# A LEAST SQUARE APPROACH FOR BIDIMENSIONAL SOURCE SEPARATION USING HIGHER ORDER STATISTICS CRITERIA

*Amir A. Khan*[†], *Valeriu Vrabie*[‡], *Jérôme I. Mars*[†], *and Alexandre Girard*[*]

[†]GIPSA-Lab–Department of Image and Signal (DIS), INP of Grenoble, BP 46, 38402 Saint Martin d'Hères Cedex, France.

[‡]Centre de Recherche en STIC (CReSTIC), University of Reims, BP 1039, 51687 Reims Cedex, France.

[*]Direction R&D, Electricité de France (EDF), 6 quai Watier, 78401 Chatou Cedex, France.

Email: {amir-ali.khan, jerome.mars}@gipsa-lab.inpg.fr, valeriu.vrabie@univ-reims.fr, alexandre.girard@edf.fr

## ABSTRACT

The anomaly detection based on processing of distributed temperature sensors data is a new research problem. The acquired data is highly influenced by the response of the ground in which the sensors are buried. It therefore becomes essential to remove the influence of this response. This response, being the most coherent factor in the acquired signal, appears as the most energetic source vector. However, its classical estimation by SVD runs the risk of taking into account energetic phenomena like precipitations. We propose to characterize such phenomena using higher order statistics thus giving a criteria of selecting only the data not influenced by such phenomena. An overlapping window approach then allows estimation of characteristic ground response source. Moreover, the corresponding ground response subspace is constructed by least squares based unmixing approach on the characteristic source. This avoids also the physically unjustifiable orthogonality condition of temporal variations of the estimated sources imposed by SVD.

## 1. INTRODUCTION

The use of fiber optic sensors has been a common practice in diverse domains covering applications as engineering structures monitoring, fault detection in electrical circuits, fire detection systems, parameter sensing in oil and gas industry, etc [1, 2]. Now-a-days, an important issue in engineering domain is the detection of anomalies, such as leakages (significant flow of water), in the dikes to avoid disaster at mass level. One of the most promising methods for this purpose is thermometric based method employing optical fiber Distributed Temperature Sensors (DTS). The major advantage of DTS is their commercial viability (low-cost telecommunications grade fiber), ability to multiplex large number of sensors along a single fiber and environmental robustness [1, 3]. The leakage detection using DTS signals is a new research problem in the signal processing domain.

The basic concept behind temperature acquisition is that a change of ground temperature is brought about by a significant flow of water through the structure due to leakages. However, this change of temperature can equally be brought about by other factors such as precipitations, seasonal effects, day/night, the existing structures (e.g. drains), etc. Moreover, since the fiber optic cable is buried in ground, the temperature signals acquired by DTS are strongly influenced by the response of the near surface (ground) where the acquisitions are made. The leakage detection thus becomes a source separation problem with sources being all of the above mentioned factors. Due to overwhelming influence of the ground response on the acquired signals, it is imperative for a leakage detection scheme to first remove its dependence.

The source separation techniques have been successfully employed in diverse domains like neural networks, biomedical engineering, telecommunications, econometrics, geophysics, image processing, audio signal separation, spatio-temporal data set analysis, etc [4, 5, 6]. More recently, they have been employed to analyze fiber sensor signals for the measurement of food color and water monitoring [7]. In multisensor signal processing (geophysics, underwater acoustic, etc.), a classical source separation technique is based on the Singular Value Decomposition (SVD). It is a useful tool to perform a separation of the initial dataset into complementary orthogonal subspaces by extracting decorrelated vectors [8, 9, 10].

In case of temperature data (a function of space and time), we estimate the sources as a function of distance to identify different phenomena in space. It is observed that the first source vector obtained by application of SVD on temperature data, representing the most significant energy, is usually related to the ground response. However, a major problem is that this estimated source vector can be influenced by energetic factors ephemeral in time (like significant precipitations). In this paper, this problem is addressed by devising a criteria based on higher order statistics (HOS), exploiting the fact that ephemeral temporal phenomena have specific statistical behavior. The data selectivity is done based on this criteria to avoid these ephemeral phenomena. The first source vector is estimated in sliding temporal windows on this selected data and a characteristic vector from amongst these windowed vectors is finally obtained using the mean operator. This characteristic vector serves as a better estimate of the ground response source. Moreover, for constructing the subspace corresponding to this estimated ground response source, the physically unjustifiable orthogonality condition imposed by SVD on the temporal variations of the estimated sources is avoided by using a least squares (LS) based approach.

## 2. SYSTEM AND DATA DESCRIPTION

The temperature data is acquired through DTS using Optical Time Domain Reflectometry (OTDR) technique based on the Raman backscattering principle [1]. A thermometric data monitoring system has been installed, by Electricité de France (EDF), at an experimental test site in the south of France (near Oraison city) with an aim to study leakages (both natural and controlled) in the dike of canal (see Fig. 1). A 2.2 km long fiber optic cable (containing 4 optic fibers, of

type multimode 50/125), is buried at the downstream toe of the canal at a depth of 1 m to intercept water leakage from the canal. Two distinct elevation levels, (Level1, from approximately 0.5 km to 1.25 km and Level2, from approximately 1.25 km to 2.2 km), will be exposed with varying intensities to direct sunlight. The cable also circumvents two drains, D1 and D2, located at 0.561 km and 0.858 km, respectively. The temperature data were recorded by the device Sensornet, Sentinel DTS-MR, with temperature and spatial resolutions of $0.01°C$ and 1-meter, respectively. To monitor temporal evolution of the anomalies, acquisitions were made over a period of five and a half weeks with a sampling interval of 2 hours. This gives a two-dimensional temperature data set, as a function of displacement along the fiber and time:

$$\mathbf{Y} = \{y(x,t) \mid 1 \leq x \leq N_x, 1 \leq t \leq N_t\}, \qquad (1)$$

where $N_x$ and $N_t$ are the number of observation points in distance and the total time samples, respectively. The data is normalized so that each acquisition has a zero mean and unity variance (Fig. 2). This removes daily and seasonal variations. Three artificial leakages, $L1$ (on day 28), $L2$ and $L3$ (on day 30), were introduced at the site with different flow rates of 5, 1 and 1 lit/min and at different positions, 1.562, 1.547 and 1.569 km, respectively.
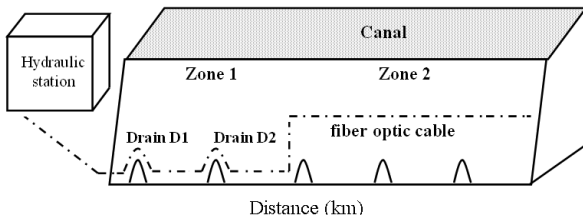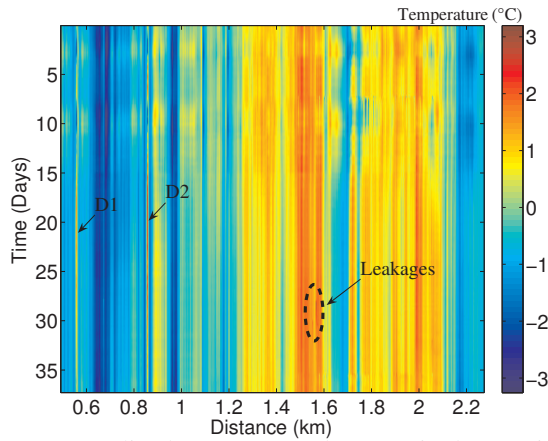


Figure 1: Data acquisition system.



Figure 2: Normalized temperature data acquired at Oraison site.

## 3. SUBSPACE DECOMPOSITION BY SVD

The SVD of the signal in Eq. (1) is defined as [9, 10]:

$$\mathbf{Y} = \mathbf{U}_N \Delta_N \mathbf{V}_N^T = \sum_{j=1}^{N} \beta_j \mathbf{u}_j \mathbf{v}_j^T, \qquad (2)$$

where $N = \min(N_x, N_t)$, $\Delta_N \in \mathbf{R}^{N \times N}$ is a matrix containing on its diagonal the singular values $\beta_j \geq 0$ arranged in a de-

scending order and $\mathbf{U}_N \in \mathbf{R}^{N_x \times N}$ and $\mathbf{V}_N \in \mathbf{R}^{N_t \times N}$ are orthogonal matrices, containing the left and right singular vectors $\mathbf{u}_j \in \mathbf{R}^{N_x}$ and $\mathbf{v}_j \in \mathbf{R}^{N_t}$ respectively. The left singular vectors $\mathbf{u}_j$ are identified as estimators of the sources defined by different factors (ground response, existing structures (drains), leakages, etc.) and are orthogonal to each other. The first vector, $\mathbf{u}_1$, being the most energetic, is linked to the ground response. This orthogonality can be justified considering that these factors are physically independent of each other. On the other hand, the vectors $\mathbf{v}_j$, representing the temporal variations of the sources, are also orthogonal by construction which is not always physically justifiable. SVD can be used to achieve separation between signal and noise subspaces [10]:

$$\mathbf{Y} = \mathbf{Y}_{sig}^{SVD} + \mathbf{Y}_{residue}^{SVD} = \sum_{j=1}^{P} \beta_j \mathbf{u}_j \mathbf{v}_j^T + \sum_{j=P+1}^{N} \beta_j \mathbf{u}_j \mathbf{v}_j^T \quad (3)$$

where $\mathbf{Y}_{sig}^{SVD}$ is given in our case by the ground response, which means that $P = 1$.

## 4. ESTIMATION OF AN AVERAGE SOURCE

The estimation of the first source, $\mathbf{u}_1$, by applying subspace decomposition over the entire data runs the risk of being influenced by ephemeral energetic factors like precipitations. To overcome this problem, SVD is thus calculated in small time blocks using overlapping temporal sliding window approach by choosing a suitable sliding window size ($\Delta T$) and an overlapping interval. It should be recalled here that since the estimated sources are a function of distance, we look to eliminate effects of temporal ephemeral phenomena and thus consider only temporal windowing and not spatial windowing. Considering the $m^{\text{th}}$ data block, $\mathbf{Y}_m = \{y(x,t) | 1 \leq x \leq N_x, t_m \leq t \leq t_m + \Delta T\}$, this decomposition can be written as:

$$\mathbf{Y}_m = \mathbf{U}_N^m \Delta_N^m \mathbf{V}_N^{mT} = \sum_{j=1}^{N} \beta_j^m \mathbf{u}_j^m \mathbf{v}_j^{mT} \qquad (4)$$

where $N = \min(N_x, \Delta T)$ and $t_m$ depends on $\Delta T$ and the overlapping interval. Amongst the first SVD sources, $\mathbf{u}_1^m$, with $m = 1, ..., M$ and $M$ the total number of blocks, there will be some uniquely linked to the ground response while others influenced by ephemeral phenomena. The goal is to select a characteristic vector $\bar{\mathbf{u}}_1$ from amongst these $\mathbf{u}_1^m$ vectors.

The application of the mean operator for the selection purpose assures that the selected vector will be adapted to all time zones. However, there will be zones where the vectors $\mathbf{u}_1^m$ are dominated by phenomena other than the ground response (like precipitation) and, if significant, they can introduce false results due to averaging. In order to avoid this situation, we employ a criteria based on higher order statistics to identify these defective transient zones (time blocks). Once these zones have been identified, their contribution can be removed from the data. The third and the fourth order statistics are considered here, namely the skewness ($\kappa_3(t)$) and the kurtosis ($\kappa_4(t)$), respectively. The estimators for skewness and kurtosis are defined using the k-statistics and their definitions can be found in [11] along with their variances, $\sigma_3$ and $\sigma_4$. It should be mentioned that the variances of skewness and kurtosis estimators do not depend on time, $t$, but

uniquely on the number of elements used to estimate them, so on the number of sensors in our application. It was found out that the skewness and kurtosis for the temporal zones containing instances of precipitation were significantly different from those not containing any precipitation. The data selection is made by considering uniquely the zones for which the skewness and kurtosis values fall within a threshold around the reference values for the respective estimators. The reference values, $\kappa_3^{ref}$ and $\kappa_4^{ref}$, are calculated here as median values of the data skewness, $\kappa_3(t)$, and kurtosis, $\kappa_4(t)$. The median was chosen here to avoid factors that act as impulsive perturbations, such as significant precipitations, with respect to the analyzed time zone. For the selection of the threshold, it was observed that placing too low a threshold does not allow the separation of ephemeral phenomena from rest of the data. A threshold of $\pm\sigma$, for example, does not permit an efficient identification of ephemeral phenomena. However, the phenomena identified with thresholds of $\pm2\sigma$ and $\pm3\sigma$ are almost identical and thus $\pm2\sigma$ can be selected as an optimum threshold. Having removed the ephemeral time zones from the data, $\mathbf{Y}$, using the above criteria, the sliding window SVD of Eq. (4) is applied on this curtailed data, $\mathbf{Y}^{sel}$. An average ground response source vector is estimated by applying the mean operator on the first SVD vectors obtained for each position of the sliding window. The adopted approach is summarized in the three step algorithm in Fig. 3, where the first step is the HOS calculation and the threshold selection, the second step is the data selection eliminating the defective zones based on results of step 1 and the third step is the final source estimation using overlapping sliding temporal windows on the selected data followed by averaging to find the characteristic source vector. Once the characteristic source vector, $\bar{\mathbf{u}}_1$ has been obtained, we move on to the subspace separation step.
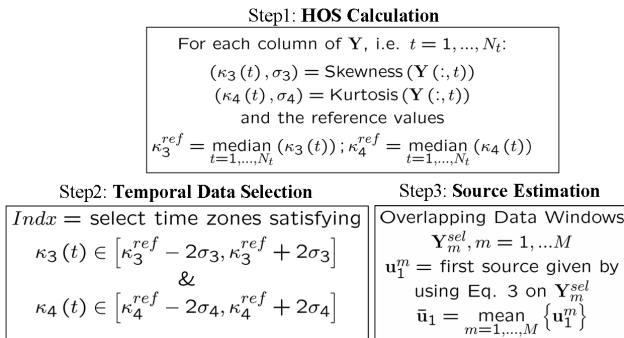


Figure 3: Algorithm for estimating an average ground response using a HOS criteria.

## 5. SUBSPACE DECOMPOSITION USING LEAST SQUARES

The construction of subspace corresponding to the average source vector, $\bar{\mathbf{u}}_1$, requires the temporal variation of this vector. One possible approach is to use the concept of source unmixing as posed in hyper spectral image processing [12]. Since, our goal is to remove the effects of ground response, characterized by $\bar{\mathbf{u}}_1$, we can rewrite this approach as follows. Let the linear model of the recorded data be given by:

$$\mathbf{Y} = \mathbf{Y}_{sig}^{LS} + \mathbf{Y}_{residue}^{LS} = \bar{\mathbf{u}}_1\alpha + \mathbf{Y}_{residue}^{LS} \qquad (5)$$

where $\alpha$ represents the temporal variations of the vector $\bar{\mathbf{u}}_1$ on the recorded data, and together with $\bar{\mathbf{u}}_1$ defines a signal subspace, $\mathbf{Y}_{sig}^{LS}$. The residue $\mathbf{Y}_{residue}^{LS}$ is defined by phenomena other than the ground response. The vector $\alpha$ can be estimated using a Least Squares (LS) procedure as:

$$\widehat{\alpha}_{LS} = (\bar{\mathbf{u}}_1^T \bar{\mathbf{u}}_1)^{-1} \bar{\mathbf{u}}_1^T \mathbf{Y} \qquad (6)$$

The final residue, $\mathbf{Y}_{residue}^{LS}$, is devoid of the effects of ground response. The residue subspace thus highlights the information linked to other important factors such as precipitations, anomalies (leakages), drains, etc. The differences between this decomposition and the one obtained by using the SVD Eq. (3) is that: (1) the source vector characterizing the ground response does not take into account ephemeral phenomena such as precipitations, and (2) there is no orthogonality condition imposed on the temporal variations of the sources.

## 6. RESULTS AND DISCUSSION

The proposed algorithm is applied to normalized real temperature data in Fig. 2. The data contains information linked to the drains, the leakages, instances of precipitation, all of whom are more or less masked by the response of the ground where DTS are buried. The normalization here has been done so as to bring each acquisition to zero mean and unity variance thus reducing the effects of seasonal variations. The ground response source estimation by processing the data in its entirety runs the risk of being influenced by phenomena like precipitation that are ephemeral in time. Thus we adopt the approach as highlighted in Fig. 3. In the first step, we calculate skewness and kurtosis of the data along with their variances. The corresponding results are presented in Fig. 4 with skewness in the top window and kurtosis in the bottom window. The thresholds, selected as $\pm2\sigma$, with $\sigma = \sigma_3$ being the variance of the skewness estimator and $\sigma = \sigma_4$ that of kurtosis estimator, are represented by the dashed lines. These results show that there are mainly two time zones (Z1 and Z2), where the values of these statistics surpass their respective thresholds. With the study of the meteorological data, it was found out that these two zones correspond primarily to instances of precipitation. The application of threshold for removing these two temporal zones results in selection of data for which the general trend observed with the higher order statistics remains more constant than with the entire data. Having curtailed the data to remove the influence of temporal ephemeral phenomena, the next step is to evaluate the ground response source vector by scanning the data using temporal sliding windows. A window size $\Delta T = 1$ day (corresponding to 12 acquisition points for the current data) was selected along with an overlapping of 25%. It should be mentioned here that other window sizes and sliding steps were also tested but no significant change in the end result was obtained. The first source vector, $\mathbf{u}_1^m$, is thus estimated in different sliding windows, with $m = 1,...,M$ and $M = 25$, the total number of overlapping windows in our application. The fact that the first sources represent about 95% of the total energy for each window and that they have approximately identical shapes means that these sources are related to the most coherent factor of the data, i.e., the ground response. The characteristic vector, $\bar{\mathbf{u}}_1$, is then obtained by taking the mean of all these overlapping vectors and is represented in Fig. 5. In this source, two separate levels can be observed, one
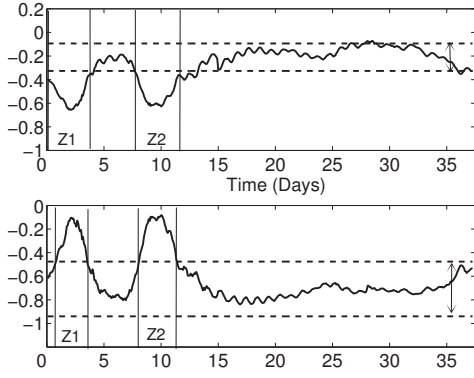
Figure 4: Higher Order Statistics: Skewness (top) and Kurtosis (bottom) with dashed lines the $\pm 2\sigma$ thresholding ranges.
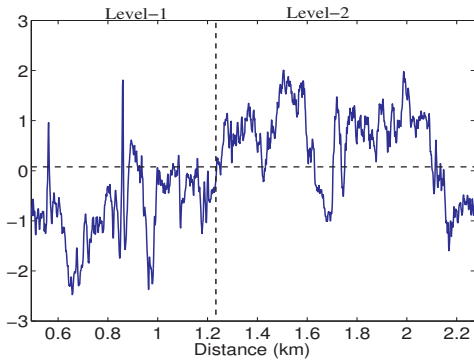


Figure 6: Residue subspace, $\mathbf{Y}_{residue}^{SVD}$, obtained with SVD.



Figure 5: Estimated source $\bar{\mathbf{u}}_1$ using the HOS criteria as highlighted by the scheme of Fig. 3.



Figure 7: Residue subspace $\mathbf{Y}_{residue}^{LS}$ using the LS approach on HOS based source $\bar{\mathbf{u}}_1$.

from approximately 0.5km to 1.25km and the other one from 1.25km to 2.2km. As evident from the fiber layout scheme of the actual test site (see Fig. 1), these levels correspond to the physical ground structure present at the actual site of data acquisition. Moreover, a strong singularity is noted in the region from $1.63 - 1.7$ km, which also corroborates well with the physical site as the material composition of the ground around this distance is different from those in other regions. This validates the fact that this source indeed represents the ground response. Moreover, due to the HOS criteria, this source is devoid of the precipitation effects while at the same time being a representative of all time zones due to averaging operation.

The signal subspace corresponding to the source obtained by application of SVD on the entire data (i.e. without applying the HOS criteria) can be constructed by using the subspace decomposition of Eq. (3). The residue, $\mathbf{Y}_{residue}^{SVD}$, obtained by subtracting this signal subspace from the normalized data is given in Fig. 6. The temporal variation of HOS based estimated source, $\bar{\mathbf{u}}_1$, is obtained using the LS based approach to avoid the orthogonality condition of SVD. The final residue, $\mathbf{Y}_{residue}^{LS}$, thus obtained using the HOS criteria and LS method is given in Fig. 7. It can be observed that the removal of the ground response results in bringing to evidence other phenomena such as the precipitations, the leakages and the drains (Fig. 6 and Fig. 7). The residue subspace obtained by applying SVD on the entire data differs from the one obtained using the proposed method. The first immediate observation is that subspaces using the proposed approach are overall smoother as compared to the approach
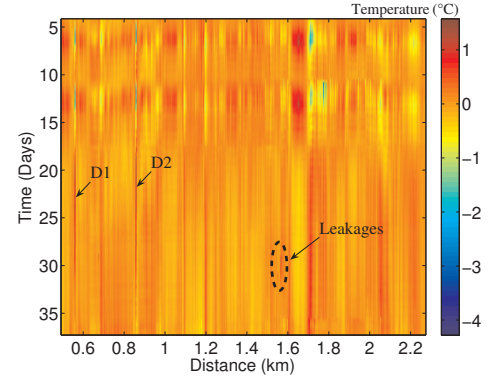
based on SVD. This is specially visible in the zones which are sparse in terms of their information content. The zoomed versions of the two residues (Fig. 8(a)-(b)) highlight this difference. Another noticeable artifact of the global SVD technique is that it compensates the time zones of more intensity with those of opposite intensity (Fig. 8(a)), which means to say that it performs a sort of positive/negative compensation in time. This is due to the orthogonality constraint imposed on the temporal variation of the sources. The overall result is that it introduces false information into the residue subspace which was not there in the first place. However, using our proposed method, these artifacts are not observed (Fig. 8(b)).



(a) Zoom of SVD residue   (b) Zoom of HOS based residue

Figure 8: Zoom of residues showing relatively smoother LS residue and undue compensation in SVD residue.

Fig. 9 shows the two residues on day 21, with $\mathbf{Y}_{residue}^{SVD}$ in the top and $\mathbf{Y}_{residue}^{LS}$ in the bottom plot. Both these plots have

been normalized by their respective maxima for ease of comparison. It can be noted that not only the LS based residue is smoother but also the singularity of the ground, in the $1.63 - 1.7$ km region (previously identified in Fig. 5, encircled here), is relatively more energetic in SVD based residue than the LS based residue. This singularity is thus better estimated by the proposed HOS and LS based approach as a part of ground response and thus eliminated in the residue.
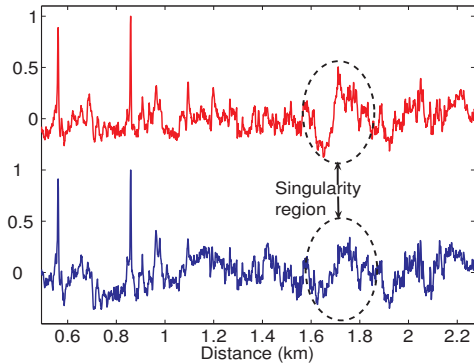


Figure 9: Residues on day 21 (top $\mathbf{Y}_{residue}^{SVD}$, bottom $\mathbf{Y}_{residue}^{LS}$) showing more pronounced singularity (encircled) in SVD residue than LS residue.

Moreover, the proposed approach reveals with relatively more intensity the useful phenomena in the residue subspace. As an example, a zoom of the two residues on day 28 (when $L1$ occurred at 1.562km) in Fig. 10 shows that the leakage $L1$ has higher SNR in $\mathbf{Y}_{residue}^{LS}$ (calculated with respect to the background in $L1$'s vicinity) than in $\mathbf{Y}_{residue}^{SVD}$. This is an important result for leakage detection objective as the leakage is better separated from the background.
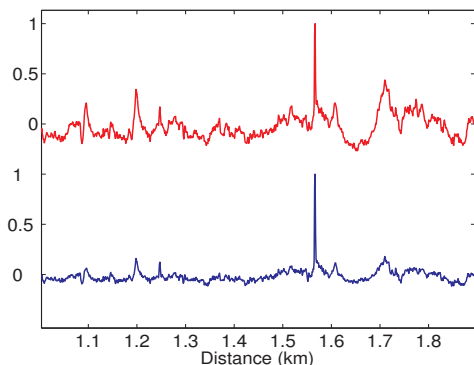


Figure 10: Zoom of residues on day 28 (top $\mathbf{Y}_{residue}^{SVD}$, bottom $\mathbf{Y}_{residue}^{LS}$) with more energetic interesting factors (leakage $L1$ here) in LS based residue.

## 7. CONCLUSION

The detection of anomalies in dikes using bidimensional temperature data acquired by distributed temperature sensors is a new research problem in the signal processing domain and there are not many priors to this work. In this paper, we presented a subspace separation technique to remove the influence of the ground response where acquisitions are made. The removal of this source is important as it masks other useful information in the data like the presence of leakages. It was observed that the first source of SVD applied on the

DTS data is linked to the ground response where the acquisitions are made. This source estimation might however be effected by transient temporal phenomena in the data. A criteria based on higher order statistics was proposed to identify them. It was shown that phenomena like precipitation can be efficiently identified by using the HOS, skewness and kurtosis. Based on this criteria, temporal data zones can be selected which are devoid of these transient phenomena. The first source vector is estimated in overlapping temporal sliding windows applied on this curtailed data. The characteristic ground response source is finally obtained as mean of the source vectors estimated in these time blocks. The temporal variation of this characteristic source (for constructing the corresponding ground response subspace) is estimated using least squares unmixing approach. This allows to avoid the unjustifiable orthogonality condition imposed by SVD on the temporal variation of the estimated source. The method was applied to a real data set and the resulting residue subspace was found not only to be devoid of the global SVD artifacts but also found out to better put to evidence the useful information.

## REFERENCES

[1] A. H. Hartog, "Progress in distributed fiber-optic temperature sensing," in *Proc. SPIE, Fiber Optic Sensor Technology and Applications, 2001*, M. A. Marcus and B. Culshaw, Eds., 2002, vol. 4578 of *0277-786X/02*, pp. 43–52.

[2] St. Grosswig, A. Graupner, and E. Hurtig, "Distributed fiber optical temperature sensing technique - a variable tool for monitoring tasks," in *Proc. 8th Intl. Symposium on Temperature and Thermal Measurements in Industry and Science*, 19-21, June 2001, pp. 9–17.

[3] B. Vogel, C. Cassens, A. Graupner, and A. Trostel, "Leakage detection systems by using distributed fiber optical temperature measurements," in *Proc. SPIE Smart Structures and Materials, 2001*, D. Inaudi E. Udd, Ed., vol. 4328, pp. 23–34.

[4] P. Comon, "Independent component analysis, A new concept ?," *Signal Processing*, vol. 36, no. 3, pp. 287–314, 1994.

[5] J. F. Cardoso, "Blind signal separation : statistical principles," *Proc. of IEEE*, vol. 86, no. 10, pp. 2009–2025, 1998.

[6] A. Hyvärinen, J. Karhunen, and E. Oja, *Independent Component Analysis*, Wiley, New York, 2001.

[7] E. Lewis, C. Sheridan, M. O'Farrell, D. King, C. Flanagan, W. B. Lyons, and C. Fitzpatrick, "Principal component analysis and artificial neural networks based approach to analysing optical fiber sensors signals," *Sensors and Actuators A*, vol. 136, pp. 28–38, 2007.

[8] V. C. Klema and A. J. Laub, "The singular value decomposition: its computation and some applications," *IEEE. Trans on Auto. Control*, vol. 25, no. 2, pp. 164–176, 1980.

[9] L. L. Scharf, *Statistical Signal Processing : Detection, Estimation, and Time Series Analysis*, Addison-Wesley, New York, 1991.

[10] V. Vrabie, J. I. Mars, and J-L . Lacoume, "Singular value decomposition by means of independent component analysis," *Signal Processing*, vol. 84, no. 3, pp. 645–652, 2004.

[11] M. G. Kendall and A. Stuart, *The Advanced Theory of Statistics, Volume1 - Distribution Theory, Second Edition*, Addison-Wesley, London, 1963.

[12] D. C. Heinz and C.-I Chang, "Fully constrained least squares linear spectral mixture analysis method for material quantification in hyperspectral imagery," *IEEE Trans. on Geoscience and Remote Sensing*, vol. 39, no. 3, pp. 529–545, March 2001.