



Provided by the author(s) and University of Galway in accordance with publisher policies. Please cite the published version when available.

Title	A statistically based fault detection and diagnosis approach for non-residential building water distribution systems
Author(s)	Hashim, Hafiz; Ryan, Paraic; Clifford, Eoghan
Publication Date	2020-11-01
Publication Information	Hashim, Hafiz, Ryan, Paraic, & Clifford, Eoghan. (2020). A statistically based fault detection and diagnosis approach for non-residential building water distribution systems. <i>Advanced Engineering Informatics</i> , 46, 101187. doi: https://doi.org/10.1016/j.aei.2020.101187
Publisher	Elsevier
Link to publisher's version	https://doi.org/10.1016/j.aei.2020.101187
Item record	http://hdl.handle.net/10379/16532
DOI	http://dx.doi.org/10.1016/j.aei.2020.101187

Downloaded 2024-04-24T15:31:29Z

Some rights reserved. For more information, please see the item record link above.



1 Article

2 **A statistically based fault detection and diagnosis approach for non-residential building water**
3 **distribution systems.**

4

5 **Hafiz Hashim^{1*}, Paraic Ryan^{2,3} and Eoghan Clifford⁴**

6 ¹School of Engineering, National University of Ireland, Galway, Ireland.

7 h.hashim1@nuigalway.ie (H.H.), paraic.ryan@ucc.ie (P.R.), eoghan.clifford@nuigalway.ie (E.C.).

8 ²Discipline of Civil, Structural and Environmental Engineering, School of Engineering, University College Cork, Ireland.

9 ³MaREI Centre, Environmental Research Institute, University College Cork, Cork, Ireland.

10 ⁴Informatics Research Unit for Sustainable Engineering, National University of Ireland, Galway, Ireland.

11 *Correspondence: h.hashim1@nuigalway.ie; Tel: 00 353 91 492219.

12 **Abstract:** Large non-residential buildings can contain complex and often inefficient water distribution systems.
13 As requirements for water increase due to water scarcity and industrialization, it has become increasingly
14 important to effectively detect and diagnose faults in water distribution systems in large buildings. In many cases,
15 if water supply is not impacted, faults in water distribution systems can go unnoticed. This can lead to unnecessary
16 increases in water usage and associated energy due to pumping, treating, and heating water. The majority of fault
17 detection and diagnosis studies in the water sector are limited to municipal water supply and leakage detection.
18 The application of detection and diagnosis for faults in building water networks remains largely unexplored and
19 the ability to identify and distinguish between routine and non-routine water usage at this scale remains a
20 challenge. This study, using case-study data, presents the application of principal component analysis and a multi-
21 class support vector machine to detect and classify faults for non-residential building water networks. In the
22 absence of a process model (which is typical for such water distribution systems), principal component analysis
23 is proposed as a data-driven fault detection technique for building water distribution systems for the first time
24 herein. Hotelling T²-statistics and Q-statistics were employed to detect abnormality within incoming data, and a
25 multi-class support vector machine was trained for fault classification. Despite the relatively limited training data
26 available from the case-study (which would reflect the situation in many buildings), meaningful faults were
27 detected, and the technique proved successful in discriminating between various types of faults in the water
28 distribution system. The effectiveness of the proposed approach is compared to a univariate threshold technique
29 by comparison of their respective performance in the detection of faults that occurred in the case-study site. The
30 results demonstrate the promising capabilities of the proposed fault detection and diagnosis approach. Such a
31 strategy could provide a robust methodology that can be applied to buildings to reduce inefficient water use,
32 reducing their life-cycle carbon footprint.

33 **Highlights:**

- 34
- 35 • Principal component analysis and support vector machines were combined to build a fault detection and
36 diagnosis method for non-residential building water networks.
 - 37 • Data from a non-residential building was used to validate the developed fault detection and diagnosis
38 model.
 - 39 • Improved performance monitoring of non-residential water distribution systems can be achieved using
40 limited process data.
 - 41 • A comparative study illustrated that the proposed approach detected and diagnosed more actual system
42 alarms than a conventional approach, while also utilizing less computational resources.

43 **Keywords:** Water distribution system, Non-routine events, Fault detection and diagnosis, Principal component
44 analysis (PCA), Support vector machine (SVM), Performance monitoring.

45 **1. Introduction:**

46 In Europe, the non-residential sector consumes about 28% of total water withdrawn (European Commission,
47 2012). Thus, the efficient use of this water in such buildings is imperative. In certain areas of Europe, overall
48 losses in water networks due to leakage are as high as 50% (Nowicki et al., 2012). For example, 41% of treated
49 water is lost through damaged pipes in Ireland (Irish Water, 2019). It has been reported that water loss in buildings

50 is primarily linked to leakage within the water distribution systems (Datta & Sarkar, 2016; Skworcow et al., 2013)
51 - in a recent study Sydney Water estimated that building leaks account for 28% of total water consumption
52 (Sydney Water, 2011). It is thus vital to reduce the wastage of water through leakage and improve the performance
53 of such water distribution systems (and thus reducing associated energy usage due to water pumping, treatment,
54 etc.) This need has resulted in a growing focus on optimising water (and associated energy) usage across all
55 building types using building management systems to detect problems (Sousa et al., 2019), particularly given the
56 water scarcities impacting many countries today (Connor & Murphy, 2017).

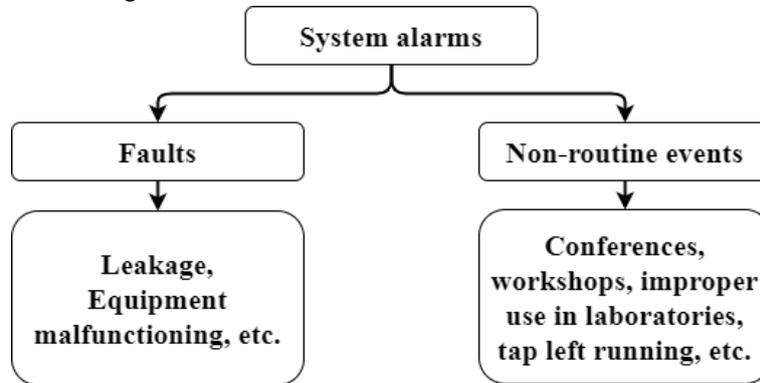
57 In many sectors, fault detection and diagnosis systems have been developed to optimise system operation and to
58 detect and diagnose system abnormalities. Examples include heating ventilation and air conditioning (HVAC)
59 systems, water and wastewater treatment processes, automotive industries, chemical and petrochemical processes,
60 wind farms and some manufacturing facilities (Hu et al., 2019; Xu et al., 2017; Zhang et al., 2017; Naderi &
61 Khorasani, 2016; Zhou et al., 2015; Yu et al., 2014; Bruton et al., 2014; Zhang et al., 2017). For example, the
62 analysis of energy consumption in buildings as a key driver of fault detection and diagnosis has been extensively
63 studied (Hu et al., 2019; Burak Gunay et al., 2019; Agostino et al., 2017; Balaras et al., 2017; Stavset & Kauko,
64 2015). Previous studies have tended to describe one fault detection and diagnosis method at a time (either principal
65 component analysis or support vector machines), often formulated as an optimization task. For instance, principal
66 component analysis was used to compress the data or to identify the anomalies using detection indices, whereas
67 the goal of using machine learning models (such as support vector machines) was to minimise the difference
68 between the measurements taken on the system and predicted values from physics-based models (Grueiro et al.,
69 2018; Escofet et al., 2016; Nasir et al., 2014). A significant challenge impacting the performance of these
70 approaches is the time varying characteristics of real time processes (such as change in mean, variance over time,
71 process noise, etc.) which can be a particular feature of water consumption.

72 Fault detection and diagnosis studies in the water sector have focused mainly on leak detection at a municipal
73 water supply level (Seyoum et al., 2017). A range of detection approaches (such as leak-noise correlation, acoustic
74 sensing, water balancing, pressure management, etc.) have been studied for these municipal systems (Robles et
75 al., 2016; Escofet et al., 2016; Pérez et al., 2015; Makaya & Hensel, 2015). Although these approaches are
76 considered the most accurate for leak detection, they are not suitable when it comes to large-scale building water
77 distribution systems due to being relatively high cost, labour-intensive and time consuming. Other approaches
78 have relied on hydraulic modelling techniques (such as transient analysis, pressure residual vector method, etc.)
79 which have been validated using modelled or synthetic data. These approaches have focused on comparisons of
80 measurements with predictions obtained from hydraulic models (Moors et al., 2018; Abdulshaheed et al., 2017;
81 Alsaydalani, 2017; Soldevila et al., 2017; Perfido et al., 2016; Moser et al., 2015; Sedki & Ouazar, 2012; Mashford
82 et al., 2012; Salam et al., 2015). Creating such models is often difficult and expensive due to the presence of non-
83 stationarity and uncertainties in complex non-residential water distribution systems.

84 Overall, the implementation of fault detection and diagnosis for water distribution systems has received limited
85 attention (Stetco et al., 2019; Geng et al., 2018). This is particularly the case for non-residential buildings
86 (Adeyeye, 2014; Cosgrove et al., 2015) where it is acknowledged significant challenges remain in developing
87 detection and diagnosis system for faults in building water distribution systems (Liu et al., 2016; Nezhad et al.,
88 2014). Furthermore, fault detection and diagnosis methods must be cognisant with challenges associated with real
89 data that comes from non-residential building water distribution systems (e.g. to distinguish between routine and
90 non-routine water usage, etc.).

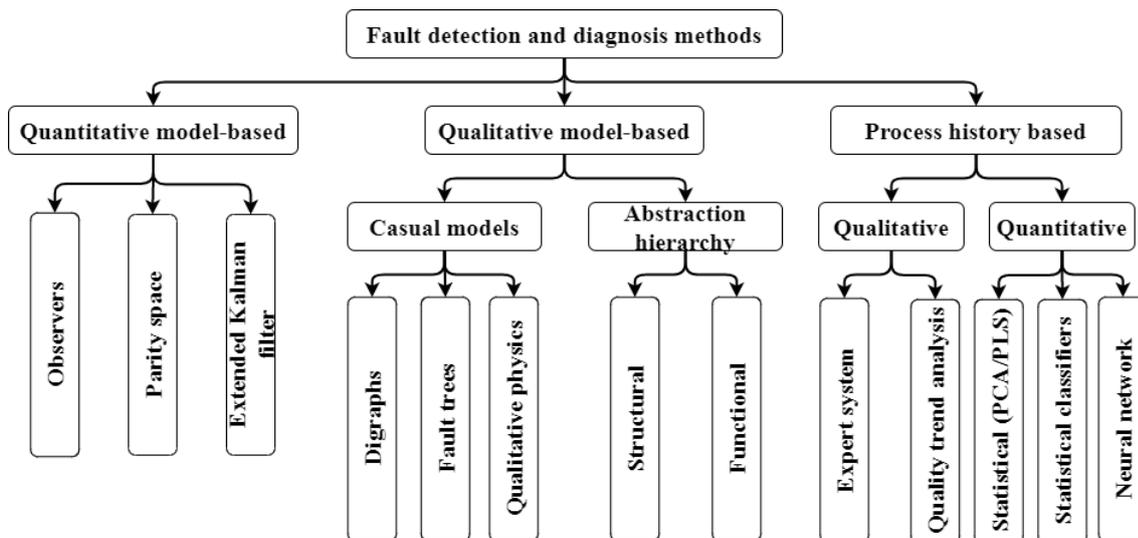
91 Currently, most building managers react to faults on an ad-hoc basis, responding to obvious faults and repairing
92 infrastructure as required (Moser et al., 2015). This can result in low-level imperceptible faults persisting within
93 a water distribution system, compromising water distribution system efficiency. The application of existing
94 approaches developed for municipal water supply systems may be limited within non-residential buildings as key
95 faults not only involve leaks but also comprise system related faults (such as equipment malfunctioning,
96 operational errors, etc.) or non-routine events caused by inefficient water usage and infrequent sudden changes in
97 demand (such as hosting a one-off event or increased production for a short period of time, taps left running, etc.
98 - Figure 1). Furthermore, it is advantageous within buildings to distinguish between faults, inefficient water usage
99 and unusual consumption which may or may not be the result of a fault.

100 Building water distribution system studies have not leveraged the more sophisticated fault detection and diagnosis
 101 approaches used in other sectors such as principal component analysis or support vector machines (Gharsellaoui
 102 et al., 2020; Zhao et al., 2019). Thus, the development and validation of robust systems that can detect system
 103 alarms, and distinguish between faults and non-routine events in non-residential water distribution systems has
 104 the potential to significantly increase the efficiency of water consumption (Pelz, 2003; Danacova et al., 2016).
 105 The present work aims to fill these gaps by developing decision models based on a real site data that can detect
 106 faults in non-residential buildings.



107
 108 Figure 1: System alarm characterization.

109 In general, fault detection and diagnosis methods can be categorized mainly as: signal-based detection, and signal
 110 and multiple model-based detection (Figure 2) (Pelz, 2003; Venkatasubramanian et al., 2003a;
 111 Venkatasubramanian et al., 2003b).



112
 113 Figure 2: Fault detection and diagnosis methods. Adapted from (Venkatasubramanian et al., 2003; Pelz, 2003).

114 Basic fault detection systems comprise alarms which trigger when relatively high levels of water are being used
 115 (Mulligan et al., 2020; Clifford et al., 2018; Quevedo et al., 2014; Perfido et al., 2016) and provide high level
 116 statistics on water consumption. Typically, in such systems when water consumption exceeds a predefined
 117 threshold, an alarm is sent automatically. While the alarm may not be designed to detect excessive water
 118 consumption due to improper use (e.g. a tap left running) or other non-routine events, it can be adapted to do so.
 119 However, such systems are unable to identify complex patterns of water usage - such as distinguishing between
 120 actual faults or non-routine events. Additionally, such systems are embedded with alarm notification features but
 121 still require assistance from experts to choose thresholds for faults or non-routine events. More effective fault
 122 detection and diagnosis approaches must be capable of detecting inevitable deviations in water time series.
 123 Furthermore, non-residential water distribution system can feature non-stationarity statistical properties over time
 124 (e.g. water consumption can vary between seasons and depend on working hours, holiday periods, etc.) which can

125 make it challenging to detect faults (Prabuchandran et al., 2019). It can also be noted (e.g. Table 1) that a
126 significant gap exists in published literature relating to fault detection and diagnosis methods that utilize real case-
127 study data. Existing studies in the literature tend to focus on the use of generated or experimental data sets (Cody
128 & Narasimhan, 2020; Grueiro et al., 2018; Moser et al., 2015; Villegas et al., 2010). The use of idealised data
129 means the effect of random variations and outliers, which undoubtedly occur in industry are not taken into account
130 in the assessment of proposed methodologies. This can lead to an over-estimation of the effectiveness of proposed
131 approaches if applied to real case-studies.

132 Principal component analysis is a versatile and well-known technique for unsupervised data analysis, which
133 projects data in a reduced space defined by orthogonal principal components (Izem et al., 2018; Wei Li et al.,
134 2018). It has been adopted in various disciplines, such as in damage detection of civil structures, data compression,
135 face recognition, image analysis, visualization as well as in fault detection in petroleum/chemical industries (Laory
136 et al., 2011; Saitta et al., 2005; Posenato et al., 2008; Harrou et al., 2013; Abdi & Williams, 2010; Li et al., 2000).
137 Principal component analysis (PCA) is a vector space transformation predominately used to reduce data
138 dimensionality by extracting correlation between variables in sets of independent variables, that explain the trend
139 of the process while optimising the variance of the original data in reduced number of dimensions (Garcia-Alvarez,
140 2014; Laory et al., 2011; Villegas et al., 2010). A PCA model leverages process history data and extracts a linear
141 combination of variables explaining the major trend of the process being analysed (Ait-Izem et al., 2018).

142 To measure the variation of samples within the PCA model and detect the abnormality of the new incoming data,
143 Hotelling T²-statistics and Q-statistics (squared prediction error) have been applied in other engineering domains
144 (Li et al., 2018; Horrigan et al., 2018; Zhang et al., 2017; Villegas et al., 2010). The T²-statistics measure the
145 variability of a score matrix and can detect abnormality within new data by comparing it to variation in the
146 parameters defined by baseline condition. Q-statistics evaluate the variability of a matrix which is the projection
147 of the original data onto a residual subspace (Benaicha et al., 2013; Ahmed et al., 2012; Zumoffen & Basualdo,
148 2008). PCA based methods have advantages over classical statistical methods (e.g. Shewhart chart, cumulative
149 sum,, exponentially weighted moving average) which while accurate in many cases, have been shown to produce
150 inaccurate results when introduced to multivariate systems (Maione et al., 2019; Xiao et al., 2017). Other
151 challenges associated with the application of classical statistical approaches to building water distribution systems
152 include the considerable variation in water consumption levels in buildings, and the possible mismatches between
153 the proposed statistical model and the process being monitored (Chen, 2010).

154 In recent times, the use of support vector machines learning has increased relative to other methods (e.g. artificial
155 neural network) when modelling the relationship between the input and the output in classification and regression
156 problems in various engineering sectors (Xu et al., 2020; Kouziokas, 2020; Gautam et al., 2020; Cody, 2020; Akil
157 et al., 2019; Grueiro et al., 2018; Liu et al., 2019; Samer et al., 2018). This is due to SVM's strong fault diagnosis
158 ability in multi-class classification problems, high training efficiency, and its effectiveness when working with
159 low-dimensional (small sample dataset) process data (Xu et al., 2017). In contrast, artificial neural network rep-
160 resents a more challenging pattern learning problem with inconsistent or unpredictable sampling of datapoints
161 arising when working with smaller datasets (Kouziokas, 2020; Gautam et al., 2020). In the context of water stud-
162 ies, comparative analysis of different machine learning approaches (such as support vector machines, artificial
163 neural network, nearest neighbours, etc.) has shown the robustness and improved performance of support vector
164 machines when compared to other approaches in identifying municipal water network leaks in the presence of
165 different uncertain conditions such as non-stationarity of water use, sensor noise, etc. (Grueiro et al., 2018; Cody,
166 2020). The major drawback of using a machine learning approach (e.g. support vector machines) is that the trained
167 classifier will only be able to deal with fault alarms that have been experienced by the water distribution system
168 previously. In practical applications, the available data can be limited to a subset of the potential fault alarms that
169 the system has experienced. However, that could be improved by utilizing knowledge of experienced faults or
170 anomalous situations to update the model and provide a more reliable performance monitoring system. As the
171 number of faults or anomalous situations in the fault detection and diagnosis model database increases, the system
172 becomes more valuable and reliable. This study utilizes support vector machines to develop fault detection and
173 diagnosis approaches for non-residential water distribution system for the first time in the literature. To this end,
174 principal component analysis (PCA) and support vector machines (SVM) have not been applied to water distri-
175 bution systems, the proposed methodology not only to detect and diagnose faults in the water distribution system,

176 but also involve robust means of identifying outliers within the water time series lead to an advancement in using
177 these computational approaches without explicit knowledge of the system.

178 Table 1 summarises relevant studies that describe the efficacy and applicability of the principal component
179 analysis and machine learning tools in various engineering domains (municipal water network, nuclear power
180 plants, rotary machines, etc.). The authors did not find any studies which applied principal component analysis
181 and support vector machines methods to building water networks. Previously, literature has focused on stationary
182 systems (and often utilized modelled data). The use of more advanced statistical techniques (e.g. principal
183 component analysis) in conjunction with support vector machines has not been widely explored in water sector.
184 Given that water distribution systems can often exhibit non-stationarity (change in statistical properties over time),
185 the combination of principal component analysis and support vector machines techniques has potential to offer
186 robust applications in water distribution system in non-residential buildings.

187 Table 1: Literature review.

References	Objective	Key points across studies
PCA-based studies		
Horrigan et al., 2018	Statistical fault detection scheme developed by utilizing commercial building energy performance data	Proposed statistical-based performance prediction model utilizing univariate data time series. Actual fault scenarios were not studied to validate the fault detection method.
Li et al., 2018	Sensor fault detection in a nuclear power plant using principal component analysis method	Fault detection methodology developed based on principal component analysis and detection indices (T^2 and Q - statistics).
Zhang et al., 2017	Assessment of T^2 and Q -statistics for detecting additive and multiplicative faults in multivariate statistical process monitoring	Statistical method (principal component analysis) impact on detecting additive and multiplicative faults have been studied in context of fault detection rate.
Ahmed et al., 2012	Fault detection and diagnosis using principal component analysis of vibration data from reciprocating compressor	T^2 - statistics performs better than Q - statistics in detecting additive and multiplicative faults. Several diagnostic features extracted, are analysed by PCA to detect leakage faults in a system.
Villegas et al., 2010	Principal component analysis for fault detection and diagnosis, experience with a pilot plant (two-tanks system)	The PCA model shows good accuracy in providing detected fault information. PCA methodology developed to identify changes across operating point, assuring model robustness for system non-linearity's
Grueiro et al., 2018	PCA-based leak detection and localization in water distribution network	Use of measured data using district metering area (DMA) technique under normal operating conditions. The PCA model shows good accuracy with an average detection rate of over 80%.
Laory et al., 2011	PCA-based damage detection during continuous static monitoring of civil structures	Use of measured sensor data from linear structures.
Posenato et al., 2006		Two conditions, damage detectability and time to detection, are evaluated with respect to changes in sensor-damage location, dynamic loading, and damage level in structures.
SVM-based studies		
Escofet et al., 2016	Model vs. data based approaches applied to fault diagnosis in potable water supply networks	SVM based methodology developed for detecting leakage and to diagnose anomalies in a system. Use of simulated data rather than measured data.
Muralidharan et al., 2014	Fault diagnosis of monoblock centrifugal pump using support vector machines	Support vector machines shows superiority over artificial neural network in terms of performance, robustness in both normal and noisy conditions.
Nasir et al., 2014	Measurement error sensitivity analysis for detecting and locating leaks in pipelines using artificial neural network and support vector machines	Support vector machines exhibit insensitiveness in detecting leaks in a noisy environment as compared to artificial neural network.
Saberi et al., 2011	Comparing performance and robustness of SVM and artificial neural network for fault diagnosis in a centrifugal pump	Measurement/modelling uncertainties were not studied.
Qu et al., 2010	SVM-based pipeline leakage detection and pre-warning system	Signal decomposition techniques along with support vector machines to extract features for classification.
Mashford et al., 2009	An approach to leak detection in pipe networks using analysis of monitored pressure values by support vector machines	Classification accuracies of discrete wavelet at different levels were calculated and compared. Feature effect on individual fault found missing.
Cody, 2020;	Compared performance and robustness of different classifiers for detecting leaks in water distribution networks	Support vector machines, artificial neural network, nearest neighbours, and Bayes classifiers were used.
Grueiro et al., 2018		Use of simulated data rather than measured data. SVM shows advantage over other classifiers in terms of performance, robustness under different uncertain conditions.
Samer et al., 2018	SVM-based single event leak detection in pressurized water network pipelines	Use of experimental data at lab-scale setup. Experiments were conducted using one inch and two-inch straight pipelines of materials ductile iron and polyvinyl chloride (PVC). Support vector machines performed with an accuracy of 98% in identifying leaks.
Akil et al., 2019	SVM-based smart monitoring of an institutional building using gas, electricity, and water use data	Use of measured water, electricity, and gas usage data. Support vector machines shows effectiveness in identifying atypical behaviour of an institutional building (such as increased electricity, water use, etc.).
Liu et al., 2019	Water pipeline leak detection based on SVM and wireless sensor network	Use of simulated and experimental data to verify the effectiveness of the proposed approach. Integrated conditional probabilities and time-frequency analysis to identify leak signals within the network.
Gautam et al., 2020	SVM-based monitoring and forecasting for leak detection in household water storage	Use of measured water use data. Compared actual values against predicted values to identify leaks.

189 This work leverages a case-study building with real-time water consumption data to illustrate the effectiveness of
190 using PCA detection indices (Hotelling T^2 -statistics and Q-statistics) in identifying faults in a non-residential
191 building water network. The proposed approach applies the PCA technique to measured time-series water data to
192 (i) identify faults and (ii) classifying the faults into various categories, represented by prediction model output
193 data, to reset and guide operational performance. The key contribution of the proposed methodology is an
194 improvement in the performance monitoring of water distribution system through the identification of faults in
195 early stage for building managers under realistic conditions. Moreover, by utilizing multi-class SVM classification
196 the faults can be classified quickly and effectively with a relatively small amount of training data. This can enable
197 building operators to target predictive maintenance and reduce the occurrence of false alarms, while also
198 developing strategies to reduce unnecessary or inefficient consumption of water.

199 The proposed fault detection and diagnosis methodology involves analysing water time series and monitoring
200 inevitable changes in a system without the use of geometrical information. In the context of performance
201 monitoring, to differentiate between features of routine and non-routine or anomalous water uses and to limit false
202 detection. Principal component analysis is based on the covariance structure of the data in the PCA space and is
203 sensitive to outlier measurements (Rousseeuw et al., 2018; Hubert et al., 2016). In the presence of outliers,
204 principal components are often attracted towards outlying data points and may not be able to capture the realistic
205 behaviour of routine variation in water distribution system. This can result in unreliable or false fault detection
206 during the process. To address this in the current study, the methodology combines principal component analysis
207 with a distance-distance approach for outlier localization (as discussed in Section 4.1) creating a robust approach
208 where principal components are not influenced by outliers. The proposed methodology is also capable of
209 considering biased uncertainties and non-stationarity of water use which are typical in water distribution system
210 modelling challenges and is in practice seen as state of things rather than process disturbances. Furthermore, in
211 this study the proposed fault detection and diagnosis approach is assessed using real site water distribution system
212 data, with associated outliers and process noise leading to a more robust assessment of the machine learning
213 techniques suitability to detect and diagnose faults in water distribution system.

214 The objective of the proposed methodology is to obtain a simpler equivalent approach than physics-based models,
215 to reduce the complexities linked with the hydraulic modelling of the non-residential water distribution systems.
216 This proposed approach is designed to be effective in identifying non-routine events and faults of different types
217 (such as equipment malfunctioning, high or low-level imperceptible process faults) not only leakage. Moreover
218 data availability for non-residential building water distribution systems represents a significant and persistent
219 challenge (Matos et al., 2020). In many such cases reliable and objective information about water use and network
220 performance is poor, lacking or otherwise unavailable (Houngbo, 2019). The dataset used in this study is typical
221 of what is available in non-residential buildings, and thus the aim was to develop a fault detection and diagnosis
222 methodology which is capable of operating under these practical data constraints and assessing the effectiveness
223 of a proposed methodology under these practical data constraints.

224 This paper is arranged as follows: In section 2 materials and methods are presented which include the case study
225 and pumping system details, measurable characteristics, followed by an introduction to principal component
226 analysis, fault detection indices (T^2 and Q-statistics), and fault detection and diagnosis methodology is outlined
227 in Section 3. Later results and discussion are presented in Section 4. Section 5 summarizes the study and provides
228 conclusions and future directions.

229 **2. Methods and materials:**

230 *2.1. Case Study:*

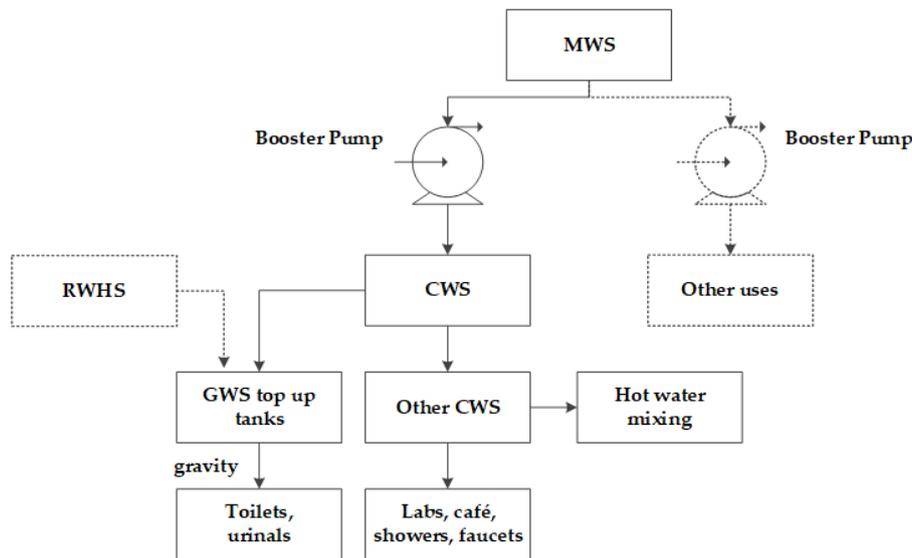
231 The case-study site was the Alice Perry Engineering Building at the National University of Ireland, Galway (NUI
232 Galway). The building is located on the west coast of Ireland in an area with a temperate marine climate. The
233 building houses approximately 1,100 students and 100 staff during teaching and exam terms (approximately 26-
234 28 weeks a year in total) and about 150 research, academic, administrative, and technical staff and 50-100
235 postgraduates during the rest of the year. The building includes lecture halls, classrooms, offices, laboratory
236 facilities, a café, and shower and toilet facilities spread across 14,000 m² of floor space on four storeys. Thus, it
237 has a variety of end-uses of its water and hot water systems. The building is managed through a building

238 management system that collects data on building performance and operational efficiency - including 11 water
 239 meters and a number of energy meters. Some of the key water uses include showers and hand wash basins, grey
 240 water from a rainwater-harvesting system for toilets and urinals and potable water for the water fountains and the
 241 café. This study deals primarily with mains water usage within the building (i.e. hot water system is not
 242 considered).

243

244 2.2. *Building Pumping System:*

245 The mains water system (MWS) in the building is divided into a cold-water system (CWS) and a potable water
 246 supply for drinking fountains and for the café. There are two sets of booster pumps (Figure 3), that deliver (i) the
 247 cold-water system and (ii) potable water supply. The cold-water system water is then divided into water for grey-
 248 water system (GWS) applications and water for other uses such as laboratories, hot water mixing in showers and
 249 faucets and non-potable uses in the café. The grey-water system supplies water to toilets and urinal flushing as
 250 required by gravity. The building has a large rainwater-harvesting system (RWHS) for supplying grey water.
 251 When the rainwater-harvesting system cannot supply sufficient grey water, two holding tanks located on the roof
 252 of the building (8m³ capacity each) are topped up by the cold-water system. During the study period, the building
 253 rainwater-harvesting system utilised the cold-water system even during periods when rainwater was available due
 254 to system faults that remained undetected prior to this study.



255

256 Figure 3: Simplified schematic of water network in the building. The dashed line indicates systems not
 257 considered in this study.

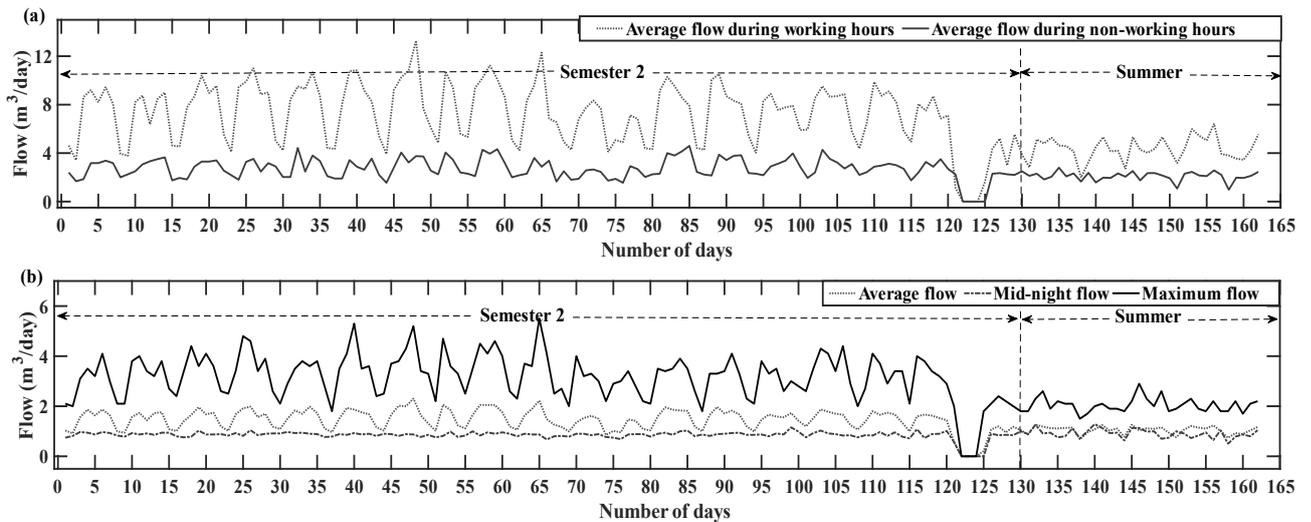
258 2.3. *Feature extraction:*

259 To perform fault detection and diagnosis in water distribution systems, measurable characteristics can play a vital
 260 role in discriminating between routine and abnormal events in the water network. The characteristics of the
 261 system, derived from meter readings, can be used to analyse the water usage baseline and diurnal patterns (or
 262 indeed patterns at any chosen time interval for which adequate data is available). For instance, analysing night
 263 time flow can assist in detecting continuous flow which can be due to leaks; analysing peak flows can assist in
 264 detecting exceptional events. Analysing water consumption at specific daily intervals (working and non-working
 265 hours) can help in identifying disruptions to normal routines. The details of the features utilized in this study are
 266 listed in Table 2. The daily flow and maximum flow features summarize the water usage over a 24-hour period
 267 from midnight to midnight. Whereas the remaining features (working hours, non-working hours, and night time
 268 flow) focus on water consumption for particular periods of the day which are linked to how the building is used.
 269 In this case, these periods comprise of four working hour time intervals, two non-working hour time intervals and
 270 a night time interval. Medium resolution data recorded at 7.5-minute intervals was collected in the case-study
 271 building from an inline displacement water meter equipped with a magnetic pulse output and analysed in
 272 developing the methodology for fault detection and diagnosis.

273 Table 2: Features (measured in m³/hr) used to characterise daily water demand.

	Feature	Description
Total	Daily flow	Average flow in 24 h
	Maximum flow	Highest flow in 24 h
Time of day	Working hours	Average flow between 6 a.m. – 9 a.m.
		Average flow between 9 a.m. – 12 p.m.
		Average flow between 12 p.m. – 3 p.m.
	Non-working hours	Average flow between 3 p.m. – 6 p.m.
		Average flow between 6 p.m. – 9 p.m.
		Average flow between 9 p.m. – 12 a.m.
Night time flow	Average flow between 12 a.m. – 6 a.m.	

274 For a given year the case-study monitoring period data was divided into Semester 1 - September to December,
 275 Semester 2 - January to May and a summer period (June-August). In general, water consumption would be higher
 276 in Semesters 1 and 2 than during the summer periods as undergraduate students would not be present during the
 277 summer periods. Figure 4a demonstrates the flow pattern at different intervals (Table 2) during Semester 2 and
 278 the summer period under the routine operating condition considered for PCA model training. Limited variation
 279 can be observed between working and non-working hour water usage in the summer due to reduced building
 280 occupancy as compared to Semester periods. In Figure 4b water consumption patterns are visualized during the
 281 summer period. The flow pattern during night time, when the building was not in use, is primarily due to the
 282 routine urinal flushing which operate irrespective of the building occupancy. These patterns derived from the
 283 water time series, can assist in analysing different faults in the water distribution system.



284

285 Figure 4: Flow characteristics within the periods of interest (a) Average 7.5 minute interval flow pattern during
 286 the 12 working hours and the 6 non-working hours, (b) Maximum, 7.5 minute interval average over 24 hours and
 287 night time flow pattern over 6 hours.

288 This study utilises multivariate statistical approaches (principal component analysis and support vector machines)
 289 to attempt to detect patterns in data, to develop a new fault detection and diagnosis approach for non-residential
 290 building water distribution systems. The theory behind the statistical methods used are outlined in the next section.
 291 It is noted that the primary advantages of the combined principal component analysis and support vector machines
 292 approaches are the ability to compress large data sets without losing data definition. In this case it can enable
 293 analysis of faults across a number of water usage features, both of which reduce the likelihood of false positives
 294 and false negatives within the fault detection and diagnosis system. In addition, support vector machine facilitates
 295 classification of faults, providing a greater level of information to building managers.

296 **3. Theory:**

297 *3.1. Principal Component Analysis:*

298 Computation of all feature values for water consumption, gives a data matrix of $X \in \mathbb{R}^{n \times m}$ whereby n rows
 299 represent the values for all features (Table 2) within a single day and m columns represent a single feature (such
 300 as maximum flow, average flow, etc.) over the observation with mean zero and unit variance.

301

$$X = \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1j} & \dots & x_{1m} \\ x_{21} & x_{22} & \dots & x_{2j} & \dots & x_{2m} \\ \dots & \dots & \dots & \dots & \dots & \dots \\ x_{i1} & x_{i2} & \dots & x_{ij} & \dots & x_{im} \\ \dots & \dots & \dots & \dots & \dots & \dots \\ x_{n1} & x_{n2} & \dots & x_{nj} & \dots & x_{nm} \end{bmatrix} = (w_1 \ w_2 \ \dots \ w_j \ \dots \ w_m) \quad (1)$$

302

303 where row vector X_i represents all feature measurements (i.e. daily flow, etc.) at specific time intervals and column
 304 vector w_j represents single feature measurements (average flow, night time flow, etc.) in a series of days. The data
 305 matrix $X \in \mathbb{R}^{n \times m}$ is then transformed by utilizing singular value decomposition as per (Jackson & Mudholkar,
 306 1979; Zenobi et al., 2011; Sheriff et al., 2017).

$$X = TP^T + E \quad (2)$$

307 where $T = [t_1, t_2, t_3, \dots, t_m] \in \mathbb{R}^{n \times m}$ is a matrix of transformed variables, where each column represents the score
 308 vectors and the i^{th} eigenvalue equals the square of the i^{th} singular value (i.e. $\lambda_i = \sigma_i^2$). $P = [p_1, p_2, p_3, \dots, p_m] \in$
 309 $\mathbb{R}^{m \times m}$ is the matrix of orthogonal vectors, where each column is populated by the eigenvectors associated with
 310 the covariance matrix of the data matrix $X \in \mathbb{R}^{n \times m}$. E is the residual matrix which ideally contain noise in the
 311 data. Typically, most of the data variance is contained in the principal components with larger eigenvalues, while
 312 the remaining principal components are considered as measurement noise which can be removed by reducing the
 313 data dimension (Grueiro et al., 2018; Laory et al., 2011). The covariance matrix S quantifies the amount of linear
 314 correlation between all possible combinations of features within the dataset and can be computed as (Sengupta &
 315 Kundu, 2016; Rosen, 2001; Jackson & Mudholkar, 1979):

$$S = \frac{1}{n-1} X^T X = P \Lambda P^T \quad (3)$$

316 where Λ is the diagonal matrix $\Lambda = \text{diag}(\lambda_1, \lambda_2, \lambda_3 \dots \lambda_m) \in \mathbb{R}^{m \times m}$ containing non-negative eigenvalues related to
 317 m principal components and magnitudes in a descending order $\lambda_1 \geq \lambda_2 \geq \dots \lambda_m \geq 0$. The model built by principal
 318 component analysis contains the same number of principal components as the original number of features in the
 319 data matrix $X \in \mathbb{R}^{n \times m}$. In the case of water distribution systems, it may contain features (such as average flow,
 320 average flow during working hours, average flow during non-working hours, etc.) that are highly correlated. In
 321 this study, the number of components is selected based on the cumulative percent variance (CPV) due to its
 322 computational simplicity and wide adoptability (Sheriff et al., 2017). It provides a good estimate of the number
 323 of principal components that need to be retained for most practical applications. Cumulative percent variance is a
 324 measure of the percentage variance $CPV(r) \geq 90\%$ captured by the first r principal components (Zumoffen &
 325 Basualdo, 2008; Johnson & Wichern, 2007):

$$CPV(r) = \frac{\sum_{i=1}^r \lambda_i}{\text{trace}(S)} 100 \quad (4)$$

326 *3.1.1. Score and orthogonal distance:*

327 To analyse the linear combination of water time series that restrain valuable information or outliers; the score
 328 distance SD_i and the orthogonal distance OD_i of each data point to the PCA subspace are given by (Hubert et al.,
 329 2005; Harris et al., 2014; Rousseeuw et al., 2018).

$$SD_i = \sqrt{\sum_{j=1}^r \frac{t_{ij}^2}{\lambda_j}} \quad (5)$$

330

$$OD_i = \|x_i - \hat{\mu} - P_{m,r} t_i'\| \quad (6)$$

331 where t_{ij} is the score value of each data point, λ_j is the non-negative eigenvalues and $P_{m,r}$ is the eigenvectors matrix
 332 related to r principal components, $\hat{\mu}$ is the mean of the covariance matrix. Assuming the data follows a multivariate
 333 normal distribution. The cut-off for the data points were estimated using Chi-square distribution $\sqrt{X_{r,0.975}^2}$. Where,
 334 $X_{r,0.975}^2$ is the 97.5% quantile of the Chi-squared distribution with r principal components (Rousseeuw & Hubert,
 335 2018).

336 3.2. *Fault detection indices:*

337 3.2.1. *Hotelling T²-statistics:*

338 T²-statistics represents the major variation in the data (Garcia-Alvarez, 2014). The T²-statistics of the i^{th} sample
 339 or observation x can be expressed by (Mujica et al., 2011; Qin, 2003) as,

$$T_i^2 = x_i^T P \Lambda_r^{-1} P^T x_i \quad (7)$$

340 where, Λ_r^{-1} is the diagonal matrix containing the eigenvalues related to retained principal components, x_i is the
 341 data vector of the i^{th} observation and P contains the loading vector associated with the r columns. Under routine
 342 process conditions, the data follow a multivariate normal distribution, the T²-statistics is related to an F -
 343 distribution considering that the population mean, and covariance are estimated from data (Qin, 2003).

344 The control limit T_α^2 can be obtained by F -distribution as follows.

$$T_\alpha^2 = \frac{r(n-1)}{n-r} F_{\alpha(r, n-r)} \quad (8)$$

345 where n is the number of observations in the data, r is the number of retained principal components, $F_{\alpha(r, n-r)}$ is the
 346 F -distribution. In the current study, the α -control limits are calculated at a confidence level of 95%. The T²-
 347 statistics can be interpreted as measuring the systematic variations of the process. Violation of routine conditions
 348 would indicate that the systematic variations are outside normal operating bounds, and thus the data set may
 349 indicate a fault or non-routine process condition (Vieira et al., 2014). The new observation is considered to be
 350 normal if it satisfies the following condition.

$$\begin{cases} T_i^2 < T_\alpha^2 & \text{--- Normal} \\ T_i^2 \geq T_\alpha^2 & \text{--- System alarm} \end{cases} \quad (9)$$

351 3.2.2. *Q-statistics or squared prediction error:*

352 The Q-statistics also known as squared prediction errors, measure the variability of the observation that violates
 353 the routine process correlation (with small or moderated magnitudes) not accounted by the principal component
 354 subspace. In case of water distribution systems, continuous flow usually does not impact on the overall water
 355 consumption in a short period and so may not be detected by conventional methods, however a low-grade fault,
 356 such as minor continuous flows, can result in undesirable water loss and accumulate to significant losses. The
 357 portion of the measurement space corresponding to the lowest $(m-r)$ eigenvalues are thus monitored (Ballabio,
 358 2015; Qin, 2012).

$$Q_i = x_i^T (I - P_i P_i^T) x_i \quad (10)$$

359 where I is the identity matrix. The control limit Q_α can be computed from its approximate distribution (Ballabio,
 360 2015; Qin, 2012).

$$Q_\alpha = \theta_1 \left[\frac{h_0 c_\alpha \sqrt{2\theta_2}}{\theta_1} + 1 \frac{\theta_2 h_0 (h_0 - 1)}{\theta_1^2} \right]^{1/h_0} \quad (51)$$

$$\theta_i = \sum_{j=r+1}^n \lambda_j^{2i} \quad (62)$$

$$h_0 = 1 - \frac{2\theta_1\theta_3}{3\theta_2^2} \quad (73)$$

361 where r is retained principal components inside the model and c_α is the standard normal deviation with the upper
 362 $(1 - \alpha)$ percentile. When an unusual event occurs, and it produces a change in the covariance structure of the
 363 model, it will be detected by a high value of Q . The new observation is considered to be normal if it satisfies the
 364 following condition.

$$\begin{cases} Q_i < Q_\alpha & \text{--- Normal} \\ Q_i \geq Q_\alpha & \text{--- System alarm} \end{cases} \quad (84)$$

365 When the Q_i of the new experimental trial violates the Q_α control limit, this is indicative of a fault or a non-routine
 366 event. The value of Q_α is defined with an assumption that the observation data is multivariate normally distributed
 367 and time-independent (Harrou et al., 2013). The value Q -statistics are small and consequently more sensitive than
 368 the T^2 -statistics. This characteristic makes the Q -statistics ideally suited to detecting minor variation within the
 369 system behaviour. It evaluates the variation of the new incoming water data which are not accounted for by the
 370 principal component subspace. On the contrary, the T^2 -statistics require significant variation in the system
 371 behaviour to be measurable (Mujica et al., 2011).

372 3.2.3. Fault alarm reduction:

373 In order to isolate actual fault alarms from system alarms, and subsequently reduce false alarms in the fault
 374 detection and diagnosis system, the time series nature of the data over periods of two days or more were
 375 considered. Non-residential public use buildings such as the one explored herein, are occasionally used for large
 376 events which result in significant increases in water consumption (i.e. conferences, seminars, etc. on a given day).
 377 Thus, fault alarms in the proposed framework are raised only when system alarms persist over two or more
 378 consecutive days. Thus, when the analysed data exceeded the α -control limits (Eq. 15) for more than two days a
 379 fault alarm is triggered. The two-day time period was selected based on a historical analysis of the events within
 380 the building (conferences, seminars, workshops, etc.). The nature of the time series alarm system could vary for
 381 different buildings based on water consumption patterns and event patterns and can be established from a high-
 382 level water audit. It is also noted that this framework step can be removed for buildings where occasional large
 383 occupancy events do not occur (i.e. food processing facilities, manufacturing facilities, etc.).

$$\begin{cases} t_{\text{acceptable}} < t_{\text{limit}} & \text{--- Normal} \\ t_{\text{acceptable}} \geq t_{\text{limit}} & \text{--- Indicative of fault} \end{cases} \quad (95)$$

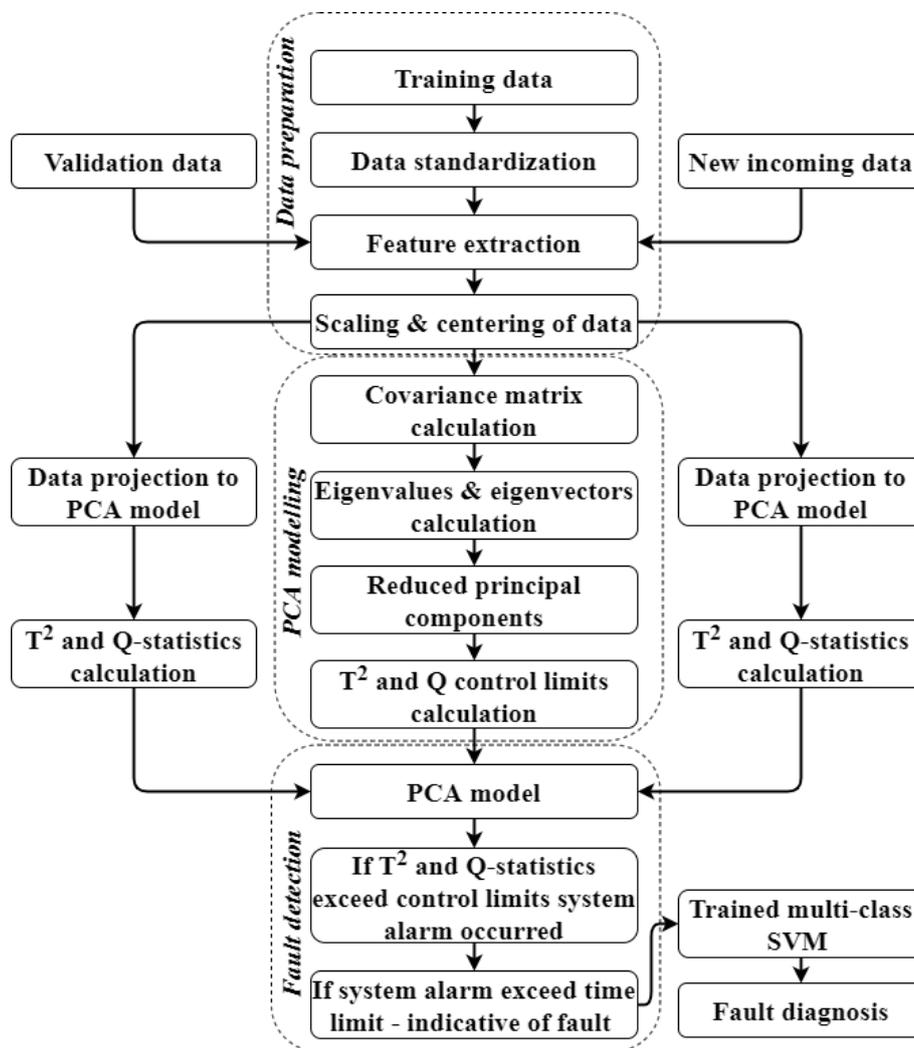
384

385 3.3. Fault diagnosis:

386 Support vector machine (SVM) is a supervised learning technique widely used for classification and regression.
 387 SVM creates a decision boundary in between classes by mapping the training data (through kernel functions) onto
 388 a higher dimensional space, and then obtaining the maximum margin hyperplane within that space. Finally, this
 389 hyperplane can be regarded as a classifier. To achieve multi-class classification, error-correcting output code
 390 (ECOC) was used. ECOC represents a powerful framework to deal with multi-class classification problems based
 391 on combining binary classifiers (such as support vector machine, neural networks, decision trees, etc.) (Bagheri
 392 et al., 2012; Lin, 2018). A detailed explanation of the support vector machine technique can be found in (Bishop,
 393 2013). The support vector machine technique has been used herein due to its wide use in different engineering
 394 sectors (such as thermal power plants, centrifugal pumps, bearing, etc.) in context of classification (Sabri et al.,
 395 2017; Xiao, 2016; Bayar et al., 2015; Chen et al., 2011, Hmeidi et al., 2008).

396 Fault diagnosis of the system alarms in this study was carried out based on the domain knowledge and process
 397 historical information of faults within the building. In this study, SVM was trained by considering the labelled
 398 flow water data. Samples with different system alarm labels (outlined in Table 3) were obtained from the training
 399 data. It is noted that fault data across various faults was somewhat limited in the training data due to lack of
 400 occurrence of these faults over the monitoring period.

401 Moreover, prior to detecting a routine or faulty condition, cross-validation was carried out to optimize classifier's
 402 hyper-parameters and assess the performance of the classifier. It is noted that SVM is less effective when the data
 403 is noisy and contains overlapping data points (Lin, 2018). The goal of cross-validation is to gauge the
 404 generalizability of an algorithm and to prevent overfitting by minimizing the influence of noisy data (Mutasa et
 405 al., 2020; Gupta, 2019; Mudry & Tjellström, 2011). In this study, the SVM classifier was trained by conducting
 406 stratified 10-fold cross validation to minimize the generalization error (minimizing the error associated with the
 407 biasness and variance of data). It also tends to provide less biased estimation of the accuracy. In this case a
 408 classification accuracy of 80.95% was achieved during training. The summary of the fault detection and diagnostic
 409 approach is outlined in Figure 5.



410

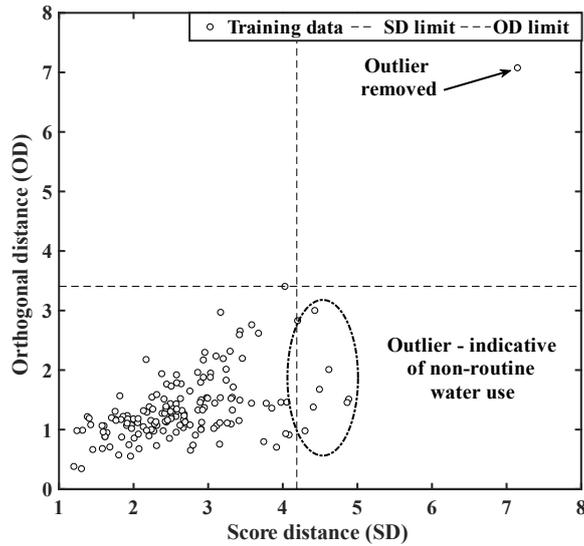
411 Figure 5: Summary of the proposed fault detection and diagnosis approach.

412 4. Results and discussion:

413 4.1. Training:

414 The water usage data was aggregated into hourly flow traces (i.e. 8 readings at 7.5-minute intervals averaged over
 415 one hour) to analyse the usage characteristics at the pilot site. In a number of cases, data points had overlapping
 416 timestamps due to metering errors and these were removed from the training dataset. Real water time series data

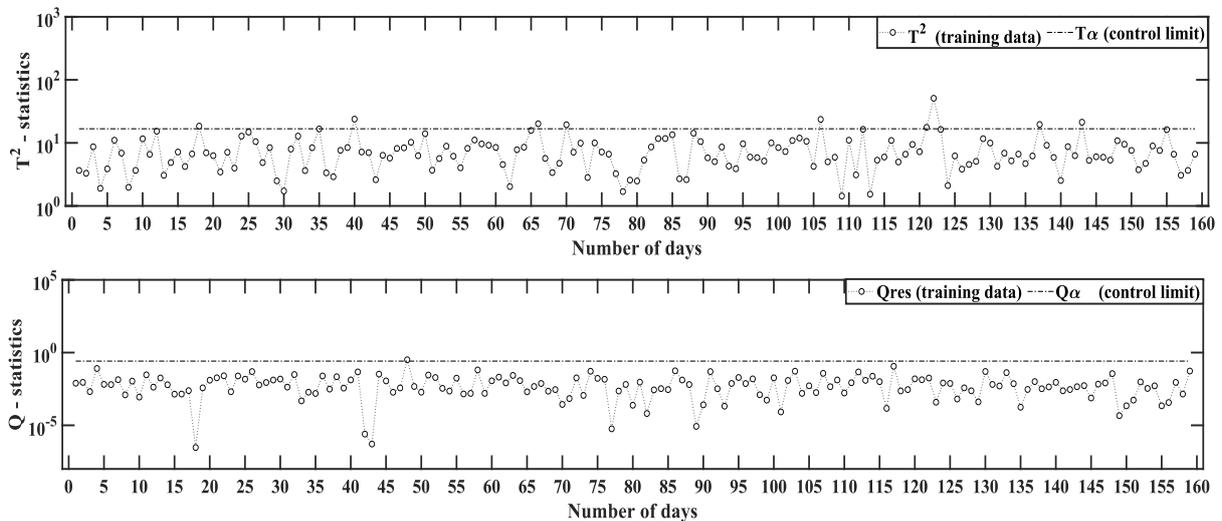
417 often contains values that exhibit atypical usage characteristics and do not follow expected or routine patterns.
 418 Such data points are normally identified as outliers (Harris et al., 2014). These outliers can be due to measurement
 419 error, variation in the water use, and process noise, but could assist in stabilizing the sensitivity of the PCA model.
 420 Thus, some outliers in the dataset may contain valuable information which could be lost if all outliers are removed,
 421 while others (i.e. measurement error and process noise) can negatively affect the model accuracy and could be
 422 removed. To localize outliers in this context, the dataset was analysed by plotting distances (i.e. score and
 423 orthogonal distances) from the centroid of the covariance data structure (Figure 6). The score distance was
 424 measured within a PCA space, while the orthogonal distance was measured between data each point and its
 425 projection in the r -dimensional space (Rousseeuw et al., 2018).



426

427 Figure 6: Distance-distance plot for the training dataset.

428 The PCA model and the SVM model were trained using six months of data from January to June (which
 429 incorporated a Semester of teaching and exams); the calculated PCA monitoring statistics (T^2 -statistics and Q -
 430 statistics) are shown in Figure 7 (dashed lines). As can be seen from Figure 7, a number of points lie above the α -
 431 control limits, constituting system alarms. These are due to a combination of non-routine events in the building
 432 (e.g. event at 122 days) and routine statistical variation related to the inherited definition of the α -control limit,
 433 which is linked to false alarm probability (Li et al., 2019). Importantly, none of these system alarms in the training
 434 data would trigger a fault alarm in the fault detection and diagnosis framework as none of the system alarms occur
 435 on consecutive days (as discussed in Section 3.2.3). It is also noted that metering fault data and one obvious fault
 436 data point were removed from the training data set, in line with standard threshold development procedure.

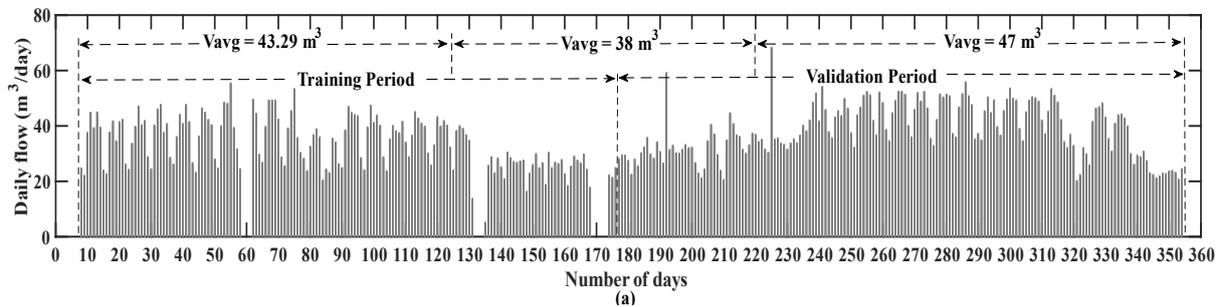


437

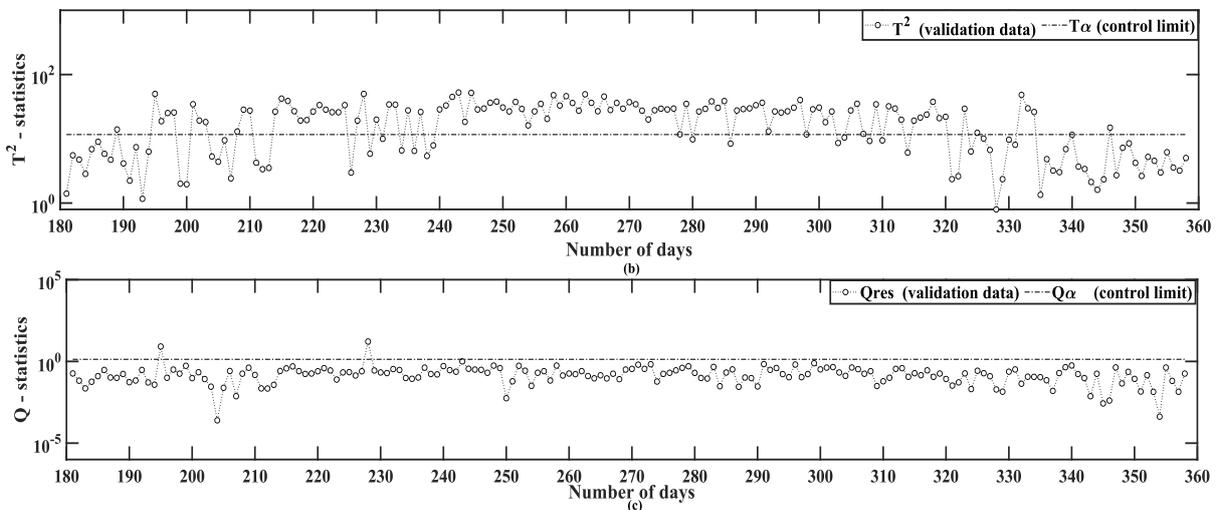
438 Figure 7: The time evolution of T^2 -statistics and Q-statistics on a semi-logarithmic scale for the training data (PCA
 439 model).

440 4.2. *Validation:*

441 Having trained the PCA model, a model validation was conducted on academic Semester 1 water consumption
 442 data (September-December) to assess the model's ability to detect faults. Importantly, it was known that this
 443 Semester 1 data contained an actual fault in the building water distribution system. A previous study (Clifford et
 444 al., 2018) focused on the use of low-resolution meter data to detect various flow signatures associated with
 445 different end-use cases in this building. This study demonstrated faults had occurred in the rainwater-harvesting
 446 system, which was eventually found to be due to a faulty valve on the inlet to one of the storage tanks on the roof
 447 top of the building. Further inspection of meter data from the period concerned indicated that the fault caused
 448 excess metered water consumption of $3.5 \text{ m}^3/\text{day}$ in the building (as the cold water system supply was engaged to
 449 fill the rainwater top-up tanks even though rainwater was available). Figure 8a illustrates the total daily usage of
 450 water over the training period and during the Semester 1 period (validation period) where the fault occurred.
 451 Figure 8b & c show the PCA model output for the validation period, which utilised the Hotelling T^2 -statistics and
 452 Q-statistics α -control limits calculated during the training phase.



453



454

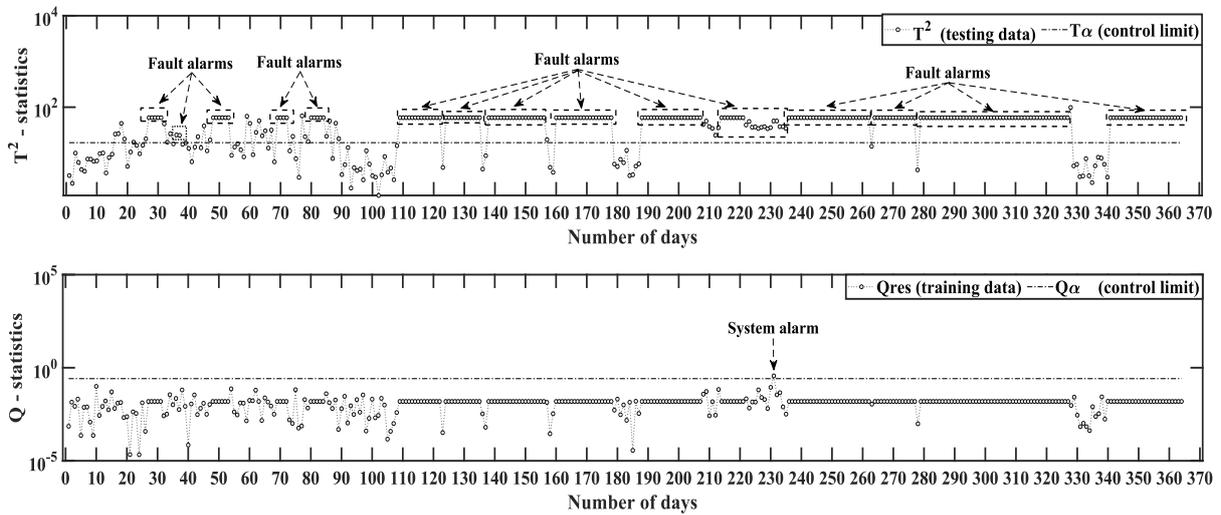
455 Figure 8: (a) Daily usage of total water use in the building, (b & c) The time evolution of T^2 -statistics and Q-
 456 statistics on a semi-logarithmic scale for the testing data (PCA model validation).

457 As can be seen from the plots, the T^2 -statistics for the Semester 1 data breach the α -control limit T_α^2 for a large
 458 portion of the monitoring period, resulting in a significant number of system alarms. As many of these system
 459 alarms occurred on consecutive days, fault alarms were triggered in the PCA fault detection and diagnosis
 460 framework in accordance with the false alarm reduction step (Section 3.2.3). While the T^2 -statistics were effective
 461 in detecting the water distribution system faults (Figure 8b), the Q-statistic did not pick up the fault in question
 462 with few Q values breaching Q_α . This is due to the fact that the T^2 -statistics is more suited to detecting faults
 463 relating to higher flows, while the Q-statistics is targeted towards detecting lower flow reading faults. To gain
 464 further insight into the PCA model performance, the results were compared to those obtained using the univariate

465 method in (Clifford et al., 2018). The PCA model outperformed the univariate model over the validation period,
 466 detecting 25% more system alarms and triggering four more fault alarms. When these additional system alarms
 467 were checked back against the historical records for the building, they correctly corresponded to known faults or
 468 known non-routine events. Comparisons between the performance of the univariate fault detection and diagnosis
 469 approach and the PCA and SVM fault detection and diagnosis approach proposed herein is explored in more detail
 470 in Section 4.3 of this paper.

471 4.3. *Evaluation:*

472 Having trained and validated the fault detection model, its performance was evaluated by considering data for the
 473 year following the validation period. As can be seen in Figure 9, a large number of data points (272 data points)
 474 exceeded the α -control limit and were thus labelled as system alarms i.e. non-routine events or due to faults in the
 475 system. Of these 272 system alarms, sixteen sets of fault alarms were raised (dashed line boxes in Figure 9). As
 476 per the Section 3.2.3 these fault alarms occur when system alarms (caused by values above the α -control limit)
 477 persisted for two consecutive days or more. Several non-routine events were also observed over the 1-year period
 478 (i.e. whereby there was a single 1 day above α -control limit). These high usage peaks raised system alarms but
 479 did not trigger a fault alarm in this fault detection and diagnosis approach.



480
 481 Figure 9: The time evolution of T^2 -statistics and Q -statistics on a semi-logarithmic scale for the new testing data
 482 (PCA fault detection).

483 Having detected the system alarms using the PCA approach, the model then used SVM to classify each system alarm.
 484 The SVM classifier checked system alarms to ensure they were not in fact routine flow events i.e. some
 485 events which are close to the α -control limit may be classified as routine flow. The remaining system alarms were
 486 classified as 1) non-routine events (non-consecutive system alarms), 2) metering error, or 3) excess usage
 487 indicatives of fault. This is shown in Table 3. It is noted that, in this case, the SVM could not be trained to classify
 488 continuous flow events for this case study as no such events occurred during the training data period.

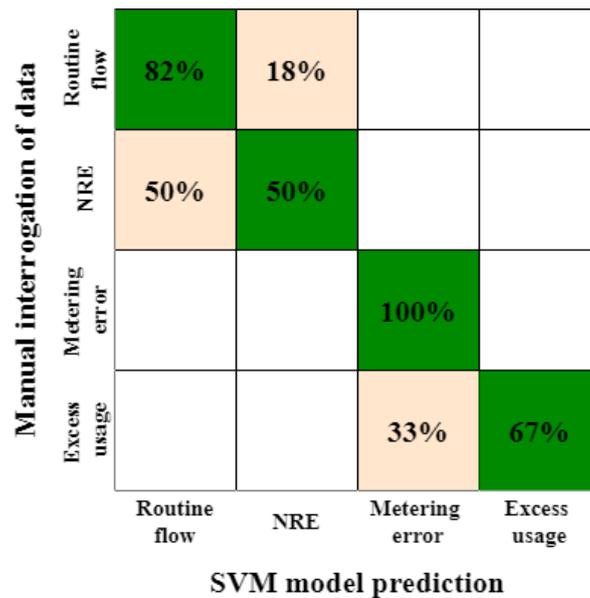
489 Table 3: System alarm classification labels.

Labels	Description	
1	Non-routine event	Unusual water uses in working or non-working hours
2	Metering error	Error due to equipment malfunctioning / zero water use during the time period
3	Excess usage	Exceptional water uses during the time period indicative of water distribution system fault

491 Importantly, the SVM classification identified the majority of the sixteen fault alarm sets (sequential α -control
492 limit breaches), fifteen were due to metering errors (the previous fault identified with the rainwater-harvesting
493 system was explored as a possible issue but was found to have been resolved). On further investigation the
494 “metering error” was in fact found to be the result of flow data not being logged and thus there was no data for
495 these periods (which in this approach are equivalent to “zero” flow days). For days with system alarms these were
496 attributed to workshops, infrequent laboratory events, conferences, etc. that occurred on the given days. In
497 retrospect it was not possible to identify the exact nature of the error.

498 These results are reflective of situations where such data sets can have a significant portion of missing data for a
499 range of reasons including: malfunctioning of monitoring equipment (the replacement of which is seldom a short-
500 term investment priority), lapses in monthly or annual contracts between building owner and third party data
501 monitoring companies, etc. It is noted that if the proposed fault detection and diagnosis system was in place over
502 the evaluation data period the faulty meter would have been identified and subsequently replaced.

503 Figure 10 represents the performance of the classifier across the three system alarm classifications via a confusion
504 matrix. Confusion matrices (aka contingency tables) reveal how the classifier mislabels (or confuses) system
505 alarms and can be used to summarise the performance of a classifier with respect to test data (Vitter & Webber,
506 2018). In this matrix, the main diagonal elements represent how often a certain system alarm label is classified
507 correctly (green), while the other boxes show the classification results for misclassification. It is noted that the
508 correct classification of the data for this comparison was determined from time consuming manual interrogation
509 of the historical test data for all system alarms over the monitoring period. In general, Figure 10 shows that the
510 classifier performs well given the relatively short training dataset of six months. The extent of error in
511 classification is a function of a) the difficulty in identification of a particular class, and b) the small sample sizes
512 of the different system alarms present in the training data. For instance, the training data contained only three
513 metering errors. The classifier can be easily trained for this error however, as metering error readings nearly
514 always show zero flow or a series of days with the same flow (which is generally highly unlikely). This results in
515 100% correct classification for metering error as shown in Figure 10. The training data also contained three non-
516 routine events; however, this data set was only sufficient to ensure 50% classification accuracy for non-routine
517 events. This is due to the fact that many non-routine events are borderline flow events, which are close to the T^2
518 and Q α -control limits and are thus highly susceptible to misclassification. For excess usage, which is perhaps the
519 most important classification as it is indicative of a fault in the water distribution system, the classifier was 67%
520 accurate. This success rate was obtained with only two excess usage data events in the training data. This is
521 partially due to the fact that excess usage events tend to have easily distinguishable characteristics such as high
522 flows and are thus easier to distinguish from routine operation than say non-routine events. While overall classifier
523 performance was good, the analysis clearly indicated that a greater number of system alarms in the training data
524 would facilitate improved classification going forward. Thus, from a practical viewpoint, a classifier can be
525 continually improved during the performance monitoring period, through continued updating of the system alarm
526 training dataset.

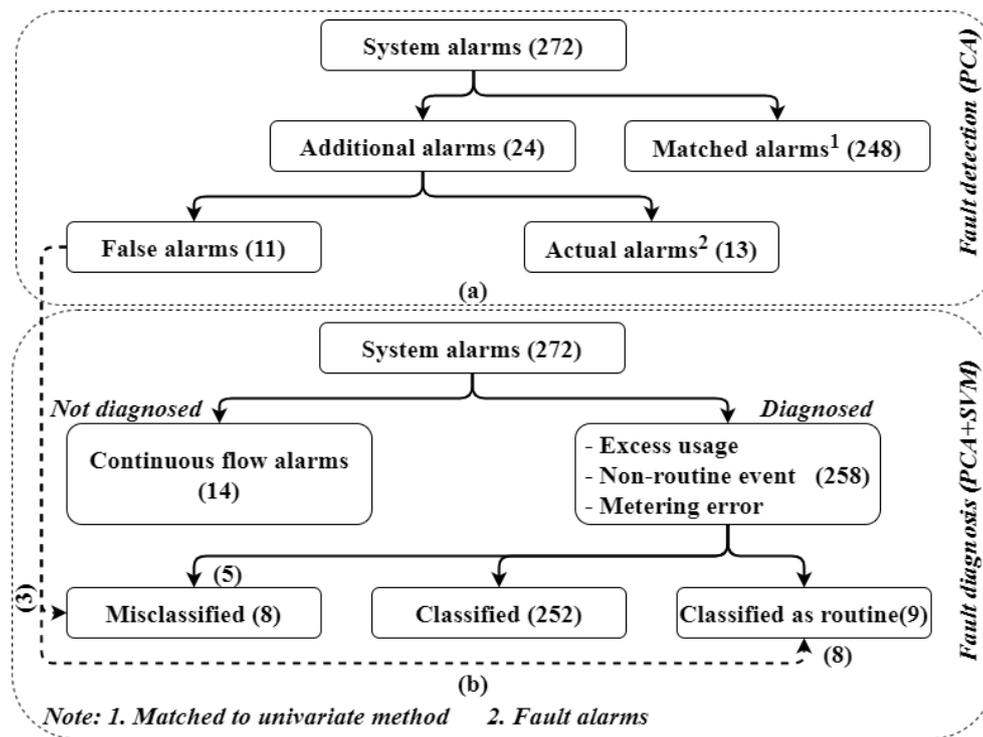


527

528 Figure 10: Confusion matrix of an actual and predicted system alarm labels.

529 *4.4. Effectiveness of fault detection and diagnosis:*

530 In order to gain further insight into the performance of the fault detection and diagnosis approach proposed herein,
 531 the results over the evaluation period were compared to those obtained from the more basic univariate approach
 532 for the same data set. The findings from this comparison are presented graphically in Figure 11a. As discussed,
 533 272 system alarms were detected by the proposed fault detection and diagnosis approach. As shown in Figure 11,
 534 248 of these mapped directly onto the univariate model results, meaning the PCA approach detected an extra 24
 535 system alarms over the univariate method. Manual interrogation of the historical data revealed that of these 24
 536 additional system alarms, 13 corresponded to non-routine events, while 11 constituted false system alarms (i.e.
 537 examination of the historical data indicated the water distribution system in the building was operating correctly
 538 at the time of the system alarm). As shown in Figure 11b however, the performance of the proposed approach was
 539 further enhanced by the SVM component of the model. The SVM classification resulted in 11 of the false system
 540 alarms being re-classified as routine flow. Consequently, overall, the method proposed herein identified all the
 541 univariate model system alarms, detected 13 additional system alarms successfully, and revealed 3 were found to
 542 be false alarms. The proposed approach also facilitated accurate classification of the system alarms (over 90%
 543 accuracy), providing a greater level of information to the building manager. Again it is noted that continuous flow
 544 alarms could not be classified due to a lack of these faults in the training data, and while the proposed approach
 545 has been shown to have advantages over the standard univariate approach, its performance would be even further
 546 enhanced through provision of larger training data sets.



547

548 Figure 11: Description of (a) system alarm detection results and (b) system alarm diagnosis results.

549 **5. Conclusion and future directions:**

550 Principal component analysis together with a multi-class classifier support vector machine were found to be useful
 551 in detecting and diagnosing faults in building water networks. Classical and advanced approaches have resolved
 552 many fault detection and diagnosis problems in water networks, but when it comes to building level multivariate
 553 systems or special cases (for instance substantial variation in the water consumption data), these approaches are
 554 not entirely effective. This study investigated a proposed fault detection and diagnosis approach for non-
 555 residential building water distribution system which combines the detection indices (T^2 and Q-statistics) and
 556 multi-class support vector machine. Hotelling T^2 -statistics and Q-statistics were used to detect system alarms in
 557 the incoming data and the latter multi-class support vector machine along with error-correcting output code was
 558 trained for system alarm classification. The results demonstrated promising capabilities of the proposed fault
 559 detection and diagnosis approach. When compared to the standard univariate approach, a greater number of system
 560 alarms were detected and found to have occurred. The multi-class support vector machine also allowed these
 561 system alarms to be classified, providing a greater level of information to building managers, which may avoid
 562 unnecessary emergency shutdown in industrial applications. Thus, the comparative study has shown that the
 563 proposed approach performs better than standard univariate approach. While the proposed approach has good
 564 capability in detecting and diagnose faults in complex non-residential water distribution systems, it also requires
 565 less computational resources compared to univariate approach.

566 Another important characteristic is adaptability. Once a new fault or anomalous situation is identified, the model
 567 can be updated to incorporate this additional information. Adaptation will make system more reliable and valuable
 568 in time and allow for detection and diagnosis of further anomalies in the water distribution systems. The approach
 569 would be further enhanced by extending training data periods beyond the six-month used in this study, especially
 570 for system alarm classification. The authors do not envisage that this would be a problem in practice.

571 Future work will aim to improve some practical aspects of the model, such as reducing false alarm probability
 572 and improving system alarm classification through enhanced training datasets. The latter is likely to be more
 573 challenging as data from known faults can be difficult to obtain as there may be limited faults in the training data
 574 or many historical faults may be of a similar classification. Furthermore, introducing faults to functional systems
 575 for the purposes of data collection is somewhat impractical. It would also be interesting and potentially very useful

576 to develop a hybrid performance monitoring system (i.e. combining model-based and data-driven approaches),
577 which could combine the strengths of each approach.

578 Acknowledgements

579 This paper has emanated from research conducted as a part of Energy Systems Integration Partnership Programme
580 (ESIPP) project with the financial support of Science Foundation Ireland under the SFI Strategic Partnership
581 Programme Grant Number SFI/15/SPP/E3125.

582 References:

- 583 Abdi, H., & Williams, L. J. (2010). Principal component analysis. *Wiley Interdisciplinary Reviews:*
584 *Computational Statistics*, 2(4), 433–459. <https://doi.org/10.1002/wics.101>
- 585 Abdulshaheed, A., Mustapha, F., & Ghavamian, A. (2017). A pressure-based method for monitoring leaks in a
586 pipe distribution system: A Review. *Renewable and Sustainable Energy Reviews*, 69(May 2015), 902–
587 911. <https://doi.org/10.1016/j.rser.2016.08.024>
- 588 Ahmed, M., Baqqar, M., Gu, F., & Ball, A. D. (2012). Fault detection and diagnosis using Principal Component
589 Analysis of vibration data from a reciprocating compressor. *Proceedings of the 2012 UKACC*
590 *International Conference on Control*, CONTROL 2012, (September 2012), 461–466.
591 <https://doi.org/10.1109/CONTROL.2012.6334674>
- 592 Ait-Izem, T., Harkat, M. F., Djeghaba, M., & Kratz, F. (2018). On the application of interval PCA to process
593 monitoring: A robust strategy for sensor FDI with new efficient control statistics. *Journal of Process*
594 *Control*, 63, 29–46. <https://doi.org/10.1016/j.jprocont.2018.01.006>
- 595 Akil, M., Tittelein, P., Defer, D., & Suard, F. (2019). Statistical indicator for the detection of anomalies in gas,
596 electricity and water consumption: Application of smart monitoring for educational buildings. *Energy and*
597 *Buildings*, 199, 512–522. <https://doi.org/10.1016/j.enbuild.2019.07.025>
- 598 Alsaydalani, M. O. A. (2017). Influence of cavitation phenomenon on the hydraulic behavior of leaks in water
599 distribution systems. *Journal of Water Supply: Research and Technology - Aqua*, 66(5), 327 LP – 339.
600 Retrieved from <http://aqua.iwaponline.com/content/66/5/327.abstract>
- 601 Bagheri, M. A., Montazer, G. A., & Escalera, S. (2012). Error correcting output codes for multiclass
602 classification: Application to two image vision problems. *AISP 2012 - 16th CSI International Symposium*
603 *on Artificial Intelligence and Signal Processing*, (Aisp), 508–513.
604 <https://doi.org/10.1109/AISP.2012.6313800>
- 605 Balaras, C., Dascalaki, E., Droutsas, K., Michas, M., Kontyiannidis, S., & Argiriou, A. (2017). Energy use
606 Intensities for Non-Residential Buildings. *Proceedings of the 48th International HVAC&R Congress*,
607 (December), 369–389. <https://doi.org/10.24094/kgkh.017.48.1.369>
- 608 Ballabio, D. (2015). A MATLAB toolbox for Principal Component Analysis and unsupervised exploration of
609 data structure. *Chemometrics and Intelligent Laboratory Systems*, 149(November), 1–9.
610 <https://doi.org/10.1016/j.chemolab.2015.10.003>
- 611 Bayar, N., Darmoul, S., Hajri-Gabouj, S., & Pierreval, H. (2015). Fault detection, diagnosis and recovery using
612 Artificial Immune Systems: A review. *Engineering Applications of Artificial Intelligence*, 46(November
613 2017), 43–57. <https://doi.org/10.1016/j.engappai.2015.08.006>
- 614 Benaicha, A., Mourot, G., Benothman, K., & Ragot, J. (2013). Determination of principal component analysis
615 models for sensor fault detection and isolation. *International Journal of Control, Automation and Systems*,
616 11(2), 296–305. <https://doi.org/10.1007/s12555-012-0142-x>
- 617 Bishop, C. M. (2013). *Pattern Recognition and Machine Learning*. *Journal of Chemical Information and*
618 *Modeling* (Vol. 53). <https://doi.org/10.1117/1.2819119>
- 619 Bruton, K., Raftery, P., Kennedy, B., Keane, M. M., & O'Sullivan, D. T. J. (2014). Review of automated fault
620 detection and diagnostic tools in air handling units. *Energy Efficiency*, 7(2), 335–351.
621 <https://doi.org/10.1007/s12053-013-9238-2>
- 622 Burak Gunay, H., Shen, W., & Newsham, G. (2019). Data analytics to improve building performance: A critical
623 review. *Automation in Construction*, 97(June 2018), 96–109. <https://doi.org/10.1016/j.autcon.2018.10.020>
- 624 Chen, K. Y., Chen, L. S., Chen, M. C., & Lee, C. L. (2011). Using SVM based method for equipment fault
625 detection in a thermal power plant. *Computers in Industry*, 62(1), 42–50.
626 <https://doi.org/10.1016/j.compind.2010.05.013>
- 627 Chen, T. (2010). On reducing false alarms in multivariate statistical process control. *Chemical Engineering*
628 *Research and Design*, 88(4), 430–436. <https://doi.org/10.1016/j.cherd.2009.09.003>
- 629 Clifford, E., Mulligan, S., Comer, J., & Hannon, L. (2018). Flow-Signature Analysis of Water Consumption in

630 Nonresidential Building Water Networks Using High-Resolution and Medium-Resolution Smart Meter
631 Data: Two Case Studies. *Water Resources Research*, 54(1), 88–106.
632 <https://doi.org/10.1002/2017WR020639>

633 Cody, R. (2020). Acoustic Monitoring for Leaks in Water Distribution Networks by.
634 Cody, R. A., & Narasimhan, S. (2020). A field implementation of linear prediction for leak-monitoring in water
635 distribution networks. *Advanced Engineering Informatics*, 45(April).
636 <https://doi.org/10.1016/j.aei.2020.101103>

637 Connor, B. O., & Murphy, C. (2017). Irish Water.
638 Cuguero Escofet, M. À., Quevedo, J., Alippi, C., Roveri, M., Puig, V., & García, D. (2016). Model- vs. data-
639 based approaches applied to fault diagnosis in potable water supply networks. *Journal of*
640 *Hydroinformatics*, (May), 1–20. <https://doi.org/10.2166/hydro.2016.218>

641 D’Agostino, D., Cuniberti, B., & Bertoldi, P. (2017). Energy consumption and efficiency technology measures
642 in European non-residential buildings. *Energy and Buildings*, 153(2017), 72–86.
643 <https://doi.org/10.1016/j.enbuild.2017.07.062>

644 Danacova, M., Fencik, R., & Nosko, R. (2016). Historical Development of the Permanent Gully Erosion - Case
645 Study Tura Luka. *Water, Resources, Forest, Marine and Ocean Ecosystems Conference Proceedings, Vol*
646 *I*, (June), 391–398. <https://doi.org/10.5593/sgem2016B31>

647 Datta, S., & Sarkar, S. (2016). A review on different pipeline fault detection methods. *Journal of Loss*
648 *Prevention in the Process Industries*, 41, 97–106. <https://doi.org/10.1016/j.jlp.2016.03.010>

649 EPA. (n.d.). Water use efficiency in buildings. In *EPA of USA* (pp. 48–71).
650 <https://doi.org/10.1002/9781118456613>

651 European Commission - DG Environment. (2012). Water Performance of Buildings, (August), 154.

652 García-Alvarez, D. (2014). Fault detection using Principal Component Analysis (PCA) in a Wastewater
653 Treatment Plant (WWTP), (January). Retrieved from
654 <http://sntk09en.guap.ru/sntk09en/main/docs/Alvarezisa.pdf>

655 Gautam, J., Chakrabarti, A., Agarwal, S., Singh, A., Gupta, S., & Singh, J. (2020). Monitoring and forecasting
656 water consumption and detecting leakage using an IoT system. *Water Supply*, 20(3), 1103–1113.
657 <https://doi.org/10.2166/ws.2020.035>

658 Geng, Y., Ji, W., Wang, Z., Lin, B., & Zhu, Y. (2018). A review of operating performance in green buildings:
659 Energy use, indoor environmental quality and occupant satisfaction. *Energy and Buildings*, 183, 500–514.
660 <https://doi.org/10.1016/j.enbuild.2018.11.017>

661 Gharsellaoui, S., Mansouri, M., Trabelsi, M., Refaat, S. S., & Messaoud, H. (2020). Fault diagnosis of heating
662 systems using multivariate feature extraction based machine learning classifiers. *Journal of Building*
663 *Engineering*, 30(September 2019), 101221. <https://doi.org/10.1016/j.jobe.2020.101221>

664 Gupta, S. (2019). 2019_Dealing with Noise Problem in Machine Learning Data-sets.pdf.

665 Harris, P., Brunson, C., Charlton, M., Juggins, S., & Clarke, A. (2014). Multivariate Spatial Outlier Detection
666 Using Robust Geographically Weighted Methods. *Mathematical Geosciences*, 46(1), 1–31.
667 <https://doi.org/10.1007/s11004-013-9491-0>

668 Harrou, F., Nounou, M. N., Nounou, H. N., & Madakyaru, M. (2013). Statistical fault detection using PCA-
669 based GLR hypothesis testing. *Journal of Loss Prevention in the Process Industries*, 26(1), 129–139.
670 <https://doi.org/10.1016/j.jlp.2012.10.003>

671 Hmeidi, I., Hawashin, B., & El-Qawasmeh, E. (2008). Performance of KNN and SVM classifiers on full word
672 Arabic articles. *Advanced Engineering Informatics*, 22(1), 106–111.
673 <https://doi.org/10.1016/j.aei.2007.12.001>

674 Hougbo, G. F. (2019). *The United Nations world development report 2019. Leaving no one behind. UNESCO*
675 *Digital Library*. <https://doi.org/10.1037/0033-2909.126.1.78>

676 Hu, R. L., Granderson, J., Auslander, D. M., & Agogino, A. (2019). Design of machine learning models with
677 domain experts for automated sensor selection for energy fault detection. *Applied Energy*, 235(May 2018),
678 117–128. <https://doi.org/10.1016/j.apenergy.2018.10.107>

679 Hubert, M., Reynkens, T., Schmitt, E., & Verdonck, T. (2016). Sparse PCA for High-Dimensional Data With
680 Outliers. *Technometrics*, 58(4), 424–434. <https://doi.org/10.1080/00401706.2015.1093962>

681 Hubert, M., Rousseeuw, P. J., & Vanden Branden, K. (2005). ROBPCA: A new approach to robust principal
682 component analysis. *Technometrics*, 47(1), 64–79. <https://doi.org/10.1198/004017004000000563>

683 Irish Water. (2019). Strategic Funding Plan 2019-2024.

684 J., P. K., Singh, N., Dayama, P., & Pandit, V. (2019). Change Point Detection for Compositional Multivariate
685 Data. Retrieved from <http://arxiv.org/abs/1901.04935>

686 Jackson, J. E., & Mudholkar, G. S. (1979). Control Procedures for Residuals Associated With Principal
687 Component Analysis. *Technometrics*, 21(3), 341–349. <https://doi.org/10.1080/00401706.1979.10489779>

688 Joe Qin, S. (2003). Statistical process monitoring: basics and beyond. *Journal of Chemometrics*, 17(8–9), 480–
689 502. <https://doi.org/10.1002/cem.800>

690 Johnson, R. A., & Wichern, D. W. (2007). Applied multivariate statistical analysis.

691 Kouziokas, G. N. (2020). SVMkernelbasedonparticleswarmoptimizedvector.pdf.

692 Laory, I., Trinh, T. N., & Smith, I. F. C. (2011). Evaluating two model-free data interpretation methods for
693 measurements that are influenced by temperature. *Advanced Engineering Informatics*, 25(3), 495–506.
694 <https://doi.org/10.1016/j.aei.2011.01.001>

695 Li, Wei, Peng, M., & Wang, Q. (2018). False alarm reducing in PCA method for sensor fault detection in a
696 nuclear power plant. *Annals of Nuclear Energy*, 118, 131–139.
697 <https://doi.org/10.1016/j.anucene.2018.04.012>

698 Li, Wei, Peng, M., & Wang, Q. (2019). Improved PCA method for sensor fault detection and isolation in a
699 nuclear power plant. *Nuclear Engineering and Technology*, 51(1), 146–154.
700 <https://doi.org/10.1016/j.net.2018.08.020>

701 Li, Weihua, Yue, H. H., Valle-Cervantes, S., & Qin, S. J. (2000). Recursive PCA for adaptive process
702 monitoring. *Journal of Process Control*, 10(5), 471–486. [https://doi.org/10.1016/S0959-1524\(00\)00022-6](https://doi.org/10.1016/S0959-1524(00)00022-6)

703 Lin, T.-K. (2018). PCA/SVM-Based Method for Pattern Detection in a Multisensor System. *Mathematical
704 Problems in Engineering*, 2018, 1–11. <https://doi.org/10.1155/2018/6486345>

705 Liu, X., Iftikhar, N., Nielsen, P. S., & Heller, A. (2016). Online Anomaly Energy Consumption Detection Using
706 Lambda Architecture. In S. Madria & T. Hara (Eds.), *Big Data Analytics and Knowledge Discovery* (pp.
707 193–209). Cham: Springer International Publishing.

708 Liu, Y., Ma, X., Li, Y., Tie, Y., Zhang, Y., & Gao, J. (2019). Water pipeline leakage detection based on
709 machine learning and wireless sensor networks. *Sensors (Switzerland)*, 19(23), 1–21.
710 <https://doi.org/10.3390/s19235086>

711 Maione, C., Barbosa, F., & Barbosa, R. M. (2019). Predicting the botanical and geographical origin of honey
712 with multivariate data analysis and machine learning techniques: A review. *Computers and Electronics in
713 Agriculture*, 157(January), 436–446. <https://doi.org/10.1016/j.compag.2019.01.020>

714 Makaya, E., & Hensel, O. (2015). Water loss management strategies for developing countries: Understanding
715 the dynamics of water leakages, 1–97. Retrieved from <https://d-nb.info/1112580042/34>

716 Mashford, J., De Silva, D., Burn, S., & Marney, D. (2012). Leak detection in simulated water pipe networks
717 using SVM. *Applied Artificial Intelligence*, 26(5), 429–444.
718 <https://doi.org/10.1080/08839514.2012.670974>

719 Matos, C., Santos, C., Pereira, S., Bentes, I., Imteaz, M., Cook, S., ... UNESCO. (2020). Environmental impact
720 of water-use in buildings: Latest developments from a life-cycle assessment perspective. *Journal of
721 Environmental Management*, 261(2), 110198. <https://doi.org/10.1016/j.jenvman.2020.110198>

722 Moors, J., Scholten, L., van der Hoek, J. P., & den Besten, J. (2018). Automated leak localization performance
723 without detailed demand distribution data. *Urban Water Journal*, 15(2), 116–123.
724 <https://doi.org/10.1080/1573062X.2017.1414272>

725 Moser, G., German Paal, S., & Smith, I. F. C. (2015). Performance comparison of reduced models for leak
726 detection in water distribution networks. *Advanced Engineering Informatics*, 29(3), 714–726.
727 <https://doi.org/10.1016/j.aei.2015.07.003>

728 Mudry, A., & Tjellström, A. (2011). Historical background of bone conduction hearing devices and bone
729 conduction hearing aids. *Advances in Oto-Rhino-Laryngology*, 71, 1–9.
730 <https://doi.org/10.1159/000323569>

731 Mujica, L. E., Rodellar, J., Fernández, A., & Güemes, A. (2011). Q-statistic and t2-statistic pca-based measures
732 for damage assessment in structures. *Structural Health Monitoring*, 10(5), 539–553.
733 <https://doi.org/10.1177/1475921710388972>

734 Mutasa, S., Sun, S., & Ha, R. (2020). Understanding artificial intelligence based radiology studies: What is
735 overfitting? *Clinical Imaging*, 65(April), 96–99. <https://doi.org/10.1016/j.clinimag.2020.04.025>

736 Naderi, E., & Khorasani, K. (2016). A Data-driven Approach to Actuator and Sensor Fault Detection, Isolation
737 and Estimation in Discrete-Time Linear Systems, 85, 165–178. Retrieved from
738 <http://arxiv.org/abs/1606.06220>

739 Nasir, M. T., Mysorewala, M., Cheded, L., Siddiqui, B., & Sabih, M. (2014). Measurement error sensitivity
740 analysis for detecting and locating leak in pipeline using ANN and SVM. *2014 IEEE 11th International
741 Multi-Conference on Systems, Signals and Devices, SSD 2014*, 7–10.
742 <https://doi.org/10.1109/SSD.2014.6808847>

743 Nezhad, A. J., Wijaya, T. K., Vasirani, M., & Aberer, K. (2014). SmartD: Smart Meter Data Analytics
744 Dashboard. *Proceedings of the 5th International Conference on Future Energy Systems*, 213–214.
745 <https://doi.org/10.1145/2602044.2602046>

- 746 Nowicki, A., Grochowski, M., & Duzinkiewicz, K. (2012). Data-driven models for fault detection using kernel
747 PCA: A water distribution system case study. *International Journal of Applied Mathematics and*
748 *Computer Science*, 22(4), 939–949. <https://doi.org/10.2478/v10006-012-0070-1>
- 749 Patabendige, S., Cardell-Oliver, R., Wang, R., & Liu, W. (2018). Detection and interpretation of anomalous
750 water use for non-residential customers. *Environmental Modelling and Software*, 100, 291–301.
751 <https://doi.org/10.1016/j.envsoft.2017.11.028>
- 752 Pelz, G. (2003). *Mechatronic Systems. Library*. <https://doi.org/10.1002/0470867906>
- 753 Pérez, R., Sanz, G., Cugueró, M. À., Ramon, P., Sanz, G., & Angel, M.-. (2015). Parameter Uncertainty
754 Modelling in Water Distribution Network Models Models, (September).
755 <https://doi.org/10.1016/j.proeng.2015.08.911>
- 756 Perfido, D., Messervey, T., Zanotti, C., Raciti, M., & Costa, A. (2016). Automated Leak Detection System for
757 the Improvement of Water Network Management. *Proceedings*, 1(2), 28. <https://doi.org/10.3390/ecsa-3->
758 [S5002](https://doi.org/10.3390/ecsa-3-S5002)
- 759 Posenato, D., Lanata, F., Inaudi, D., & Smith, I. F. C. (2006). Model free interpretation of monitoring data.
760 *Lecture Notes in Computer Science (Including Subseries Lecture Notes in Artificial Intelligence and*
761 *Lecture Notes in Bioinformatics)*, 4200 LNAI, 529–533. https://doi.org/10.1007/11888598_47
- 762 Posenato, D., Lanata, F., Inaudi, D., & Smith, I. F. C. (2008). Model-free data interpretation for continuous
763 monitoring of complex structures. *Advanced Engineering Informatics*, 22(1), 135–144.
764 <https://doi.org/10.1016/j.aei.2007.02.002>
- 765 Qin, S. J. (2012). Annual Reviews in Control Survey on data-driven industrial process monitoring and
766 diagnosis, 36, 220–234.
- 767 Qu, Z., Feng, H., Zeng, Z., Zhuge, J., & Jin, S. (2010). A SVM-based pipeline leakage detection and pre-
768 warning system. *Measurement: Journal of the International Measurement Confederation*, 43(4), 513–519.
769 <https://doi.org/10.1016/j.measurement.2009.12.022>
- 770 Quevedo, J., Chen, H., Cugueró, M. À., Tino, P., Puig, V., García, D., ... Yao, X. (2014). Engineering
771 Applications of Artificial Intelligence Combining learning in model space fault diagnosis with data
772 validation / reconstruction : Application to the Barcelona water network. *Engineering Applications of*
773 *Artificial Intelligence*, 30, 18–29. <https://doi.org/10.1016/j.engappai.2014.01.008>
- 774 Quiñones-Grueiro, M., Bernal-de Lázaro, J. M., Verde, C., Prieto-Moreno, A., & Llanes-Santiago, O. (2018).
775 Comparison of Classifiers for Leak Location in Water Distribution Networks *. *IFAC-PapersOnLine*,
776 51(24), 407–413. <https://doi.org/10.1016/j.ifacol.2018.09.609>
- 777 Quiñones-Grueiro, M., Verde, C., Prieto-Moreno, A., & Llanes-Santiago, O. (2018). An unsupervised approach
778 to leak detection and location in water distribution networks. *International Journal of Applied*
779 *Mathematics and Computer Science*, 28(2), 283–295. <https://doi.org/10.2478/amcs-2018-0020>
- 780 Robles, D., Puig, V., Ocampo-Martinez, C., & Garza-Castañón, L. E. (2016). Reliable fault-tolerant model
781 predictive control of drinking water transport networks. *Control Engineering Practice*, 55, 197–211.
782 <https://doi.org/10.1016/j.conengprac.2016.06.014>
- 783 Rosen, C. (2001). *A chemometric approach to process monitoring and control with application to wastewater*
784 *treatment operation- . Department of Industrial Electrical Engineering and Automation.*
- 785 Rousseeuw, P. J., & Hubert, M. (2018). Anomaly detection by robust statistics. *Wiley Interdisciplinary Reviews:*
786 *Data Mining and Knowledge Discovery*, 8(2), 1–14. <https://doi.org/10.1002/widm.1236>
- 787 Rousseeuw, P. J., Raymaekers, J., & Hubert, M. (2018a). A Measure of Directional Outlyingness With
788 Applications to Image Data and Video. *Journal of Computational and Graphical Statistics*, 27(2), 345–
789 359. <https://doi.org/10.1080/10618600.2017.1366912>
- 790 Rousseeuw, P. J., Raymaekers, J., & Hubert, M. (2018b). A Measure of Directional Outlyingness With
791 Applications to Image Data and Video. *Journal of Computational and Graphical Statistics*, 27(2), 345–
792 359. <https://doi.org/10.1080/10618600.2017.1366912>
- 793 Sabri, M., Sinal, B., & Kamioka, E. (2017). ADAPTIVE THRESHOLD BASED APPROACH TO
794 PERFECTLY DETECT HEART CYCLE IN ECG DATA, (44), 492–498.
- 795 Saitta, S., Raphael, B., & Smith, I. F. C. (2005). Data mining techniques for improving the reliability of system
796 identification. *Advanced Engineering Informatics*, 19(4), 289–298.
797 <https://doi.org/10.1016/j.aei.2005.07.005>
- 798 Salam, A. E. U., Tola, M., Selintung, M., & Maricar, F. (2015). Application of SVM and ELM Methods to
799 Predict Location and Magnitude Leakage of Pipelines on Water Distribution Network, (19), 7970.
- 800 Samer El-Zahab. (2018). An accelerometer-based leak detection system _ Elsevier Enhanced Reader.pdf.
- 801 Sean Mulligan, Louise Hannon, Paraic Ryan, Sudeep Nair, E. C. (2020). Jo ur n re of. *Journal of Building*
802 *Engineering*, 102248. <https://doi.org/10.1016/j.janxdis.2020.102248>
- 803 Sedki, A., & Ouazar, D. (2012). Hybrid particle swarm optimization and differential evolution for optimal

804 design of water distribution systems. *Advanced Engineering Informatics*, 26(3), 582–591.
805 <https://doi.org/10.1016/j.aei.2012.03.007>

806 Sengupta, R. N., & Kundu, D. (2016). Statistical Methods, 413–520.

807 Seyoum, S., Alfonso, L., Andel, S. J. Van, Koole, W., Groenewegen, A., & Van De Giesen, N. (2017). A
808 Shazam-like Household Water Leakage Detection Method. *Procedia Engineering*, 186, 452–459.
809 <https://doi.org/10.1016/j.proeng.2017.03.253>

810 Sheriff, M. Z., Botre, C., Mansouri, M., Nounou, H., Nounou, M., & Karim, M. N. (2017). Process Monitoring
811 Using Data-Based Fault Detection Techniques: Comparative Studies. *Fault Diagnosis and Detection*,
812 (December). <https://doi.org/10.5772/67347>

813 Skworcow, P., Paluszczyszyn, D., Ulanicki, B., Jung, B. S., Boulos, P. F., Wood, D. J., ... Giustolisi, O. (2013).
814 Pressure , Leakage and Energy Management in Water Distribution Systems. *Water Resources*
815 *Management*, i(1), 266. <https://doi.org/10.1016/j.proeng.2015.08.855>

816 Sliskovic, D., Grbic, R., & Hocenski, Z. (2012). Multivariate Statistical Process Monitoring. *Tehnicki Vjesnik-*
817 *Technical Gazette*, 19(1), 33–41. <https://doi.org/10.1016/j.patcog.2009.11.008>

818 Soldevila, A., Fernandez-Canti, R. M., Blesa, J., Tornil-Sin, S., & Puig, V. (2017). Leak localization in water
819 distribution networks using Bayesian classifiers. *Journal of Process Control*, 55, 1–9.
820 <https://doi.org/10.1016/j.jprocont.2017.03.015>

821 Sousa, V., Silva, C. M., & Meireles, I. (2019). Performance of water efficiency measures in commercial
822 buildings. *Resources, Conservation and Recycling*, 143(October 2018), 251–259.
823 <https://doi.org/10.1016/j.resconrec.2019.01.013>

824 Stavset, O., & Kauko, H. (2015). *Report -possibilities for smart energy solutions*.

825 Stetco, A., Dinmohammadi, F., Zhao, X., Robu, V., Flynn, D., Barnes, M., ... Nenadic, G. (2019). Machine
826 learning methods for wind turbine condition monitoring: A review. *Renewable Energy*, 133, 620–635.
827 <https://doi.org/10.1016/j.renene.2018.10.047>

828 Sydney Water. (2011). Best practice guidelines for water management in aquatic leisure centres.

829 Venkatasubramanian, V, Rengaswamy, R., Kavuri, S. N., & Yin, K. (2003). A review of process fault detection
830 and diagnosis Part III: Process history based methods. *Computers & Chemical Engineering*, 27(3), 293–
831 311. [https://doi.org/10.1016/S0098-1354\(02\)00160-6](https://doi.org/10.1016/S0098-1354(02)00160-6)

832 Venkatasubramanian, Venkat, Rengaswamy, R., & Yin, K. (2003). A re v iew of process fault detection and
833 diagnosis Part I : Quantitati v e model-based methods, 27, 293–311.

834 Vieira, A. S., Beal, C. D., Ghisi, E., & Stewart, R. A. (2014). Energy intensity of rainwater harvesting systems:
835 A review. *Renewable and Sustainable Energy Reviews*, 34, 225–242.
836 <https://doi.org/10.1016/j.rser.2014.03.012>

837 Villegas, T., Fuente, M. J., & Rodríguez, M. (2010). Principal component analysis for fault detection and
838 diagnosis. experience with a pilot plant. *Proceedings of the 9th WSEAS International Conference on*
839 *Computational Intelligence, Man-Machine Systems and Cybernetics*, 147–152.

840 Vitter, J. S., & Webber, M. E. (2018). A non-intrusive approach for classifying residential water events using
841 coincident electricity data. *Environmental Modelling and Software*, 100, 302–313.
842 <https://doi.org/10.1016/j.envsoft.2017.11.029>

843 Xiao, S., Lu, Z., & Xu, L. (2017). Multivariate sensitivity analysis based on the direction of eigen space through
844 principal component analysis. *Reliability Engineering and System Safety*, 165(March), 1–10.
845 <https://doi.org/10.1016/j.res.2017.03.011>

846 Xiao, W. (2016). A Probabilistic Machine Learning Approach to Detect Industrial Plant Faults. *International*
847 *Journal of Prognostics and Health Management*, (c), 1–11.

848 Xu, L., Wang, X., Bai, L., Xiao, J., Liu, Q., Chen, E., ... Luo, B. (2020). Probabilistic SVM classifier ensemble
849 selection based on GMDH-type neural network. *Pattern Recognition*, 106.
850 <https://doi.org/10.1016/j.patcog.2020.107373>

851 Xu, X., Wang, H., Zhang, N., Liu, Z., & Wang, X. (2017). Review of the Fault Mechanism and Diagnostic
852 Techniques for the Range Extender Hybrid Electric Vehicle. *IEEE Access*, 5, 14234–14244.
853 <https://doi.org/10.1109/ACCESS.2017.2725298>

854 Yu, Y., Woradechjumroen, D., & Yu, D. (2014). A review of fault detection and diagnosis methodologies on
855 air-handling units. *Energy and Buildings*, 82, 550–562. <https://doi.org/10.1016/j.enbuild.2014.06.042>

856 Zenobi, R., Knochenmuss, R., Yamashitat, M., Fenn, J. B., Wong, A., Jime, B., ... Agency, E. (2011).
857 Introduction to multivariate analysis. *J Am Chem Soc*. <https://doi.org/10.1002/jms>

858 Zhang, H., Qi, Y., Wang, L., Gao, X., & Wang, X. (2017). Fault detection and diagnosis of chemical process
859 using enhanced KECA. *Chemometrics and Intelligent Laboratory Systems*, 161(December 2016), 61–69.
860 <https://doi.org/10.1016/j.chemolab.2016.12.013>

861 Zhao, Y., Li, T., Zhang, X., & Zhang, C. (2019). Artificial intelligence-based fault detection and diagnosis

862 methods for building energy systems: Advantages, challenges and the future. *Renewable and Sustainable*
863 *Energy Reviews*, 109(February), 85–101. <https://doi.org/10.1016/j.rser.2019.04.021>
864 Zhou, A., Yu, D., & Zhang, W. (2015). A research on intelligent fault diagnosis of wind turbines based on
865 ontology and FMECA. *Advanced Engineering Informatics*, 29(1), 115–125.
866 <https://doi.org/10.1016/j.aei.2014.10.001>
867 Zumoffen, D., & Basualdo, M. (2008). From large chemical plant data to fault diagnosis integrated to
868 decentralized fault-tolerant control: Pulp mill process application. *Industrial and Engineering Chemistry*
869 *Research*, 47(4), 1201–1220. <https://doi.org/10.1021/ie071064m>
870
871