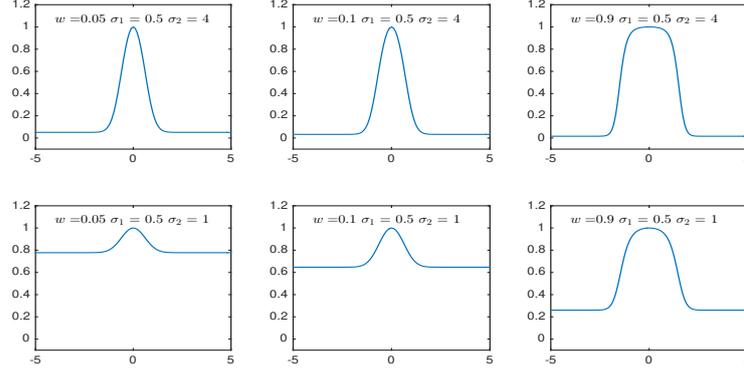


Figure 4: Influence function

For  $w = 1$  we naturally obtain a constant weight; thus all the data are equally weighted and, in particular, the optimisation criterion will be sensitive to large magnitude of data, i.e. to outliers. Taking  $w = 0.1, \sigma_1 = 0.5, \sigma_2 = 1$  reduces the influence of outliers since the weight decreases from 1 for data around the origine to 0.63 for data with large magnitude. Ultimately, what is important is to have separate sensitivities for low and high values of the  $x$  variable. Of course, the situation is even clearer with the choice  $(w = 0.1, \sigma_1 = 0.5, \sigma_2 = 4)$  where this time the weight of the large values of  $x$  is negligible, which greatly reduces the influence of outliers.

The maximum of the likelihood function  $\mathcal{V}$  with respect to  $m$  is obtained for :

$$\sum_{i=1}^N w_i (x_i - m) = 0$$

$$w_i = \mu \frac{p_{1,i}}{\sigma_1^2} + (1 - \mu) \frac{p_{2,i}}{\sigma_2^2}$$
(16)

Given the expressions of  $p_{1,i}$  and  $p_{2,i}$  (15) which depend on  $m$ , the non-linear equation (16) is solved iteratively with respect to  $m$ , for example according to the scheme (17) initialized with weights  $w_i^0$  equal to unity :

$$m^{iter+1} = \frac{\sum_{i=1}^N w_i^{iter} x_i}{\sum_{i=1}^N w_i^{iter}}$$

$$w_i^{iter} = \mu \frac{p_{1,i}^{iter}}{\sigma_1^2} + (1 - \mu) \frac{p_{2,i}^{iter}}{\sigma_2^2}$$

$$p_{1,i}^{iter} = \frac{1}{\sqrt{2\pi}\sigma_1} \exp\left(-\frac{(x_i - m^{iter})^2}{2\sigma_1^2}\right)$$

$$p_{2,i}^{iter} = \frac{1}{\sqrt{2\pi}\sigma_2} \exp\left(-\frac{(x_i - m^{iter})^2}{2\sigma_2^2}\right)$$
(17)

Table 4 shows an estimation result of the mean of a sample of 50 values (with 4 outliers) and by extension its standard deviation. The first two rows of this table relate to the standard valuation without and with outliers. The robust evaluation in the third row shows estimates close to those obtained in the absence of outliers, thus demonstrating the robustness of the method.

The figure (5) shows the 50 values of the sample with its outliers and the weights  $w_i$  used for the calculation of the mean, which also makes it easy to locate the outliers. Of course, the robustness of the estimate with respect to outliers is related to their proportion to healthy values. The setting of the parameters  $\sigma_1, \sigma_2$  and  $\mu$  conditions this robustness; it can be realized heuristically by learning but also by a more analytical optimization procedure, directly from the likelihood function.

- M-estimators [12, 38, 72]

M-estimators were introduced as a generalization of the maximum likelihood minimization estimate of a  $\rho$  function over the available data set  $z_i, i = 1, \dots, N$ . Thus, the M-estimator(s) associated with the

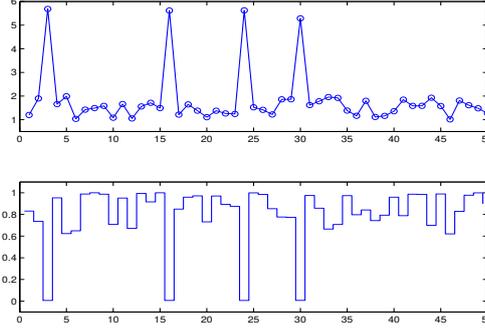


Figure 5: Series  $x_i$  with abnormal values and weights  $w_i$

	Mean	Standard deviation
Conventional evaluation without abnormal value	1.501	0.277
Conventional evaluation with abnormal value	1.821	1.145
Robust evaluation with abnormal value	1.501	0.282

Table 4: Robust estimation of a mean and standard deviation

data and the function  $\rho$  is given by :

$$\hat{\theta} = \operatorname{argmin}_{\theta} \left( \sum_{i=1}^N \rho(z_i, \theta) \right) \quad (18)$$

A well-known example is the function of Cauchy or Lorenz :

$$\rho(z_i, a, b) = \frac{c^2}{2} \log \left( 1 + \left( \frac{\varepsilon_i}{c} \right)^2 \right) \quad (19)$$

260 where  $z_i = \{x_i, y_i\}$  and where  $\varepsilon_i = y_i - ax_i - b$  is an image of model errors when representing data with a straight line. The sensitivity of  $\rho$  to the errors  $\varepsilon$  becomes explicit :

$$\frac{\partial \rho}{\partial \varepsilon_i} = \frac{\varepsilon_i}{1 + \left( \frac{\varepsilon_i}{c} \right)^2} \quad (20)$$

and concludes that a large  $\varepsilon_i \gg c$  error produces a small insensitivity on  $\hat{\theta}$  and a small  $\varepsilon_i \ll c$  error produces a sensitivity in the order of  $\varepsilon_i$ . Thus, the choice of the  $c$  threshold, at the user's discretion, determines the robustness of the estimator with respect to outliers. Starting from (19) and (20), the  
 265 reader will be able to establish the optimality equations of the parameters  $a$  and  $b$  of the regression model by himself:

$$\begin{aligned} \sum_{i=1}^N w_i(a, b)(y_k - ax_k - b)x_i &= 0 \\ \sum_{i=1}^N w_i(a, b)(y_k - ax_k - b) &= 0 \end{aligned} \quad (21)$$

with the following expression of weights:

$$w_i(a, b) = \frac{1}{1 + \left( \frac{y_i - ax_i - b}{c} \right)^2} \quad (22)$$

The non-linear system (21, 22) can be solved by a simple iteration mechanism from an initial choice of weights  $w_i$  for example to the unit value. This procedure is to be compared with that obtained by an  
 270 ordinary least squares method, the only difference being the use of adapted weights [33], i.e. function of outliers.

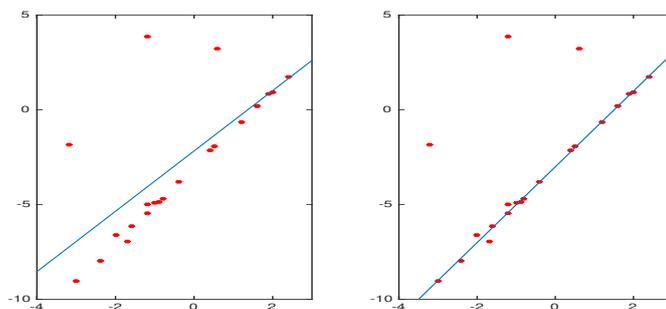


Figure 6: Ordinary and robust regression

The figure (6) shows the arrangement of 22 pairs of points  $\{x_i, y_i\}$  of which three are obviously outliers, as well as the regression lines obtained by an ordinary regression and a robust regression. The true parameters  $a$  and  $b$  of the system are 2 and  $-3$ , those obtained by ordinary least squares 1.594 and  $-2.169$ , and those obtained by the robust procedure are 2.001 and  $-3.008$ , which highlights the merits **275** of the robust estimator.

### 3. Reconciliation of data in the presence of outliers

This section is devoted to a technique for dealing with measurements potentially contaminated by outliers. It has a twofold objective, on the one hand to reconcile the measurements against a model and, on the other hand, to detect/locate outliers. Successively, the principle of measurement reconciliation is **280** given, then its refinement by complementing it with a procedure for robust estimation of the variables. To illustrate its implementation, the case of a non-linear model system is treated numerically, then a comparison with three more classical techniques is proposed and discussed.

#### 3.1. Basic principle of data reconciliation

The purpose of data reconciliation,[42, 64], is to make the measurements made on a system compatible **285** with its model, which is assumed to be accurate because it is based on the laws of matter or energy conservation [32]. As such, the reconciliation methods are close to the state estimation methods established in a much more general framework. An important consequence of reconciliation is the detection of outliers. Indeed, the reconciled values can be compared to the measurements; the discrepancies found can be analysed, the largest of them being able to testify to the presence of outliers. The simplest **290** formulation in the context of a linear model system linking the true quantities  $x^*$  is discussed in this section:

$$M x^* = 0, \quad x^* \in \mathcal{R}^v, \quad M \in \mathcal{R}^{n \times v} \quad (23)$$

whose measurements  $x_m$  are defined in the additive form with respect to noises  $\varepsilon$  related to the instrumentation and the measurement procedure:

$$x_m = x^* + \varepsilon + F d \quad (24)$$

where the matrix  $F$  gives the directions of influence of the faults  $d$  on the measurements. Typically the **295**  $x_m$  measures do not exactly verify the model (23) of the system and the reconciliation principle aims at correcting them to satisfy this model. Since the measures are assumed not to be totally outliers, the corrected quantities must remain close to them and for this reason of proximity, the estimated  $\hat{x}$  which minimizes the criterion here chosen quadratically with a weighting matrix  $W$ :

$$\Phi = \| x_m - x^* \|_{W^{-1}}^2 \quad (25)$$

is defined by :

$$\begin{aligned} \hat{x} &= \arg \min_{x^*} \Phi \quad \text{under } M x^* = 0 \\ &= (I - W M^T (M W M^T)^{-1} M) x_m \end{aligned} \quad (26)$$

**300**

which is explicit:

$$\hat{x} = (I - WM^T(MWM^T)^{-1}M) x_m \quad (27)$$

This expression provides, on the one hand, estimates of true quantities that are consistent in the sense of the satisfaction of the system model and, on the other hand, estimates of the corrections that have been made to the measurements :

$$\begin{aligned} \tilde{x} &= x_m - \hat{x} \\ &= WM^T(MWM^T)^{-1}M x_m \end{aligned} \quad (28)$$

305 The analysis of the magnitudes of the corrective terms  $\tilde{x}$  provides information on the magnitude of the errors and their distribution. To clarify this point it is then possible to specify the expression of the corrective terms and the reader will be able to check from (24, 28) the following expression of the corrective terms :

$$\tilde{x} = WM^T(MWM^T)^{-1}M \varepsilon + WM^T(MWM^T)^{-1}MF d \quad (29)$$

where the influence of "small" random  $\varepsilon$  errors and the more important influence of  $d$  faults appears.

310 The interpretation of (29) in terms of performance of detection and isolation of faults  $d$  is to be found in [50]. As the  $\varepsilon$  errors are small in magnitude, the approximation can be adopted:

$$\tilde{x} \sim WM^T(MWM^T)^{-1}MF d \quad (30)$$

which makes explicit the corrective terms according to the faults. Therefore the fault  $d_i$  (the  $i$ th component of  $d$ ) is detectable if the  $i$ th column of the matrix  $WM^T(MWM^T)^{-1}MF$  is non-zero. The isolation of the fault(s) obviously depends on the structure of this same matrix. In a certain way, this reconciliation procedure does not have the desired robustness in the sense that the faults are corrected but to the detriment of their dissemination on all the variables. This justifies the following paragraph which presents a robust approach to the problem of reconciliation of measurements, the robustness allowing not to disseminate the measurement faults on all the variables during the reconciliation procedure. As an important result, the correction of the measurements essentially will concern those subject to faults.

320 **Remark 2.** *It is also interesting to consider the problem of reconciliation when using only part of the available measures. To ignore the measure of the  $p$ -th variable, an elegant way of dealing with this case is to choose a diagonal weighting matrix whose  $(p,p)$  element takes an infinite value.* ■

325 **Remark 3.** *The previous data reconciliation principle extends to non-linear and dynamic systems. For the sake of brevity, let us consider only the case of bilinear systems for which the models involve products of variables. This is very frequent in the chemical or mineralurgical field when total and partial flow material balances are established. In this case, if  $x$  and  $y$  designate, for example, the vector of flows and the vector of concentrations in a chemical or mineral species, the constraints (23) are extended as:  $Mx^* = 0$  and  $Mx^* \otimes y^* = 0$ , where the  $\otimes$  operator makes the term to term product of two vectors. The measurement equations then become  $x_m = x^* + \varepsilon_x$  and  $y_m = y^* + \varepsilon_y$ . Criterion (25) is amended as :*

$$\Phi = \|x_m - x^*\|_{W_x^{-1}}^2 + \|y_m - y^*\|_{W_y^{-1}}^2 \quad (31)$$

and the estimation of the reconciled variables  $\hat{x}$  and  $\hat{y}$  results from the Lagrangian optimization:

$$\mathcal{L} = \|x_m - x^*\|_{W_x^{-1}}^2 + \|y_m - y^*\|_{W_y^{-1}}^2 + \lambda^T Mx^* + \mu^T Mx^* \otimes y^* \quad (32)$$

with respect to the variables  $x^*, y^*, \lambda$  and  $\mu$ , thus leading to the estimates  $\hat{x}$  and  $\hat{y}$ . We leave it to the reader to carry out this rather classic calculation and to repeat the analysis of the residuals of the resulting  $\tilde{x}$  and  $\tilde{y}$  estimates.

335

■

### 3.2. Robust data reconciliation

Numerous developments are complementing the basic principle of data reconciliation that we have just recalled. Thus, extensions made it possible to deal with dynamic systems [16], non-linear systems [13], the presence of poorly known parameters [71], the localisation of measurement faults [47], the taking into account of missing measures [39].

In this paragraph, only the specific point of the robustness of reconciliation against outliers is addressed. To introduce this issue, it should be remembered that reconciliation is based on the minimization of a criterion formed from the discrepancies between the variables and their respective measures. The validity and optimal character of this approach are eminently linked to the strong assumption of normality of measurement errors. In practice, this assumption can be defeated in the presence of large errors that constitute outliers, which can hardly be considered as realizations of normal random variables.

We thus bring more realism by posing the reconciliation problem in the following way: from  $x_m$  measurements estimate the true quantities  $x^*$  of a linear model system (23). Starting from the assumption that the number of large errors is low, the first technique proceeds by reconciling the measurements by weighted ordinary least squares (*OLS*), then detects and locates the large errors (analysis of the corrective terms), and finally repeats the reconciliation procedure cited above by assigning a very low weight to the measurements for which large errors have been located. The major drawback of this approach is that the first reconciliation can be highly erroneous due to the presence of the large errors; this can then make it difficult to locate the large errors by analyzing the corrective terms.

To directly take into account the presence of gross errors, a more appropriate error distribution law is used. Recall that the class of M-estimators provides estimates that are robust to large errors. Let us consider the measurement case  $x_m = x^* + \varepsilon$ , the components of  $\varepsilon$  being noted  $\varepsilon_i$ . The estimation criterion (25) is now taken as :

$$\Phi = \frac{c^2}{2} \sum_{i=1}^v \log \left( 1 + \left( \frac{\varepsilon_i}{c} \right)^2 \right) \quad (33)$$

Let us recall the important role played by the constant  $c$  in this criterion. It appears clearly that the errors of magnitude higher than  $c$  are more taken into account in  $\Phi$  than those of lower magnitude, i.e. lower than  $c$ . Therefore, minimizing  $\Phi$  tends to reduce the influence of large errors on the estimates. The previous formulation can then be repeated, taking into account the model (23) and the objective function (34). We leave it to the reader to find the following estimate:

$$\begin{aligned} \hat{x} &= (I - WM^T(MWM^T)^{-1}M) x_m \\ W &= I_v + \frac{1}{c^2} \text{diag}(\tilde{x} \otimes \tilde{x}) \\ \tilde{x} &= x_m - \hat{x} \end{aligned} \quad (34)$$

The non-linear system (34) is solved with respect to  $\hat{x}$  in a iterative manner from an initial choice of the matrix of weights  $W$ , for example the identity matrix.

**Remark 4.** *As before, the case of bilinear systems extends the robust M-estimator formulation. To do this, we consider again the variables  $x$  and  $y$  defined in remark 2 and the estimation criterion to be considered is now:*

$$\Phi = \frac{c_x^2}{2} \sum_{i=1}^v \log \left( 1 + \left( \frac{\varepsilon_{x,i}}{c_x} \right)^2 \right) + \frac{c_y^2}{2} \sum_{i=1}^v \log \left( 1 + \left( \frac{\varepsilon_{y,i}}{c_y} \right)^2 \right) \quad (35)$$

the  $c_x$  and  $c_y$  parameters setting the outlier insensitivity. The reconciliation procedure is then based on finding the minimum  $\Phi$  under the respect of the equations  $Mx^* = 0$  and  $Mx^* \otimes y^* = 0$ . Since this is a standard procedure, the establishment of the optimality equations is not explained here. ■

### 3.3. Application

- 375 The following example, although small in size, illustrates the benefits of the robust approach for reconciling data, some of which are polluted by outliers, that need to be detected and located. The data used will also be processed (section 3.4) by three other approaches that also allow outliers to be detected and located. In order not to burden the presentation with an abundance of numerical results, the example uses only one set of outliers.
- 380 Figure (3.3) shows a material transport network in the chemical industry (but it could also be in the mineral industry, water, gas, oil products distribution circuits, etc.) made up of nine production units and sixteen connecting routes between these units.

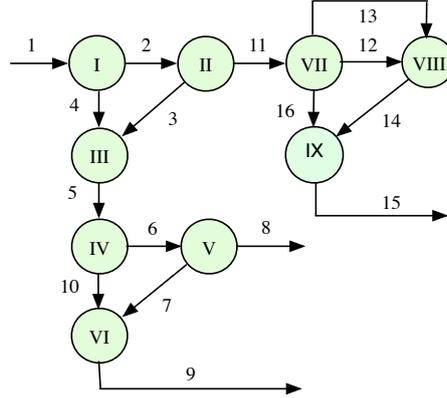


Figure 7: Material transport network

- The law of conservation of the flows of material makes it possible to write, in permanent regime, the system of equations (36, 37) linking the flows  $x_i^*$  and the concentrations  $y_i^*$  in a chemical species. These
- 385 equations only reflect the conservation of total and partial flows and could in a more general case take into account the variations of material stocks in the production units but also of other constituents.

$$\begin{cases} x_1^* - x_2^* - x_4^* & = 0 \\ x_2^* - x_3^* - x_{11}^* & = 0 \\ x_3^* - x_4^* - x_5^* & = 0 \\ x_5^* - x_6^* - x_{10}^* & = 0 \\ x_6^* - x_8^* - x_7^* & = 0 \\ x_7^* + x_{10}^* - x_9^* & = 0 \\ x_{11}^* - x_{12}^* - x_{13}^* - x_{16}^* & = 0 \\ x_{12}^* + x_{13}^* - x_{14}^* & = 0 \\ x_{16}^* + x_{14}^* - x_{15}^* & = 0 \end{cases} \quad (36)$$

$$\begin{cases} x_1^* y_1^* - x_2^* y_1^* - x_4^* y_4^* & = 0 \\ x_2^* y_2^* - x_3^* y_3^* - x_{11}^* y_{11}^* & = 0 \\ x_3^* y_3^* - x_4^* y_4^* - x_5^* y_5^* & = 0 \\ x_5^* y_5^* - x_6^* y_6^* - x_{10}^* y_{10}^* & = 0 \\ x_6^* y_6^* - x_8^* y_8^* - x_7^* y_7^* & = 0 \\ x_7^* y_7^* + x_{10}^* y_{10}^* - x_9^* y_9^* & = 0 \\ x_{11}^* y_{11}^* - x_{12}^* y_{12}^* - x_{13}^* y_{13}^* - x_{16}^* y_{16}^* & = 0 \\ x_{12}^* y_{12}^* + x_{13}^* y_{13}^* - x_{14}^* y_{14}^* & = 0 \\ x_{16}^* y_{16}^* + x_{14}^* y_{14}^* - x_{15}^* y_{15}^* & = 0 \end{cases} \quad (37)$$

The measures  $(x_m, y_m)$  of these sixteen pairs of variables over a given time period are recorded in table (5). The purpose of data validation is twofold: to detect outliers (here, two biases of respective

amplitudes 12 and 6 affect the measures of the second component of  $x$  and the eleventh component of  $y$ ), and to propose replacement values for the outliers.

	1	2	3	4	5	6	7	8
$x_m$	54.23	80.39	59.33	10.52	47.04	59.61	33.18	24.76
$\hat{x}_{RLS}$	56.00	67.81	58.31	11.81	46.49	59.39	34.12	25.27
$\tilde{x}_{RLS}$	-1.77	12.58	1.02	-1.29	0.55	0.22	-0.95	-0.51
$\hat{x}_{OLS}$	57.88	72.98	62.17	15.10	47.08	59.67	34.10	25.57
$\tilde{x}_{OLS}$	-3.65	7.41	-2.84	-4.57	-0.04	-0.06	-0.92	-0.81
$y_m$	8.26	6.64	6.20	1.16	7.47	6.81	4.22	10.92
$\hat{y}_{RLS}$	8.14	6.92	6.11	1.13	7.38	7.08	4.29	10.85
$\tilde{y}_{RLS}$	0.12	-0.28	0.09	0.03	0.09	-0.27	-0.07	0.06
$\hat{y}_{OLS}$	8.56	7.03	5.83	1.16	7.32	7.05	4.25	10.78
$\tilde{y}_{OLS}$	-0.31	-0.40	0.37	-0.01	0.15	-0.23	-0.03	0.14

	9	10	11	12	13	14	15	16
$x_m$	21.58	13.58	10.08	16.60	2.81	19.99	9.30	9.01
$\hat{x}_{RLS}$	21.22	12.90	9.51	16.53	2.73	19.26	9.51	9.76
$\tilde{x}_{RLS}$	0.36	0.68	0.58	0.07	0.08	0.73	-0.20	-0.75
$\hat{x}_{OLS}$	21.51	12.60	10.81	17.16	2.98	20.13	10.81	9.32
$\tilde{x}_{OLS}$	0.07	0.99	-0.73	-0.56	-0.16	-0.14	-1.50	-0.31
$y_m$	3.34	6.10	17.62	6.52	2.52	6.25	11.75	0.20
$\hat{y}_{RLS}$	3.25	6.02	11.85	6.57	2.53	5.99	11.85	0.29
$\tilde{y}_{RLS}$	0.10	0.08	5.77	-0.05	-0.01	0.25	-0.09	-0.09
$\hat{y}_{OLS}$	3.21	6.04	13.99	7.84	2.75	7.09	13.99	-0.91
$\tilde{y}_{OLS}$	0.14	0.06	3.63	-1.32	-0.23	-0.84	-2.23	1.11

Table 5: Measured, estimated and corrective terms for flows and concentrations

In order to highlight the contribution of the robust technique, this table shows the results of estimating  $\hat{x}$  and  $\hat{y}$  of variables  $x$  and  $y$  on the one hand by robust least squares ( $RLS$ ) on the other hand by ordinary least squares ( $OLS$ ), as well as the corrective terms  $\tilde{x}$  and  $\tilde{y}$ . Concerning the estimates of the variable  $x$  through  $RLS$ , it can be seen that its second component is the most corrected (12.58), the other components being only slightly corrected. The results are quite different with  $OLS$  where not only the second component of  $x$  is adjusted but also the 3, 4 and 15 components. The same is true for the variable  $y$  where the eleventh component of  $y$  is corrected by 5.77 for  $OLS$ , the others being only slightly corrected. Robustness is reflected by the fact that a faulty measure is corrected without correcting other measures. It should also be noted that the correction terms with  $RLS$  (12.58 and 5.77) are completely related to the magnitude of the simulated faults (12 and 6), which is not the case with  $OLS$ .

Figure (8) gives an overview of the correction terms (CT) in absolute value of the variables  $x$  and  $y$  for both techniques ( $\tilde{x}_{RLS}$  and  $\tilde{y}_{RLS}$  for  $RLS$  and  $\tilde{x}_{OLS}$  and  $\tilde{y}_{OLS}$  for  $OLS$ ). This figure only graphically translates the results of the table (5) and shows the correct location of the corrections of the two measurement faults. To judge an average effect of the procedures  $RLS$  and  $OLS$ , 12 simulations were made with the same variables  $x$  and  $y$  in fault but by generating random measurement noise of small amplitude. Figure (9) visualizes the mean corrective terms from these 12 simulations and confirms if need be the relevance of the robust approach for the detection of measurement errors where variables not subject to errors have not been corrected in a sensitive way. Moreover, with respect to the magnitude of the errors,  $RLS$  gives a fairly accurate estimate, which is not the case with  $OLS$ . Of course, in practice, this technique is repeated over time in order to continuously monitor the system in question.

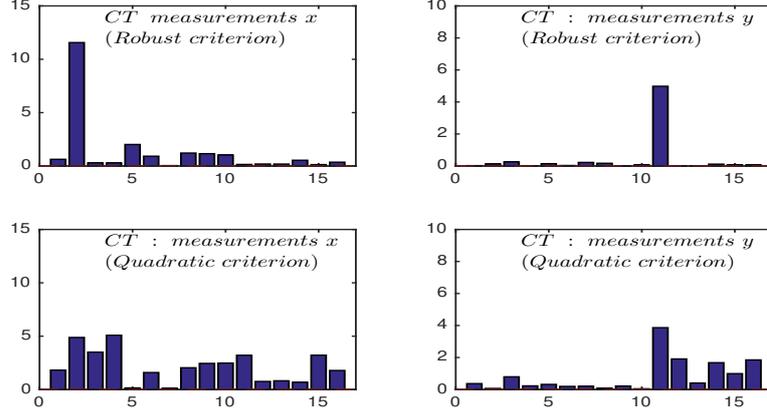


Figure 8: Corrective terms CT for flows and concentrations

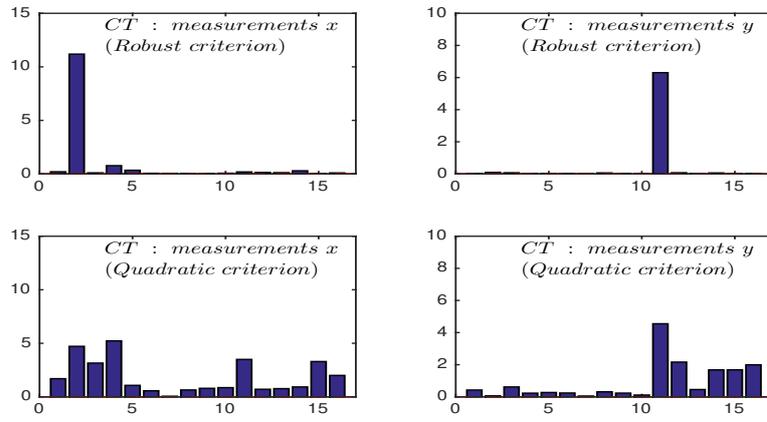


Figure 9: Corrective terms CT for flows and concentrations. Mean over 12 cases.

### 3.4. Comparisons and discussion

Roughly speaking, it can be said that the techniques for detecting and locating outliers are based  
 415 on two strategies, respectively described as global or sequential, depending on whether they seek to  
 determine them simultaneously or one by one. In all cases, and this also appears in our approach which  
 can be qualified as global, the choice of a detection threshold (or even several) remains necessary. In  
 the previous section, our approach has shown its capacities of detection/localization/identification of  
 outliers. It should now be compared with approaches widely used in the literature. Of course, we are  
 420 constrained to some limitations and have chosen to give some results from a sequential approach, an  
 approach using brute force technique and an approach using the signature of the residues with respect  
 to outliers.

#### 3.4.1. Sequential detection and localization

As its name suggests, the sequential procedure seeks to detect and locate one after the other suspicious  
 425 variables. Since the variable  $x_1$  is first suspected, the aim is to estimate all the variables from all the  
 measures except that of  $x_1$ ; this estimation is directly related to the reconciliation procedure of section  
 3.1 by taking advantage of the remark 2. The consistency of the estimates  $\hat{x}_i, i = 1, \dots, v$  obtained  
 is then tested by evaluating the vector of residuals  $M \hat{x}$  and then globally its norm  $\Phi_{R,1} = \| M \hat{x} \|$ .  
 This evaluation is then restarted by suspecting the other variables  $x_i, i = 2, \dots, v$ . The criteria  $\Phi_{R,i} = \|$   
 430  $M \hat{x} \|, i = 1, \dots, v$  obtained are then analyzed, the one whose value is judged below a threshold close  
 to zero corresponding to a suspicious variable. The procedure can be carried out again to detect on the  
 one hand a possible suspicious  $x$  variable and on the other hand one or more suspicious  $y$  variables.

Of course, the search for another outlier can benefit from the previous outlier estimate. This estimation allows the correction of the affected measurement and thus the search for another outlier taking into account the previous correction. The table (6) gathers the results obtained by this sequential technique, which results from 16 eliminations for a variable  $x$  and as much for a variable  $y$ , each elimination being followed by a reconstruction of the variables, the quality of each reconstruction being evaluated by the norm of the residual vector of the redundancy models (36, 37) computed with the estimated variables. The first two lines of this table relate to the elimination of a variable  $x_i$ , the next two lines to that of a variable  $y_i$ . Examination of the norms for the residuals of the redundancy equations unambiguously points to  $x_1$  and  $y_{12}$  as variables with aberrant measurements. In the reconstruction, let us indicate that the corrective terms affecting these two measures were respectively 12.1 and 5.9 values in accordance with the biases that had been created. For this example, let us note that the detection and localization of outliers is done without ambiguity, the contrast between the minimum values of the two residual criteria with respect to the other values being significant: 1 compared to 12 and 9 compared to 74. 435 440 445

**Remark 5.** *The above sequential procedure does not claim to provide the best solution for isolating outliers. Indeed it proceeds by decoupling, a first outlier is detected, localized and corrected by the reconciliation technique. This reconciled value is substituted for the measurement of the variable concerned, the procedure is resumed to process a second outlier and so on until the procedure is stopped. It is clear that in the presence of several simultaneous outliers, the estimate of the first value may be partially contaminated by the second. A more advanced version of this sequential technique can be used to iteratively refine the different estimates.* ■ 450

deleted variable	$x_1$	$x_2$	$x_3$	$x_4$	$x_5$	$x_6$	$x_7$	$x_8$	$x_9$	$x_{10}$	$x_{11}$	$x_{12}$	$x_{13}$	$x_{14}$	$x_{15}$	$x_{16}$
$\  M x \parallel$	12	1	15	15	17	17	17	17	17	17	15	17	17	17	17	17
deleted variable	$y_1$	$y_2$	$y_3$	$y_4$	$y_5$	$y_6$	$y_7$	$y_8$	$y_9$	$y_{10}$	$y_{11}$	$y_{12}$	$y_{13}$	$y_{14}$	$y_{15}$	$y_{16}$
$\  M (x \otimes y) \parallel$	88	74	76	88	88	88	88	88	88	88	9	77	77	88	88	76

Table 6: Sequential outlier search technique

### 3.4.2. Detection and localization using brute force technique

In computer science, brute-force search or exhaustive search, also known as generate and test, is a very general problem-solving technique and algorithmic paradigm that consists of systematically enumerating all possible candidates for the solution and checking whether each candidate satisfies the problem's statement. 455

Often not very efficient in terms of calculation time, this procedure is nevertheless easy to implement. In an optimization problem, as it analyzes all possible solutions, one is sure to highlight the optimal solution. For the example we are interested in, we have to search  $n_f = 2$  (this number could be modified) simultaneously aberrant measures among 32 (16 variables  $x_i$  and 16 variables  $y_i$ , which makes 496 different situations to examine. For each situation, only two faulty measures are set. The set of variables is then reconciled with respect to the redundancy equations according to the same technique as previously described, this reconciliation using all the measures except those of the two faulty variables. We thus have 496 reconciliation results, 460 465

the best situation is the one where the corrections essentially affect  $n_f = 2$  variables among the 32 and where the residual vector has the smaller norm. For this situation, the table (7) shows the correction terms  $\tilde{x}$  and  $\tilde{y}$  made to the measurements, which can be compared to the terms in the table (5). Variables  $x_2$  and  $y_{11}$  have the largest corrections, the other variables are only slightly corrected. Thus, this procedure essentially causes the corrections to be carried over to the outliers, which is indeed the desired effect. In the end, the results obtained by this technique are quite similar to those obtained with the proposed method. However, as indicated above, this method is time consuming and is unsuitable for large systems. 470

	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16
$\tilde{x}$	0.15	12.1	0.27	0.27	0.06	0.23	0.28	0.41	0.14	0.04	0.17	0.23	0.23	0.05	0.30	0.28
$\tilde{y}$	0.12	0.10	0.21	0.02	0.15	0.11	0.22	0.04	0.17	0.06	5.80	0.06	0.01	0.29	0.07	0.11

Table 7: Corrective terms. Brute force technique algorithm

### 3.4.3. Detection and localization of outliers from the signature of model residuals

475 As previously mentioned, the diagnostic technique requires the synthesis of indicators that can reveal the presence of anomalies that may affect the measurement. Here, these indicators are simply the result of the adequacy of the measurements with respect to the redundancy equations.

**Definition 2.** (*Fault Signatures*).

480 The signature of a fault  $f_j$  is the binary vector  $FS(f_j) = [s_{1,j}, \dots, s_{n,j}]^T$  where  $s_{i,j} = 1$  if  $f_j$  is a variable of the equation used to form the redundancy equation  $r_i$ , otherwise  $s_{i,j} = 0$ .

This definition implicitly assumes that the occurrence of  $f_j$  is observable on the result of  $r_i$ , i.e., if  $r_i$  is satisfied, then fault  $f_j$  did not occur. This is known as a simple fault exoneration assumption. ■

**Definition 3.** (*Fault Isolability*). A surveillance system does not respect the isolability property, if two theoretical signatures are identical, i.e. if the Hamming distance between two signatures is zero. To 485 allow an unambiguous recognition of the fault by tolerating  $k$  degradations on the real signature, the theoretical signatures must be at least  $2.k + 1$  bits apart. Indeed, if the distance between two theoretical signatures is only  $2.k$  bits,  $k$  degradations of one of them bring the real signature at an equal distance from both. The failure is thus not isolable. ■

490 Looking again at the example in section 3.3, the equations (36, 37) allow us to establish the table (8) of occurrence (with the symbols "1" and "." to translate the presence or absence) of the variables  $x_i, y_i, i = 1, \dots, 16$ . This table which consists of 32 columns for the variables and 18 rows for the equations reflects the influence of outliers in the different equations. It is important to note that the columns of this table are distinct from each other (except for the 12 and 13 columns).

**Remark 6** (Multiple outliers). The synthesis of the signature of multiple outliers, i.e. occurring con- 495 comitantly, is a generalization of the previous case. As an example, the simultaneous presence of outliers affecting variables  $x_2$  and  $x_{10}$  is characterized by a signature that is deduced by applying a logical operator to columns 2 and 10 of the table (8). ■

**Remark 7** (Structuring the residues). Signatures for multiple outliers may become identical, limiting 500 the ability to isolate them. For example, as simultaneous outliers on variables  $x_6$  and  $x_{10}$  have the same signature as simultaneous outliers on variables  $x_8$  and  $x_{10}$ , i.e.  $[\dots 1 1 1 \dots]^T$ , the joint outliers  $x_6, x_{10}$  and  $x_8, x_{10}$  are detectable but cannot be isolated. Subject to satisfying some structural conditions of the equations, this isolation can be solved by structuring the residuals by combining the equations so as to eliminate some variables [50]. For the example shown, merging the 4 and 5 equations 505 of the system (36) generates the redundancy equation  $x_5^* - x_8^* - x_9^* = 0$ , which allows to complete the table (8) with a tenth line translating the occurrences of the 16 variables  $x_i$  in this equation. As a consequence of this addition, the signatures of the  $x_6, x_{10}$  and  $x_8, x_{10}$  variable pairs become distinct, which solves the previous isolation problem. ■

How to use this signature table for detection/location of outliers affecting measurements?

510 With measurements of the variables to be analyzed, the residual of the redundancy models is evaluated, generating a signature. The detection and location of faults are carried out by comparing this signature with those in the table (8). As the latter are in binary form, it is therefore appropriate that the actual signature is also in binary form : this is done by comparing each residual component to a threshold. Having normalized the experimental signature and the theoretical signatures, it is then necessary to choose a criterion to compare them, and here the Hamming distance has been chosen.

	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16
	<i>x</i>																<i>y</i>															
1	1	1	.	1	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.
2	.	1	1	.	.	.	.	.	.	.	1	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	
3	.	.	1	1	1	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	
4	.	.	.	.	1	1	.	.	.	1	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	
5	.	.	.	.	.	1	1	1	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	
6	.	.	.	.	.	.	1	.	1	1	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	
7	.	.	.	.	.	.	.	.	.	.	1	1	1	.	1	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	
8	.	.	.	.	.	.	.	.	.	.	.	1	1	1	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	
9	.	.	.	.	.	.	.	.	.	.	.	.	1	1	1	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	
10	1	1	.	1	.	.	.	.	.	.	.	.	.	.	.	1	1	.	1	.	.	.	.	.	.	.	.	.	.	.	.	
11	.	1	1	.	.	.	.	.	.	.	1	.	.	.	.	.	1	1	.	.	.	.	.	.	1	.	.	.	.	.	.	
12	.	.	1	1	1	.	.	.	.	.	.	.	.	.	.	.	1	1	1	.	.	.	.	.	.	.	.	.	.	.	.	
13	.	.	.	.	1	1	.	.	.	1	.	.	.	.	.	.	.	1	1	.	.	1	.	1	.	.	.	.	.	.	.	
14	.	.	.	.	.	1	1	1	.	.	.	.	.	.	.	.	.	.	1	1	1	.	.	.	.	.	.	.	.	.	.	
15	.	.	.	.	.	.	1	.	1	1	.	.	.	.	.	.	.	.	.	1	.	1	1	.	.	.	.	.	.	.	.	
16	.	.	.	.	.	.	.	.	.	.	1	1	1	.	1	.	.	.	.	.	.	.	.	.	1	1	1	.	.	.	1	
17	.	.	.	.	.	.	.	.	.	.	.	1	1	1	.	.	.	.	.	.	.	.	.	.	1	1	1	.	.	.	.	
18	.	.	.	.	.	.	.	.	.	.	.	.	1	1	1	.	.	.	.	.	.	.	.	.	.	1	1	1	.	.	.	

Table 8: Fault signature table for  $x, y$  variables

The example in section (3.3) is again considered with the measurements recorded in the table (9), where the outliers are always  $x_2$  and  $y_{11}$ . The table (10) shows a result of outlier detection/localization. The  $r$  line is the residuals from the measurements and the 18 redundancy equations (36, 37) that should be analyzed against the theoretical outlier signatures. Examination of the values taken by  $r$  shows a group of values close to 0 and another behaviour of significantly larger values without being able to clearly establish a threshold separating these two groups. The line  $r_N$  represents a normalization of the respective residuals of the redundancy equations (36, 37) by their maximum values. The line  $s_t^{(2,11)}$  corresponds to the theoretical signature of the fault pair on  $\{x_2, y_{11}\}$  and its comparison with the theoretical signatures of the table (8) confirms the presence of faults on  $\{x_2, y_{11}\}$ . The reader will easily notice the difficulty in choosing a threshold to apply to the values of  $r_N$  in order to binarize the  $r_N$  residual to find the theoretical signature  $s_t^{(2,11)}$  see for example the values 0.36 and 0.23 for residuals 11 and 12). Of course this is a particular conclusion related to this example, for different measurements with lower noise levels the outliers were perfectly localized.

In conclusion, it is clear that in this technique, the calculation of the residuals of the redundancy equations is easy to implement, but the adequate setting of a threshold to binarize them often remains problematic.

	1	2	3	4	5	6	7	8
$x$	56.67	78.76	57.1	9.32	49.08	60.60	32.39	26.63
$y$	7.73	6.43	6.09	1.02	7.18	7.14	3.80	10.30
	9	10	11	12	13	14	15	16
$x$	21.04	11.56	9.80	16.90	2.58	22.05	10.21	10.83
$y$	3.05	5.89	18.30	7.16	2.60	6.22	12.61	0.08

Table 9: Measures of  $x$  and  $y$  variables

530

To conclude this triple comparative analysis, with the exception of the redundancy equation residuals analysis approach which suffers from a significant difficulty in threshold setting, the so-called sequential and brute force approaches give detection/location results that are quite comparable to the one we

	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18
$r$	12	12	1.2	0.1	1.6	0.2	1	3	1	58	21	13	12	35	9	52	10	7
$r_N$	1	0.92	0.1	0	0.1	0.02	0.1	0.2	0.1	1	0.36	0.23	0.2	0.6	0.16	0.9	0.16	0.13
$s_t^{(2,11)}$	1	1	0	0	0	0	0	0	0	1	1	0	0	0	0	1	0	0

Table 10: Residuals of redundancy equations

propose. However, at the computational level, the brute force approach quickly finds its limitations when  
535 the size of the system increases. The sequential approach remains competitive but has a methodological  
disadvantage, as the estimation of the variables is done in a decoupled way and therefore presents a  
certain pessimism. On the other hand, the proposed approach carries out a global estimation of the  
variables, but, to remain critical, it is necessary to adjust the hyper-parameters  $c_x$  and  $c_z$  of the objective  
function  $\Phi$  (36), in order to satisfy the attenuation of the effects of outliers. However, practice shows  
540 that the choice of these parameters can be made in a wide range of operation.

#### 4. Robust Principal Component Analysis

Principal Component Analysis (*PCA*) is a widely used statistical tool for analyzing data collected from  
a running system to monitor its behavior. However, one of the major drawback of the "ordinary"  
*PCA* approach results from the use of least squares estimation techniques, which often fail to overcome  
545 the bias of accidental measurements, which is unfortunately quite frequent in practice. However, a  
*PCA* model can be constructed from the data without prior filtering, this construction being robust  
to the presence of large errors. The obtained model *PCA* being healthy, i.e. not (or only slightly)  
contaminated by outliers, its use for diagnosis (detection and localization of measurement errors) is  
then efficient [61, 29, 46].

550 After a reminder of its classical formulation and its robust version, we focus on a much less known aspect  
of *PCA* which concerns the generation of residues suitable for the detection/localisation of outliers. The  
performance of these residuals comes from their structuring, which combines the reconstruction of vari-  
ables and their projection in an ad-hoc subspace. The numerical example proposed remains modest in  
size, but allows us to show the implementation of this approach and the results of detection/localisation.

##### 4.1. Recalls on Principal Component Analysis

In practice, we have a matrix of  $X \in \mathcal{R}^N$  data, row vectors  $x_i^T$ , which brings together the  $N$  measure-  
ments made on the  $n$  variables of the system. To search for the set of principal axes, we proceed as  
follows:

- evaluate the matrix of experimental variances and covariances of the centered data:

$$\Sigma = X^T X \quad (38)$$

- 560 • solve, with respect to the  $P$  and  $\Lambda$ , equation :

$$\Sigma P = P \Lambda \quad (39)$$

$P \in \mathcal{R}^{n,n}$  being the [orthogonal matrix](#) of eigenvectors  $p_i$  of  $\Sigma$  and  $\Lambda \in \mathcal{R}^{n,n}$  that, diagonal, of its  
eigenvalues  $\lambda_i$ .

It can also be shown that:

$$X = T P^T \quad T = X P \quad (40)$$

- 565 Relationships (40) essentially find their interest when you decrease the size of the representational space  
(the number of main components used). Once the number  $\ell$  of components to be retained is determined,

the  $X$  matrix of the data can be approximated. To do this, the eigenvector matrix is partitioned into the form :

$$P = \begin{pmatrix} \hat{P} & \tilde{P} \end{pmatrix} \quad \hat{P} \in \mathcal{R}^{n \times \ell} \quad (41)$$

From the decomposition (41), we can then explicite the part  $\hat{X}$  of the data explained by the  $\ell$  first eigenvectors and the residual  $\tilde{X}$  explained by the remaining components :

$$\hat{X} = X\hat{P}\hat{P}^T \quad (42)$$

$$\tilde{X} = X(I - \hat{P}\hat{P}^T) \quad (43)$$

The following sections show how to apply these relationships to a new observation whose consistency you want to test. 570

#### 4.2. Robust Principal Component Analysis Formulation

A major difficulty with the *PCA* is its sensitivity to outliers. To reduce this sensitivity, various techniques can be used, notably the one that consists in performing the *PCA* directly on the data that may be contaminated by outliers using an algorithm that is robust to these outliers. In [23] the authors define a "local" matrix of variances and covariances in the sense that the proposed form tends to emphasize the contribution of close observations to the detriment of distant observations due to the presence of outliers. This matrix noted  $\Sigma^r$  is defined in the following form according to the observations  $x_i$  : 575

$$\Sigma^r = \frac{\sum_{i=1}^{N-1} \sum_{j=i+1}^N w_{i,j} (x_i - x_j)(x_i - x_j)^T}{\sum_{i=1}^{N-1} \sum_{j=i+1}^N w_{i,j}} \quad (44)$$

$$w_{i,j} = \exp\left(-\frac{\beta}{2}(x_i - x_j)^T V^{-1}(x_i - x_j)\right) \quad (45)$$

$\beta$  being a parameter to be adjusted to effectively obtain a reduction in the influence of distant observations, the authors recommend a value close to 2. Decompositions (40) to (43) are then made using  $\Sigma^r$  instead of  $\Sigma$  in (39). 580

#### 4.3. Principle of variable reconstruction

Knowing the robust *PCA* model, the consistency of a new  $x$  measurement vector can now be tested. Considering the previous results (42, 43), we can write the following decomposition:

$$x = \hat{x} + \tilde{x} \quad (46a)$$

$$\hat{x} = C^{(\ell)}x \quad (46b)$$

$$\tilde{x} = (I - C^{(\ell)})x \quad (46c)$$

$$C^{(\ell)} = \hat{P}\hat{P}^T \quad (46d)$$

where  $\hat{x}$  and  $\tilde{x}$  are respectively the projection of  $x$  on the spaces generated by the  $\ell$  main components and the  $n - \ell$  remaining components (residual space). The analysis of the magnitude of the components of  $\tilde{x}$ , or even those of  $\hat{x}$ , can reveal the presence of measurement faults. However, note that  $\hat{x}$  is obtained from all the components of the  $x$  measurement vector. Consequently, the presence of an outlier in the observation vector  $x$  makes the estimate  $\hat{x}$  sensitive to this value and to avoid this we can try to express this estimate using only a part of the observation vector  $x$ . 585

Let's try to estimate the  $r$ -th component of the  $x$  vector. By noting  $c_{ij}$  the elements of the matrix  $C^{(\ell)}$ , the  $r$ -th component of  $\hat{x}$  (46b) becomes explicit: 590

$$\hat{x}_r = \sum_{j=1, j \neq r}^n c_{rj}x_j + c_{rr}x_r \quad (47)$$

where we have particularized the contribution of the  $r$ th component  $x_r$  of the  $x$  measure for reasons that are going to make sense. For the estimation (47), if one wishes not to use the  $r$ -th component  $x_r$  of the measure  $x$ , one can replace, in the right-hand side of the equation (47),  $x_r$  by  $\hat{z}_r$ , which gives (under the condition  $c_{rr} \neq 1$ ) the desired estimate :

$$\hat{x}_r = \frac{[c_{-r}^T \quad 0 \quad c_{+r}^T]}{1 - c_{rr}} x \quad (48)$$

595 where the indices  $-r$  and  $+r$  are respectively used to construct a vector formed by the first  $r-1$  and the last  $n-r$  elements of the vector  $c_r$ . Thus, the  $r$ -th component of  $x$  is estimated using all its components except the  $r$ th. If only the  $r$ -th component of  $x$  is subject to error, then the resulting estimate is not sensitive to this error. Consequently, this partial reconstruction of the measurement vector is noted:

$$\hat{x}^{(r)} = [x_{-r}^T \quad \hat{x}_r \quad x_{+r}^T]^T, \quad \hat{x}^{(r)} \in \mathcal{R}^n \quad (49)$$

600 where  $(\cdot)^{(r)}$  recalls  $r$  index is that of the variable not used in the reconstruction. For the purpose of diagnosis, the estimate (49) should be analysed, e.g. by comparing it with measurements. In fact, it is more interesting to analyze the projection of this estimate in the residual space and this for the  $n$  possible reconstructions according to the value of the  $r$  index. These projections are an indicator of the presence of a fault and can be explained as follows:

$$\tilde{x}^{(r)} = (I - C^{(\ell)})\hat{x}^{(r)} \quad (50)$$

After postponing (49) in (50), we can show that:

$$\tilde{x}^{(r)} = P_r^{(\ell)} x \quad (51a)$$

$$P_r^{(\ell)} = (I - C^{(\ell)}) \left( I + \frac{\xi_r \xi_r^T C^{(\ell)}}{1 - \xi_r^T C^{(\ell)} \xi_r} \right) (I - \xi_r \xi_r^T) \quad (51b)$$

where the vector  $\xi_r \in \mathcal{R}^n$  has all its components equal to one except the  $r$  rank which is equal to zero.

605 **Remark 8.** *It is important to note again that the previous reconstruction  $\hat{x}^{(r)}$  is done using all available measures except the one of rank  $r$ . There are thus  $n$  reconstruction possibilities and this remark will be used later during the phase of isolating the outlier(s). The same applies to the projection  $\tilde{x}^{(r)}$ . ■*

**Remark 9.** *The matrix  $P_r^{(\ell)}$  (51b) has two important special properties. Given its definition the reader can check that :*

$$P_r^{(\ell)} \xi_r = 0 \quad (52)$$

$$\xi_r^T P_r^{(\ell)} = 0 \quad (53)$$

610 *which highlights the peculiar structure of the  $P_r^{(\ell)}$  matrix, namely that the  $r$ -th column and the  $r$ -th row of this matrix all have zero components. Obviously, the examination of the expression (51a) shows that the  $\tilde{x}^{(r)}$  projection has a structure particularly adapted to the detection and localization of aberrant measurements. ■*

**Remark 10.** *The proposed reconstruction with (47, 48, 51a) is concerned with a single component of the measurement vector. Using a classical hypothesis of observability, it is possible to reconstruct several variables at the same time from the same measurement vector. This extension is particularly useful in the presence of several measurements in fault simultaneously [8], [9]. ■*

#### 4.4. Detection and localization of abnormal measurements

To specify the way to detect measurement faults, let us consider the case of a healthy data  $x^*$  corrupted by a noise of zero mean value  $\epsilon$  and a fault of magnitude  $d$  acting in the direction  $\xi_f$  ( $d$  and  $\xi_f$  not being known):

$$x_m = x^* + \epsilon + \xi_f d \quad (54)$$

In this expression  $x^*$  is the true value (and thus satisfies the *PCA* model),  $\xi_f$  is the null vector except its  $f$  component equal to the unit and  $x_m$  is an available observation of  $x^*$ . Under (51a), the residual calculated by reconstructing only the  $r$  rank component of  $x$  is explicit :

$$\begin{aligned}\tilde{x}^{(r)} &= P_r^{(\ell)}(x^* + \epsilon + \xi_f d) \\ &= P_r^{(\ell)}(\epsilon + \xi_f d)\end{aligned}\quad (55)$$

whose mathematical expectation is:

$$\mathcal{E}(\tilde{x}^{(r)}) = P_r^{(\ell)} \xi_f d \quad (56)$$

which highlights the role played by the  $r$ -th row and  $r$ -th column of the  $P_r^{(\ell)}$  projection matrix. One can generalize this analysis by calculating the projection matrices for the various possible directions of fault  $\xi_r$ ,  $r = 1, \dots, n$  as well as the resulting residuals. The analysis of these residuals, thanks to the properties (52) and (53), then makes it possible to detect and localize the fault if it exists. Indeed, let us consider all the possible reconstructions  $r = 1, \dots, n$ . Starting from (56) for the different projection matrices we can state the two rules :

- $R_1$  : if the direction of reconstruction  $\xi_r$  is that of the fault, i.e. if  $r = f$ , then all the components of the vector  $P_r^{(\ell)} \xi_f$  are zero
- $R_2$  : if the direction of reconstruction  $\xi_r$  is different from that of the fault, then the components of the vector  $P_r^{(\ell)} \xi_f$  are not a priori null, except the component of rank  $r$ .

The implementation of this isolation technique is systematic. It requires, at each time instant, the calculation of the projection of the reconstructions according to a set of  $n$  directions, but the projection matrices can be calculated once and for all and applied to the analysis of all new observations acquired on the system.

#### 4.5. Example

A simple example has been constructed with  $n = 7$  variables and  $N = 120$  measures. The  $X$  matrix collecting the data is expressed :

$$X = \begin{pmatrix} x_1^T \\ x_2^T \\ \dots \\ x_N^T \end{pmatrix} \quad (57)$$

where the components of  $x_i$  ( $i = 1, \dots, N$ ) are :

$$\begin{aligned}x_{i,1} &= \sin^2(0.2i)(1 + \cos(0.33i)) \\ x_{i,2} &= 2x_{i,1}(1 + x_{i,1}) \\ x_{i,3} &= \sin(0.5i)(1 + \cos(0.16i)) \\ x_{i,4} &= 3.5x_{i,1} - x_{i,2} \\ x_{i,5} &= x_{i,1} + 0.5x_{i,2} \\ x_{i,6} &= x_{i,1} + x_{i,3} \\ x_{i,7} &= 0.5x_{i,2} + x_{i,3}\end{aligned}\quad (58)$$

To these seven variables are added realizations of variables distributed according to centered normal laws of the same standard deviation equal to 0.02.

A constant amplitude bias 1.5 simulates the presence of outliers  $d_j$  affecting variables  $x_j$ ,  $j = \{1, 2, 3, 5, 7\}$ : observations 14 to 20 for  $x_1$ , from 29 to 35 for  $x_2$ , from 44 to 59 for  $x_3$ , from 74 to 80 for  $x_5$ , from 104 to 110 for  $x_7$ . The objective is to detect them and especially to locate them.

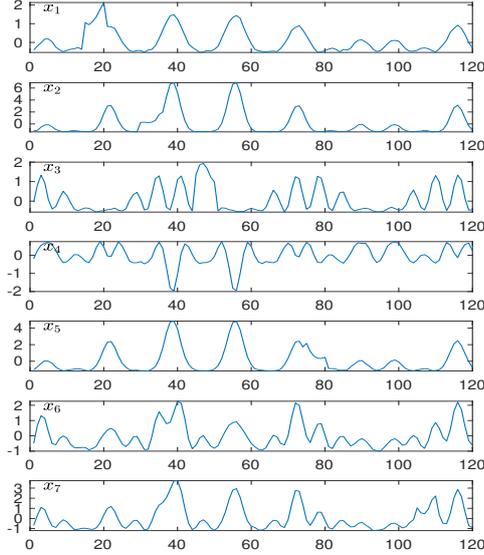


Figure 10: Raw data

Figure (10) represents the temporal evolution of these variables. In accordance with what was said in section 4.4, table 11 indicates the theoretical sensitivity of the projections (51a) with respect to the faults  $d_i, i = 1, \dots, n$  where  $\tilde{x}_{i,j}^{(r)}$  designates the  $j$ th component of the residue evaluated at time  $i$  without using the measure of the variable  $r$ .

For each variable, seven reconstructions can be proposed, which justifies the presence of the seven projections  $\tilde{x}_{i,j}^{(r)}$  in table 11 established thank to the properties (52, 53) (as a reminder,  $r$  is the index of the variable not used for the reconstruction,  $i$  the time,  $j$  the number of one of the seven variables). The presence of a cross indicates the structural influence of an outlier on a residual projection, as opposed to the 0 symbol.

Thus, the analysis of the magnitudes of the residuals  $\tilde{x}_{i,j}^{(r)}$  for  $r = 1 \dots n$  reveals the presence of faults and makes it possible to determine the component of the measurement affected by this fault. Note, however, that the sensitivity table only reflects the occurrence of faults in the residuals independently of their numerical values. Be aware that certain operating conditions can lead to very low numerical sensitivities that do not allow for any conclusion.

Using the raw data contaminated by outliers, we determined the robust *PCA* model by applying the propositions in sections 4.2 and 4.3. The analysis of the decay of the normalized eigenvalues of the variance and covariance matrix, allows us to limit to for the number of principal components to be retained for the reconstruction of the variables using the *PCA* model. Given this model, the procedure in section 4.3 can then be applied to reconstruct the variables and to project he reconstruction errors.

Figures (11) to (16) display the results of detection and isolation of outliers. Each figure has seven graphs, each relating to one of the seven system variables. For reasons of space, only the figures relating to the reconstructions in directions  $\xi_1, \xi_2, \xi_4$  are given.

The 1 to 7 graphs in figure 11 visualize the estimates (in red color) of the seven variables  $\hat{x}_{i,j}^{(1)}, j = 1, \dots, 7$  obtained by reconstruction without using the measurement of the first variable. They can be compared to the measurements (in blue color) and thus highlight the corrections made.

The graphs 1 to 7 in figure 12 concern the residuals  $\tilde{x}_{i,j}^{(1)}, j = 1, \dots, 7$  obtained by projecting the previous reconstructions of the variables obtained without using the 1 variable measure. For this, (51a) was used with the projection matrix  $P_1^{(\ell)}$  elaborated with the direction  $\xi_1 = [1 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0]^T$ .

In a similar fashion, Figures 13 and 14 on the one hand and Figures 15 and 16 on the other have been constructed with the respective directions of projection.  $\xi_2 = [0 \ 1 \ 0 \ 0 \ 0 \ 0 \ 0]^T$  and  $\xi_4 = [0 \ 0 \ 0 \ 1 \ 0 \ 0 \ 0]^T$ . To analyse the different residues obtained, reference should be made to the properties of the projection

		$d_1$	$d_2$	$d_3$	$d_4$	$d_5$	$d_6$	$d_7$
Residues, without using the measure of variable 1	$\tilde{x}_{.,1}^{(1)}$	0	0	0	0	0	0	0
	$\tilde{x}_{.,2}^{(1)}$	0	×	×	×	×	×	×
	$\tilde{x}_{.,3}^{(1)}$	0	×	×	×	×	×	×
	$\tilde{x}_{.,4}^{(1)}$	0	×	×	×	×	×	×
	$\tilde{x}_{.,5}^{(1)}$	0	×	×	×	×	×	×
	$\tilde{x}_{.,6}^{(1)}$	0	×	×	×	×	×	×
	$\tilde{x}_{.,7}^{(1)}$	0	×	×	×	×	×	×
Residues, without using the measure of variable 2	$\tilde{x}_{.,1}^{(2)}$	×	0	×	×	×	×	×
	$\tilde{x}_{.,2}^{(2)}$	0	0	0	0	0	0	0
	$\tilde{x}_{.,3}^{(2)}$	×	0	×	×	×	×	×
	$\tilde{x}_{.,4}^{(2)}$	×	0	×	×	×	×	×
	$\tilde{x}_{.,5}^{(2)}$	×	0	×	×	×	×	×
	$\tilde{x}_{.,6}^{(2)}$	×	0	×	×	×	×	×
	$\tilde{x}_{.,7}^{(2)}$	×	0	×	×	×	×	×
⋮								
Residues, without using the measure of variable 7	$\tilde{x}_{.,1}^{(7)}$	×	×	×	×	×	×	0
	$\tilde{x}_{.,2}^{(7)}$	×	×	×	×	×	×	0
	$\tilde{x}_{.,3}^{(7)}$	×	×	×	×	×	×	0
	$\tilde{x}_{.,4}^{(7)}$	×	×	×	×	×	×	0
	$\tilde{x}_{.,5}^{(7)}$	×	×	×	×	×	×	0
	$\tilde{x}_{.,6}^{(7)}$	×	×	×	×	×	×	0
	$\tilde{x}_{.,7}^{(7)}$	0	0	0	0	0	0	0

Table 11: Sensitivity of projected residues to faults

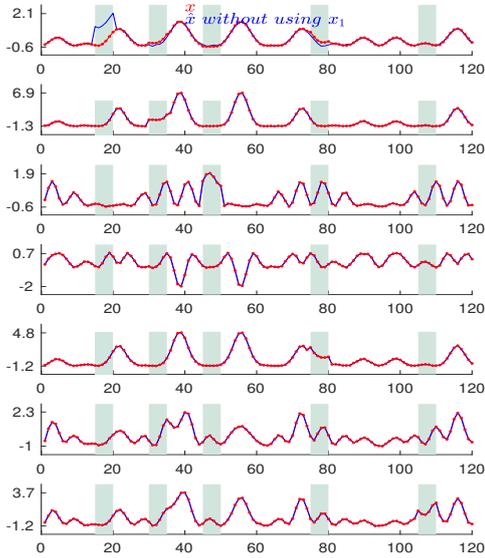


Figure 11: Reconstructions without using variable 1

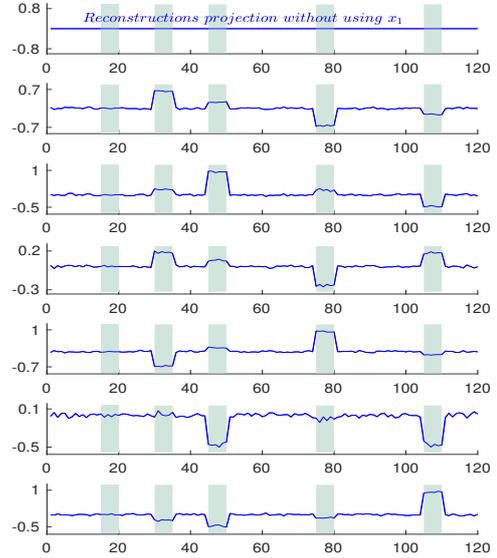


Figure 12: Projection without using variable 1

matrix. Let us consider the first time interval 14 - 20 where a fault has been applied on the variable  $x_1$ . **675** The seven projections are substantially null when the reconstruction is done without the measurements of the first variable. Two hypotheses can be stated: absence of fault or presence of a fault in the direction

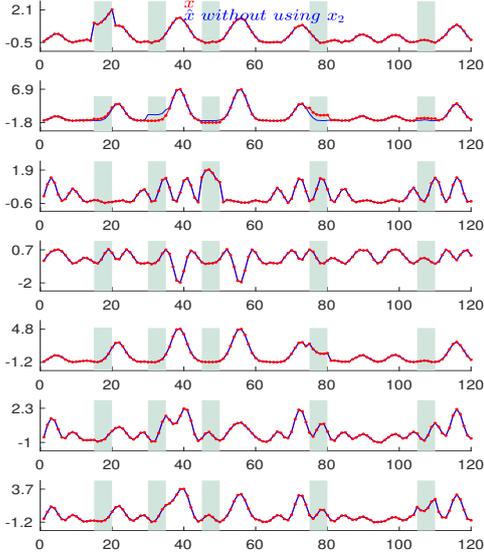


Figure 13: Reconstruction without using variable 2

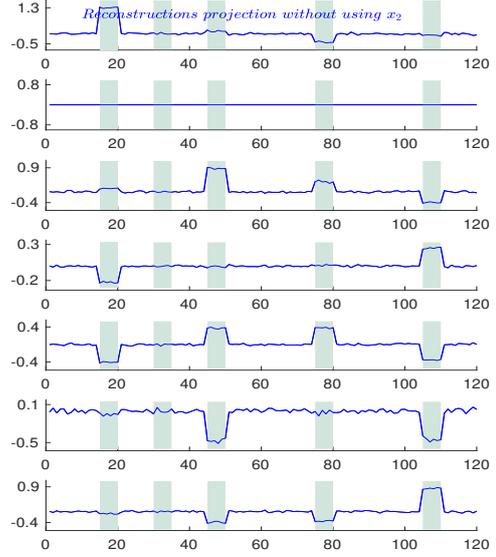


Figure 14: Projection without using variable 2

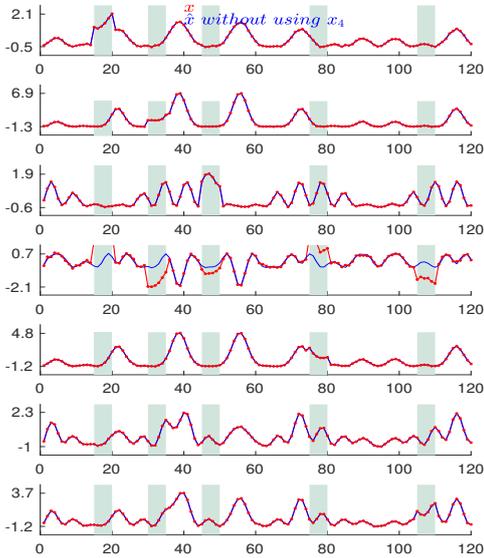


Figure 15: Reconstruction without using variable 4

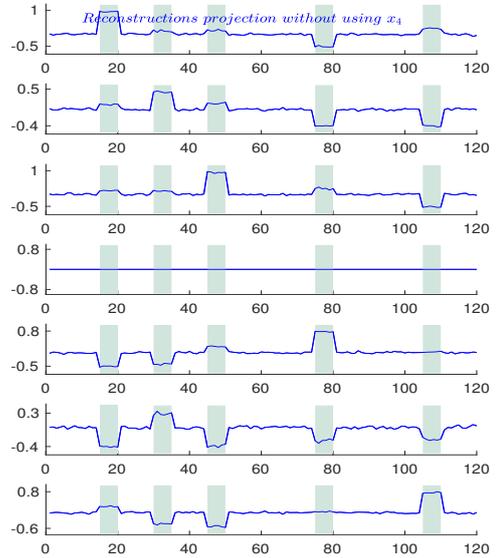


Figure 16: Projection without using variable 4

$\xi_1$  that is to say affecting the variable  $x_1$ . This figure alone is not enough to resolve the ambiguity. In a similar fashion examining the figure 14 (constructed without using the measure of the second variable) always for the time interval 14 - 20 reveals non-zero values for several projections and therefore the existence of a fault. Since two of the projections for the first variable are non-zero, it is this variable that is in fault. Still for this time interval 14 - 20, the reader will be able to interpret in a similar way the seven graphs of figure 16 in order to confirm the previous conclusion. The detection and isolation of a measurement fault for this interval is thus achieved. They can be repeated for the other time intervals and conclude to the presence of faults on the other two variables. The graphical results are in perfect agreement with those of the detection/location signatures in Table 10.

This result can be further reinforced from the figures 17, 18 and 19. The graph number  $j$  in figure 17 shows the sums of the projections taken in absolute values when the variable  $j$  is not used. As an example, the first graph of this figure is obtained from the sum of the signals of figure 11 taken in absolute values. Clearly, the figure 17 highlights the signatures of the faults and allows to find the

presence of fault on the variables  $\{1, 2, 3, 5, 7\}$ . The instants of occurrence of the faults are also obtained from this figure by applying a classical jump detection technique.

More synthetically, this jump detection can be performed from the signal represented in figure 19 which is none other than the sum of the signals in figure 17. The chosen magnitude threshold equal to 2 makes it possible to find the times of occurrence of the five faults which are respectively defined by the 695 intervals :  $[14\ 20]$ ,  $[29\ 35]$ ,  $[44\ 50]$ ,  $[74\ 80]$ ,  $[104\ 110]$ .

Finally, figure 18 partially reproduces those dedicated to reconstruction errors. However, it is limited to visualizing the reconstruction error of each variable  $j, j = 1, \dots, n$  when the measure of this variable is not used in the reconstruction procedure. From each graph, an estimate of the bias can be extracted, the bias corresponding, with the exception of noise, to the reconstruction error. Thus, for the first graph 700 relating to the reconstruction error of the first variable, the magnitude of the reconstruction error in the interval  $[14\ 20]$  indicates a value of 1.5 which effectively corresponds to the bias affecting this variable. The reader will make the same kind of observation for variables 2, 3, 5, 7 from the graphs associated with these variables.

This succinct presentation of the use of *PCA* for the detection of outliers would require much further 705 development. The previous example presents a relatively simple situation where the magnitude of the faults is large enough to allow their detection. The magnitude of the detectable faults is to be compared with the magnitude of the measurement noise but also with the quality of the *PCA* model. The reader will easily understand that in order to remain concise, we have omitted to present a systematic analysis of the hyper-parameters of the proposed technique and their influences on the results of fault detec- 710 tion/location: number of measurements, richness of information in the measurements used, robustness factor  $\beta$  in the calculation of the variance-covariance matrix, signal-to-noise ratio.

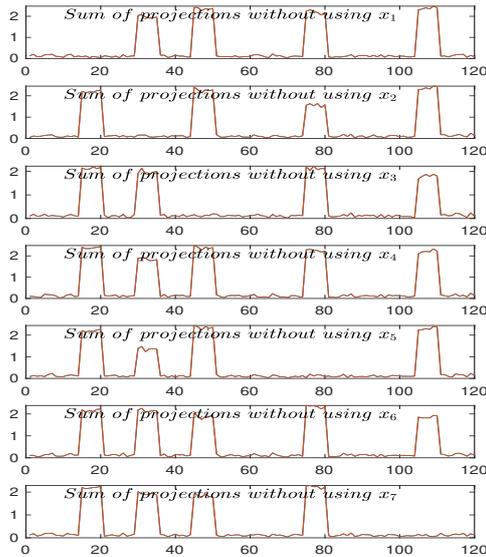


Figure 17: Sum of projections

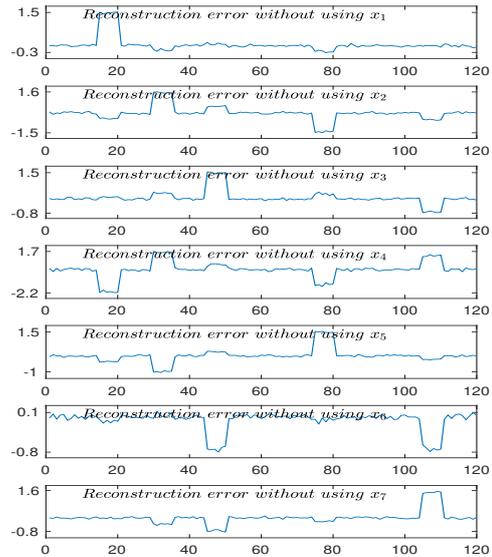


Figure 18: Reconstruction errors

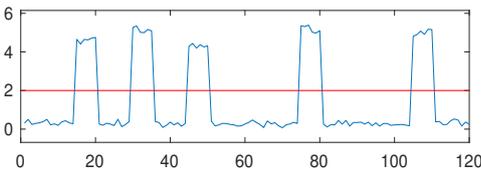


Figure 19: Global fault detection indicator

## 5. Conclusion

715 Although partial and with restrictive assumptions regarding the static nature of the systems consid-  
ered, this presentation has attempted to draw the reader's attention to the presence of outliers in the  
measurements and how to take them into account. Two points of view are presented, the first relating  
to the location of these outliers and their replacement by substitute values, the second relating to their  
accommodation, i.e. their acceptance by minimising their influence in the estimation procedure. This  
720 second point of view has been detailed on the one hand in a robust measurement validation procedure  
and on the other hand in the use of robust principal component analysis, the robustness being to be  
understood in terms of reducing the influence of outliers. The two techniques presented in sections 3  
and 4 were presented in the somewhat reductive framework of static systems, but their extension to the  
case of dynamic systems does not pose any methodological problems.

## Remerciements

725 Some of the results presented in this communication have been established by a research team. The  
author wishes to express his gratitude to Dr. Gilles Mourot. Special thanks to former PhD students  
Mohamed-Faouzi Harkat and Yvon Tharault.

## 6. Bibliography

- 730 [1] H. Alighardashi, N.M. Jan, B. Huang. Expectation Maximization Approach for Simultaneous Gross  
Error Detection and Data Reconciliation Using Gaussian Mixture Distribution. *Industrial & Engi-  
neering Chemistry Research*, 56 (49), 14530-14544, 2017.
- [2] M. Antolin, E. Del Barrio, J-M. Loubes. A data driven trimming procedure for robust classification.  
2017. *hal* – 01437147, <https://arxiv.org/pdf/1701.05065.pdf>.
- 735 [3] A. Archimbaud. Détection non-supervisée d'observations atypiques en contrôle de qualité : un  
survol. *Journal de la Société Française de Statistique*, 159 (3), 2018.
- [4] C.C. Aggarwal. *Outlier Analysis* (2 ed.). Springer, 2016.
- [5] V. Barnett, T. Lewis. *Outliers in Statistical Data*, 3rd edition, USA: Wiley & Sons, 1994.
- [6] C. Barreyre. *Statistiques en grande dimension pour la détection d'anomalies dans les données  
fonctionnelles issues des satellites*. Thèse INSA de Toulouse, 2018.
- 740 [7] R.J. Beckman, R.D. Cook, R. D. Outlier.....s. *Technometrics*, 25 (2), 119-149, 1983.
- [8] A. Ben Aicha, G. Mourot, K. Benothman, J. Ragot. Détermination de modèles ACP pour la  
détection et la localisation de défauts de capteurs. *Journal Européen des Systèmes Automatisés*, 46  
(1), 9-32, 2012.
- 745 [9] A. Ben Aicha, G. Mourot, K. Benothman, J. Ragot. Determination of Principal Component Anal-  
ysis models for sensor fault detection and isolation. *International Journal of Control, Automation  
and Systems*, 11 (2), 296-305, 2013.
- [10] S.E. Benkabou. *Détection d'anomalies dans les séries temporelles : application aux masses de  
données sur les pneumatiques*. Thèse de l'Université de Lyon, 2018.
- 750 [11] A. Blazquez-Garcia, A. Conde, U. Mori, J.A. Lozano. A review on outlier/anomaly detection in  
time series data. *arXiv:2002.04236*, 2020.

- [12] C.L. Brown, R.F. Brcich, C. Debes. Adaptive M-estimators for robust covariance estimation, In Proceedings of the International Symposium on Applied Stochastic Models and Data Analysis, Brest, France, 2005.
- [13] O. Cencic, R. Frühwirth. Data reconciliation of non-normal observations with nonlinear constraints, *Journal of Applied Statistics*, 45 (13), 2411-2428, 2018. **755**
- [14] V. Chandola, V. Kumar. Outlier Detection : A Survey. *ACM Computing Surveys*, 41, 2009.
- [15] N. Chèze, J.M. Poggi. Détection par boosting de données aberrantes en régression. *Revue des Nouvelles Technologies de l'Information*, 159-171, 2008.
- [16] A.S. Chinen, J.C. Morgan, B.P. Omell, D. Bhattacharyya, D.C. Miller. Dynamic Data Reconciliation and Model Validation of a MEA-Based CO<sub>2</sub> Capture System using Pilot Plant Data, *IFAC-PapersOnLine*, 49 (7) 639-644, 2016. **760**
- [17] P. Čížek, S. Sadıkoğlu, Robust nonparametric regression: A review. Retrieved from <https://doi.org/10.1002/wics.1492>, 2020.
- [18] A. Deleforge, F. Forbes, R. Horaud. High-dimensional regression with Gaussian mixtures and partially-latent response variables. *Statistics and Computing*, 25 (5), 893-911, 2015. **765**
- [19] W.J. Dixon. *The Annals of Mathematical Statistics*. 21 (4), 488-506, 1950.
- [20] F.Z. Dogru, O. Arslan. Robust mixture regression modeling using the least trimmed squares (LTS)-estimation method. *Communications in Statistics - Simulation and Computation*, 2017.
- [21] R. Domingues, M. Filippone, P. Michiardi, J. Zouaoui. A comparative evaluation of outlier detection algorithms: Experiments and analyses, *Pattern Recognition*, 74, 406-421, 2018. **770**
- [22] S.S. Du, Y. Wang, S. Balakrishnan, P. Ravikumar, A. Singh. Robust Nonparametric Regression under Huber's contamination Model, arXiv:1805.10406, 2018.
- [23] M. Fekri, A. Ruiz-Gazen. Robust weighted orthogonal regression in the errors-in-variables model. *Journal of Multivariate Analysis* 88, 89-108, 2003.
- [24] S. Fellaou, T. Bounahmidi. Mass Balance Reconciliation for Bilinear Systems: A Case Study of a Raw Mill Separator in a Typical Moroccan Cement Plant. *Engineering, Technology & Applied Science Research*, 6 (3), 2016. **775**
- [25] A.J. Fox. Outliers in Time Series. *Journal of the Royal Statistical Society: Series B (Methodological)*, 34 (3), 350-363, 1972.
- [26] P. Fearnhead, G. Rigall. Change point Detection in the Presence of Outliers. *Journal of the American Statistical Association*, 2017. **780**
- [27] F.E. Grubbs. Procedures of Detection Outlying Observations in Samples, *Technometrics*, 4 (1), 1-21, 1969.
- [28] F.R. Hampel. A general qualitative definition of robustness. *The Annals of Mathematical Statistics*, 42 (6), 1887-1896, 1971. **785**
- [29] M.F. Harkat, G. Mourot, J. Ragot. An improved PCA scheme for sensor FDI: application to an air quality monitoring network. *Journal of Process Control*, 16 (6), p. 625-634, 2006.
- [30] F. Harle. Détection de ruptures multiples dans des séries temporelles multivariées : application à l'inférence de réseaux de dépendance. *Traitement du signal et de l'image*. Thèse de l'Université Grenoble Alpes, 2016. **790**

- [31] D.M. Hawkins. Identification of Outliers. Chapman and Hall: New York, 1980.
- [32] D. Hodouin. Process Observers and Data Reconciliation Using Mass and Energy Balance Equations. In: Sbárbaro D., del Villar R. (eds) Advanced Control and Supervision of Mineral Processing Plants. Advances in Industrial Control, Springer, London, 2010.
- 795 [33] P.W. Holland, R.E. Welsch. Robust regression using iteratively reweighted least squares. Commun. Statistics-Theory and Methods, A6: 813-827, 1977.
- [34] C. Hong, M. Hausrecht. Multivariate Conditional Outlier Detection: Identifying Unusual Input-Output Associations. International Florida Artificial Intelligence Research Society Conference, 2018.
- 800 [35] B Hoppenstedt. Towards a Hierarchical Approach for Outlier Detection in Industrial Production Settings. EDBT/ICDT Joint Conference, Lisbon, Portugal, 2019.
- [36] J. Horton, R.L. Stuart. Multiple imputation in practice: comparison of software packages for regression models with missing variables. The American Statistician, 55, 244-254, 2001.
- [37] P.J. Huber. Robust estimation of location parameter. The Annals of Mathematical Statistics, 35  
805 (1), 73-101, 1964.
- [38] M. Hubert, P.J. Rousseeuw, K. Van den Branden. ROBPCA: a New Approach to Robust Principal Component Analysis. Technometrics, 47, 64-79, 2005.
- [39] S.A. Imtiaz, S.L. Shah, S. Narasimhan. Missing Data Treatment Using Iterative PCA and Data Reconciliation, IFAC Proceedings Volumes, 37 (9), 2004.
- 810 [40] R. Isermann, P. Ballé: Trends in the Application of Model-Based Fault Detection and Diagnosis of Technical Processes. Control Engineering Practice, 5 (5), 709-719, 1997.
- [41] M. Kallas, G. Mourot, D. Maquin, J. Ragot. Data driven approach for fault detection and isolation in nonlinear system. International Journal of Adaptive Control and Signal Processing, 32 (11), 1569-1590, 2018.
- 815 [42] M. King. Data Reconciliation. in Statistics for Process Control Engineers: A Practical Approach, John Wiley & Sons, Ltd, Chichester, UK, 2017.
- [43] N. Kolokas, T. Vafeiadis, D. Ioannidis, D. Tzovaras. A generic fault prognostics algorithm for manufacturing industries using unsupervised machine learning classifiers. Simulation Modelling Practice and Theory, 103, 2020.
- 820 [44] T. Korpela, O. Suominen, Y. Majanne, V. Laukkanen, P. Lautala. Robust data reconciliation of combustion variables in multi-fuel fired industrial boilers. Control Engineering Practice, 55, 101-115, 2016.
- [45] Y.H. Kuo, Z. Li, D. Kifer. Detecting Outliers in Data with Correlated Measures. <https://arxiv.org/pdf/1808.08640>, 2018.
- 825 [46] W. Li, M. Peng, Q. Wang. Improved PCA method for sensor fault detection and isolation in a nuclear power plant. Nuclear Engineering and Technology, 51 (1), 146-154, 2019.
- [47] C.E. Llanos, M.C. Sánchez, R.A. Maronna. A robust methodology for the sensor fault detection and classification of systematic observation errors. Computer Aided Chemical Engineering, 40, 1525-1530, 2017.

- [48] D. Lim, B. Park, D. Nott, X. Wang, T. Choi. Sparse signal shrinkage and outlier detection in high-dimensional quantile regression with variational Bayes. *Statistics end its interface*, 13, 237-249, 2020. **830**
- [49] A. Mami, A. Jaber, O. Almagrouk. Applying Bootstrap Robust Regression Method on Data with Outliers. *International Journal of Sciences: Basic and Applied Research (IJSBAR)*. 143-160, 2020.
- [50] D. Maquin, J. Ragot. Comparison of gross errors detection methods in process data. 30th IEEE Conference on Decision and Control, 2253-2261, Brighton, 1991. **835**
- [51] J. Marzat, H. Piet-Lahanier, S. Bertrand. Cooperative fault detection and isolation in a surveillance sensor network: a case study. *IFAC-PapersOnLine*, Elsevier, 251 (24), 790-797, 2018.
- [52] M.A. Moussa. Data gathering and anomaly detection in wireless sensors networks. Université Paris-Est, 2017. **840**
- [53] M. Nikulin, A. Zerbet. Détection des observations aberrantes par des méthodes statistiques. *Revue de statistique appliquée*, 50 (3), 25-51, 2002.
- [54] J.C. Ondo, T.B.M.J. Ouarda, V. Fortin, B. Bobée. Procédures bayésiennes pour la détection d'observations singulières : synthèse bibliographique. *Journal de la société française de statistique*, 142 (2), 41-74, 2001. **845**
- [55] R. Pothina, R. Ganguli. Detection of Subtle Sensor Errors in Mineral Processing Circuits Using Data-Mining Techniques. *Mining, Metallurgy & Exploration*, 37, 399-414, 2020.
- [56] J.E. Potter, M.C. Suman. Thresholdless redundancy management with arrays of skewed instruments. *Integrity in electronic flight control systems. AGARDograph*, 224:15-25, 1997.
- [57] V. Planchon. Traitement des valeurs aberrantes : concepts actuels et tendances générales. *Biotechnology, Agronomy, Society and Environment*, 9 (1), 19-34, 2005. **850**
- [58] P.J. Rousseeuw, A.M. Leroy. *Wiley Series in Probability and Statistics. Wiley Series in Probability and Statistics*, 1987.
- [59] P.J. Rousseeuw, M. Hubert. Anomaly Detection by Robust Statistics, "WIREs Data Mining and Knowledge Discovery, e1236, 1-14, 2018. **855**
- [60] N.N.R. Ranga Suri, M. Narasimha Murty, G. Athithan. *Outlier detection : techniques and applications, a data mining perspective*. Springer, 2019.
- [61] P. Saha, N. Roy, D. Mukherjee, A.K. Sarkar. *Application of Principal Component Analysis for Outlier Detection in Heterogeneous Traffic Data, Procedia Computer Science*, 83, 2016.
- [62] M.A. Sayed. Représentations pour la détection d'anomalies : Application aux données vibratoires des moteurs d'avions. Université Paris-Saclay, 2018 (in french). **860**
- [63] O. Shetta, M. Niranjana. Robust subspace methods for outlier detection in genomic data circumvents the curse of dimensionality. *Royal Society Open Science*, 7: 190714. <http://dx.doi.org/10.1098/rsos.190714>, 2020.
- [64] J.E. Serth, W.A. Heenan. Gross error detection and data reconciliation in steam-metering system. *AIChE J.*, 32 (5), 733-742, 1986. **865**
- [65] F. Tahir et al. Process Monitoring and Fault Detection on a Hot-Melt Extrusion Process Using in-Line Raman Spectroscopy and a Hybrid Soft Sensor. *Computers and Chemical Engineering*, 125, 400-414, 2019.

- 870 [66] R.S. Tsay, D. Pena, A.E. Pankratz. Outliers in multivariate time series. *Biometrika*, 87 (4), 789-804, 2000.
- [67] E.C. de Valle, R. de Arajo Kalid, A.R. Secchi, A. Kiperstok. Collection of benchmark test problems for data reconciliation and gross error detection and identification. *Computers & Chemical Engineering*, 111, 134-148, 2018.
- 875 [68] A. Virouleau, A. Guilloux, S. Gaïffas, M. Bogdan. High-dimensional robust regression and outliers detection with slope. hal-01798400, <https://hal.archives-ouvertes.fr/hal-01798400>, 2018.
- [69] X. Xu, H. Liu, M. Yao. Recent Progress of Anomaly Detection. *Complexity*, Article ID 2686378, <https://doi.org/10.1155/2019/2686378>, 2019.
- [70] S. Zair. Détection de données aberrantes appliquée à la localisation GPS. Thèse de doctorat Université Paris Sud, 2016.
- 880 [71] Z. Zhang, Y.Y. Chuang, J.Chen. Using clustering based logical equation set to decompose large scale chemical processes for parallel solving data reconciliation and parameter estimation problem. *Chemical Engineering Research and Design*, 120, 396-409, 2017.
- [72] Y. Zhang and X. Wang. Dynamic Data Reconciliation Based on an Improved Robust M-Estimator. Chinese Automation Congress, Xi'an, China, 2018.
- 885

Prof. José RAGOT  
Centre de Recherche en Automatique de Nancy  
UMR 7039 - Université de Lorraine - CNRS  
2, Avenue de la forêt de Haye, TSA 60 604  
54 518 Vandoeuvre-lès-Nancy Cedex, FRANCE

Measurement journal

This proposal falls within the general theme of sensor diagnostics and data validation. More specifically, it deals with how to detect and locate abnormal values in a given time series by a sensor or set of sensors. Thus, to be clear, this proposal does not deal with the design of sensors but focuses exclusively on the coherence analysis of data provided by sensors.

Highlights of our paper :

- Localization of abnormal values in time series
- Robust approach in data reconciliation
- Structuration of fault indicator in Principal Component Analysis

Prof. José RAGOT  
Centre de Recherche en Automatique de Nancy  
UMR 7039 - Université de Lorraine - CNRS  
2, Avenue de la forêt de Haye, TSA 60 604  
54 518 Vandoeuvre-lès-Nancy Cedex, FRANCE

Measurement journal

The authors whose names are listed immediately below certify that they have NO affiliations with or involvement in any organization or entity with any financial interest (such as honoraria; educational grants; participation in speakers' bureaus; membership, employment, consultancies, stock ownership, or other equity interest; and expert testimony or patent-licensing arrangements), or non-financial interest (such as personal or professional relationships, affiliations, knowledge or beliefs) in the subject matter or materials discussed in this manuscript.

Author's name : José RAGOT

Date : Mai 2020, 14