

Leak Detection and Isolation in Water Distribution Networks using Principal Component Analysis and Structured Residuals

J. Gertler, J. Romera, V. Puig and J. Quevedo

Abstract— Leaks are present to some extent in all water-distribution systems. This paper proposes a leakage localization method based on pressure measurements and the application of principal component analysis to the fault diagnosis in water distribution systems. First, some theoretical basics are introduced, from model building and modeling the fault effects to monitoring. Then a simple hydraulic case study is presented to illustrate the proposed methodology, its particularities and the detection results.

I. INTRODUCTION

Water loss in distribution system networks is an issue of great concern for water utilities, strongly linked with operational costs and water resources savings. Continuous improvements in water loss management are being applied and new technologies are developed to achieve higher levels of efficiency. Usually a leakage detection method in a DMA (District Metered Area) starts by analyzing input flow data, such as minimum night flows and consumer metering data [1]. Once the water distribution district is identified to have a leakage, various techniques are used to locate the leakage for pipe replacement or repair. Methods for locating leaks range from ground-penetrating radar to acoustic listening devices or physical inspection [2][3]. Some of these techniques require isolating and shutting down part of the system. The whole process could take weeks or months with an significant volume of water wasted. Techniques based on locating leaks from pressure monitoring devices allow a more effective and less costly search in situ.

A methodology to detect and pre-localize leaks is being developed within a project carried out by Aguas Barcelona, Water Technological Centre CETaqua, and the Technical University of Catalonia (UPC) [4]. The objective of this project is to develop and apply an efficient system to detect and locate leaks in a water distribution network. It integrates methods and technologies available and in use by water companies, including DMA and flow/pressure sensor

data, in conjunction with mathematical hydraulic models. The method is based on the analysis of pressure variations produced by a leakage in the water distribution network [5]. This technique differs from others in the literature, such as the reflection method (LRM) or the inverse transient analysis (ITA), since it is not based on the transient analysis of pressure waves [6][7][8][9][10]. Alternatively, the leakage detection procedure is performed by comparing real pressure and flow data with their estimation using the simulation of the mathematical network model as suggested by Pudar [11]. In order to develop this methodology, the project includes a characterization of district metered areas and consumers, considered a critical issue for a correct model calibration. The project also proposes a methodology to place pressure sensors within a district metered area network, to optimize leakage detection using a minimum number of sensors. Finally, the leakage detection methodology proposed will be tested with sensors installed in three DMA's used as case studies.

In this article, as an alternative to the work previously developed in [5], the application of principal component analysis (PCA) to fault diagnosis in water distribution systems is proposed [14] [15]. The technique is composed of two phases. In the first phase, a model of the fault-free system is built off-line using PCA. Then, the fault effects are characterized through an experimental fault sensitivity analysis. Once the fault sensitivity matrix has been obtained, it is used to generate from the PCA model a set of structured residuals [12]. The second, on-line phase involves the computation of these residuals and their evaluation against a threshold, as usually done in model based FDI [13]. When an inconsistency is detected (some residual violates its threshold), then the fault can be isolated by using the fault-to-residual mapping established when designing the structured residuals.

The structure of the reminder of the paper is the following: First, in Section II, water distribution systems mathematical modeling is recalled. Section III introduces the theoretical basis of the proposed leak detection and isolation methodology is presented. Section IV the proposed methodology is illustrated using a simple hydraulic case study. Finally, in Section V, the main conclusions of the paper are provided.

J. Gertler is with Department of Electrical and Computer Engineering at George Mason University, 4400 University Drive, Fairfax, VA 22030 USA (e-mail: jgertler@gmu.edu).

J. Romera, V. Puig, and J. Quevedo are with the Advanced Control Systems Group, Universitat Politècnica de Catalunya (UPC), Rambla Sant Nebridi 10, 08222 Terrassa (Spain) (e-mail: {juli.romera, vicenc.puig, joseba.quevedo}@upc.edu).

II. WATER DISTRIBUTION SYSTEMS

A. Introduction

A water distribution system consists of three major components: pumps, distribution storage, and distribution piping network. Most systems require pumps to supply lift to overcome differences in elevation, and energy losses caused by friction. Pipes may contain flow-control devices, such as regulating or pressure-reducing valves.

The purpose of a distribution system is to supply the system's users with the amount of water demanded, under adequate pressure for various loading conditions. A loading condition is a spatial pattern of demands that defines the users' flow requirements.

In the analysis of water-distribution systems, the assumption is usually made that all the demands occur at the nodes. For the purposes of this paper, we will assume that leaks as well occur at the nodes.

B. Mathematical model

The governing laws for flow in pipe systems under steady conditions are *conservation of mass* and *energy*. The *law of conservation of mass* states that the rate of storage in a system is equal to the difference between the inflow to and outflow from the system. In pressurized water distribution networks, no storage can occur within the pipe network, although tank storage may change over time. Therefore, in a pipe, or a junction node, the inflow and outflow must balance. For a junction node,

$$\sum q_{in} - \sum q_{out} = q_{ext}$$

where q_{in} and q_{out} are the pipe flow rates into and out of the node and q_{ext} is the external demand or supply. Conservation of energy states that the difference in energy between two points is equal to the energy added to the flow in components between these points minus the frictional and minor losses. An energy balance can be written for paths between the two end points of a single pipe, between two fixed graded nodes (a node for which the total energy is known, such as a tank) through a series of pipes, valves, and pumps, or around a loop that begins and ends at the same point. In a general form for any path,

$$\sum_{i \in I_p} h_{P,j} - \sum_{i \in I_p} h_{L,i} = \Delta E$$

where: $h_{L,i}$ is the headloss across component i along the path, $h_{P,j}$ is the head added by pump j , and ΔE is the difference in energy between the end points of the path.

The primary network component is a pipe. The relationship between pipe flow (q) and energy loss caused by friction (h_L) in individual pipes can be represented by a number of equations, including the Darcy-Weisbach and Hazen-Williams equations. The general relationship is of the form

$$h_L = Kq^r$$

where K is a pipe coefficient that depends on the pipe's diameter, length, and material and r is an exponent in the range of 2.

III. PROPOSED LEAK DETECTION AND ISOLATION METHODOLOGY USING PCA STRUCTURED RESIDUALS

A. Model building for the fault-free system

Let us consider that we have n measurements $\mathbf{x} = [x_1, \dots, x_n]'$ (pressures or flows in the water system). Normal water consumption in the k nodes are considered disturbances $v_1 \dots v_k$ and are assumed to be random. Faults $f_1 \dots f_k$ are leaks in the same k nodes, assumed to be deterministic steps or ramps. Due to the network topology there will be m relations on the variables \mathbf{x} . These relations around a given operating will point will be assumed linear.

To build a model for the fault-free system, we need a training data-set $\mathbf{X} = [\mathbf{x}(1) \dots \mathbf{x}(N)]$. All variables in the training set must be "centered" (their average over the training set deducted). (To simplify notation, we will just use x for the centered data.)

From the training data-set, we form the covariance matrix

$$\mathbf{R} = \mathbf{X} \mathbf{X}' / N \quad (1)$$

and find its eigenvalues $\sigma_1^2 \dots \sigma_n^2$ and eigenvectors $\mathbf{q}_1 \dots \mathbf{q}_n$. Observe the eigenvalues; according to the theory of PCA, $m-k$ of the eigenvalues will be zero (or near zero). Any (centered) $\mathbf{x}(t)$ vector in the training set can then be described as

$$\mathbf{x}(t) = \sum \mathbf{q}_i p_i(t) \quad i = 1 \dots n - m + k \quad (2)$$

where the eigenvectors in the sum are those belonging to nonzero eigenvalues (these span the representation space) and where

$$p_i(t) = \mathbf{q}_i' \mathbf{x}(t) \quad (3)$$

is the projection of the $\mathbf{x}(t)$ vector on the \mathbf{q}_i direction. The n -dimensional $\mathbf{x}(t)$ is thus projected into an $n-m+k$ dimensional space. In that space, it may be expressed with its projections:

$$\mathbf{x}(t) \Rightarrow \mathbf{p}(t) = [p_1(t) \dots p_{n-m+k}(t)]' \quad (4)$$

A variable $x_i(t)$ that is constant over the training set becomes identically zero after centering. Then $x_i(t)=0$ needs to be considered as an additional equation (or x_i ignored as a variable).

In a typical water distribution network, pressure is measured at each distribution node while pressure and flow are measured at the source node. With ν distribution nodes, and assuming the source pressure is constant, $n=\nu+1$. There are ν node equations, plus a branch equation for the source, yielding $m=\nu+1$. The number of disturbances is $k=\nu$. Thus the dimension of the representation space is $n-m+k=\nu$; clearly, the reduction in dimensionality is not significant.

B. Modeling the fault effects

While sensor and actuator faults may be handled by a simple extension of the fault-free system model, process faults require explicit modeling. This may be done on the physical system, by emulating faults, or on a simulator. In either case, a fault f_j is introduced (with known location and size) and the resulting $\mathbf{x}(t|f_j)$ vector is observed (centered with the average of the fault-free training data). This may then be utilized in two different ways:

- As the direction of the response, by normalizing $\mathbf{x}(t|f_j)$ into a unit-length vector:
- As the gain (sensitivity) of the response, by scaling $\mathbf{x}(t|f_j)$ to unit fault-size.

Usually the fault responses are not completely in the representations space (cannot be described by equation (2)) but reach, at least partially, into its complement, the residual space. This may be advantageous in fault isolation. In the water flow system, however, the faults are co-linear with the disturbances (act in the same nodes) therefore the fault responses are entirely in the representation space. To utilize this, the normalized or scaled responses to the emulated (simulated) faults are transformed into the representation space, by (3) and (4). In the following, we will make use of the transformed sensitivity vectors $\mathbf{s}_{\bullet j} = [s_{1j} \dots s_{(n-m+k)j}]'$, $j=1 \dots k$, where

$$s_{ij} = \mathbf{q}_i' \mathbf{x}(t|f_j) / f_j \quad i=1 \dots n-m+k, j=1 \dots k \quad (5)$$

Note that it is necessary to simulate each fault with a number of different sizes, to check the linearity of the response or assess its deviation from linearity, due to the nonlinearity of the water system.

C. Monitoring

In the monitoring phase, measurements $\mathbf{x}(t)$ are taken from the physical (or simulated) system (and centered with the average of the fault-free training data-set). These centered measurements are transformed into the representation space, as $\mathbf{p}(t)$. These latter are then used as residuals in fault detection and isolation.

The residuals may be evaluated in a directional or a structural framework. The directional approach is rather complex and is not discussed here. In a structured framework, the fact that the dimension of the residual vector $\mathbf{p}(t)$ equals the number of faults allows for the design of a diagonal structure. Define $\mathbf{S} = [\mathbf{s}_{\bullet 1} \dots \mathbf{s}_{\bullet k}]$, then the transformed residual

$$\mathbf{r}(t) = \mathbf{S}^{-1} \mathbf{p}(t) \quad (6)$$

obeys a diagonal structure: each residual responds to one specific fault. Such structure allows for the isolation of multiple simultaneous faults.

The residuals first need to be filtered to reduce the effect of noise (in the case of the water distribution system, this "noise" is primarily the random water consumption). A first-order recursive filter

$$\mathbf{r}_F(t) = \alpha \mathbf{r}_F(t-1) + (1-\alpha) \mathbf{r}(t) \quad (7)$$

is usually sufficient, with $\alpha=0.9$ (or similar value). Note that there is a tradeoff between noise filtering and fault-response speed.

Even with residual filtering, the elements of $\mathbf{r}_F(t)$ need to be subjected to statistical testing. The thresholds are best established by observing the variation of the training data-set, after centering, transformation into the representation space, and filtering (7). (The thresholds may be determined theoretically as well. The eigenvalue σ_i^2 is the variance of the training data in the \mathbf{q}_i direction. However, the distribution of the disturbances is uni-directional; after centering the training data, it becomes bi-directional, truncated in the negative direction. Thus the usual assumption of normal distribution could not be applied).

Due to the diagonal structure of $\mathbf{r}(t)$, the elements of the residual vector are estimates of the individual fault sizes (fault identification). These estimates, however, beyond being noisy, are also subject to an error due to the linearization of a nonlinear relationship.

IV. APPLICATION EXAMPLE

A. Introduction

To illustrate the proposed PCA based leak detection and isolation methodology, the network presented in Figure 1 is used as case study. The system has four nodes ($v=4$), and a tank (source node). Pressure (P_i) is measured at all nodes but the source pressure is constant so it is ignored as a variable. Also, flow is measured at the source (q_0). So, there are $n=v+1=5$ measurements. There is a disturbance (demand) d_i and a possible leak f_i in all ordinary nodes, $k=v=4$. There is a node equation for all ordinary (not source) nodes and a branch equation for the source, so $m=v+1=5$:

$$\begin{aligned} d_1 &= q_0 - q_1 - q_2 \\ d_2 &= q_1 - q_2 \\ d_3 &= q_3 + q_4 \\ d_4 &= q_2 - q_4 \\ P_0 - P_1 &= (R_0 q_0)^2 \end{aligned}$$

Thus, $n-m+k = n-1=4$. Note that the node equations are non-linear, since pressure is measured instead of flow, so the above considerations may not strictly apply.

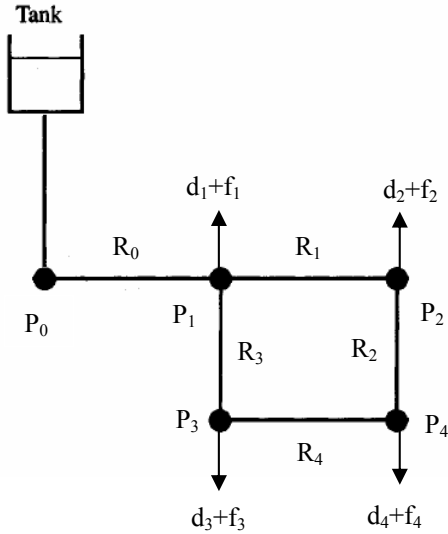


Fig. 1. Case study

The whole non linear model is the following:

$$\begin{bmatrix} -\frac{1}{R_1} & 0 & -\frac{1}{R_3} & 0 \\ \frac{1}{R_1} & -\frac{1}{R_2} & 0 & 0 \\ 0 & 0 & \frac{1}{R_3} & \frac{1}{R_4} \\ 0 & \frac{1}{R_2} & 0 & -\frac{1}{R_4} \end{bmatrix} \begin{bmatrix} \sqrt{P_1 - P_2} \\ \sqrt{P_2 - P_4} \\ \sqrt{P_1 - P_3} \\ \sqrt{P_3 - P_4} \end{bmatrix} = \begin{bmatrix} d_1 \\ d_2 \\ d_3 \\ d_4 \end{bmatrix} - \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix} q_0$$

The water distribution system is simulated using the non-linear model in SIMULINK (Figure 2).

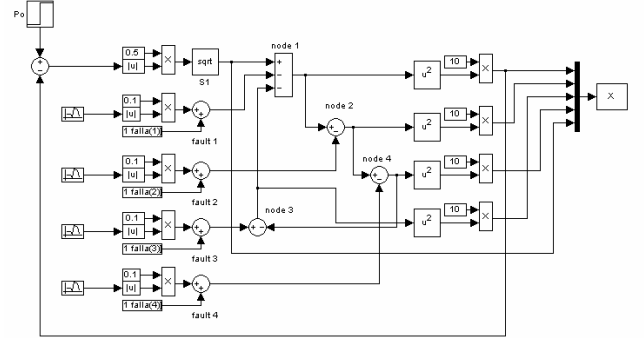


Fig. 2. SIMULINK model

B. Fault-free system model using PCA

Using the water network simulator, a training data-set of a length $N=150$ is generated driving the system with random loads (demands) at the nodes in the fault-free situation. The random demands were generated using the absolute value of normally distributed variables (with zero mean and a standard deviation of $0.1 \text{ m}^3/\text{s}$). (Note that this captures the uni-directional nature of the load but places the maximum likelihood at zero load that is not really correct.)

The covariance matrix (1) from the normalized recorded measurements produces the following eigenvalues

$$[0, 0.1430, 0.7965, 2.3515, 7.1432]$$

Notice that the first eigenvalue is 0. This is in concordance with the PCA theory that establishes that $m-k$ of the eigenvalues will be zero (so, in this case, $5 - 4 = 1$).

C. Fault system model using PCA

In order to see how linearly the projection $p_i(t|f_j)$ of the $\mathbf{x}(t|f_j)$ vector on the \mathbf{q}_i direction behaves, a ramp fault varying from $0.1 \text{ m}^3/\text{s}$ to $1.5 \text{ m}^3/\text{s}$ is applied to every node (Figure 3).

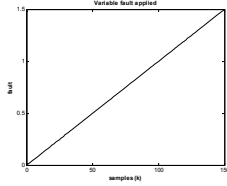


Fig. 3. Ramp fault

Figure 4 shows the variation of the projections $p_i(t|f_j)$ when the ramp fault is applied to only one node each time. From these plots, it can be noticed that the behavior of these projections is linear for small fault sizes but tends to be more non-linear when the fault size increases. Approximating the evolution of the projections $p_i(t|f_j)$ by a straight line will allow to derive the fault sensitivities presented in (5), yielding the sensitivity matrix S that is used in (6) to obtain diagonally structured residuals.

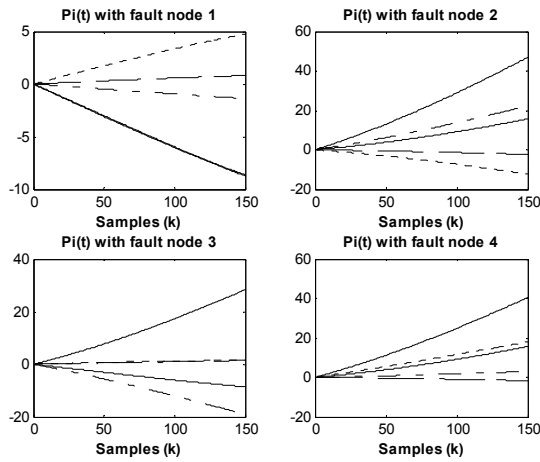


Fig. 4. Non-linearity behavior

Note that due to the non-linearity of the system, the validity of the sensitivity matrix S is also restricted to the neighborhood of the operating point (average load situation) in which it is computed. In the remainder of this section, the S matrix will be used to detect small leaks (around $0.2 \text{ m}^3/\text{s}$). The S matrix is obtained, in the selected operating point, by applying faults of size $0.2 \text{ m}^3/\text{s}$, and then normalizing the resulting fault responses to unit fault size (dividing by 5). This yields

$$S = \begin{bmatrix} 0,6275 & -1,0502 & 0,9567 & -0,7022 \\ -0,9625 & 11,4656 & -10,6152 & 1,2668 \\ 3,4881 & -5,5405 & 1,6238 & 10,5610 \\ -6,0044 & 24,1213 & 14,2013 & 20,7291 \end{bmatrix}$$

D. Monitoring phase

To show the performance of the proposed approach, a fault scenario consisting of two simultaneous leaks (one of

magnitude $0.3 \text{ m}^3/\text{s}$ introduced at the first node and a second one of magnitude $0.1 \text{ m}^3/\text{s}$ introduced at the second node) is created with the simulator. Figure 5 shows the elements of $\mathbf{x}(t)$, in the fault-free and faulty scenario. Figure 6 presents the diagonally structured residuals $\mathbf{r}(t)$, obtained using (6), with the sensitivity matrix S as presented in the previous section. In the top picture, there is the system response without faults, while at the bottom these residuals are presented in the faulty scenario. In Figure 7, the diagonal residuals are shown after filtering.

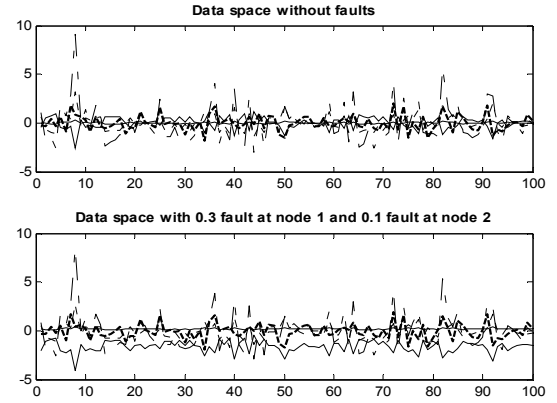


Fig. 5. Data space $\mathbf{x}(t)$ in non-faulty and faulty situation

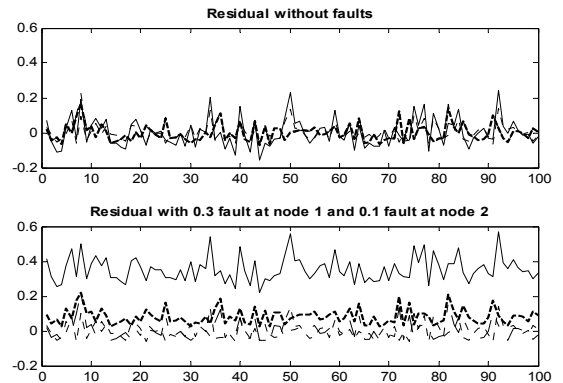


Fig. 6. Structured residuals $\mathbf{r}(t)$ in non-faulty and faulty situation

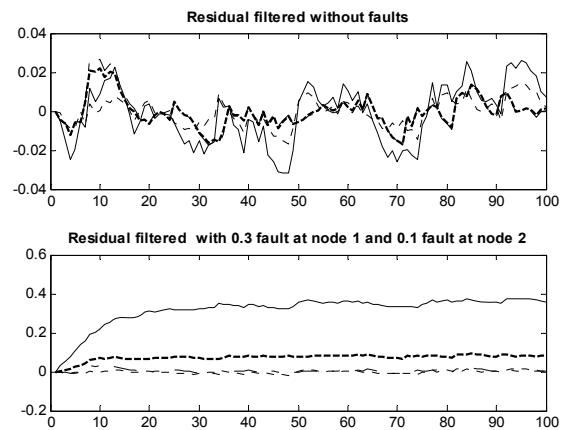


Fig. 7. Filtered structured residuals $\mathbf{r}_f(t)$ in non-faulty and faulty situation

The fault indication is presented in the same figure. From this figure it can be read that the evaluation of the residuals leads to

r_1	r_2	r_3	r_4
1	1	0	0

that clearly indicates the presence of one leak each in node 1 and node 2. Finally, in Figure 8, the filtered diagonally structured residuals are compared to the respective thresholds. The result of this comparison leads to 0 (no fault indicated) or 1 (fault indicated) (the detection is shown as 0.5 to compare it easily).

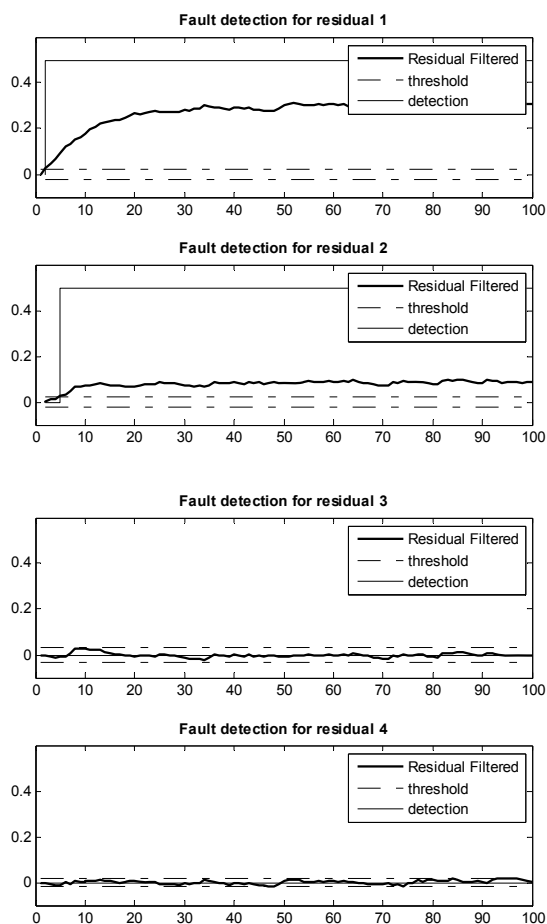


Fig. 8. Detection results

V. CONCLUSIONS

A PCA-based approach is proposed for the detection and isolation of leaks in water distribution systems. A PCA model of the fault-free system is obtained from operational (or simulated) data. Response to the specific leaks (faults) is modeled by emulating (simulating) each fault. These responses may then be utilized in fault diagnosis either in a directional or a structured (diagonal) framework. In the monitoring phase, the residuals are subjected to filtering

and threshold testing. The water distribution problem has two special characteristics: First, there is a disturbance (normal usage) in each node. This limits the reduction of dimensionality normally offered by PCA (subnetworks, rather than the full network, might be considered). Second, the faults are co-linear with the disturbances (act in the same nodes). This hinders the spatial separation of faults from normal data. While the above characteristics limit the utilization of PCA's full potential, it is still a powerful diagnostic tool and a strong alternative to other approaches.

REFERENCES

- [1] MacDonald G. "DMA Design and Implementation", a North American Context. Leakage Conference, IWA. 2005.
- [2] Farley M., Trow. "Losses in Water Distribution Networks". IWA Publishing UK. 2003.
- [3] Colombo, A.F., Lee, P., Karney, B.W. "A selective literature review of transient-based leak detection methods". *Journal of Hydro-environment Research*, 2, pp. 212-227, 2009.
- [4] Landeros E., Pérez R., Peralta A., Cembrano G. "Leakage detection using pressure sensors and mathematical models", *Water Loss 2009*, Cape Town, South Africa. 2009.
- [5] Pérez, R., Puig, V., Pascual, J., Peralta, A., Landeros, E. and Jordanas, LI. "Pressure sensor distribution for leak detection in Barcelona water distribution network". *Water Science & Technology*, Vol 9 No 6 pp 715-721, 2009
- [6] Ferrante, M., Brunone, B. "Pipe system diagnosis and leak detection by unsteady-state tests. Part 1. Harmonic analysis". *Advances in Water Resources*, Volume 26, Issue 1, January 2003, Pages 95-105.
- [7] Ferrante, M., Brunone, B. Pipe system diagnosis and leak detection by unsteady-state tests. Part 2. Wavelet analysis. *Advances in Water Resources*, Volume 26, Issue 1, January 2003, Pages 107-116.
- [8] Misiunas, D., Vítkovský, J., Olsson, G., Simpson, A., Lambert, M. "Pipeline Break Detection Using Pressure Transient Monitoring". *J. Water Resour. Plng. and Mgmt.* Volume 131, Issue 4, pp. 316-325, 2005.
- [9] Verde, C., Visairo, N., Gentil, S. "Two leaks isolation in a pipeline by transient response". *Advances in Water Resources*, Volume 30, Issue 8, pp. 1711-1721, 2007.
- [10] Misiunas, D., Vítkovský, J., Olsson, G., Simpson, A., Lambert, M. (2005). "Pipeline Break Detection Using Pressure Transient Monitoring". *J. Water Resour. Plng. and Mgmt.* Volume 131, Issue 4, pp. 316-325, 2005.
- [11] Pudar R.S., Liggett J. A. "Leaks in Pipe Networks". *Journal of Hydraulic Engineering*, Vol. 118, No. 7, July 1992, pp. 1031-1046.
- [12] Gertler, J., Li, W., Huang, Y. and McAvoy, T., "Isolation Enhanced Principal Component Analysis." *AIChE Journal*, **45**, pp. 323-334, 1999.
- [13] Gertler J.J., *Fault Detection and Diagnosis in Engineering Systems*, Marcel Dekker, 1998.
- [14] Kresta, J.V., J.F. MacGregor and T.E. Marlin, "Multivariate Statistical Monitoring of Processes," *Can. J. Chem. Eng.*, 69, 35 (1991).
- [15] B. M. Wise and N. B. Gallagher, "The Process Chemometrics Approach to Chemical Process Fault Detection and Supervision," *Journal of Process Control*, 6(6), pps 329-348, (1996).