

Leak localization in water distribution networks using classifiers with cosenoidal features

I. Santos-Ruiz ^{*,***,1} J. Blesa ^{**,**,*} V. Puig ^{***}
F. R. López-Estrada ^{*}

^{*} *Tecnológico Nacional de México, Instituto Tecnológico de Tuxtla Gutiérrez, TURIX Dynamics - Diagnosis and Control Group. Carr. Panamericana km 1080, 29050 Tuxtla Gutiérrez, Chiapas, Mexico (e-mail: {idelossantos,frlopez}@ittg.edu.mx)*

^{**} *Supervision, Safety and Automatic Control Research Center (CS2AC) of the Universitat Politècnica de Catalunya, Campus de Terrassa, Gaia Building, Rambla Sant Nebridi, 22, 08222 Terrassa, Spain. (e-mail: {vicenc.puig,joaquin.blesa}@upc.edu)*

^{***} *Institut de Robòtica i Informàtica Industrial (CSIC-UPC). Carrer Llorens i Artigas 4-6, 08028 Barcelona, Spain.*

^{****} *Serra Hünter Fellow Department of Automatic Control (ESAI), Universitat Politècnica de Catalunya (UPC), Avda. Eduard Maristany, 16 08019 Barcelona, Spain*

Abstract: This paper presents a leak localization approach for water distribution networks using classifiers with pressure residuals as input features. This approach is based on applying a non-linear transformation to the residuals of the node pressures to increase the separability of the leak classes. The transformed features can be interpreted as the direction cosines in the subspace spanned by the residuals of the measured pressures. In order to illustrate the method, different tests were performed with MATLAB[®] applying four different classification algorithms on a synthetic dataset obtained from an EPANET model of the Hanoi network. Then, by considering the cosenoidal features, a significant improvement in the leak location error was achieved. In this way, the leak location error decreases by more than 97% compared to the use of residual features when accurate measurements are used, and about 50% when noisy measurements with 60 dB SNR are used.

Keywords: Fault Diagnosis, Water Distribution Network, Leak Localization, Machine Learning, Feature Extraction, Direction Cosines.

1. INTRODUCTION

Loss due to leaks in water distribution systems is one of the main problems in managing drinking water. The percentage of chemically treated water loss in pipelines before reaching final consumers is around 30% worldwide, but in some cities, it exceeds 60% (OECD, 2016). Faults due to broken pipes, even small ones, generate important losses of water volumes when they remain unrepaired for a long time. Therefore, it is necessary to detect in the short term and localize them as accurately as possible for their rapid repair.

Leaks are not always visible, because leaking water can drain down the pipe instead of emerging towards the surface. Therefore, the precise localization requires the use of specialized instrumentation (vibration meters, geophones, etc.). Also, computational systems that, by monitoring the hydraulic variables of the network (pressures and flow rates), allow alerting about the existence of leaks and delimit them in an area of few meters to facilitate the work

of the maintenance staff. In Puig et al. (2017), the main monitoring and diagnostic techniques in pipelines and water distribution networks are described as well as the control strategies frequently used to minimize the effect of leaks. The localization of leaks in water distribution networks (WDNs) is a challenging problem due to the uncertainty that affects the performance of the localization algorithms. The most relevant is the uncertainty in the leak size, in the measurements (sensors noise), in the users' demand, and some parameters of the hydraulic model of the network such as the roughness and the actual diameter of the pipes, among others.

In large cities, a WDN is divided into different sectors, also known as District Metered Areas (DMAs). For control and billing purposes, DMAs are usually equipped with pressure and flow sensors at their inlets. Some of the recently proposed leak localization methods use these sensors and additional pressure sensors in inner nodes of the DMA that are cheaper and easier to install than flow sensors. In Pérez et al. (2011), a model-based method that relies on the pressure measurements and leak sensitivity analysis

¹ Corresponding author: idelossantos@ittg.edu.mx

was proposed. In this methodology, pressure residuals, i.e., differences between pressure measurements provided by sensors and the corresponding estimations obtained by using the hydraulic network model, are used. These residuals are computed on-line and compared against associated thresholds that take into account the effects of modeling uncertainty and noise. When a residual exceeds the thresholds, this is matched against a leak sensitivity matrix in order to identify which leak was presented. Several further works (Casillas et al., 2013; Pérez et al., 2014; Pérez et al., 2017) have proof that this method can provide reasonable results in real cases where the performance is affected by uncertainty in measurements, sensor noises, and mismatches between estimated and real demand users and between estimated and real hydraulic parameters (roughness and the actual diameter of the pipes, among others). For instance, see Cugueró-Escofet et al. (2015); Blesa and Pérez (2018) for discussions about the effect of uncertainties in residual correlation methods.

In the last years, artificial intelligence methods have been applied for leak localization purposes using pressure measurements. For example, in Mashford et al. (2009) was proposed a method to localize leaks using Support Vector Machines (SVM) that analyzes data obtained by a set of pressure sensors of a pipeline network allowing to localize and estimate the size of the leak. In a similar way, the use of k -Nearest Neighbors, Bayesian classifiers, Fisher discriminant analysis, and convolutional neural network for leak location have also been proposed in Soldevila et al. (2016), Soldevila et al. (2017), Romero-Tapia et al. (2018) and Javadiha et al. (2019), respectively. The performance in leak location of some of these methods has been assessed in Quiñones-Grueiro et al. (2018).

The main advantage of these artificial intelligence methods is that they are data-driven methods that formulate the problem of leak localization as a supervised multiclass classification problem. They use a matrix of node pressures or their residuals (differences between actual pressures and leak-free pressures) to train a classifier, which will then be used to predict the network nodes closest to the leaks. Therefore, if enough real pressure data were available from the network, it would not be necessary for any hydraulic model. However, in practice is in general not possible to have real data considering all leak scenarios and operating conditions in the DMA. Then, artificial data can be generated with a hydraulic model considering model uncertainties that can be extracted from real leak-free data (Blesa and Pérez, 2018). In this way, model uncertainties are considered in the design of the leak localization method.

To the best of the authors' knowledge, all the leak localization methods that have been formulated as a multiclass classification problem considering inner pressure measurements, use the pressure values or residual pressures as features. An exploratory analysis of the residuals suggests that leaks in the same node tend to show a characteristic direction in the sensor subspace (see Figure 1). In fact, this is a fundamental hypothesis in model-based leak localization methods based on leak sensitivity matrix (Pérez et al., 2011). However, if data-driven methods use residuals as input characteristics in a raw Cartesian form where both magnitudes and directions of residual vectors appear

implicitly combined, makes classification more difficult. Therefore, to improve the performance of classifiers for leak localization, in this paper it is proposed to take the features of a subspace derived from the original residual subspace by means of a non-linear transformation that summarizes only the information on the direction of leaks, discarding the information about its magnitude.

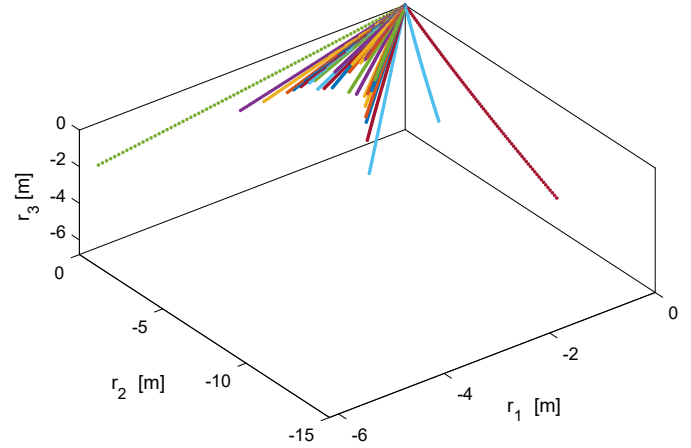


Fig. 1. Leaks of different magnitudes located in different nodes, plotted on the residuals subspace. Different colors are used for each node.

The structure of the reminder of the paper is as follows: In section 2, the proposed leak localization methodology is described. In Section 3, results obtained from the application a well-known case study are presented and discussed. Finally, Section 4 summarizes main conclusions and suggests future research paths.

2. METHODOLOGY

Consider a network consisting of n nodes, and suppose that pressure measurements on s inner nodes are available. These measurements at any time, that are sensitive to leaks, are denoted by $\mathbf{x} = [x_1, x_2, \dots, x_s]$, and the corresponding leak-free pressures are represented by $\mathbf{x}^* = [x_1^*, x_2^*, \dots, x_s^*]$. It is assumed that the s sensors have been placed on selected nodes under some optimality criteria (Casillas et al., 2015; Blesa et al., 2016).

From simulations with the network model, a matrix of nodal pressures $\mathbf{X} \in \mathbb{R}^{m \times s}$ is constructed considering m different leakage scenarios (different leakage magnitudes and different leaky nodes):

$$\mathbf{X} = \begin{matrix} \begin{matrix} \text{Number of pressure sensor} \rightarrow \\ \begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1s} \\ x_{21} & x_{22} & \cdots & x_{2s} \\ \vdots & \vdots & \ddots & \vdots \\ x_{m1} & x_{m2} & \cdots & x_{ms} \end{bmatrix} \\ \downarrow \\ \text{Leak scenarios} \end{matrix} \end{matrix} \quad (1)$$

Also, a vector of targets $\mathbf{y} \in \mathbb{N}^m$ is constructed containing the desired output classes (node numbers) associated with each leakage scenario in \mathbf{X} . A matrix of leak-free pressures \mathbf{X}^* , the same size as \mathbf{X} , is also generated, considering the nominal operating conditions of the network, according to the expected users' demand at the time of each leakage scenario. Each row \mathbf{x}_i of \mathbf{X} is a sample of node pressures for which a residual $\mathbf{r}_i = \mathbf{x}_i - \mathbf{x}_i^*$ can be calculated.

As mentioned in the introduction, pressure residuals provide relevant information to locate leaks. The previous works cited are based directly on the classical Cartesian components of these residuals to estimate the location of leaks. Nevertheless, according to Vector Analysis, vectors can also be expressed in a form where information about magnitude and direction is decoupled (Young, 2017). Thus, for any residual vector $\mathbf{r} = [r_1, r_2, \dots, r_s]$, the decoupled expression is

$$\mathbf{r} = M [\cos \theta_1, \cos \theta_2, \dots, \cos \theta_s] \quad (2)$$

where M is the residual magnitude, and

$$\cos \theta_k = f_k(r_1, r_2, \dots, r_s) = \frac{r_k}{\sqrt{r_1^2 + r_2^2 + \dots + r_s^2}} \quad (3)$$

are the so-called “direction cosines” and uniquely describe the direction of vector \mathbf{r} in the s -dimensional subspace, as shown in Figure 2.

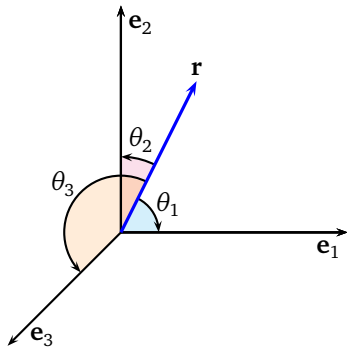


Fig. 2. Geometric interpretation of θ_i in a three-dimensional case. The vectors \mathbf{e}_i are unitary and they indicate the direction of each residual, e.g. $\mathbf{e}_1 = [1, 0, 0]$.

Although s direction cosines can be calculated for each residual, only $s - 1$ are independent of each other, because as can be deduced from (3), they satisfy the relationship

$$\sum_{k=1}^s \cos^2 \theta_k = 1. \quad (4)$$

Considering (4), the number of features used by the classifier can be reduced by one, because any one of the direction cosines is determined by the others.

Regarding the training procedure and the predictive use of the classifiers in the location of leaks, these have been previously described by Ferrandez-Gamot et al. (2015). The modification proposed in this work consists in the use of a new type of input features. As will be shown later, the replacement of the original features r_k by the new features $\cos \theta_k$ improves the performance of classifiers in leak location, facilitating the class separability. In the diagram of Figure 3, the gray box shows where the feature transformation is applied to improve the classifier performance in the leak localization process.

In a way, the proposed feature transformation can be seen as an ad-hoc kernelization, because the original features are projected into another subspace through a non-linear transformation. However, in this proposal, the subspace dimension does not increase as is usually the case when the kernel trick is applied.

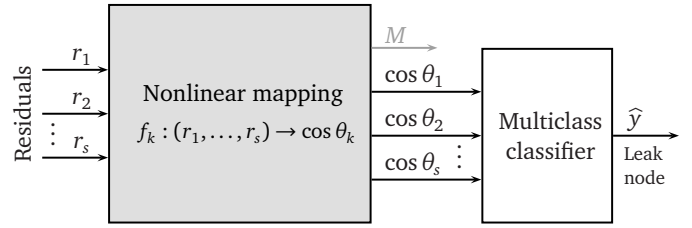


Fig. 3. Feature transformation to improve the classifier performance in leak location.

3. RESULTS AND DISCUSSION

3.1 Experimental setup

Different algorithms for data-driven leak location using classifiers were implemented in MATLAB®, and tested with a synthetic database of leaks in the Hanoi network (Fujiwara and Khang, 1990), which is shown in Figure 4. This network consists of 32 nodes (31 junction nodes and one reservoir) and 34 pipelines with a total length of 39 420 m.

The database used was generated through a steady-state simulation with the EPANET software (Rossman, 2000) using a hydraulic model that considers the pressure available in the reservoir, the geometry of the pipes and their roughness, as well as the demands on the consumption nodes. The leaks were simulated by manipulating the demands on the consumption nodes using the EPANET/MATLAB Toolkit interface (Eliades et al., 2016), increasing the base demand by an amount equal to the simulated leakage rate. The procedure for using the EPANET solver from MATLAB has been described by Vegas Niño et al. (2018).

In each node of the network, leaks of different magnitudes were simulated, considering flows $Q_{\text{leak}} = \{1, 2, \dots, 50\} 1/s$. Only single leaks (non-concurrent leaks) have been considered. In this way, a matrix of node pressures with 1 550 hypothetical leak scenarios was built, corresponding to the 51 different leakage magnitudes for each of the 31 junction nodes. This dataset was partitioned into 2 subsets, one for training and one for testing: Half of data, corresponding to $Q_{\text{leak}} = \{1, 3, \dots, 49\} 1/s$, were used for training; the other half, corresponding to $Q_{\text{leak}} = \{2, 4, \dots, 50\} 1/s$, were used for testing. Using a test dataset other than the training dataset will assess the predictability of the classification models used to locate leaks. The entire dataset is available in binary MATLAB® format (MAT-file) for download at <http://github.com/isantosruiz/direction>.

From the database, four different classification algorithms were tested: k -Nearest neighbors, Naïve Bayes, Decision tree, and Linear discriminant. The training and prediction tests of the classification models were performed using the MATLAB® Statistics and Machine Learning Toolbox.

To assess the performance of classifiers in leak location using different input features, the following error measure (often called “classification loss”) is used:

$$E = 1 - \frac{\sum_i c_{ii}}{\sum_i \sum_j c_{ij}}, \quad (5)$$

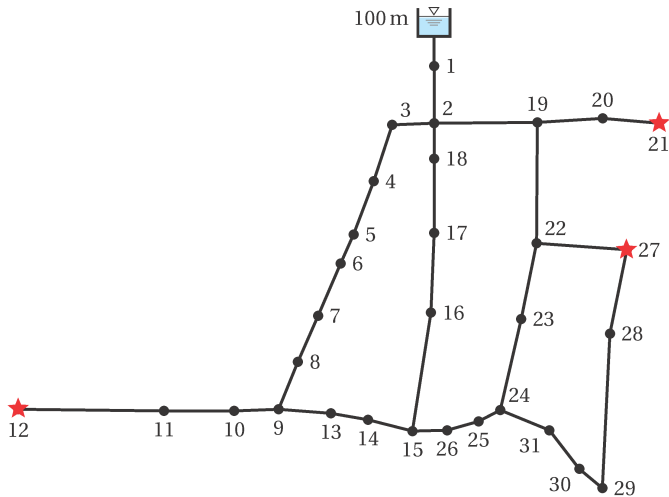


Fig. 4. Hanoi Network. Pressure sensors are located at the starred nodes.

where $[c_{ij}]$ is the confusion matrix. The localization error is calculated both in the training subset itself and in the subset selected specifically for testing. Then, to quantify the improvement obtained by the use of cosenoidal features, instead of unprocessed residuals, the following improvement index is used:

$$I = \frac{E_{\text{res}} - E_{\text{cos}}}{E_{\text{res}}}, \quad (6)$$

where E_{res} is the leak localization error obtained using residual features, and E_{cos} is the corresponding error when using cosenoidal features.

Additionally, since the classification loss (5) only provides a reference of the classification goodness and not how good it is the leak localization, the *Average Topological Distance* (ATD) is also calculated. This indicator was proposed by Soldevila et al. (2016) to assess the overall performance of a leak localization method in a real DMA. The ATD is defined as the average value of the minimum distance in nodes between the node with the leak and the node predicted by the leak localization method, and is computed as follows:

$$\text{ATD} = \frac{\sum_i \sum_j c_{ij} d_{ij}}{\sum_i \sum_j c_{ij}}, \quad (7)$$

where $[d_{ij}]$ is a symmetric matrix such that each element d_{ij} contains the minimum topological distance in nodes between the nodes referred by i and j .

3.2 Simulation results

From the dataset described above, classifiers of four different machine learning methods, using both residual and cosenoidal features, were trained and tested to compare performance with each type of feature. In all cases, only the pressure measurements corresponding to nodes 12, 21, and 27 were used. For classification, the node numbers $\{1, 2, \dots, 31\}$ where leaks occur were used as class labels. The results of the performance test are summarized in Tables 1 and 2. Based on these results, it was determined that the classification error with k -NN had decreased 99.3% when using cosines, with respect to when residuals are used directly as features. When Naïve Bayes, Decision Tree, and

Linear Discriminant classifiers are used, the classification error is reduced by 99.7%, 99.6%, and 97.0%, respectively.

Table 1. Classification error in training data (resubstitution loss) using different feature sets.

Classification Method	Features	
	Residuals	Cosines
k -Nearest Neighbors [†]	0.4813	0.0013
Naïve Bayes	0.7665	0.0026
Decision Tree	0.4516	0.0013
Linear discriminant	0.7729	0.0232

[†] Using Euclidean distance with $k = 5$.

Table 2. Classification error in testing data using different feature sets.

Classification Method	Features	
	Residuals	Cosines
k -Nearest Neighbors [†]	0.3458	0.0026
Naïve Bayes	0.7613	0.0026
Decision Tree	0.5936	0.0026
Linear discriminant	0.7665	0.0232

[†] Using Euclidean distance with $k = 5$.

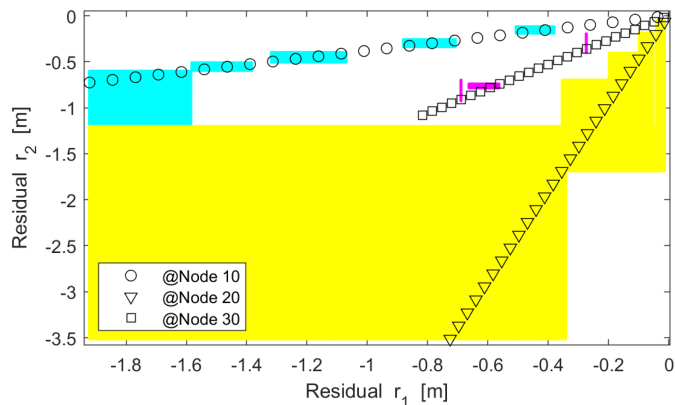
The lower classification error obtained when using cosines as features lead to a better class separability. Therefore, with the four classifiers tested, cosines better capture the directionality of leaks in the residual subspace. This is shown in Figure 5(b), where it is observed that the convergence region towards three of the 31 classes in the Hanoi network, using cosenoidal features, are better defined than those corresponding to Figure 5(a) where unprocessed residual features are used. Both figures correspond to a classification by decision tree but with a different types of features.

In order to analyze the robustness of the classifiers fed by cosenoidal features, leak localization tests were performed considering measurement noise at node pressures. The noise was assumed Gaussian and characterized by the signal-to-noise ratio:

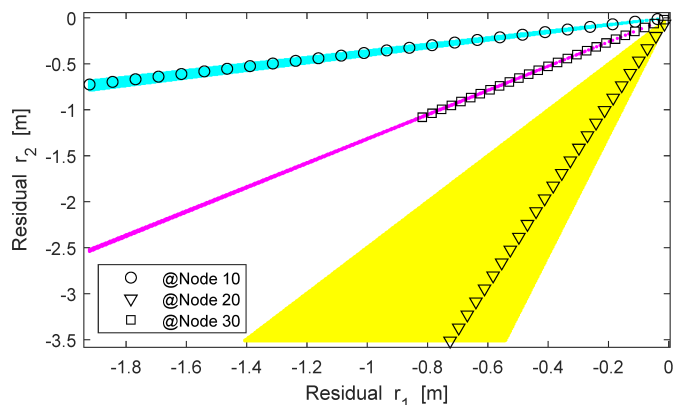
$$\text{SNR} = 20 \log_{10} \left(\frac{\text{True pressure}}{\text{Pressure noise}} \right). \quad (8)$$

The results presented in Table 3 show that the percentage of improvement when using the cosenoidal features with respect to the residual features is higher than 50% for a wide noise margin. The rate of improvement decreases as the proportion of the noise in the signal increases, because when the noise is considerable (SNR less than 40 dB), the leak directions captured by the cosines become irrelevant to locate the leak. However, in the worst case, the improvement percentages remain close to zero, which means that the use of cosenoidal features is not decreasing performance.

To assess the overall performance of the classifiers in the leak location task, the ATD was calculated as defined in (7). The results are presented in Table 4 for noise-free measurements, confirming the best performance when using the cosenoidal features. With noisy measurements, the ATD increases (see Table 5), but its variation is consistent with the increase in the classification error.



(a) Using residual features



(b) Usign cosenoidal features

Fig. 5. Convergence regions of three classes (leaks in three different nodes) of a leak locator by decision tree. Each color corresponds to predictions of different classes. The geometric figure symbols show the data of each class used for training.

Table 3. Improvement index in classifier performance when using cosenoidal features for different noise magnitude in measurements, according to (6).

Classification Method	Signal/Noise Ratio (SNR)					
	∞^\dagger	80 dB	60 dB	40 dB	20 dB	0 dB ‡
k -Nearest Neighbors	0.993	0.873	0.394	0.006	0.013	0.005
Naïve Bayes	0.997	0.939	0.405	0.009	-0.003	0.003
Decision Tree	0.996	0.917	0.473	0.021	-0.018	-0.007
Linear discriminant	0.970	0.929	0.550	0.016	-0.040	-0.018

† Noise-free measurements.

‡ Noise and measurements of the same magnitude.

Table 4. Average topological distance on noise-free testing data using different feature sets.

Classification Method	Features	
	Residuals	Cosines
k -Nearest Neighbors †	0.6555	0.0026
Naïve Bayes	2.4090	0.0026
Decision Tree	1.3239	0.0026
Linear discriminant	2.1974	0.0232

† Using Euclidean distance with $k = 5$.

Table 5. Average topological distance on noisy testing data, SNR = 60 dB, using different feature sets.

Classification Method	Features	
	Residuals	Cosines
k -Nearest Neighbors †	0.9523	0.5948
Naïve Bayes	2.4387	0.9742
Decision Tree	1.3935	0.7303
Linear discriminant	2.2090	0.7432

† Using Euclidean distance with $k = 5$.

4. CONCLUSION

The use of cosenoidal features, instead of raw residuals, showed a good performance in leak location using different classifiers, considering accuracy and robustness. The non-linear transformation to obtain the direction cosines implies a low computational cost but significantly improves performance because it reduces most of the localization error in the four machine learning techniques tested. In future work, it is expected to test this methodology using physical measurements in water distribution networks with a greater number of nodes. In addition, considering that the current approach requires a well-calibrated hydraulic network model to obtain training data, in the future, it is intended to develop methodologies less dependent on the model in the search for leak localization techniques completely based on data.

REFERENCES

- Blesa, J., Nejjari, F., and Sarrate, R. (2016). Robust sensor placement for leak location: analysis and design. *Journal of Hydroinformatics*, 18(1), 136–148.
- Blesa, J. and Pérez, R. (2018). Modelling uncertainty for leak localization in water networks. *IFAC-PapersOnLine*, 51(24), 730 – 735.
- Casillas, M.V., Garza-Castañón, L.E., and Puig, V. (2013). Model-based leak detection and location in water distribution networks considering an extended-horizon analysis of pressure sensitivities. *Journal of Hydroinformatics*.
- Casillas, M., Garza-Castañón, L., and Puig, V. (2015). Optimal sensor placement for leak location in water distribution networks using evolutionary algorithms. *Water*, 7(11), 6496–6515.
- Cugueró-Escofet, P., Blesa, J., Pérez, R., Cugueró-Escofet, M.A., and Sanz, G. (2015). Assessment of a leak localization algorithm in water networks under demand uncertainty. *IFAC-PapersOnLine*, 48(21), 226 – 231.
- Eliades, D.G., Kyriakou, M., Vrachimis, S., and Polycarpou, M.M. (2016). EPANET-MATLAB Toolkit: An Open-Source Software for Interfacing EPANET with MATLAB. In *Proc. 14th International Conference on Computing and Control for the Water Industry (CCWI)*, 8. The Netherlands. doi:10.5281/zenodo.831493.
- Ferrandez-Gamot, L., Busson, P., Blesa, J., Tornil-Sin, S., Puig, V., Duviella, E., and Soldevila, A. (2015). Leak localization in water distribution networks using pressure residuals and classifiers. *IFAC-PapersOnLine*, 48(21), 220–225.
- Fujiwara, O. and Khang, D.B. (1990). A two-phase decomposition method for optimal design of looped

- water distribution networks. *Water resources research*, 26(4), 539–549.
- Javadiha, M., Blesa, J., Soldevila, A., and Puig, V. (2019). Leak localization in water distribution networks using deep learning. In *2019 6th International Conference on Control, Decision and Information Technologies (CoDIT)*, 1426–1431.
- Mashford, J., de Silva, D., Marney, D., and Burn, S. (2009). An approach to leak detection in pipe networks using analysis of monitored pressure values by support vector machine. In *Third International Conference on Network and System Security*, 534–539.
- OECD (2016). Water Governance in Cities. *OECD Studies on Water*.
- Pérez, R., Puig, V., Pascual, J., Quevedo, J., Landeros, E., and Peralta, A. (2011). Methodology for leakage isolation using pressure sensitivity analysis in water distribution networks. *Control Engineering Practice*, 19(10), 1157–1167.
- Pérez, R., Sanz, G., Puig, V., Quevedo, J., Nejjari, F., Meseguer, J., Cembrano, G., Mirats, J.J., and Sarrate, R. (2014). Leak localization in water networks. *IEEE Control Systems Magazine*, August, 24–36.
- Pérez, R., Cugueró, J., Sanz, G., Cugueró, M.A., and Blesa, J. (2017). Leak monitoring. In V. Puig, C. Ocampo-Martínez, R. Pérez, G. Cembrano, J. Quevedo, and T. Escobet (eds.), *Real-time Monitoring and Operational Control of Drinking-Water Systems*, 115–130. Springer International Publishing, Cham.
- Puig, V., Ocampo-Martínez, C., Pérez, R., Cembrano, G., Quevedo, J., and Escobet, T. (eds.) (2017). *Real-time Monitoring and Operational Control of Drinking-Water Systems*. Springer International Publishing. doi: 10.1007/978-3-319-50751-4.
- Quiñones-Grueiro, M., de Lázaro, J.M.B., Verde, C., Prieto-Moreno, A., and Llanes-Santiago, O. (2018). Comparison of classifiers for leak location in water distribution networks. *IFAC-PapersOnLine*, 51(24), 407 – 413.
- Romero-Tapia, G., Fuente, M., and Puig, V. (2018). Leak localization in water distribution networks using fisher discriminant analysis. *IFAC-PapersOnLine*, 51(24), 929–934.
- Rossman, L.A. (2000). EPANET 2: Users manual. Technical Report EPA/600/R-00/057, US Environmental Protection Agency.
- Soldevila, A., Blesa, J., Tornil-Sin, S., Duviella, E., Fernandez-Canti, R.M., and Puig, V. (2016). Leak localization in water distribution networks using a mixed model-based/data-driven approach. *Control Engineering Practice*, 55, 162–173.
- Soldevila, A., Fernandez-Canti, R.M., Blesa, J., Tornil-Sin, S., and Puig, V. (2017). Leak localization in water distribution networks using bayesian classifiers. *Journal of Process Control*, 55, 1–9.
- Vegas Niño, O.T., Martínez Alzamora, F., Alonso Campos, J.C., and Tzatchkov, V. (2018). Using the EPANET Toolkit v2. 00.12 with different programming environments.
- Young, E.C. (2017). *Vector and tensor analysis*. CRC Press.