



ELSEVIER

Contents lists available at [ScienceDirect](http://www.sciencedirect.com)

## Journal of Sound and Vibration

journal homepage: [www.elsevier.com/locate/jsvi](http://www.elsevier.com/locate/jsvi)

# On robust regression analysis as a means of exploring environmental and operational conditions for SHM data



N. Dervilis\*, K. Worden, E.J. Cross

Dynamics Research Group, Department of Mechanical Engineering, University of Sheffield, Mappin Street, Sheffield S1 3JD, England

## ARTICLE INFO

## Article history:

Received 29 May 2014

Received in revised form

25 February 2015

Accepted 27 February 2015

Handling Editor: K. Shin

Available online 20 March 2015

## ABSTRACT

In the data-based approach to structural health monitoring (SHM), the absence of data from damaged structures in many cases forces a dependence on novelty detection as a means of diagnosis. Unfortunately, this means that benign variations in the operating or environmental conditions of the structure must be handled very carefully, lest they lead to false alarms. If novelty detection is implemented in terms of outlier detection, the outliers may arise in the data as the result of both benign and malign causes and it is important to understand their sources. Comparatively recent developments in the field of robust regression have the potential to provide ways of exploring and visualising SHM data as a means of shedding light on the different origins of outliers. The current paper will illustrate the use of robust regression for SHM data analysis through experimental data acquired from the Z24 and Tamar Bridges, although the methods are general and not restricted to SHM or civil infrastructure.

© 2015 The Authors. Published by Elsevier Ltd. This is an open access article under the CC BY-NC-ND license (<http://creativecommons.org/licenses/by-nc-nd/4.0/>).

## 1. Introduction

The effect of changing environmental and operational conditions on a structure is gaining significant attention in the fields of structural health monitoring (SHM) and system identification. This issue is a key concern, as the measured responses from a structure and the extracted features that are sensitive to damage or structural degradation are usually also sensitive to any change in operational and environmental conditions. This effect is especially obvious in monitoring procedures for civil infrastructure as the measured responses of structures in operation, like bridges or wind turbines, are subject to continuous variations due to temperature, humidity, ice, wind loading or even traffic loading for bridges. In such cases the structures will react by introducing into their time responses effects which can mask any indication or sensitivity in structural response that would reveal the presence of damage or that could affect the prediction of certain parameters through a regression analysis.

In such cases, the effects of the environmental and operational variation must be considered and identified before using a reliable feature for revealing any structural condition. This is one of the dominant factors that has to be considered when SHM technology is adopted by industry if it is a system that has to run continuously and on-line.

In the SHM literature, one of the most challenging tasks is centred around the influence of temperature on structural response. Especially for bridges, which provide the illustrative material for this study, temperature is generally considered to

\* Corresponding author.

E-mail address: [N.Dervilis@sheffield.ac.uk](mailto:N.Dervilis@sheffield.ac.uk) (N. Dervilis).

be an important environmental factor which affects the dynamic response of the structure, due to the influence on the stiffness of structural parameters and on the boundary conditions of a structure [1–5].

Besides temperature, the importance of other operational conditions such as wind conditions on the bridge structure will be addressed in this study (a concern for long-span bridges).

In the context of civil monitoring campaigns, the authors need to state the importance of monitoring such large structures even if the immediate concern is not the diagnosis of failure. As clearly stated in [6,7], health and performance monitoring can offer certain long-term advantages that can help when key management decisions need to be made. This includes issues like material degradation and fatigue assessment, structural evaluation and knowledge enhancement after earthquake events or explosions, designing modifications to the current structure and for future reference as well as better understanding of measured structural responses influenced by daily and seasonal variation. In summary, the aim of monitoring systems regarding large static structures could be in the short-term detection of fault or health novelty/prognosis or in the long-term to develop robust tools that can lead to effective management, integration and understanding of structural performance with the minimum cost.

Various methods and algorithms have been proposed in order to counteract and remove the influence of external variations such as principal component analysis, auto-associative neural networks or, more recently, cointegration [1,8]. Although these methods exhibit a series of advantages and disadvantages in terms of removing the influence of operational and environmental conditions, very little effort has been carried out in terms of characterising and distinguishing which of the outliers indicated in data are “good” in terms of representing environmental/operational variations and which are “bad” in terms of representing damage or structural degradation. This distinction would be a significant step forward as someone could remove by this classification procedure any external benign variations, leaving the algorithm to detect alone the true structural (health or performance) anomalies.

In this paper, an approach of robust regression and robust multivariate statistics is exploited as a means of characterising and distinguishing the influence of environmental and operational conditions on the structural response. It will be shown via the successful implementation of robust regression analysis that environmental and operational conditions can manifest themselves differently compared to damage condition. For more details regarding the different nature of outliers in multivariate statistics the reader can consult [9,10].

There are many definitions regarding what is an outlier. Hawkins [11] and Barnett and Lewis [12] give two general definitions. Barnett and Lewis indicate that “an outlying observation, or outlier, is one that appears to deviate markedly from other members of the sample in which it occurs”. Hawkins defines an outlier “as an observation that deviates so much from other observations as to arouse suspicion that it was generated by a different mechanism”. But which are the main mechanisms that produce outliers? According to Hawkins [11] there are two basic mechanisms that are responsible for the contamination of data with outliers. The first mechanism assumes that the data is coming from heavy-tailed distributions (like the  $t$ -distribution). The second mechanism assumes that data is coming from different kinds of distributions. One of these may be the general distribution that generates the “good” measurements and the other one may be the “infecting” distribution that gives the contaminated measurements. If the contaminating distribution has tails that are heavier than those of the good, then there will be a possibility for the observations coming from the infecting distribution to appear as outliers.

The algorithms that will be described here are the least trimmed squares (LTS) regression algorithm and the minimum covariance determinant (MCD) estimator. Outlier detection methods have been utilised for many different applications in the past and cover a broad range of fields of research such as econometrics, computer sciences, medical and biological sciences, meteorology and even political science. Some recent advances regarding robust outlier detection can be found in [13–20].

Generally, in linear regression analysis the studied examples of data are of the type  $(\{x_i\}, \{y_i\})$  where  $\{x_i\}$  is an  $m$ -dimensional input vector and the response  $\{y_i\}$  is often a one-dimensional vector [9,10]. Cases for which an  $\{x_i\}$  is far away from the majority of the  $\{x_i\}$  observations are called leverage points.

To distinguish between leverage points (good and bad as is explained later) and outliers, someone must take into account the regression model and response  $\{y_i\}$  in parallel with  $\{x_i\}$  as the linear pattern of the multivariate space dictated by the majority of the observations will lead to the best result. This is the key reason that a high-breakdown (very robust) regression tool such as the least trimmed squares (LTS) regression algorithm is critical.

To summarise the central objective of the paper is straightforward, it is to introduce a new means of analysing and visualising SHM data. This work aims to reveal the structural “DNA” of the data by distinguishing environmental and operational variations from damage. Robust residuals (LTS) and discordancy distances (MCD) are a critical and vital starting point of the analysis. This can lead not only to robustly visualise and analyse how differently the observations manifest themselves in space in terms of different sources of outliers but also can lead to the removal of the external benign variations. The results may tell one to change the model due to alternations of the physical characteristics of the system due to outlier presence or go back to the original data and explain the source of outliers. The main novel element of this work for SHM is the high level estimation of the difference between leverage points and outliers. Via this analysis one can identify and detect benign variability at an early stage and establish a normal condition clear from external influences offering a vital advantage over other multivariate methodologies. At the end of the day the strategy steps focus on explaining that different forms of outliers give distinct and different characteristics with respect to environmental and operational variations and damage. Characteristics not only can detect different

forms of outliers and leverage points but also and most importantly relate and correlate it with the different sources of these outliers.

To demonstrate the effectiveness of these methods, two experimental applications to the Z24 Bridge and Tamar Bridge are presented. First the theoretical background of the algorithms is explained and then the methods are applied to the two experimental examples.

## 2. Multiple outlier detection: the minimum covariance determinant (MCD) estimator

The classic discordancy measure for indicating outliers, as used in many of the previous studies, is the Mahalanobis squared-distance (MSD), which is given by the following equation:

$$D_i^2 = (\{x_i\} - \{\mu_x\})^T [\Sigma]^{-1} (\{x_i\} - \{\mu_x\}) \quad (1)$$

where  $\{x_i\}$  is the potential outlier,  $\{\mu_x\}$  is the mean of the sample observations and  $[\Sigma]$  is the sample covariance matrix. The mean and covariance matrix could be inclusive or exclusive measures; that is to say that the statistics may or may not have been computed from data where outliers are already present. Generally, in many different fields the test set (outlier) is not known *a priori* and an inclusive approach is a necessity. However, in the context of SHM or condition monitoring this situation presents a series of drawbacks regarding the use of multivariate statistics.

The MSD tells one how far away a specific measurement is from the centre of the training data cloud, relative to the size of the cloud.

The main disadvantage of the classical distance measures (like the (MSD)) is that they can suffer from a multiple outlier “masking effect”. If there were groups of outliers already present in the training data, they would have a critical influence on the sample mean and covariance in such a way that they would subsequently indicate small distances on new observations or outlying data and thus cause the outliers to become invisible. The arithmetic mean and unbiased covariance matrix are statistics that suffer heavily from multiple outliers present in the data. Specifically, when outliers from a cluster cloud that lie inside the data are present then they will directly move the arithmetic mean towards them and even expand the classical tolerance ellipsoid in their direction [9].

The application of robust computation of location and covariance estimation of multivariate data is of significant interest in the investigation for inclusive outliers. This is the reason that the method that is introduced here is the minimum covariance determinant (MCD) estimator which is much more robust against outliers in the training data. This algorithm has already been used in an SHM context in [21].

To make it clear to the reader, the significant importance of the MCD tool against the MSD index, in simple terms, searches inclusively for multiple outliers by removing the “masking effect” and revealing in multivariate data their infectious presence.

The computation of the MCD estimator is not a trivial procedure and requires an extensive calculation. In the current study, the FAST-MCD algorithm is implemented [9,10,22–26]. The algorithm is given in detail in the references [9,10,22–25], and the code was provided via a statistical Matlab library called LIBRA [23]. A brief description of the algorithm is provided by explaining the basic steps of the FAST-MCD technique. The description is given in order to make the present paper a little more self-contained.

A multivariate data matrix  $[X] = (\{x_1\}, \dots, \{x_m\})^T$  is assumed of  $m$  points in  $n$  dimensional space ( $n \times m$ ) where  $\{x_i\} = (x_{i1}, \dots, x_{in})^T$  is a single observation of the feature vector of interest. Robust estimates of the centre  $\{\mu\}$  and the scatter matrix  $[\sigma]$  of  $X$  can be calculated by the MCD estimator. The MCD tool looks for the  $h$  ( $> m/2$ ) observations out of  $m$  whose classical covariance matrix has the lowest possible determinant. The raw MCD estimate of location (arithmetic mean) is then computed from the average of these  $h$  points and the raw MCD estimation of scatter is the covariance matrix multiplied by a consistency factor.

The calculation of the lowest determinant is critical as the algorithm moves from one approximation of the MCD to another one with lower determinant in order to find the optimal estimates of the centre  $\{\mu\}$  and the scatter matrix  $[\sigma]$  that are free from inclusive outliers. This theorem and the proof that follows it are not obvious and can be found in the appendix of [22].

Based on these raw MCD estimates, a reweighting step can be added in order to increase the finite sampling efficiency. The advantage is that MCD estimates can resist up to  $(m-h)$  outliers and in turn, the number  $h$  (or equally  $a = h/m$ ) controls the robustness of the estimator. The highest resistance compared to contamination is achieved by calculating  $h = (n+m+1)/2$  [9,10,22–26]. It is proposed that when a large proportion of contamination is assumed then  $h = an$  with  $a = 0.5$ . Detecting outliers can be challenging when  $m/n$  is small because some data points can become coplanar. This is a general problem in the machine learning community called the “curse of dimensionality”. It is recommended [23] that when  $m/n > 5$ ,  $a$  should be 0.5. Generally, the MCD estimates of location and scatter are affine equivariant which means that they are invariant under affine transformation behaviour.

This last property is crucial as the underlying model is then immune to different variable scales and data rotations. Rousseeuw and Van Driessen [22] developed the FAST-MCD algorithm based on a Concentration step (C-step). C-steps select the  $h$  observations with the smallest distances and the scatter matrix with the lowest determinant [22] and the main details are given here. Assume  $[X] = (\{x_1\}, \dots, \{x_m\})$  and let an  $h$ -subset  $H_1 \rightarrow (1, \dots, n)$  which has  $|H_1| = h$ . Then  $\{\hat{\mu}_1\} = (1/h) \sum_{i \in H_1} \{x_i\}$

and  $[\widehat{\Sigma}_1] = (1/h) \sum_{i \in H_1} (\{x_i\} - \{\widehat{\mu}_1\})(\{x_i\} - \{\widehat{\mu}_1\})'$  and when  $[\widehat{\Sigma}_1] \neq 0$  relative distances can be defined as

$$\{d_1(i)\}^2 = (\{x_i\} - \{\widehat{\mu}_1\})' [\widehat{\Sigma}_1]^{-1} (\{x_i\} - \{\widehat{\mu}_1\}) \quad \text{for } i = 1, \dots, m \quad (2)$$

If in the initial subset  $\det([\widehat{\Sigma}_1]) = 0$  then the initial subset can be further extended by adding observations until  $\det([\widehat{\Sigma}_1]) > 0$ . In this way there is avoidance of the algorithm becoming stuck in the initial full dimensional space.

The procedure continues by selecting an appropriate  $H_2 \rightarrow \{\{d_1(i)\}; i \in H_2\} = \{(d_1)_{1:m}, \dots, (d_1)_{h:m}\}$  where  $(d_1)_{1:m} \leq (d_1)_{2:m} \leq \dots \leq (d_1)_{m:m}$  are the ordered distances and then  $\{\widehat{\mu}_2\}$  and  $[\widehat{\Sigma}_2]$  are calculated based on  $H_2$ . In turn, one should have  $\det([\widehat{\Sigma}_2]) \leq \det([\widehat{\Sigma}_1])$ . If  $\det([\widehat{\Sigma}_1]) > 0$ , the C-step leads to  $[\widehat{\Sigma}_2]$  with  $\det([\widehat{\Sigma}_2]) \leq \det([\widehat{\Sigma}_1])$ . This C-step repeated condition is followed in the algorithm until a stopping criterion is fulfilled. The stopping criterion is when  $\det([\widehat{\Sigma}_{\text{new}}] = 0)$  or when  $\det([\widehat{\Sigma}_{\text{new}}]) = \det([\widehat{\Sigma}_{\text{old}}])$ . Each step increases the density of the data or “concentrates” it.

In order to avoid any confusion regarding the stopping criterion if the new determinant is equal to zero which means that the covariance matrix is singular then the new candidate subset is rejected and the exact previous one is restored.

The calculation chain of determinants is converged in finite steps as one has a finite number of  $h$ -subsets. But the final calculation of  $\det([\widehat{\Sigma}_{\text{new}}])$  may not converge to the global minimum of the MCD objective function. This is the reason why an approximation of the MCD solution is obtained by introducing a large number of random initial conditions for  $H_1$  and via the C-step the lowest determinant is kept. In practical terms, a resampling technique is followed.

### 3. Threshold calculation

Setting an appropriate threshold in the absence of any damage-state data, as is the case in this study, is a non-trivial task. In many studies presented in the published literature, the assumption made is that the multivariate data are normally distributed, with the MSD subsequently approximated by a chi-squared distribution in  $n$ -dimensional space. For the purposes of this study, another method for setting the threshold was followed; a Monte Carlo simulation based on extreme value statistics was used regarding the MCD index. The procedure that was conducted in order to calculate the threshold is as follows:

- An  $n \times m$  (number of dimensions  $\times$  number of observations) matrix is constructed with each individual element a randomly generated number from a normal distribution with zero mean and unit standard deviation.
- The discordancy value (MCD index) as described in previous sections is evaluated for all matrix values, where the robust mean and robust covariance are computed. The largest (i.e. extreme) value recorded for each trial matrix is stored.
- The process is repeated for a large number of trials (10 000) in order to generate an array of “extreme” distance calculations. Next, all the values are ordered in terms of magnitude. The critical values (alpha value,  $\alpha$ ) can take different values such as 5 per cent or 1 per cent for a test of discordancy. In this paper,  $\alpha$  is set equal to 1 per cent giving a 99 per cent confidence limit.

The warning levels for the LTS regression error were defined here in a similar way as the MCD threshold using extreme value statistics.

### 4. Robust regression: least trimmed square (LTS) estimator

Many regression estimators break down in the presence of outliers. Generally, there are different kinds of outliers. As described in the introduction, a point could be either a regular measurement, an outlier or either a good or bad leverage point.

A leverage point is a pair  $(\{x_i\}, \{y_i\})$  where the input  $(\{x_i\})$  is far from the majority of input.

A point  $(\{x_i\}, \{y_i\})$  is considered as a good leverage point if it follows the pattern of the majority. Specifically, in the usual multiple linear regression model given by [10]

$$\{y\} = [\mathbf{X}]\{\theta\} + \{\varepsilon\} \quad (3)$$

where  $\{\theta\}$  are the regression coefficients and  $\{\varepsilon\}$  is a regression error.

$\{x_i\}$  is a bad leverage point if the pair  $(\{x_i\}, \{y_i\})$  does not solve the regression relationship of the majority.

The hat matrix  $[\mathbf{H}] = [\mathbf{X}][(\mathbf{X})^t[\mathbf{X}]]^{-1}[\mathbf{X}]^t$  and its diagonal elements are often used as diagnostics in order to identify leverage points. Unfortunately, the hat matrix as well as the classic Mahalanobis distance suffers from the “masking effect” that was described previously [10]. This follows the fact that there exists a relationship between the hat matrix diagonal elements and the classic MSD via

$$h_{ii} = \frac{(\text{MSD}_i)}{n-1} + \frac{1}{n} \quad (4)$$

where  $n$  is the number of observations. As a result, the  $h_{ii}$  does not necessarily detect the leverage points.

Via this robust data mining methodology (with a simple re-arrangement) one could uncover hidden patterns of the collected databases and in this current study this leads to a means to identify the difference between environmental and operational abnormalities and damage outliers.

The fast least trimmed squares (LTS) estimator as proposed by Rousseeuw et al. [9,10,16,22,24,27,28] will be briefly described. Generally, in order to fit a linear regression model one assumes that

$$y_i = x_{i1}\theta_1 + \dots + x_{in}\theta_n + \theta_0 \quad \text{for } i = 1 \dots m \tag{5}$$

where  $\theta_i$  are the regression coefficients and  $(\{x_i\}, y_i)$  are the data point coordinates. Assume for now that there is a single output. The basic objective of this algorithm is to minimise the function:

$$\sum_{i=1}^h (r^2)_{i:n} \quad \text{for } i = 1 \dots m \tag{6}$$

where  $r^2$  is the square of the residual which is the difference between the observed values and the predicted value  $(y - \hat{y})$ . The objective in this work is to find in a similar fashion with the MCD estimator  $h$ -subsets with the smallest least squares function (5). The LTS regression line is the least square model of these  $h$ -points.

Rousseeuw et al. [9,10,16,22,24,27,28] state the major advantages of the algorithm against the classic least median squares (LMS) such as a smooth objective function, less sensitivity in the presence of local effects as well as statistical efficiency as the LTS estimator is asymptotically normal. As a major step further, the fast-LTS can deal with large data sets and is computationally inexpensive and much faster than least median squares (LMS).

The basic concept behind the C-step was described in the MCD estimator section, to connect it with the LTS estimator the basic properties are given and the proofs can be found in [9,10,16,22,24,27,28].

The novel element of the LTS algorithm is that any approximation considered for calculating the LTS coefficients leads to another approximation with lower (optimum) objective function (see [9,10,16,22,24,27,28] for the analytic proof).

Briefly, the core of the LTS algorithm is explained as follows. Assume that there is a data set consisting of  $\{x_i\}$  input vectors and  $y_i$  responses. Create an initial subset  $H_1$  and calculate according to Eq. (6):

$$Q_1 = \sum_{i \in H_1} (r_1(i))^2$$

$$r_1(i) = y_i - \left( \hat{\theta}_1^1 x_{i1} + \hat{\theta}_2^1 x_{i2} + \dots + \hat{\theta}_n^1 x_{in} \right) \tag{7}$$

For  $i = 1, \dots, m$  where  $n$  is the dimension and  $m$  is the number of observations. Then a new subset  $H_2$  is created such as  $\{|r(i)|; i \in H_2\} = \{|r_1|_{1:n}, \dots, |r_1|_{h:n}\}$  where the ordered values of the residuals are

$$|r_1|_{1:n} \leq |r_1|_{2:n} \leq \dots \leq |r_1|_{h:n} \tag{8}$$

From this the coefficients of the least squares estimates are calculated for the  $h$  observations of  $H_2$ . As a chain this leads to a new calculation:

$$Q_2 = \sum_{i \in H_2} (r_2(i))^2$$

such that

$$Q_2 \leq Q_1 \tag{9}$$

Regarding the algorithm, the first step is to create an initial number of subsets  $h$ . The number of subsets is given by default as  $h = (n + m + 1)/2$  where  $m$  is observations in  $n$  dimensional space ( $n \times m$ ). But any number of subsets between  $(n + m + 1)/2 \leq h \leq n$  can be introduced.

Next if the number of variables is one (univariate) the LTS estimator is computed by giving the  $\{\hat{\theta}\}$  best parameters of Eq. (5) by executing the exact algorithm of [9,10,16,22,24,27,28].

If there is multivariate data set then the steps are similar to that of the previous concept but slightly re-arranged. Draw  $h$ -subsets  $H_1$ , carry out the calculations and identify  $\hat{\theta}_{\text{sub}}$  via C-steps and keep the optimal solution with best parameters  $\hat{\theta}_{\text{full}}$  that gives the lowest  $Q$ .

The results are presented by combining the FAST-LTS algorithm for revealing the leverage points and an outlier map introduced by the MCD index. These descriptive plots between the residuals of the LTS regression and robust distances provide a robust estimate of multivariate location and nature of predicted outliers.

### 5. Description of the results map

The results that follow (for both the Tamar and Z24 Bridges) are presented in the form of a residual outlier map and it can be used in order to classify the observations according to the robust regression model. This visualisation plot displays the standardised LTS residuals (the residuals divided by a robust estimate of their scale) versus the robust distances calculated by applying the MCD estimator on the input variable. In more detail and to make more clear the usage of the residual outlier map that will be generated, it displays the LTS residuals versus the robust score distances, and can be used in order to classify and characterise the observations as between outliers and leverage points. The two red horizontal lines in Fig. 1 correspond to the LTS threshold and the vertical red line to the MCD threshold.

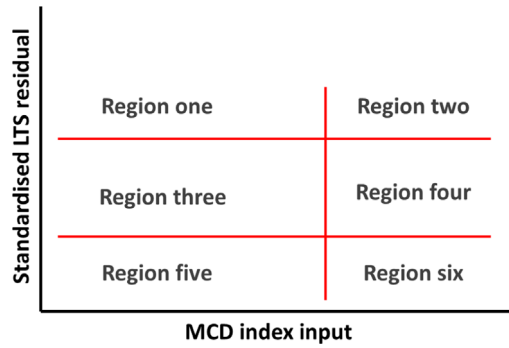


Fig. 1. Residual outlier map. (For interpretation of the references to colour in this figure caption, the reader is referred to the web version of this paper.)

Table 1  
Residual outlier map description.

Region	Classification	Description
One	Vertical outlier	Outside horizontal thresholds but within vertical threshold
Two	Bad leverage points	Outside horizontal thresholds and outside vertical threshold
Three	Normal points	Within horizontal thresholds and vertical threshold
Four	Horizontal outlier–good leverage points	Within horizontal thresholds but outside vertical threshold
Five	Vertical outlier	Outside horizontal thresholds but within vertical threshold
Six	Bad leverage points	Outside horizontal thresholds and outside vertical threshold

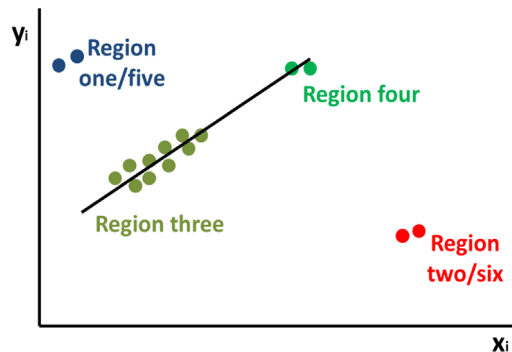


Fig. 2. Regression example in accordance to Table 1 including normal data (light green), vertical outliers (blue), good leverage points (dark green) and bad leverage points (red). (For interpretation of the references to colour in this figure caption, the reader is referred to the web version of this paper.)

As can be seen in Fig. 1, the results will be presented in a map-plot similar to this one. The plots are divided into six regions which are summarised in Table 1 [9,10,22–25].

If  $(\{x_i\}, \{y_i\})$  is a leverage point then this indicates the “outlyingness” of  $\{x_i\}$ , but does not take into account the regression response  $\{y_i\}$  [10].

A bad leverage point lies far from the bulk/mass corresponding to the majority of the data. Such a point can be proven disastrous as it attracts or even shifts the classic least squares regression (consequently the word “leverage” is used).

On the other hand, if  $(\{x_i\}, \{y_i\})$  follows the linear relation it can be called a good leverage point, because as a result it will improve the performance of the regression model.

To make it even clearer, Fig. 2 that summarises the terminology that is described in Table 1 is added [10]. As can be seen the majority of the observations are normal data which is region three. Blue and red points deviate from the linear model and therefore are called regression outliers, but region four with bright green colour is not. Region four and region two/six are leverage points, because the  $\{x_i\}$  point is outlying. As a result region four (bright green) has good leverage points or horizontal outliers and region two/six (red colour) has bad leverage points. The data points in the blue region one/five are called vertical outliers, because they are regression outliers but not leverage points.

The importance of this map for SHM is significant, as for the first time there is a distinction between the way that environmental/operational conditions and damage can manifest themselves physically. It has to be pointed out that one cannot simply remove the outliers as they could be critical points during regression. In some cases the plots may indicate a change to the model. One could be able to go back to the original data and extract information about leverage points/outliers



Fig. 3. The Tamar Bridge.

and discover where they come from. Practically, if during regression the algorithm detects vertical outliers one can go back and change the model and re-check if by fixing the model, these points are then becoming regular data or good leverage points. This will indicate that if the training set was clear from damage these vertical outliers are environmental or operational variations and should be included in the analysis.

It is obvious that, at the end of the day, someone can explain if the novelty detection that occurs is coming from a mutation that changes the physical system or from outliers away from the normal condition cloud.

## 6. Case study: the Tamar Bridge

The Tamar Bridge (Fig. 3) is situated in the south-west of the United Kingdom and connects Saltash in the county of Cornwall with the city of Plymouth in Devon. The bridge is a major road across the River Tamar and plays a significant economic role. The construction belongs to the family of suspension bridges and used to be the longest in England.

Initially, in 1961 the bridge had a main span of 335 m and side spans of 114 m. If the anchorage and approach are included, the overall length of the structure is 643 m. The bridge stands on two concrete towers with a height of 73 m with the bridge deck suspended at mid-height.

In the late 1990s, an upgrade was performed regarding the structure after an EU directive. The main aspects of these additions were to strengthen and widen the bridge in order to allow heavier haulage vehicles. This upgrade was significant and raised questions regarding the performance of the bridge. In this direction various sensors systems were installed to extract data for parameters such as tensions on stays added during the upgrade, wind velocity, temperature, deflection and tilt sensors.

For the purposes of this paper, the vibration-based data is used. The sensor systems were installed by members of the Vibration Engineering Section of the Department of Civil and Structural Engineering of University of Sheffield in 2006.<sup>1</sup> Eight accelerometers are implemented in orthogonal pairs to four stay cables and three sensors measure deck accelerations. The time-data were stored with a sampling frequency of 64 Hz at 10 min intervals. The data were then passed directly to a computer-based system and via automated modal analysis, the natural frequencies were calculated. For more details the reader is referred to [1,2].

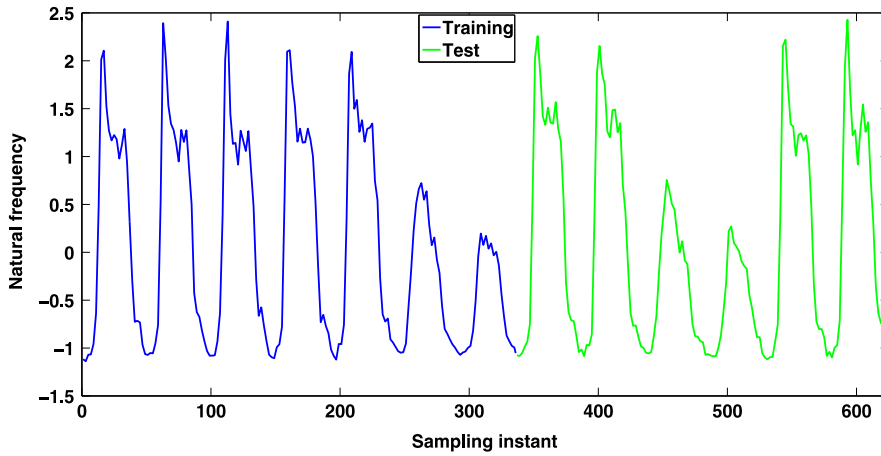
### 6.1. Results

As stated in the introduction, confounding influences can cause discontinuous changes in the features used for regression analysis due to external operational influences. This changes then can be assumed wrongly as false alarms and be characterised as outliers. Regarding the Tamar Bridge it was found [1,2] that a small number of data regions caused regression models to significantly fail where wind speed or deck acceleration was not used as a model input.

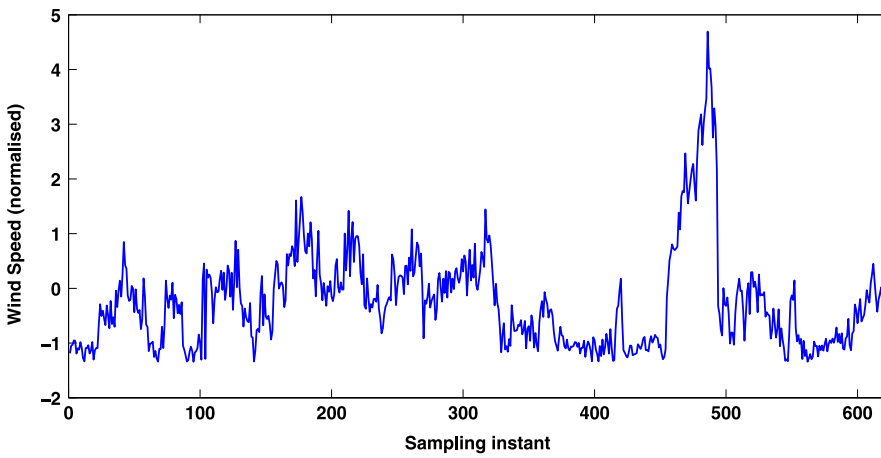
This critical failure was connected with times when there were high wind speeds, with the wind crossing the bridge deck in the transverse direction to the span. Specifically, when the wind is from the east or west, there is no significant influence on the first natural frequency, conversely when the wind is from the north or south, i.e. normal to the bridge span, above 25 mph the first natural frequency is affected.

Using the robust regression method it will be shown that it classifies this “abnormal” behaviour as vertical outliers for the regression model when using only the traffic loading and the first natural frequency. This could be proved as a significant result as someone could see that there is a component missing from the regression input and not classify this strange behaviour as a false alarm. In summary it was found in [1,2] that traffic loading is introducing a dominant operational factor that affects the modal frequencies of the deck and this is the reason that traffic loading is used as an input.

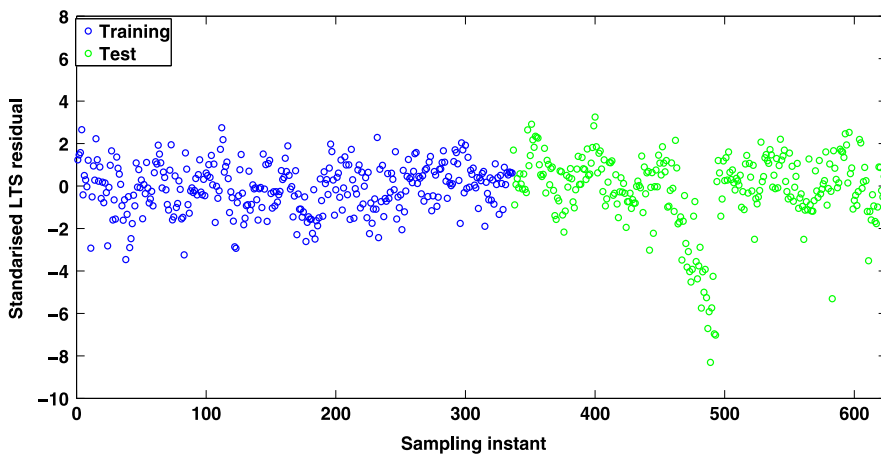
<sup>1</sup> The researchers are now at the University of Exeter, UK.



**Fig. 4.** First natural frequency against time.



**Fig. 5.** Wind speed variability (period of interest is between observations 465–495).

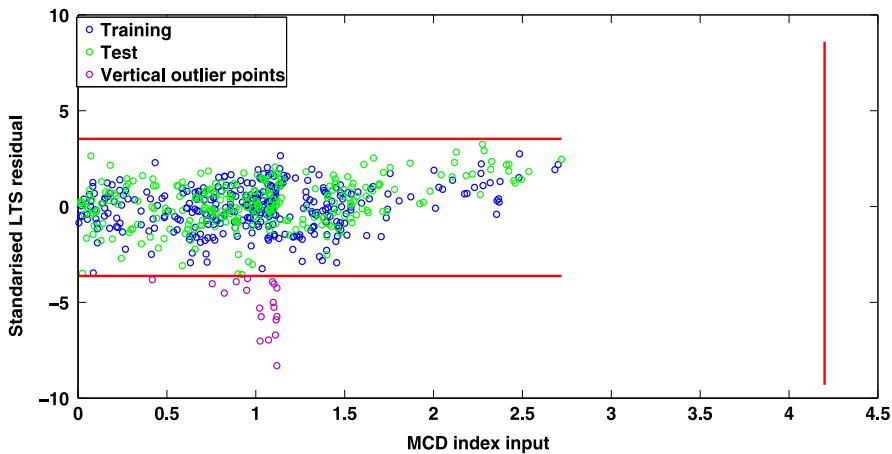


**Fig. 6.** Plot of LTS residual against time.

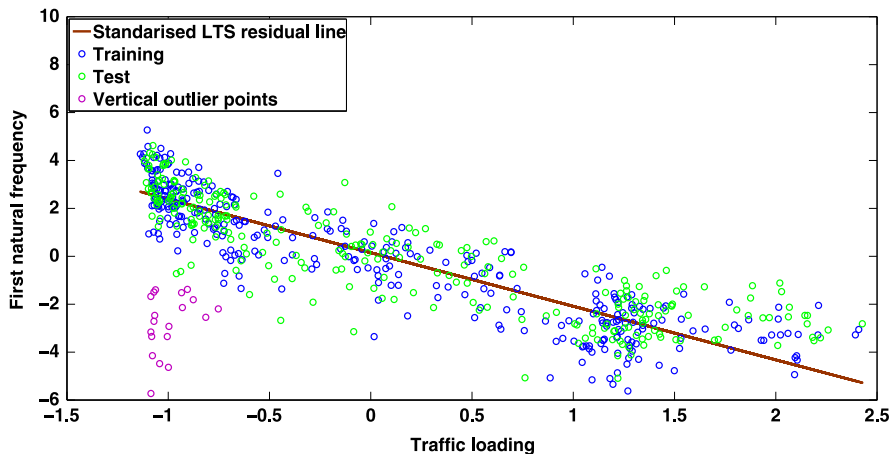
As can be seen from Fig. 7 one can distinguish regular data which are characterised by small residual and MCD distances from the data that are characterised as vertical outliers. These observations belong to region five.

Indeed as can be seen in Figs. 7 and 8 the LTS visualisation captures and reveals the ill points without indicating any bad leverage points. This visualisation corresponds to the wind behaviour that was described earlier. The first 336 points were





**Fig. 7.** Plot of LTS residual versus MCD robust distance for regression between traffic loading and first natural frequency. The vertical outlier points corresponds to high wind values as seen in Fig. 5. (For interpretation of the references to colour in this figure caption, the reader is referred to the web version of this paper.)



**Fig. 8.** Plot of traffic loading versus first natural frequency and the actual robust regression output line. The pink points correspond to the vertical outlier points as found in Fig. 7. (For interpretation of the references to colour in this figure caption, the reader is referred to the web version of this paper.)



**Fig. 9.** View of Z24 Bridge.

used in order to establish a training set and the next 288 points were projected against the extracted parameters of the training set (see Fig. 4). This is more clear by looking at Figs. 5 and 6. This high wind variation between values 465 and 495 is reflected in the LTS detection method as vertical outliers during the regression run (pink points in Figs. 7 and 8). It has to be clear that the experimental data is free from any damage introduction in the Tamar Bridge.

This is vital information and a novel approach to the SHM community. To clarify these conclusions at a second level another more advanced example from civil engineering is presented next; a complicated behaviour of the Z24 Bridge is

described. It states that environmental or operational components (again as vertical outliers) appear with a completely different nature compared to damage data during regression analysis.

### 7. The Z24 Bridge

The Z24 Bridge was a concrete highway structure in Switzerland connecting Koppigen and Utzenstorf, and in the late 1990s; before its demolition procedure, it was used for SHM purposes under the “SIMCES” project [1,29]. The Z24 is a pre-stressed bridge consisting of three spans and two lanes with an overall length of 60 m (Fig. 9). During a whole year of monitoring of the bridge, a series of sensor systems captured modal parameter measurements, as well as a family of

**Table 2**  
Progressive damage scenarios.

Sequence	Damage scenarios
1	Settlement of foundation
2	Tilt of foundation
3	Spalling of concrete at soffit
4	Landslide
5	Failure of concrete hinges
6	Failure of anchor heads
7	Number of post-tensioning tendon failures

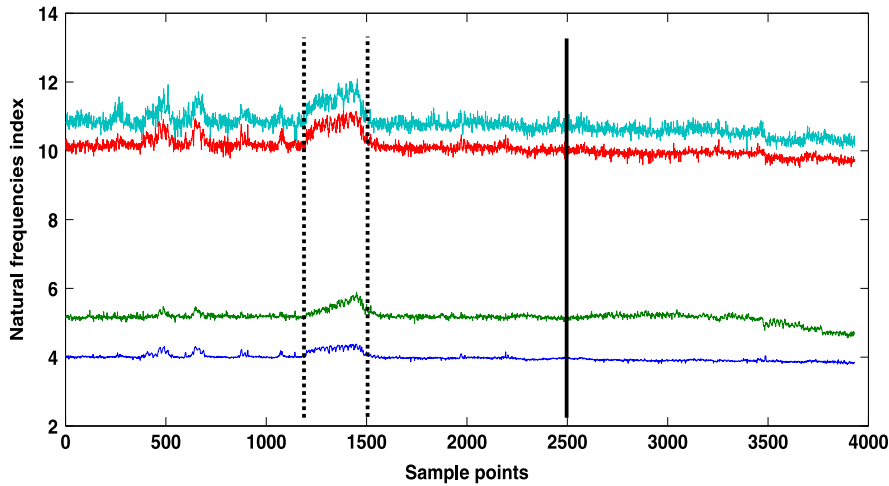


Fig. 10. Time history of frequencies.

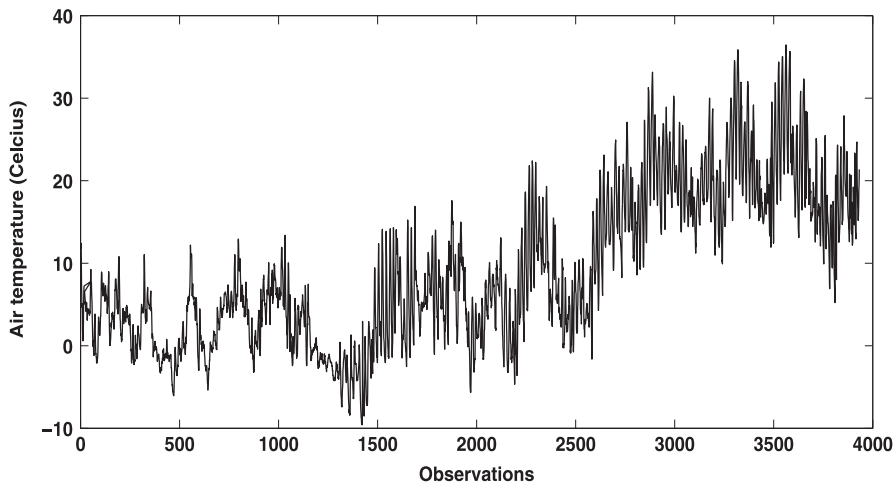


Fig. 11. Time history of air temperature.

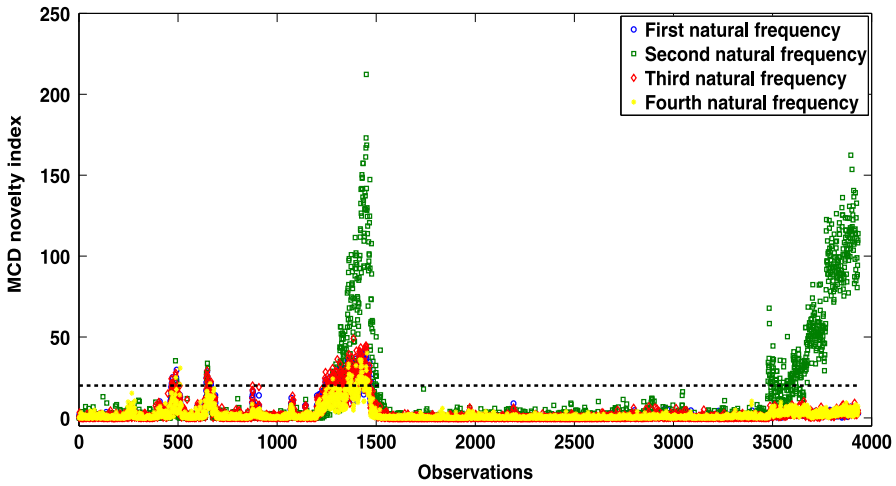


Fig. 12. MCD univariate robust distances of four natural frequencies.

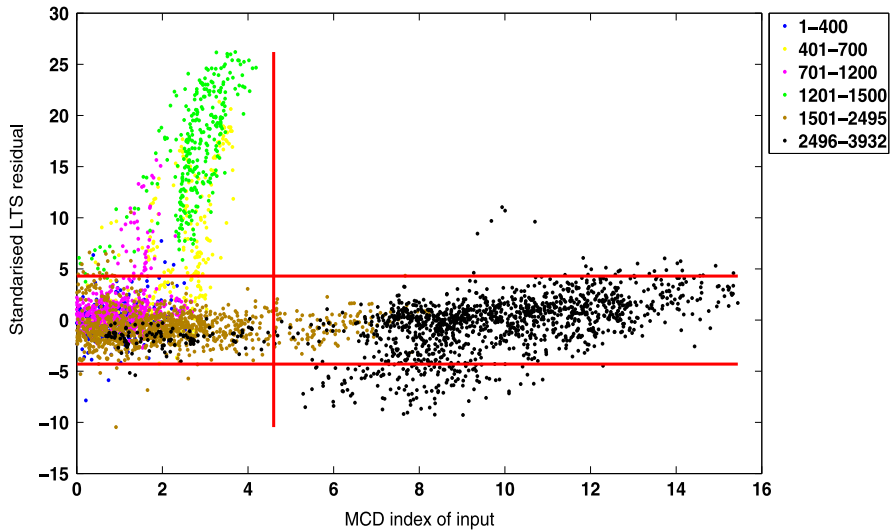
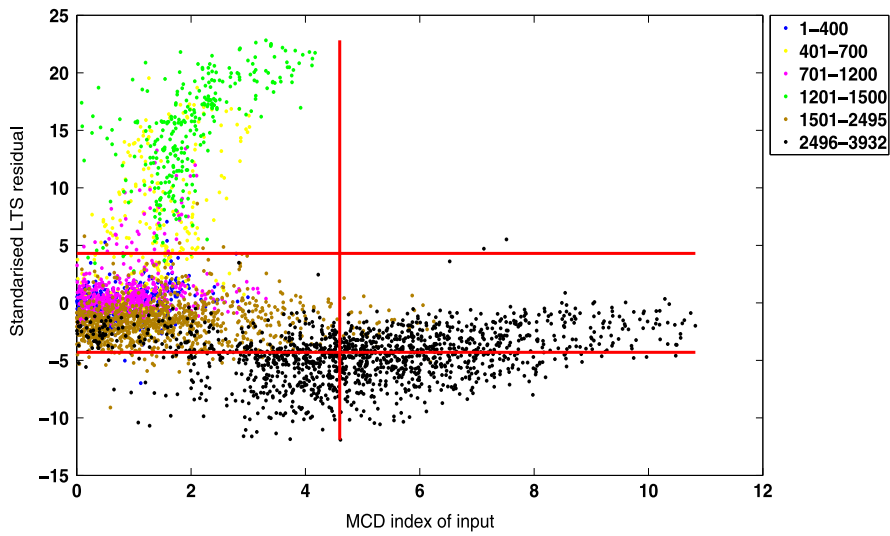


Fig. 13. Plot of LTS residual versus MCD robust distance for regression between air temperature and first natural frequency (top) and temperature in deck soffit and first natural frequency (bottom).

environmental measurements such as air temperature, soil temperature, humidity, and wind speed. The critical point in this benchmark project was the introduction of different types of real progressive damage scenarios towards the end of the monitoring year (Table 2).

For the purposes of this study, the four natural frequencies that were extracted over a period of year, including the period of structural failure of the bridge, are used. Fig. 10 shows the four natural frequencies with values between 0 and 12 Hz. The beginning of the introduced failure occurs at observation 2496. The modal parameters regarding the Z24 Bridge were estimated via Stochastic Subspace Identification by taking into account the ambient vibration.

In the first instance, it can be noted from Fig. 10 that there are some visible fluctuations between observations 1250 and 1460 but there are no dramatic visible fluctuations after the introduction of damage, making the frequency sequences nonstationary. The fluctuations are highly related to periods of very cold temperatures under 0 °C and there is a direct connection with increased stiffness based on the freezing of the asphalt layer of the bridge deck, see Fig. 11. These large temperature fluctuations are suitable candidates in order to check the sensitivity of the LTS method in characterising the environmental fluctuation differently compared to damage data sets.

Once the outliers are revealed, a decision for a suitable feature has to be made and a more complicated algorithm has to be applied in order to remove external influences like environmental factors and create a suitable normal condition training set. This extra step is investigated in the next session by combining the robust distances with LTS algorithms.

The univariate MCD robust distance was calculated for each of the four natural frequencies in order to reveal each frequency's internal fluctuations, as shown in Fig. 12. This analysis is critical, as important conclusions can be derived. The obvious observation could be that the natural frequencies fluctuate as is revealed by looking at the dramatic index increase

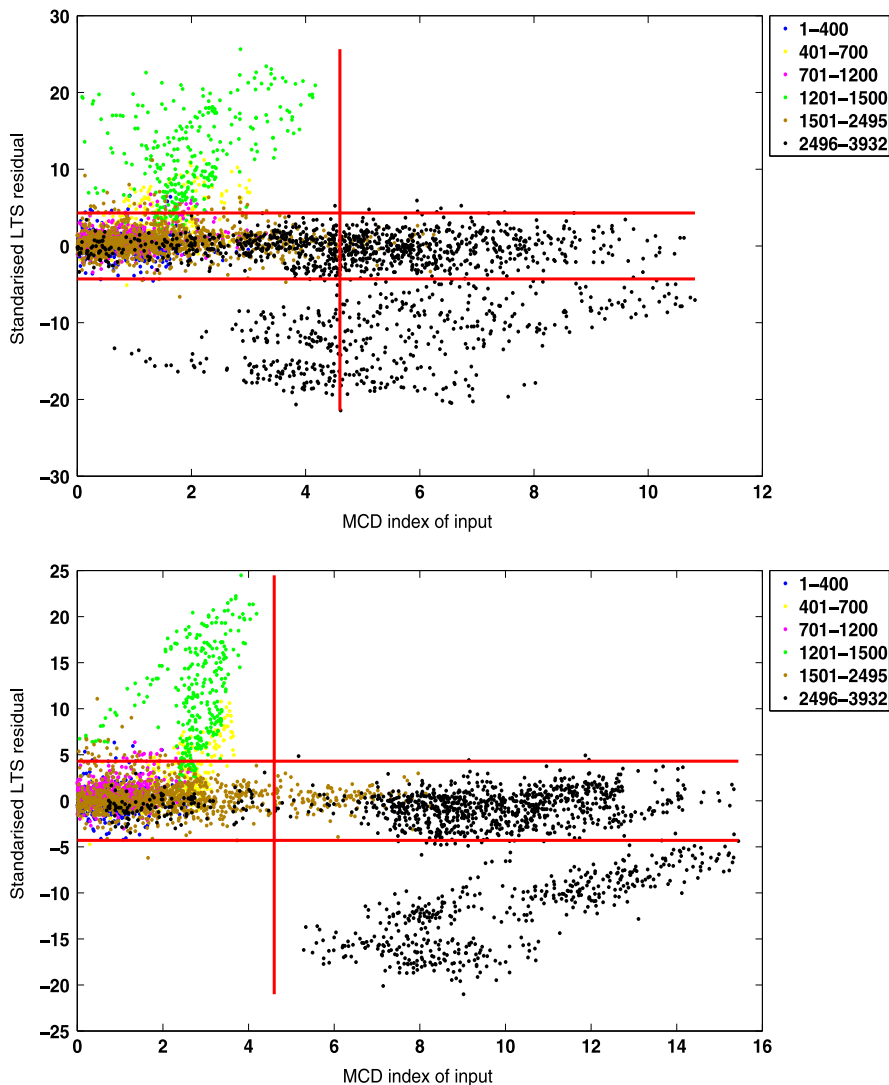


Fig. 14. Plot of LTS residual versus MCD robust distance for regression between air temperature and second natural frequency (top) and temperature in deck soffit and second natural frequency (bottom).

between the observations 1250 and 1460 where the cold temperatures are operating. Also, the index value of the damaged condition increases after observation 3500 (mainly for the second natural frequency). The frequencies present a common pattern where again the MCD index increases during the cold period, but after observation 3500 the damage is not clearly noticeable (for the first, third and fourth frequencies). Important conclusions can be found by looking at the presented patterns of Fig. 12 more carefully. Between observations 430–520 and 590–690 two distinctive peaks can be seen that correspond also to temperatures below zero. An analogous index increase appears again when cold temperatures occur.

7.1. Results

The plots of LTS residual versus MCD robust distance are presented in Figs. 13–16. The regressed variables are between the natural frequencies and two different temperature data sets. Specifically, in the plots can be seen a comparison when ambient (top) and bridge deck (bottom) temperatures are used. The first 400 measurements are used as a training set and the rest of the data was used as a testing set.

As can be seen from Figs. 13–16, one can categorise the points in all six different regions that were mentioned before. Most importantly the difference between vertical outliers/region one (cold temperature influences), horizontal outliers (damage) and bad leverage points (damage) can be seen which is most harmful for calibration methods as they disturb any regression relationship.

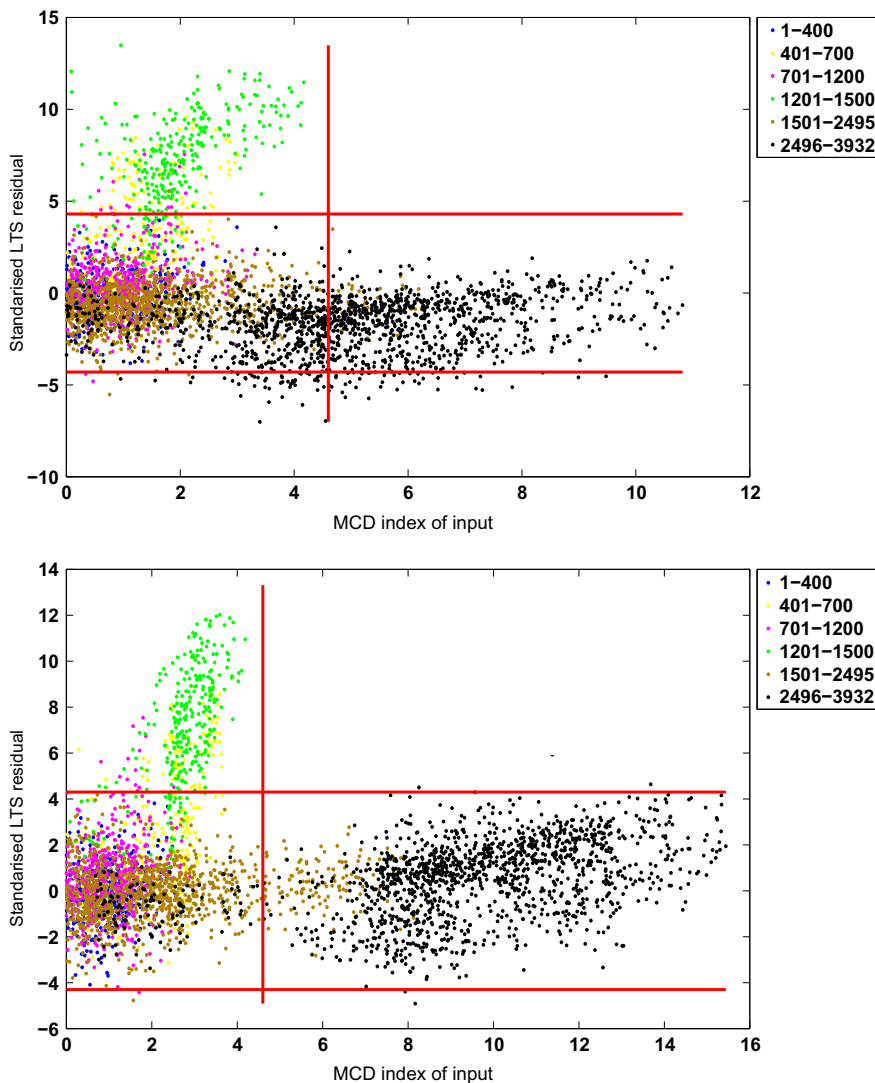
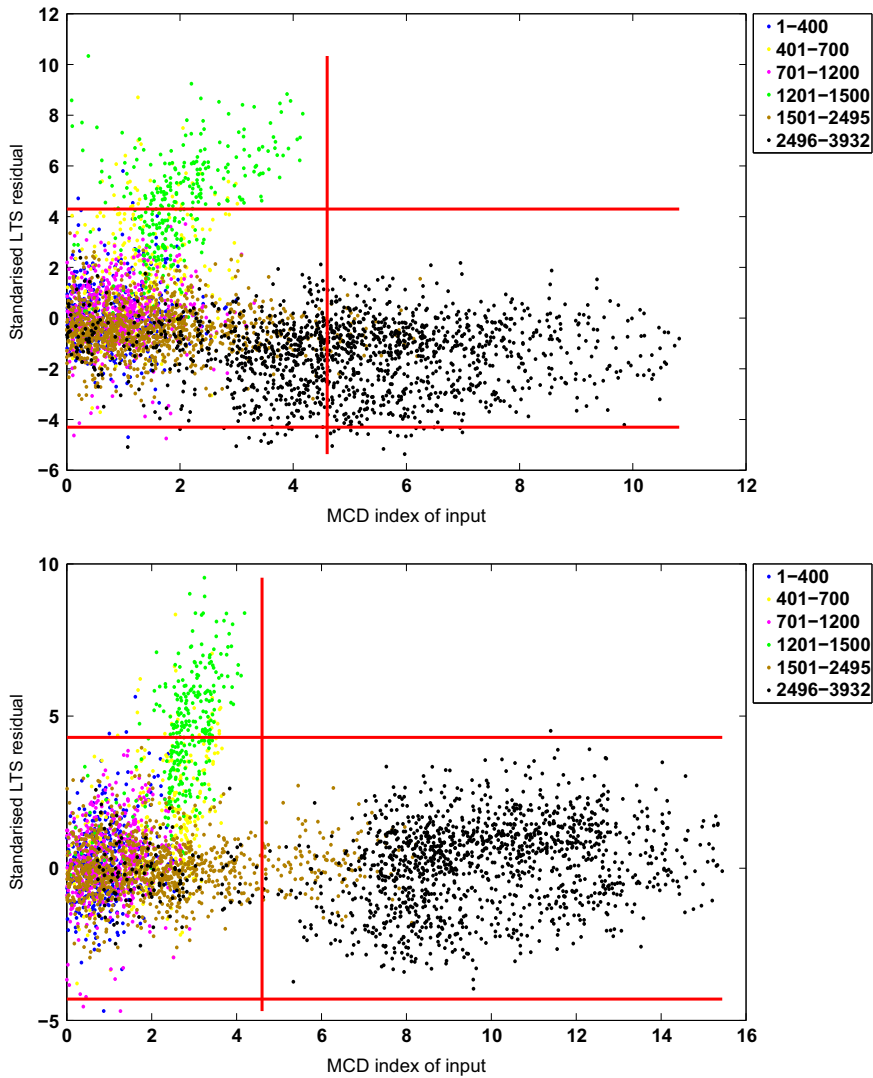


Fig. 15. Plot of LTS residual versus MCD robust distance for regression between air temperature and third natural frequency (top) and temperature in deck soffit and third natural frequency (bottom).



**Fig. 16.** Plot of LTS residual versus MCD robust distance for regression between air temperature and fourth natural frequency (top) and temperature in deck soffit and fourth natural frequency (bottom).

**Table 3**

Description of data sets as they appear in Figs. 13–16.

Observation	Condition
1–400	Undamaged
401–700	Undamaged (with some cold temperature variations)
701–1200	Undamaged (with some cold temperature variations)
1201–1500	Very cold temperature
1501–2496	Undamaged
2497–3932	Damaged

The results clearly indicate what was also found in the case of the Tamar Bridge that the environmental fluctuations are distinctly different in nature than the damaged condition (region one or five). The abnormal behaviour coming from the influence of cold temperature variation appears as vertical outliers. Table 3 that follows explains the different region points as they appear in the graphs.

To make it easier for the reader, the large fluctuation appeared in the graph of Fig. 12 between observations 400 and 700 with the two small distinctive peaks and between observations 1200 and 1500 when the temperature reaches the coldest values appear in the LTS results as vertical outliers.

There is some difference in performance of the method regarding if the ambient or deck temperature is used as a variable but this is insignificant regarding the characterisation of cold temperature influence as vertical points.

When the deck temperature is utilised, the extra resolution added to the results classifies the damaged condition after point 2496 as bad outliers/leverage points and clears any misclassification created by using the (less proper) ambient temperature, mainly regarding the second natural frequency where nonlinear effects appear.

### 8. The importance of robust measures

As stated before, the robust measures are much more efficient and effective compared to the classic least-squares (LS) regression and Mahalanobis squared-distance (MSD) (see MCD section). Specifically, when outliers from a cluster cloud that lie inside the data are present then they will directly move the arithmetic mean towards them and even expand the classical tolerance ellipsoid in their direction [9]. In the classic LS regression model like the classical MSD distance, it suffers from the same masking effect and the regression line is heavily influenced by the physical presence of bad leverage points or vertical outliers [9,10,22–25].

In order to clearly point-out their importance, the training data was chosen as to include the environmental fluctuations of Z24 Bridge in order to include in the training set some bad points. As can be seen in Fig. 17 the regression line is heavily influenced when non-robust measures are used. Furthermore, in Figs. 18–21 the reader can clearly see that when LTS and MCD indices are used their effectiveness in revealing different outlier characteristics is superior. The environmental fluctuations (green points) appear again as vertical outliers. When the classic measures of LS and MSD are used, they are not able to reveal any vertical outliers and they mask their critical presence.

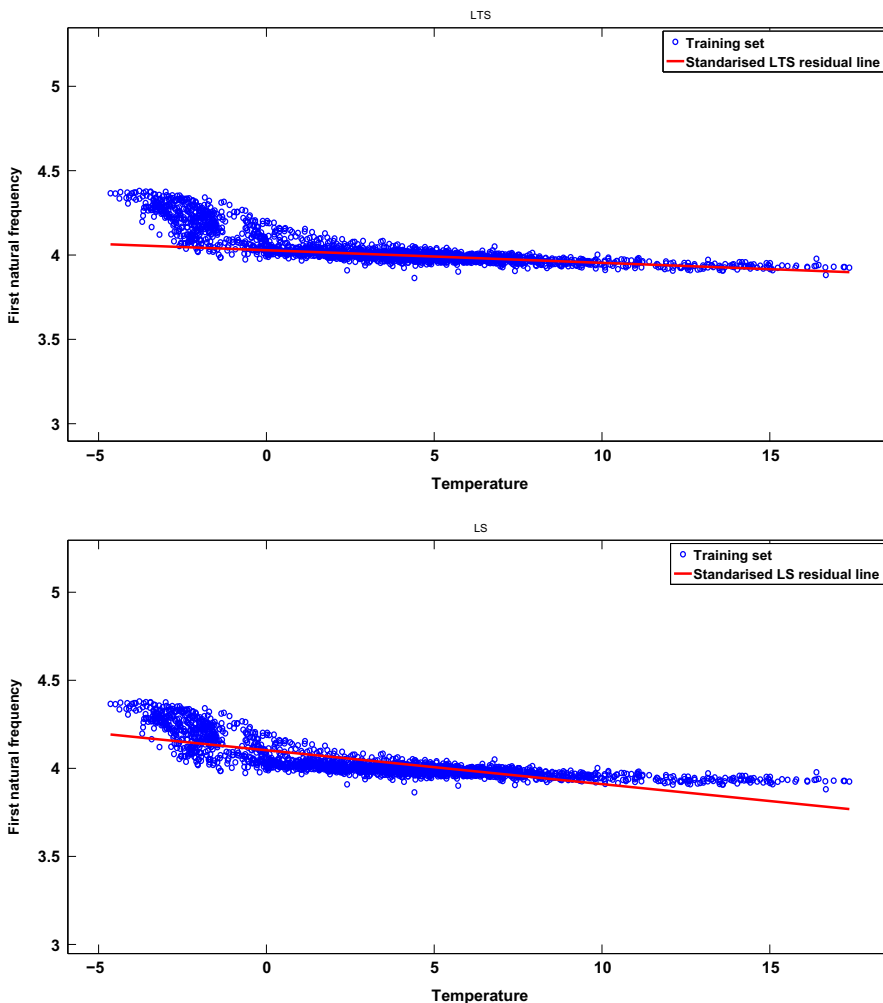
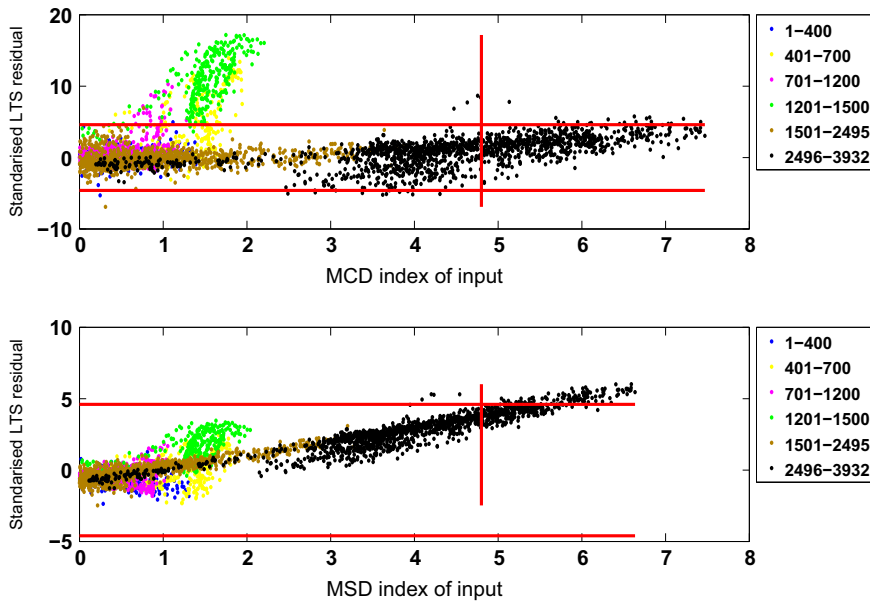
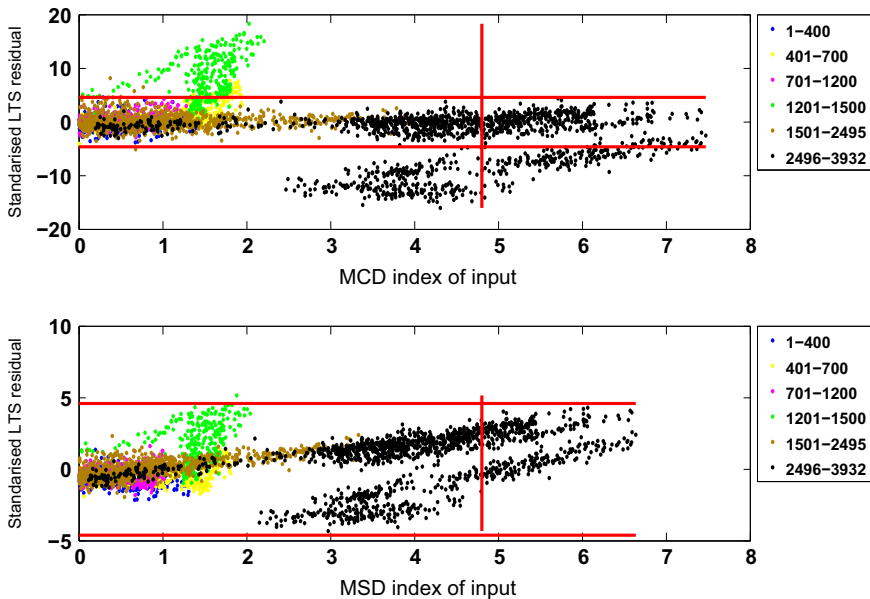


Fig. 17. Plot of temperature versus first natural frequency and the actual regression output line of robust measures (top) and classic measures (bottom).



**Fig. 18.** Plot of LTS residual versus MCD robust distance for regression between deck soffit temperature and first natural frequency (top) and plot of classic LS residual versus MSD distance (bottom). (For interpretation of the references to colour in this figure caption, the reader is referred to the web version of this paper.)

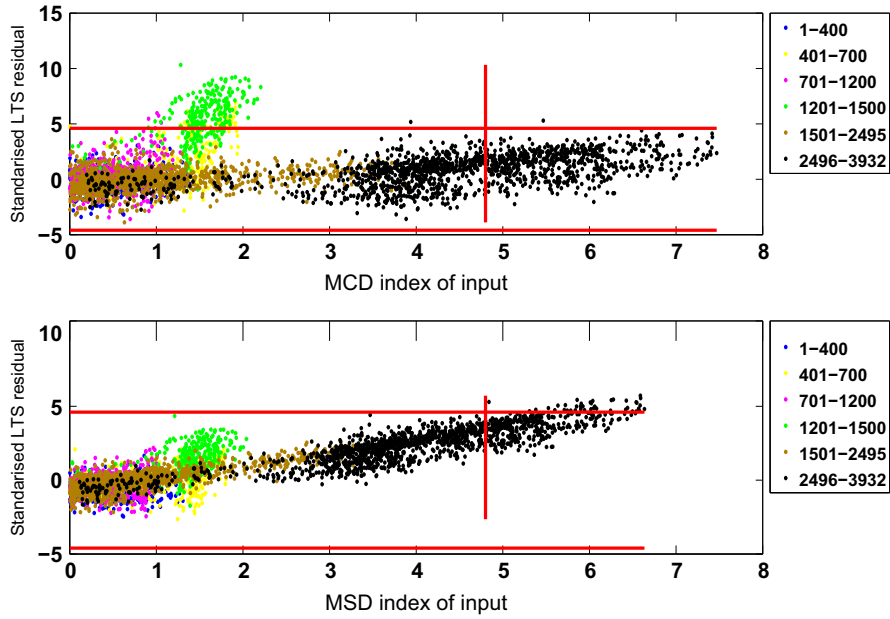


**Fig. 19.** Plot of LTS residual versus MCD robust distance for regression between deck soffit temperature and first natural frequency (top) and plot of classic LS residual versus MSD distance (bottom). (For interpretation of the references to colour in this figure caption, the reader is referred to the web version of this paper.)

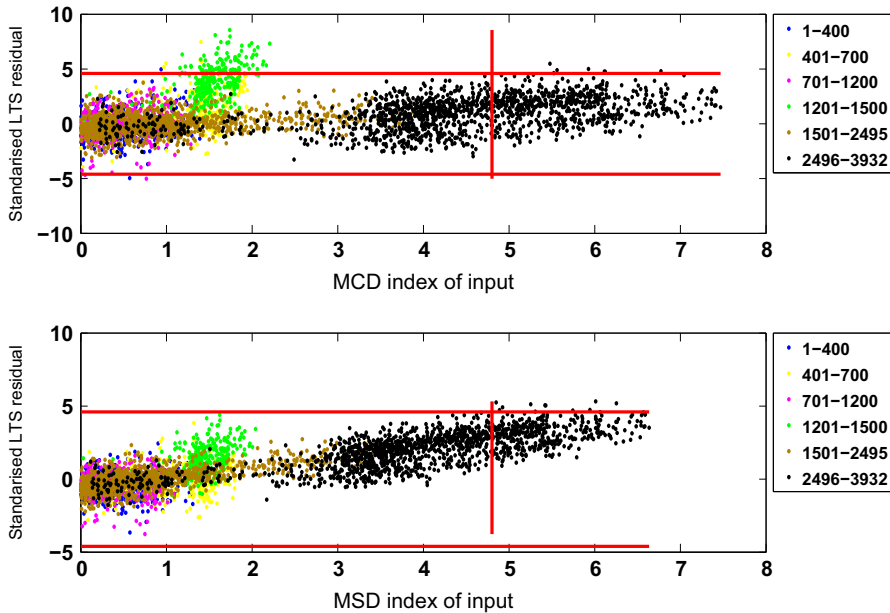
## 9. Conclusion

The approach proposed in this paper is a novel one in the SHM community. It reveals that environmental and operational variations (as feature vectors or points) can manifest themselves completely differently in physical appearance compared to fault observations. Since different outlier detection techniques are based on disjoint sets of assumptions and different technical bases, a direct comparison between them is not always fair and possible. In a lot of cases, the data structure and the outlier generating mechanisms dictate which method will outperform the others and reveal the internal variability.





**Fig. 20.** Plot of LTS residual versus MCD robust distance for regression between deck soffit temperature and first natural frequency (top) and plot of classic LS residual versus MSD distance (bottom). (For interpretation of the references to colour in this figure caption, the reader is referred to the web version of this paper.)



**Fig. 21.** Plot of LTS residual versus MCD robust distance for regression between deck soffit temperature and first natural frequency (top) and plot of classic LS residual versus MSD distance (bottom). (For interpretation of the references to colour in this figure caption, the reader is referred to the web version of this paper.)

This is the reason that this proposed method breaks new ground in identifying, during monitoring operation, the differences in the physical space of the suspicious data when sourced from damaged or undamaged condition or due to the influence of external factors. Outliers are often regarded as “garbage” points that have to be removed, but a lot of times they carry critical information. Furthermore, outliers could lead to the adoption of a misleading model which can result in biased parameter estimation and wrong results. Consequently, it is important to identify and detect them before proceeding to any analysis.

Once the outliers are revealed, a decision for a suitable feature has to be made and a more complicated algorithm has to be applied in order to remove external influences like environmental factors and create a suitable normal condition training set. Such an example is the cointegration method for removal of environmental trends in SHM [8] or auto-associative neural networks [30].

The authors do not claim that the introduction into the SHM field of these methods is a “panacea” but it has worked very well for characterising outliers in complicated-data examples. Further research is carried out in order to move from linear LTS to nonlinear robust regression analysis.

## Acknowledgements

The authors gratefully acknowledge Professor James Brownjohn and Dr Ki-Young Koo of the University of Exeter for providing access to the Tamar Bridge data.

Also, the authors gratefully acknowledge the support of the UK Engineering and Physical Sciences Research Council (EPSRC) through Grant reference number EP/J016942/1.

## References

- [1] E.J. Cross, On Structural Health Monitoring in Changing Environmental and Operational Conditions, PhD Thesis, University of Sheffield, 2012.
- [2] E.J. Cross, K.Y. Koo, J.M.W. Brownjohn, K. Worden, Long-term monitoring and data analysis of the Tamar bridge, *Mechanical Systems and Signal Processing* 35 (1) (2013) 16–34.
- [3] B. Peeters, J. Maeck, G. De Roeck, Vibration-based damage detection in civil engineering: excitation sources and temperature effects, *Smart materials and Structures* 10 (3) (2001) 518.
- [4] S. Alampalli, Effects of testing, analysis, damage, and environment on modal parameters, *Mechanical Systems and Signal Processing* 14 (1) (2000) 63–74.
- [5] P. Cornwell, C.R. Farrar, S.W. Doebling, H. Sohn, Environmental variability of modal properties, *Experimental Techniques* 23 (6) (1999) 45–48.
- [6] J.M. Brownjohn, Structural health monitoring of civil infrastructure, *Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences* 365 (1851) (2007) 589–622.
- [7] R. Moss, S. Matthews, In-service structural monitoring, a state of the art review, *Structural Engineer* 73 (2) (1995).
- [8] E.J. Cross, K. Worden, Q. Chen, Cointegration: a novel approach for the removal of environmental trends in structural health monitoring data, *Proceedings of the Royal Society A: Mathematical, Physical and Engineering Science* 467 (2133) (2011) 2712–2732.
- [9] A.M. Leroy, P.J. Rousseeuw, *Robust Regression and Outlier Detection*, Wiley Series in Probability and Mathematical Statistics, Vol. 1, Wiley, New York, 1987.
- [10] P.J. Rousseeuw, B.C. Van Zomeren, Unmasking multivariate outliers and leverage points, *Journal of the American Statistical Association* 85 (411) (1990) 633–639.
- [11] D.M. Hawkins, *Identification of Outliers*, Vol. 11, Chapman and Hall, London, 1980.
- [12] Barnett, Vic, and Toby Lewis. *Outliers in statistical data*. Vol. 3. New York: Wiley, 1994.
- [13] M. Hubert, M. Debruyne, Minimum covariance determinant, *Wiley Interdisciplinary Reviews: Computational Statistics* 2 (1) (2010) 36–43.
- [14] M. Schyns, G. Haesbroeck, F. Critchley, Relaxmcd: smooth optimisation for the minimum covariance determinant estimator, *Computational Statistics & Data Analysis* 54 (4) (2010) 843–857.
- [15] A.E. Attar, R. Khatoun, M. Lemercier, Diagnosing smartphone's abnormal behavior through robust outlier detection methods, *Global Information Infrastructure Symposium*, 2013, IEEE, 2013, pp. 1–3.
- [16] P. Rousseeuw, M. Hubert, High-breakdown estimators of multivariate location and scatter, *Robustness and Complex Data Structures*, Springer, 2013, pp. 49–66.
- [17] A. Nurunnabi, D. Belton, G. West, Robust segmentation in laser scanning 3d point cloud data, *2012 International Conference on Digital Image Computing Techniques and Applications (DICTA)*, IEEE, 2012, pp. 1–8.
- [18] T. Verdonck, M. Hubert, P. Rousseeuw, Robust covariance estimation for financial applications, *EURANDOM-ISI Workshop on Actuarial and Financial Statistics*, Eindhoven, 29–30 August 2011.
- [19] V. Fritsch, G. Varoquaux, B. Thyreau, J.-B. Poline, B. Thirion, Detecting outlying subjects in high-dimensional neuroimaging datasets with regularized minimum covariance determinant, *Medical Image Computing and Computer-Assisted Intervention–MICCAI 2011*, Springer, 2011, pp. 264–271.
- [20] A.M. Variyath, J. Vattathoor, Robust control charts for monitoring process mean of phase-I multivariate individual observations, *Journal of Quality and Reliability Engineering* (2013), <http://dx.doi.org/10.1155/2013/542305>.
- [21] N. Dervilis, E.J. Cross, R.J. Barthorpe, K. Worden, Robust methods of inclusive outlier analysis for structural health monitoring, *Journal of Sound and Vibration* 333 (20) (2014) 5181–5195.
- [22] P.J. Rousseeuw, K. Van Driessen, A fast algorithm for the minimum covariance determinant estimator, *Technometrics* 41 (3) (1999) 212–223.
- [23] S. Verboven, M. Hubert, *Libra: a matlab library for robust analysis*, *Chemometrics and Intelligent Laboratory Systems* 75 (2) (2005) 127–136.
- [24] M. Hubert, P.J. Rousseeuw, S. Van Aelst, High-breakdown robust multivariate methods, *Statistical Science* (2008) 92–119.
- [25] D.A. Jackson, Y. Chen, Robust principal component analysis and outlier detection with ecological data, *Environmetrics* 15 (2) (2004) 129–139.
- [26] S. Serneels, T. Verdonck, Principal component analysis for data containing outliers and missing elements, *Computational Statistics & Data Analysis* 52 (3) (2008) 1712–1727.
- [27] P.J. Rousseeuw, K. Van Driessen, Computing LTS regression for large data sets, *Data Mining and Knowledge Discovery* 12 (1) (2006) 29–45.
- [28] P.J. Rousseeuw, Least median of squares regression, *Journal of the American statistical association* 79 (388) (1984) 871–880.
- [29] G.D. Roeck, The state-of-the-art of damage detection by vibration monitoring: the SIMCES experience, *Journal of Structural Control* 10 (2) (2003) 127–134.
- [30] N. Dervilis, M. Choi, S. Taylor, R. Barthorpe, G. Park, C. Farrar, K. Worden, On damage diagnosis for a wind turbine blade using pattern recognition, *Journal of Sound and Vibration* 333 (6) (2014) 1833–1850.