

# Detecting Technical Anomalies in High-Frequency Water-Quality Data Using Artificial Neural Networks

Javier Rodriguez-Perez, Catherine Leigh, Benoit Liquet, Claire Kermorvant, Erin Peterson, Damien Sous, and Kerrie Mengersen\*



Cite This: <https://dx.doi.org/10.1021/acs.est.0c04069>



Read Online

ACCESS |



Metrics & More

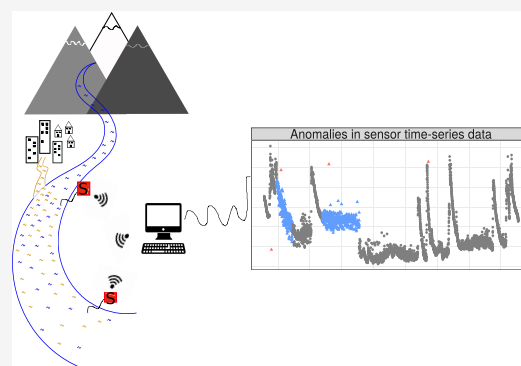


Article Recommendations



Supporting Information

**ABSTRACT:** Anomaly detection (AD) in high-volume environmental data requires one to tackle a series of challenges associated with the typical low frequency of anomalous events, the broad-range of possible anomaly types, and local nonstationary environmental conditions, suggesting the need for flexible statistical methods that are able to cope with unbalanced high-volume data problems. Here, we aimed to detect anomalies caused by technical errors in water-quality (turbidity and conductivity) data collected by automated in situ sensors deployed in contrasting riverine and estuarine environments. We first applied a range of artificial neural networks that differed in both learning method and hyperparameter values, then calibrated models using a Bayesian multiobjective optimization procedure, and selected and evaluated the “best” model for each water-quality variable, environment, and anomaly type. We found that semi-supervised classification was better able to detect sudden spikes, sudden shifts, and small sudden spikes, whereas supervised classification had higher accuracy for predicting long-term anomalies associated with drifts and periods of otherwise unexplained high variability.



## INTRODUCTION

Monitoring the water quality of rivers is becoming increasingly relevant in the Anthropocene as the need for management strategies to safeguard water resources in human-dominated landscapes accelerates. With the advent of new technologies, data on water properties can be obtained with high frequency in near-real time, which can facilitate fast and adaptive management strategies.<sup>1,2</sup> However, the generation of high-volume, high-velocity data can create problems of quality control because of a combination of (1) technical issues with the water-quality sensors that result in technical anomalies, and (2) operational problems with traditional, often manual and periodic, methods of quality control that are no longer feasible. In order to improve the quality and utility of high-volumes of data, we therefore need to explore efficient methods to complement or replace traditional methods of quality control, including manual anomaly detection (AD) and correction and quality coding.<sup>3,4</sup> Improving water-quality monitoring will also require robust and flexible statistical methods that can handle data streaming, while also providing big data tools for data-driven decision making.<sup>5–7</sup>

The accurate detection of anomalies in high-frequency, high-volume water-quality data is additionally challenged by several factors. For example, the common low frequency of anomalies in water-quality time series constrains the ability to not only detect anomalies but also to evaluate the performance of detection methods.<sup>8,9</sup> There is also a broad range of

possible anomaly types, ranging from sudden changes in water-quality values to anomalously constant values maintained for long periods.<sup>10,11</sup> In addition, the differentiation of technical anomalies, for example, related to sensor malfunction or database failure, from the real deviation of the natural system from its typical behavior is not always straightforward.<sup>1</sup> This is because natural fluctuations under environmental conditions at a given site can influence the level of detectability of sensor-related anomalies from the normal signal: the greater the fluctuation, the harder the detection of technical errors. All of these issues can constrain the generality and robustness of detection methods under a broad set of environmental conditions.<sup>10,12</sup>

A combination of methods may be required to address the challenges raised above. For instance, sudden spikes and shifts in values and impossible and/or out-of-sensor-range values on single or just a few observation points can be easily differentiated from adjacent, nonanomalous values. Despite the extremely low frequency of such point-based anomalies in time-series data, rule-based, regression-based, and feature-

Received: June 22, 2020

Revised: August 28, 2020

Accepted: August 28, 2020

Published: August 28, 2020



ACS Publications

© XXXX American Chemical Society

A

<https://dx.doi.org/10.1021/acs.est.0c04069>  
Environ. Sci. Technol. XXXX, XXX, XXX–XXX

based methods may provide optimal performance in terms of their detection.<sup>1,11</sup> Common anomalies in chemical and biochemical data from in situ sensors also comprise multiple observations, such as sensor drift and periods of unusually high or low variability, which may indicate the need for sensor calibration or maintenance.<sup>13</sup> The detection of these multiple-points anomalies can be tackled by various methods,<sup>7,11</sup> but remains challenging and often still requires user intervention. Several methods can be combined to extend the range of detected anomalies. Recent work by Leigh et al.<sup>11</sup> showed that a combination of rule-, feature-, and regression-based methods facilitated the correct classification of impossible values, sudden isolated spikes and level shifts, although drift and periods of high variability still tended to be associated with high rates of false positives.

Considering the current global effort to produce real-time, high-frequency, and long-term monitoring of aquatic ecosystems, and the limitations of existing quality control approaches,<sup>1,4,13</sup> there is a pressing need to explore new methods. Artificial neural networks (ANNs) are a promising alternative.<sup>6,7,14</sup> One basic strength of ANNs in analyzing high-frequency time-series data is that they do not require a priori knowledge of the underlying physical and environmental processes. ANNs have therefore been used to tackle multiple problems of data streaming including predicting temporal patterns in both marine<sup>15,16</sup> and freshwater systems.<sup>6,17,18</sup> They have the potential to detect anomalies and sudden changes that occur over short-time intervals in water-quality data,<sup>6,18,19</sup> suggesting that they can be readily applied to optimize AD in complex and fluctuating environmental systems. Previous work on AD applied to water-quality variables compared ANNs to other machine learning methods and found that ANNs have comparable performance to other methods in terms of AD.<sup>12,18,20,21</sup> However, authors also acknowledge that remains understudied the optimization techniques (especially on those machine learning methods whose learning process depends on a large number of hyperparameters) that aim to improve detection capacity under real-world conditions.<sup>7,20</sup>

The aim of this study is to test the ability of ANNs to detect technical anomalies in water-quality time-series data in contrasting riverine and estuarine environments. We were especially interested in testing the ability of ANNs to detect multiple-point anomalies, for example, those that may result from sensor drift, given that methods previously examined were limited in their ability to accurately detect such context-dependent events. The studied variables are turbidity and conductivity collected at high-frequency by autonomous in situ sensors. The monitored sites in subtropical Australia and temperate France provide contrasting ranges of environmental conditions, particularly in terms of the studied variables and their magnitudes and dynamics. The comparison and evaluation of detection methods in regions that vary in environmental conditions and across a common range of variables is essential to determine the methods' broad suitability for AD.<sup>7</sup> Given that the environmental fluctuations in each river (freshwater and estuarine sites) may influence the ability to detect the different types of anomalies, we calibrated ANNs using models that differed in learning method and hyperparameter values. To do so, we implemented a Bayesian optimization method using a multiobjective approach (i.e., based on scores especially suited to unbalanced classification problems), in order to find the

combination of ANN hyperparameters best suited for AD at each location. The performance of the calibrated ANNs in detecting different types of anomalies was also evaluated and compared to regression-based methods developed for and applied to similar data.<sup>11</sup>

## MATERIALS AND METHODS

**Study Sites.** The studied water-quality data derive from two freshwater rivers in north eastern Australia, the Pioneer River (PR) and Sandy Creek (SC), and an estuarine river in south western France, the Adour River estuary. Both study areas have seasonality in climate and therefore provide nonstationary environmental conditions of water properties throughout the year.

The Australian study area is characterized by humid subtropical climate and strong seasonality: the wet season (typically occurring between December and April) has higher rainfall and air temperatures than the dry season that has lower rainfall and is associated with low to zero river flows.<sup>22</sup> PR is in the Mackay Whitsunday region of northeast Australia, with a length of 120 km and a monitored catchment area of 1466 km<sup>2</sup>, with the upper reaches flowing predominantly through National or State Parks and its middle and lower reaches flowing through land dominated by sugarcane farming. SC is a low-lying coastal-plain stream, 72 km long with a monitored catchment area of 326 km<sup>2</sup> and with a similar land-use and land-cover profile to that of the lower PR. Both study sites are in the freshwater reaches of these rivers.

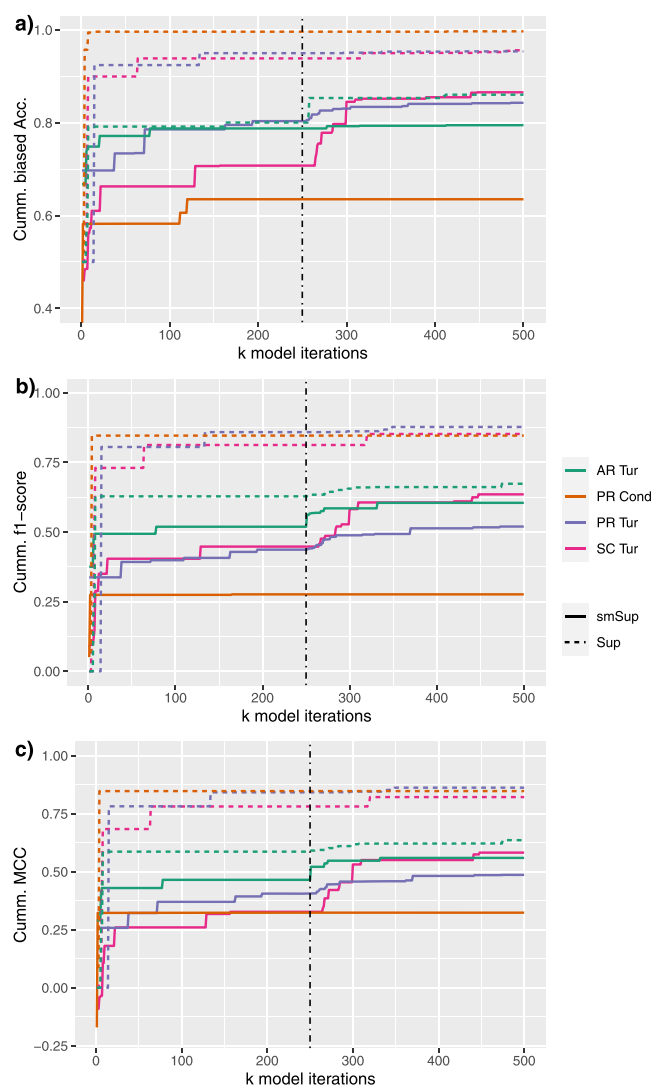
In contrast, south-western France is characterized by an oceanic temperate climate, with drier and warmer conditions from April to October. The Adour River (AR) in the Aquitaine region of southwestern France is 330 km in length and has a catchment area of 16,880 km<sup>2</sup>, with reaches flowing throughout a mosaic of agricultural and forested lands. The study site is in the lower estuary between Bayonne city and the river mouth about 3 km downstream. The tidal range varies from 1 to 4.5 m. Under the combined influence of rainfall and snowmelt during late spring, river flow is highly variable, with minima and maxima of about 80 and 3000 m<sup>3</sup>/s, respectively.<sup>23</sup> Driven by such strong fluctuations in tidal range and river discharge, the estuary is classified as a highly variable, time-dependent salt-wedge.<sup>23</sup>

**Water-Quality Data.** Autonomous multiparameter in situ sensors have been installed at each site. At PR and SC, YSI EXO2 sondes with YSI Smart Sensors are housed in flow cells in water-quality monitoring stations on the rivers' banks. At predefined time intervals (see below), monitoring systems transport water via a pumping system from the river up to the flow cell. Sensors are equipped with wipers to minimize biofouling and all equipment undergo regular maintenance and calibration, with sensors calibrated and equipment checked approximately every six weeks following manufacturer guidelines. At AR, a YSI 6920 sonde is installed 1 m below the water surface on a floating pier, providing direct measurement of water quality in the river. Data are retrieved during maintenance operations (sensor cleaning, memory erasing, battery replacement) approximately every eight weeks.

At each site, the in situ sensors record turbidity (NTU) and electrical conductivity (conductivity;  $\mu\text{S}/\text{cm}$ ). The temporal resolution of data differs among rivers according to local variability under environmental conditions; measurements are taken every 60 and 90 min in PR and SC, respectively, and every 10 min in the AR. At PR and SC, measurements are

sometimes also taken at more frequent intervals during high flow events. The data sequences are one year (12 March 2017 to 12 March 2018) and five months (11 March to 13 August 2019) for the Australian and French sites, respectively, totaling 6280, 5402, and 14,185 observations at PR, SC, and AR, respectively.

Water-quality parameters were strongly affected by the nonstationary flow regimes of each site (see Figure 1).



**Figure 1.** Observed trends for each water-quality variable and site, including the type of anomaly. Shapes correspond to the different data values (i.e., circles for “normal” values and triangles for anomalous values) and colors to the different anomaly types. Anomaly types were classified by local water-quality experts. For instance, Class 1 anomalies are defined as large sudden spikes (type A), sudden shifts (type D), small sudden spikes (type J), and Class 3 anomalies as drift (type H), high variability (type E), and untrustworthy data not defined by other types (type L).

Turbidity is a visual property of water clarity with higher values indicating lower clarity. High turbidity often occurs during flood events, when rivers become loaded with sediment eroded from the catchment, and during local re-suspension events that can result from salt-wedge arrival, flow instability, or navigational and dredging activities. In contrast, low values tend to occur during low discharge periods, for example,

during the dry season, when the slow currents promote sediment deposition or, in estuarine environments, when low-turbidity marine waters have entered the system. These overall trends can become much more complex, and sometimes reversed, under the influence of specific local features. Conductivity reflects the concentration of ions in the water. The dominant effect is generally due to salts, such as sodium chloride, with marine waters typically having much higher conductivity than inland waters. Other dissolved ionic compounds such as nutrients also contribute to conductivity.

In the Australian sites, the water regime is driven by rainfall patterns, generating strong positive or negative phases for either turbidity and conductivity values, coinciding with new and sudden inputs of fresh water or periods of low and zero flow (see upper and mid panels in Figure 1). By contrast, AR is strongly conditioned by the interplay between river discharge and tide forcing. The basic trend is that during high-discharge events, the turbidity reaches high values corresponding to massive seaward fluxes of suspended sediments, and the conductivity is minimal because the marine waters are flushed out from the estuary (see lower panels in Figure 1). During moderate discharge events, the tidal influence becomes dominant and drives tidal fluctuations of conductivity and turbidity inside the estuary: at high tide, the estuary is filled by salty (high conductivity) and clearer (low turbidity) marine waters while the reverse behavior is observed at low tide. This pattern is further affected by strong vertical gradients, the salty marine waters being heavier than fresh riverine waters, which induces complex time-dependent dynamics at the tidal scale.<sup>23</sup>

We focused on turbidity and conductivity data, two common water quality-variables measured in the majority of monitoring programs worldwide and commonly used as proxies of nutrient and/or sediment pollution in rivers, estuaries and marine environments.<sup>24</sup> As a global problem, sediment pollution also causes major economic issues related drinking water and dredging.<sup>25</sup> In addition, turbidity and conductivity are typically more stable in rivers through time than other properties such as dissolved oxygen and water temperature, which fluctuate daily as well as seasonally.

**Definitions and Types of Anomalies.** Environmental sensors generating high-frequency data sets are usually subject to technical anomalies that can derive from fouling of sensors, sensor calibration shifts, power supply problems or unforeseen environmental conditions that adversely affect the sensor equipment.<sup>1,26</sup> In general, such technical anomalies strongly depart from the expected pattern of the nonanomalous observations.<sup>10</sup> In addition, strong inter- and intraseasonal variations in environmental conditions from site to site may create local differences in the occurrence and types of anomalies, which can constrain their transferability to other sites.<sup>27</sup> We thus need to develop methods for detecting these technical anomalies, that can overcome these challenges,<sup>1,11,27</sup> particularly in the context of data streaming.

The present analysis focuses on technical AD (AD) in water-quality time series collected by in situ sensors. Types of such anomalies have been described<sup>1,11</sup> and previous work has indicated the need to explore new statistical methods to better detect the full suite of these types, notably those occurring over multiple, continuous observations (e.g., as a result of sensor drift), given that such anomalies commonly occur in water-quality time series yet are not as successfully detected as



other types of anomalies (e.g., single-point anomalies, out-of-range values).<sup>11</sup>

Following the framework provided by,<sup>11</sup> a range of technical anomalies in the time-series was first identified by local water-quality experts (i.e., for SC and PR by<sup>11</sup> and for AR by DS). For each site and variable (i.e., turbidity and conductivity), the identified anomalies were labelled along with their types based on anomaly classes defined by.<sup>11</sup> In the present work, we first focused on anomaly types of Class 1 defined as a single observation generating sudden changes in value from the previous observation. Class 1 anomalies were also considered a high priority in terms of detection, given that (true) sudden changes in turbidity and conductivity may be used as early warning signals of water quality by local environmental agencies and that they can strongly influence water-quality assessments and consequent management decisions. We additionally focused on Class 3 anomalies, which include technical anomalies such as long-term calibration offsets and changes comprised of multiple dependent observations. Class 3 was considered lower priority than Class 1 and may require a posteriori user intervention (i.e., after data collection rather than in real time) to confirm observations as anomalous. Specifically, using the nomenclature introduced by ref 11, we based our analyses of AD on those types of Class 1 anomalies defined as (a) large sudden spikes (type A), sudden shifts (type D), small sudden spikes (type J), and Class 3 anomalies defined as drift (type H), high variability (type E), and untrustworthy data not defined by other types (type L). For more information about definitions of each anomaly type and class see.<sup>11</sup>

The present analysis does not focus on AD of those types associated with impossible, out-of-sensor-range and missing values (Class 2) given that they can be easily detected by automated, hard-coded classification rules.<sup>11</sup> For each site, we therefore removed all Class 2 anomalies from the data prior to analysis. The filtered data were then log-transformed to remove exponential data variance (see Figure 1) to produce the time-series used for the presented analysis.

**Learning Methods and Data Processing.** One of the most challenging issues for AD in time-series data is that most data are usually “normal” or nonanomalous, while anomalous values are rare.<sup>8,9</sup> One of the most common discriminating approaches to solve such a problem is to compare the similarity between two sequences of time-dependent variables, specifically between the sequence of observed and predicted values.<sup>11,28</sup> When the anomalous values are prelabeled, we can apply semi-supervised classification based on training the learning process with the nonanomalous data, fitting models with prediction errors and predicting anomalous events as those observed values falling outside prediction intervals.<sup>11,28</sup> Supervised classification feeds the learning process with both a sequence of labeled values for AD (including both anomalous and nonanomalous values, tagged accordingly). The learning process can then generate a sequence of probabilities which can be binary-classified as anomalous or nonanomalous according to a predefined threshold parameter.

For each water-quality variable and site, time-series data were partitioned according to each learning process (i.e., semi-supervised or supervised classification). For semi-supervised classification, we retained the “normal” values (i.e., we discarded the anomalous values) and divided the time-series into four contiguous sequences of equal size: namely, two adjacent sequences of values for training and two other

adjacent sequences for validation (for details, see Figure S1 in the Supporting Information). Each pair of contiguous sequences for training or validation comprised one sequence for prediction (predictor variable) followed by an adjacent sequence for the outcome (outcome variable), respectively. For supervised classification, the sequence of values (including both anomalous and “normal” values) for prediction and the sequence of labeled values for AD were each divided into two adjacent sequences of equal length, for training and validation (for details, see Figure S1 in the Supporting Information). Thus, predictor and outcome variables were composed of one sequence for training followed by an adjacent sequence for validation, respectively. Given the proportionally low number of anomalies in our data set, we did not use a “test” data set during the optimization process.

Water quality in each of the studied sites undergoes intra- and interannual variation associated with local, seasonal and annual cycles and stochastic environmental events. We therefore incorporated within our analysis a “sliding window” or moving sequence of values of a defined length along the time series. For each site, the sliding window was defined according to the temporal resolution of data, and was 60, 90, and 10 min for PR, SC, and AR, respectively (see above). For both predictor and outcome variables, we constructed  $n \times p$  matrices with various time spans, defining the temporal resolution of time-series data and covering meaningful and regular environmental processes occurring in rivers (e.g., 24, 12, 6 h). For example, the matrices of  $n \times 1$  reflected the time span of 60, 90, and 10 min for PR, SC, and AR, respectively, whereas the matrices covering time spans of 24 h were defined by matrices of  $n \times 24$ ,  $n \times 16$ , and  $n \times 96$  for PR, SC, and AR, respectively. The differences in matrix  $n \times$  columns of each site are a consequence of their differences in the temporal resolutions defined above. For more information about the code and functions to allow data processing and matrix construction, see [https://github.com/benoit-liquet/AD\\_ANN](https://github.com/benoit-liquet/AD_ANN).

**Metrics for Model Performance and Evaluation.** To evaluate and compare the classification performance of AD, we calculated the four categories of the confusion matrix (i.e., true and false positives and true and false negatives; TP, FP, TN, FN, respectively) based on the discrimination threshold value  $t$  (typically 0.5; see below for detail on the threshold values we used). From these, we calculated accuracy (Acc), sensitivity (sn) and specificity (sp), and positive and negative predictive values (PPV and NPV, respectively), which allowed us to compare results directly with those of Leigh et al.<sup>11</sup>

$$\text{Acc} = \frac{\text{TP} + \text{TN}}{\text{TP} + \text{FP} + \text{TN} + \text{FN}}$$

$$\text{sn} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

$$\text{sp} = \frac{\text{TN}}{\text{TN} + \text{FP}}$$

$$\text{PPV} = \frac{\text{TP}}{\text{TP} + \text{FP}}$$

$$\text{NPV} = \frac{\text{TN}}{\text{TN} + \text{FN}}$$

Acc, sn, sp, NPV, and PPV range between 0 and 1. sn assesses the probability of determining TP correctly, whereas sp of

determining TN correctly. Acc assesses the ability to differentiate TP and TN correctly. Finally, PPV and NPV define the proportion of anomalous versus “normal” observations, specifically the negative and positive predictive values, respectively. When dealing with unbalanced classification problems (i.e., when there are far fewer anomalous than non-anomalous observations), many evaluation metrics are biased towards the majority class, maximizing TN classification while minimizing TP classification.<sup>29,30</sup> We therefore also calculated evaluation metrics specially formulated to provide an optimal classification for both positive and negative values in unbalanced data sets.<sup>30</sup> The first was balanced accuracy (b.Acc), which is defined as the arithmetic mean between the sn and sp values

$$\text{b.Acc} = \frac{\text{TP}/(\text{TP} + \text{FN}) + \text{TN}/(\text{TN} + \text{FP})}{2}$$

The second and third were the  $f_1$  and the Matthews correlation coefficient (MCC), defined as

$$f_1 = \frac{2 \times \text{TP}/(\text{TP} + \text{FP}) \times \text{TN}/(\text{TN} + \text{FP})}{\text{TP}/(\text{TP} + \text{FP}) + \text{TN}/(\text{TN} + \text{FP})}$$

$$\text{MCC} = \frac{\text{TP} \times \text{TN} - \text{FP} \times \text{FN}}{\sqrt{(\text{TP} + \text{FP}) \times (\text{TP} + \text{FN}) \times (\text{FP} + \text{TN}) \times (\text{TN} + \text{FN})}}$$

Specifically,  $f_1$  is the harmonic mean between the sn and PPV, ranging between 1 (i.e., perfect precision and recall) and 0. MCC ranges between −1 (i.e., total disagreement between predictions and observations), 0 (i.e., random prediction) and +1 (i.e., perfect prediction and balanced ratios of the four confusion matrix categories).

#### ANNs and Their Implementation in Time-Series AD.

ANN are machine learning methods, in which neurons (the basic unit of the learning process) are typically aggregated in layers that generate connections between input and output data. The typical structure of an ANN comprises input, hidden and output layers, usually connected sequentially. The input layer contains the observed data as variables (input neurons) and the output layer the predicted values (output or target neurons).<sup>31</sup> The input and output layers are connected by hidden layers, and the relationships between the neurons are set by nonlinear activation functions. For building and checking the performance of different types of model structure (i.e., the network structure typically associated with the number of layers and neurons), the ANN must be trained, by which the weights associated with connections between neurons are optimized using various methods and training algorithms.

A special case of ANNs adapted to time series applications are recurrent neural networks (RNNs), which can be considered a special case of auto-regressive integrated moving average (ARIMA) and nonlinear autoregressive moving average models.<sup>32</sup> RNN architecture mimics the cyclical connectivity of neurons, making them well suited for the analysis and prediction of nonstationary time series.<sup>33–35</sup> Long short-term memory (LSTM) networks are a redesign of the RNNs capable of learning long-term correlations in a sequence.<sup>36</sup> LSTMs are structured in network units known as memory blocks composed of self-connected memory cells and three multiplicative units, namely the “input”, “output” and “forget gates” connected to all cells within each memory

block.<sup>34</sup> Unlike RNNs, LSTM avoids the “vanishing gradients problem”, that is, when the error signal is used to train, meaning that the network exponentially decreases the further one goes backwards in the network, by means of the gates of the network units.<sup>37</sup> As a result, LSTMs allow the model to be trained successfully using backpropagation through time, which is key for accounting for long-term dependent time-series sequences. Although we applied LSTM in subsequent analyses, we use the more generic term ANN from this point onward.

In this study, ANNs were computed and fitted with “Keras”, a model-level library which provides high-level building for programming for developing deep-learning models. “Keras” allows the implementation of a wide variety of neural-network building blocks (e.g., layers, activation functions, optimizers) and supports the latest and most effective advances in deep network training (including RNNs).<sup>38</sup> “Keras” is a high-level wrapper of “TensorFlow” and helps to provide a simplified way of building neural networks from standard types of layers while facilitating a reproducible platform for developing deep learning approaches in computational environmental sciences.<sup>39</sup> “Keras” is written in “Python” and the “Keras” package provides an “R” interface to the deep-learning native functions.<sup>40</sup>

In each model run, we compiled a model with a predefined set of hyperparameters (see below), in which the internal model parameters were iteratively updated throughout the training steps (epochs; for definitions of ANNs hyperparameters, see Text S1 in the [Supporting Information](#)). During each model run, iterations were computed until the error from the model was minimized or reached a predefined value. In order to assess the performance of the learning process, we used a predefined set of standard metrics of “Keras” commonly used for classification and regression problems. Specifically, for both training and validation data sets, we calculated the loss function for semi-supervised classification and the mean-square error and the accuracy for supervised classification. For each model run, we accounted for overfitting using learning curves which showed how error changes as the training set size increases.<sup>37</sup> To do so, we compared the shape and dynamics of the learning curves of the training and validation data, which can be used to diagnose the bias and variance of the learning process. For instance, the “training learning curve” is calculated from the training data and provided information on the learning process, whereas the “validation learning curve” is calculated from the hold-out data and gives information on the model generality. For details of the dynamics of the learning curves for either semi-supervised and supervised classification with our data sets, see the [Supporting Information](#). SI contains the code and R functions to allow the implementation of ANNs using “Keras” in the “R” statistical language.<sup>41</sup>

**Optimization of HyperParameters and Their Influence on AD Performance.** ANNs include a broad suite of hyperparameters that affect the ability to learn patterns from the data and the performance of model predictions. Before training, we selected ten hyperparameters that affect (i) the network structure and (ii) the training algorithm; see the [Supporting Information](#) for details on hyperparameter definitions. Furthermore, we defined two additional hyperparameters related to (iii) the “sliding window” defined by  $n \times p$  matrices delimited at regular temporal intervals (see above) and (iv) the “threshold classification” as a measure of the

discrimination threshold value to compute the four categories of the confusion matrix (see above). For each hyperparameter we then defined a range of values (see details in the [Supporting Information](#)) that affect the performance of ANNs. The range of values depended on the type of the hyperparameter and varied from a continuous searching space (for those hyperparameters characterized by double precision, such as the learning rate) and discrete values (for those defined by integer values, such as number of layers or units, and functions or algorithms, such as the optimization algorithm). In the case of the “sliding window” we decided to use discrete values, defining the temporal resolution of the recorded time-series sequence and covering meaningful and regular environmental processes occurring in rivers (e.g., 24, 12, 6 h). For more information about this issue, see the [Supporting Information](#).

Given the large number of possible models to be tested with all of the combinations of hyperparameter values, we tuned the ANN models using a Bayesian optimization method, which is a common class of optimization methods, especially in deep-learning networks.<sup>42</sup> The Bayesian optimization method works by constructing a posterior distribution of functions (assuming a Gaussian process) associated with the variability of each hyperparameter, which best describes the “objective function”, defined as the “cost” associated with the optimization problem. With each iteration of the algorithm, the posterior distribution improves and the algorithm becomes more accurate in those regions of the parameter space with higher likelihood to maximize or minimize the objective function.<sup>42</sup> The Bayesian statistical model comprises two components: (i) a Bayesian statistical model for modelling the objective function, and (ii) an acquisition function for deciding where to sample next. In our case, we applied the mlrMBO toolbox implemented in the R statistical language. Compared with other black-box benchmark optimizers, the mlrMBO toolbox performs well for expensive optimization scenarios for single- and multiobjective optimization tasks, with continuous or mixed parameter spaces.<sup>43</sup>

Our aim for the optimization procedure was to find combinations of hyperparameters that resulted in the best performance for AD classification (see above). The optimization procedure followed two steps, firstly generating a random search space and secondly focusing on search shrinks, based on the results generated during the random search. First, we started the algorithm by generating a random design ( $n = 250$ ), which included a varying number of hyperparameters with the aim to generate enough variability to detect the most promising values. In our case we generated iterations by varying randomly the parameter combinations of the 12 hyperparameters. After computing the model and the optimization scores of each iteration, we secondly focused on search shrinks of the search space following a Bayesian optimization method with  $n = 250$  iterations based on maximizing performance metrics or objective functions. In our case, we followed a multiobjective Bayesian optimization method based on maximizing, in each  $k$  iteration run, the value of b.Acc,  $f_1$ , and MCC scores, defined above. Such an approach allowed us to maximize the different and complementary properties summarized by each score for unbalanced classification. Following the multiobjective Bayesian method, the optimization resulted in a total of  $n = 500$  iterations (i.e.,  $n = 250$  for random search and  $n = 250$  for

Bayesian optimization) for each combination of water-quality variables, sites, and learning procedures.

To examine the effects of hyperparameters predictor variables (hyperparameters) on the dependent variables (performance metrics), we applied methods for causal inference using random forests.<sup>44</sup> Specifically, we determined the statistical importance of each hyperparameter as a predictor on the dependent variables or optimization scores (e.g., b.Acc,  $f_1$ , or MCC) by calculating the “Variable Importance” (VI). VI reflects the model performance across the entire range of predictor and response variables, converted into a set of ordinal ranks. VI can be further used to test how the model response changes as the value of any of the predictor variables is changed, meaning that VI is similar to a standard “one-parameter-at-a-time” sensitivity analysis.<sup>45</sup> For each hyperparameter, we computed VI to measure the relative importance (or dependence, or contribution) of such hyperparameter predictor variables in terms of their effect on optimization scores. In our case, we computed the out-of-bag error, which is an error estimation technique used to evaluate the accuracy of a random forest after permuting each predictor variable. We used the R statistical language and the “randomForest” library.<sup>46</sup>

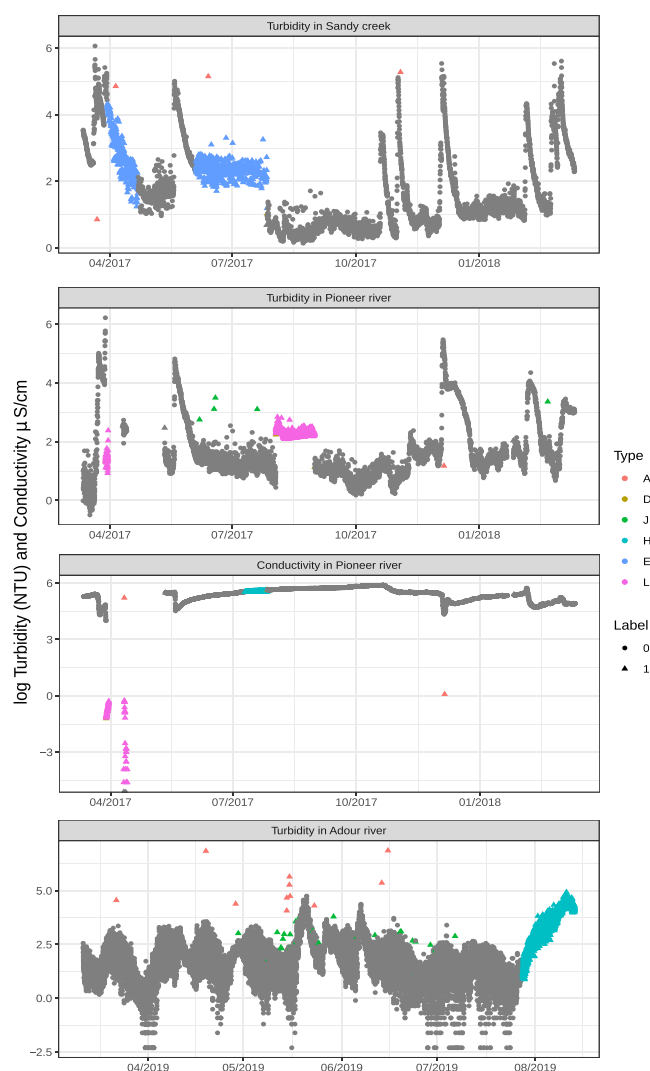
For each water-quality variable and site, we retained the “best” model as that which maximized the optimization scores. Specifically, for the complete set of candidate models, we averaged the value of b.Acc,  $f_1$ , or MCC and we retained the “best” model as that maximizing the averaged optimization scores. We compared the shape and dynamics of the learning curves from the “best” models (a) to diagnose whether or not the training and validation data sets were sufficiently representative (i.e., one data set could capture the statistical characteristics relative to other data sets) and (b) to test the behavior of the learning process (i.e., underfit, overfit, good fit).

## RESULTS

**Hyper-Parameters Optimization for AD.** The learning rate for AD stabilized early in the optimization of ANN models, with few improvements on the performance beyond  $n = 200$  model iterations (although see the ANN model for conductivity at SC). For semi-supervised classification, the costs of computing ranged from 0.312 h for turbidity in SC and 3.53 h for turbidity in AR, whereas for supervised classification the cost ranged from 1.11 h for conductivity in PR and 99.5 h turbidity in AR after  $n = 500$  model iterations (see Table S2 in the [Supporting Information](#)). Comparing learning methods, supervised classification had better performance and generated consistently higher values for b.Acc,  $f_1$ , or MCC, showing a similar and consistent pattern of the accumulative curve along the optimization process ([Figure 2](#)). Compared to supervised classification, semi-supervised classification required a larger number of model iterations for maximizing any of the three performance metrics, notably for turbidity in SC and AR.

Considering the whole suite of model iterations, we found a consistent pattern with respect to VI of those hyperparameters affecting optimization scores. Overall, the “Learning hyperparameters” had higher VI values than “Model hyperparameters” for model performance (see Tables S2–S4 in the [Supporting Information](#)). In addition, th.class had higher VI in semi-supervised classification, whereas s.win had minor VI for both supervised and semi-supervised classification.





**Figure 2.** Cumulative optimization scores for hyperparameter optimization along  $k$  model iterations. Each panel showed the procedure of multiobjective optimization procedure of (a) balanced accuracy, (b)  $f_1$ -score, and (c) MCC occurring in each  $k$  iteration run. The optimization begins searching with 250 random iterations, and then with 250 iterations of the Bayesian optimization procedure. Colors define each water-quality variable and site, whereas line patterns indicate the learning method. Abbreviations: Sandy creek (SC), Pioneer River (PR), and Adour River (AR); turbidity (Tur) and conductivity (Cond); semi-supervised classification (smSup) and supervised classification (Sup).

Specifically, the hyperparameters with higher VI for b.Acc,  $f_1$  and MCC, respectively, were th.class (78.2, 78.4, and 76.5, respectively; VI values averaged across sites, learning methods, and water-quality variables), dropout (46.2, 47.8, and 49.0), b.size (43.5, 48.4, and 51.8), momen (31.8, 34.3, and 36.2), lrate (31.3, 33.2, and 35.6). Comparing learning methods, we found that VI of hyperparameters differed depending on the combination of types of anomalies and water-quality variables. For instance, semi-supervised classification had higher VI values for th.class, dropout, momen, and lrate, whereas supervised classification for b.size and activ. Comparing sites, th.class had the highest VI values in semi-supervised classification for conductivity in PR, dropout and lrate in semi-supervised classification for turbidity in SC, b.size for supervised classification for turbidity in PR, momen for semi-

supervised classification for turbidity in AR and activ for supervised classification for turbidity in AR. For details of the VI of those hyperparameters independently affecting b.Acc,  $f_1$ , and MCC, see Tables S3–S5 in the [Supporting Information](#).

After hyperparameter optimization, the “best” models also had performed well in terms of b.Acc,  $f_1$ , and MCC scores, for all water-quality variables and learning methods; [Table 1](#); for details of hyperparameters and values of the best models, see the [Supporting Information](#). However semi-supervised classification had less balanced predictions (i.e., lower values for b.Acc,  $f_1$  score and MCC) than supervised classification, meaning that anomaly cases (positives) were proportionally less-correctly predicted than “normal” cases (negatives). Specifically, semi-supervised classification had lower rates of TP that we classified as true (sn = 0.652; values averaged across sites and water-quality variables) and lower proportions of positives and negatives that were true (PPV = 0.512), and that generated a moderately balanced detection rate for either TP and TN results (b.Acc = 0.757,  $f_1$  = 0.490 and MCC = 0.665). By contrast, supervised classification provided a higher and balanced detection rate for both TP and TN (b.Acc = 0.822,  $f_1$  = 0.622 and MCC = 0.762), and thus anomaly cases (positives) were as predicted correctly as “normal” cases (negatives). In contrast, supervised classification had higher detection rates of TP (sn = 0.704) and higher rates of correct classification of TP and TN (PPV = 0.643). For detailed information about the performance of each model, see [Table 1](#).

For “best” models, we additionally checked that our training and validation data sets were sufficiently representative, based on the learning process of either training or validation data. Overall, we found strong variation for each water-quality variable, site and learning method, suggesting that the presence of certain types of anomalies were not always consistent between training and validation data sets. Both for semi-supervised and supervised classification, we found that the validation data were easier to predict than training data (i.e., the “training learning curve” had poorer performance than the “validation learning curve”), suggesting that the validation data had lower complexity of anomaly types; for details of learning curves, see Figures S2–S5 in the [Supporting Information](#). For supervised classification, we additionally found that the validation data were unable to produce a good fit (i.e., both training and validation learning curves were almost flat), probably as a combination of the low number of anomaly cases (conductivity in PR; see [Figures 1](#) and S8 in the [Supporting Information](#)) and the presence of a long-term anomaly event at the end of the data (turbidity in AR; see [Figures 1](#) and S9 in the [Supporting Information](#)); the latter pattern did not happen when fitting semi-supervised classification in both data sets, for which the learning process produced a good fit for both training and validation processes. Overall, we did not detect overfitting of the training data, a result confirmed by the relatively medium-to-low values of dropout (i.e. <0.5) in most of the “best” models (see Tables S2–S4 in the [Supporting Information](#) for details of the values of hyperparameters of the “best” models).

#### AD among Types, Learning Methods, and Sites.

Anomaly types of Class 1 were present, but at low abundances, in all water-quality data sets, comprising, on average, 0.137% of cases ([Table 2](#)). By contrast, the majority of anomalies were classified as Class 3, which provided 11.6% of cases ([Table 2](#)). Anomalies of Class 3 were context-

**Table 1. Performance of the Best Model for AD<sup>a,b</sup>**

| site | Var  | Train | TN     | FN   | FP   | TP  | Acc   | sn    | sp    | PPV   | NPV   | b.Acc | $f_1$ | MCC  |
|------|------|-------|--------|------|------|-----|-------|-------|-------|-------|-------|-------|-------|------|
| SC   | Tur  | ArAD  | 4348   | 829  | 134  | 91  | 0.822 | 0.099 | 0.970 | 0.404 | 0.840 | 0.535 | 0.16  | 0.22 |
| SC   | Tur  | smSup | 1997   | 514  | 2485 | 405 | 0.445 | 0.441 | 0.446 | 0.140 | 0.795 | 0.443 | 0.21  | 0.40 |
| SC   | Tur  | Sup   | 4353   | 540  | 117  | 374 | 0.878 | 0.409 | 0.974 | 0.762 | 0.890 | 0.692 | 0.53  | 0.64 |
| PR   | Tur  | ArAD  | 5405   | 711  | 144  | 20  | 0.864 | 0.027 | 0.974 | 0.122 | 0.884 | 0.501 | 0.04  | 0.06 |
| PR   | Tur  | smSup | 3302   | 43   | 2180 | 684 | 0.642 | 0.941 | 0.602 | 0.239 | 0.987 | 0.772 | 0.38  | 0.69 |
| PR   | Tur  | Sup   | 5330   | 49   | 208  | 669 | 0.959 | 0.932 | 0.962 | 0.763 | 0.991 | 0.947 | 0.84  | 0.97 |
| PR   | Cond | ArAD  | 5705   | 448  | 56   | 71  | 0.920 | 0.137 | 0.990 | 0.559 | 0.927 | 0.564 | 0.22  | 0.29 |
| PR   | Cond | smSup | 2254   | 1    | 3410 | 480 | 0.445 | 0.998 | 0.398 | 0.123 | 1.000 | 0.698 | 0.22  | 0.60 |
| PR   | Cond | Sup   | 6091   | 0    | 69   | 48  | 0.989 | 1.000 | 0.989 | 0.410 | 1.000 | 0.994 | 0.58  | 0.65 |
| AR   | Tur  | ArAD  | 12,744 | 1573 | 452  | 38  | 0.863 | 0.024 | 0.966 | 0.078 | 0.890 | 0.495 | 0.04  | 0.05 |
| AR   | Tur  | smSup | 9540   | 1469 | 3654 | 142 | 0.654 | 0.088 | 0.723 | 0.037 | 0.867 | 0.406 | 0.05  | 0.07 |
| AR   | Tur  | Sup   | 13,017 | 925  | 179  | 495 | 0.924 | 0.349 | 0.986 | 0.734 | 0.934 | 0.668 | 0.47  | 0.55 |

<sup>a</sup>For each site and water-quality variable, details of the hyperparameters and their values of the best model are shown in the [Supporting Information](#). We calculated performance scores for ARIMA with AD (ArAD), semi-supervised ANNs (smSup), and supervised ANNs classification (Sup). For abbreviations of sites and variables see [Figure 2](#). <sup>b</sup>Abbreviations: True negatives (TN), false negatives (FN), false positives (FP), true positives (TP), accuracy (Acc), sensitivity (sn), specificity (sp), negative proportion of values (NPV), positive proportion of values (PPV), balanced accuracy (b.Acc),  $f_1$  score ( $f_1$ ), and Matthew's Correlation Coefficient (MCC).

**Table 2. Number of Anomalous Events According to Each Type (Columns), Site, and Water-Quality Variable (Rows)<sup>a,b</sup>**

| site | Var  | A  | D | J  | H    | E   | L   |
|------|------|----|---|----|------|-----|-----|
| SC   | Tur  | 4  | 1 | 0  | 0    | 914 | 0   |
| PR   | Tur  | 1  | 3 | 5  | 0    | 0   | 718 |
| PR   | Cond | 2  | 2 | 0  | 397  | 0   | 80  |
| AR   | Tur  | 11 | 0 | 26 | 1574 | 0   | 0   |

<sup>a</sup>The anomaly types shown here were classified by local water-quality experts, and their classification is detailed in the [Material and Methods](#) and in Leigh et al.<sup>11</sup> For abbreviations of sites, variables, and training see [Figure 2](#). <sup>b</sup>Abbreviations: large sudden spikes (A), sudden shifts (D), small sudden spikes (J), and Class 3 anomalies defined as drift (H), high variability (E), and untrustworthy data not defined by other types (L).

dependent with respect to each site and water-quality variable. Specifically, SC had two anomalous periods of high variability (type E) occurring during the first four months of monitoring (see [Figure 1](#)). PR had two drift sequences (type H) for turbidity and one period of untrustworthy data (type L) and one drift sequence (type H) for conductivity during the first

period of monitoring. Finally, AR had a long-drift sequence (type H) at the end of the monitoring period.

Comparing the performance of each “best” model with respect to detecting of the different types of anomalies, we found, on average, that semi-supervised classification had higher capacity for detecting Class 1 anomalies (45.7% vs 23.6% for semi-supervised and supervised classification, respectively; averaged values across sites and water-quality variables), whereas supervised classification had proportionally higher capacity for detecting Class 3 anomalies (72.6% vs 68.8%, respectively) ([Table 3](#)). However, such differences between learning methods were particularly context-dependent and based on the different combinations of types of anomalies at each site. Semi-supervised classification better detected large-sudden spikes (type A) for turbidity in AR, and small sudden spikes (J) for turbidity in PR and AR. Supervised classification, by contrast, had better performance for detecting drift (type H). Both semi-supervised and supervised classification performed well for detecting high variability (E) and untrustworthy anomalies for turbidity in SC and PR and conductivity in PR. Sudden shifts (type D) for turbidity in SC and drift (H) for conductivity in PR were not detected by any learning method.

**Table 3. Performance of the “Best” Model for AD by Anomaly Type<sup>a</sup>**

|    | site | Var  | Train | A     | D     | J     | H     | E     | L     |
|----|------|------|-------|-------|-------|-------|-------|-------|-------|
| 1  | SC   | Tur  | ArAD  | 1.000 | 0.000 |       |       | 0.094 |       |
| 2  | SC   | Tur  | smSup | 0.750 | 0.000 |       |       | 0.975 |       |
| 3  | SC   | Tur  | Sup   | 0.750 | 0.000 |       |       | 0.986 |       |
| 4  | PR   | Tur  | ArAD  | 1.000 | 1.000 | 1.000 |       |       | 0.010 |
| 5  | PR   | Tur  | smSup | 0.000 | 0.333 | 0.800 |       |       | 0.925 |
| 6  | PR   | Tur  | Sup   | 0.000 | 0.333 | 0.000 |       |       | 0.942 |
| 7  | PR   | Cond | ArAD  | 1.000 | 1.000 |       | 0.000 |       | 0.362 |
| 8  | PR   | Cond | smSup | 0.500 | 0.500 |       | 0.000 |       | 1.000 |
| 9  | PR   | Cond | Sup   | 0.500 | 0.500 |       | 0.000 |       | 0.975 |
| 10 | AR   | Tur  | ArAD  | 1.000 |       | 0.846 | 0.003 |       |       |
| 11 | AR   | Tur  | smSup | 1.000 |       | 0.231 | 0.541 |       |       |
| 12 | AR   | Tur  | Sup   | 0.000 |       | 0.038 | 0.728 |       |       |

<sup>a</sup>For each site and water-quality variable, values represent the percentage of data values detected relative to the total number of respective anomalies labeled in the data set (see [Table 3](#) for details). For abbreviations of sites, variables, and training see [Figure 2](#) and of anomaly types see [Table 2](#).



## ■ DISCUSSION

In this work, we found that ANN provided good classification for AD in high-frequency water-quality data, given their ability to deal flexibly with the challenges associated with local environmental variability. We tested the capacity of ANNs for AD under a broad range of variables, anomaly types, and real-world conditions. Our data come from separate monitoring programs in different parts of the world and under contrasting environmental systems (i.e., estuarine, freshwater), so our work presents an opportunity to test the robustness and performance of ANNs for AD. We used turbidity and conductivity data, which are commonly measured by water management agencies and monitoring programs, providing an avenue to test other water quality and quantity variables in the future, which may lead to an overall increase in the performance of water-quality monitoring systems. Results of the AD showed that semi-supervised classification was able to cope better with short-term anomalies associated with a single observation or time point than supervised classification, which showed improved performance for detecting anomalies dependent on multiple context-dependent observations.

Previous work has shown that regression time-series methods are useful for AD in stationary and nonstationary time-series data sets,<sup>47,48</sup> but the detection of certain anomaly types, such as sensor drift and periods of anomalously high variability, remain challenging. ANNs are a versatile method that can train models using different learning methods, and as shown by our study, can provide the flexibility required to detect a broad suite of anomaly types, including improved performance for detecting extended periods of untrustworthy data. In our case, we applied regression-based and semi-supervised ANNs similarly for AD; that is, models first predict data sequences based on “normal” cases and then classify as anomalous cases those departing from a given threshold value. Despite the good performance of ANNs for AD, our findings demonstrate that ANNs could have limitations regardless of the underlying statistical method used and their applicability in near-real time AD. In our study, ca. 100 model iterations, which were necessary to obtain acceptable model performances, required three orders of magnitude of computing time longer than regression-based ARIMA. Compared to these regression-based time-series methods, the Bayesian optimization of hyperparameter values allowed us to optimize classification of anomalies, at the expense of increased computing costs for hyperparameter optimization.

For the Australian sites, we found that the performance of semi-supervised ANNs (proposed here) and regression-based ARIMA (proposed by ref 11) were equivalent in terms of correct classification, notably by providing high false-detection rates (falses) of both anomalies (positives) and “normal” cases (negatives): semi-supervised classification had higher rates of FP and FN which resulted in low values of sn and PPV. We also found that supervised ANNs considerably minimized false detection rates (i.e., low values of FP and FN and maximized sn and PPV values), when compared with the regression-based models and semi-supervised ANNs. Semi-supervised ANNs and regression-based ARIMA generated higher rates of false alarms (i.e., both FP and FN), which overestimate anomalous events, than supervised ANNs.

We did not detect over-training in the “best” models, but after checking the shape and dynamics of the learning curves we found that the training and validation data sets had

inconsistent numbers and presence of anomaly types. For instance, the validation data set in turbidity in AR had a long-term anomaly event which is not present in the training data (see Figure 1). As a result, there is inconsistency between training and validation data sets in the performance, typically the validation data had higher performance than the training data (see Figures S2–S5 in the Supporting Information). The difference in complexity was mainly a consequence of the presence of Class 3 anomalies in the data sets, usually occurring as single long-term anomalous events in each time series. Cross-validation could solve this problem, but it is also true that multiple partitioning of the data sets could have divided the time-series into multiple smaller independent time-series sequences with inconsistent representation of anomaly types among data sets. Multiple partitioning of the time-series sequence could be especially critical during the optimization process, notably on those hyperparameters tuning the size of the subsamples processed during the learning process (i.e., *b.size* and *momen*, and also for *s.win*). All of the above issues were a consequence of the proportionally lower abundance of anomalous events compared to “normal” events in the data sets, a phenomenon that makes the detection of anomalies challenging.

We found that semi-supervised ANNs were especially suited to AD of short-term anomalous events (i.e., sudden spikes, sudden shifts and small sudden spikes defined as Class 1 anomalies, following the terminology of ref 11). Although semi-supervised and supervised ANNs were comparable regarding their ability to detect long-term anomalous events (i.e., drift and periods of high variability or Class 3 anomalies), the latter had better performance in terms of correct AD (see above). At least for AD of long-term events, our results with ANNs thus outperformed those of regression-based ARIMA,<sup>11</sup> which is an important step forward for AD in water-quality data given such anomalies have consistently proved challenging to detect by other automated AD methods. This is because such anomalous, often very context-dependent events can behave similarly to natural, nonstationary water-quality events such that they are only detected manually by a trained eye very familiar with the local conditions. Our novel use of ANNs and Bayesian optimization for hyperparameter selection therefore holds much promise for AD in high-frequency water-quality data from a broad range of environments and ecosystems, including rivers, estuaries, and marine waters.

Although our study is based on the analysis of data from three sites only, it provides a methodological framework for AD in high-frequency water-quality data collected from sites with contrasting nonstationary environmental processes.<sup>49,50</sup> We found that the nonstationary environmental conditions of each site played a substantial role in both (i) determining the types of anomalies present and (ii) increasing the uncertainty around the accurate detection of anomalies. The water-quality variables analyzed here are affected by long-term environmental processes occurring at each site (i.e., seasonal precipitation patterns in both Australian and French sites), as well as regular and short-term events (e.g., cyclone floods in Australian sites and tidal regimes in the French site). Both short- and long-term environmental processes may interact and this could affect the AD performance. For Australian sites, the detection of sudden spikes and shifts (i.e., short-term anomalies) and of drifts or periods with high variability (i.e., long-term anomalies) were rarely masked by natural environmental processes (i.e., seasonal rainfall patterns or cyclone

floods), and that resulted in the optimal classification of short- and long-term anomalous events by using either semi-supervised and supervised classification, respectively (Table 3). At the French site, by contrast, strong tidal regimes occurred alongside medium-to-high seasonal discharge events, and that conditioned and limited the accuracy of detecting both short- and long-term anomalies (Table 3).

In this work, we calibrated ANNs models using Bayesian optimization of hyperparameter values, by means of multi-objective optimization. Although optimization methodology is well-established for the detection of anomalies in water-quality data,<sup>7</sup> multi-objective optimization procedures have rarely been applied for AD in time-series analysis for either prediction or detection anomalies using ANNs methods; but see Perelman et al.<sup>19</sup> for an application of detecting anomalies using a single-objective optimization. The application of multiobjective methodology allows us to get a balanced performance of models for detecting anomalies in a broad set of environmental conditions. Notwithstanding the potential utility of ANNs for AD demonstrated above, there is substantial room for improvement in the performance under real-world conditions.<sup>12</sup> First, we found that the ability for AD is context-dependent, meaning that accuracy is conditioned on the spatiotemporal environmental variability of the data set available. Anomalies are rare events, meaning that we need large data sets spanning the entire environmental variability of each locality to ensure their inclusion in a data set.<sup>11</sup> Second, we performed our detections based on the analysis of independent environmental variables (i.e., turbidity or conductivity), but the application of ANNs with multivariate time-series could enable us to account for the temporal correlation of multiple variables monitored at the same time.<sup>19,24,51</sup> Finally, methodological improvements provide additional avenues to increase the performance of ANNs in AD, such as Bayesian RNNs,<sup>52,53</sup> which allow quantification of the uncertainty and the use of ensemble averaging for ANNs, combining unsupervised and supervised classification<sup>54</sup> or the combination of ANNs with other methodologies.<sup>6,55</sup>

The study of AD in high-frequency data has numerous applications in environmental and health monitoring, and fault and fraud detection, where there is a need to provide near-real time solutions and optimal performance for monitoring and to ensure the quality of data streaming. We have demonstrated that ANNs are a flexible method for providing optimal performance for AD, given they are able to cope with both long- and short-term nonstationary processes that condition high-frequency data. However, ANNs and their hyperparameter optimization have been understudied in the context of AD in water-quality. Given the promising results from our study, we therefore recommend further investigation and development of such methods to improve the accurate detection of a broad suite using machine learning or deep-learning methods of anomalies detection under a wide range of environmental conditions.<sup>7</sup> Environmental data are intrinsically variable though time and space and thus our approach is transferable for AD in complex spatiotemporal applications and ecosystems. Our findings will therefore be of relevance to water scientists and managers throughout the world in order to broaden the applicability of ANNs for efficient water monitoring.

## ■ ASSOCIATED CONTENT

### Supporting Information

The Supporting Information is available free of charge at <https://pubs.acs.org/doi/10.1021/acs.est.0c04069>.

Scheme of the data processing for time-series AD for each learning method; brief definition of the hyperparameters and range of hyperparameter values for fitting and optimizing ANN models; costs of computing time during the learning process; VI of values of hyperparameters to maximize the biased accuracy, the  $f_1$ -score, and the MCC; learning curves for semi-supervised and supervised learning processes for the “best” ANN models after Bayesian optimization; and data of observed trend and the probability of AD for the “best” ANN models after Bayesian optimization (PDF)

## ■ AUTHOR INFORMATION

### Corresponding Author

**Kerrie Mengersen** – Univ. Pau & Pays de l'Adour E2S UPPA, Laboratoire des Mathématiques et de leurs applications, CNRS, 64600 Anglet, France; Australian Research Council Centre of Excellence for Mathematical and Statistical Frontiers (ACEMS), 4000 Brisbane, Australia; School of Mathematical Sciences, Queensland University of Technology, 4000 Brisbane, Australia; Email: [k.mengersen@qut.edu.au](mailto:k.mengersen@qut.edu.au)

### Authors

**Javier Rodriguez-Perez** – Univ. Pau & Pays de l'Adour E2S UPPA, Laboratoire des Mathématiques et de leurs applications, CNRS, 64600 Anglet, France

**Catherine Leigh** – Biosciences and Food Technology Discipline, School of Science, RMIT University, 3000 Bundoora, Australia; Australian Research Council Centre of Excellence for Mathematical and Statistical Frontiers (ACEMS), 4000 Brisbane, Australia; Institute for Future Environments, Queensland University of Technology, 4000 Brisbane, Australia

**Benoit Liquet** – Univ. Pau & Pays de l'Adour E2S UPPA, Laboratoire des Mathématiques et de leurs applications, CNRS, 64600 Anglet, France; Department of Mathematics and Statistics, Macquarie University, 2109 Sydney, Australia

**Claire Kermorvant** – Univ. Pau & Pays de l'Adour E2S UPPA, Laboratoire des Mathématiques et de leurs applications, CNRS, 64600 Anglet, France; Australian Research Council Centre of Excellence for Mathematical and Statistical Frontiers (ACEMS), 4000 Brisbane, Australia; [orcid.org/0000-0002-1972-8937](https://orcid.org/0000-0002-1972-8937)

**Erin Peterson** – Australian Research Council Centre of Excellence for Mathematical and Statistical Frontiers (ACEMS), 4000 Brisbane, Australia; Institute for Future Environments and School of Mathematical Sciences, Queensland University of Technology, 4000 Brisbane, Australia

**Damien Sous** – Université de Toulon, Aix Marseille Université, CNRS, IRD, Mediterranean Institute of Oceanography (MIO), 83062 La Garde, France; Univ. Pau & Pays de l'Adour E2S UPPA, Chaire HPC Waves, Laboratoire des Sciences de l'Ingénieur Appliquées à la Mécanique et au Génie Electrique – Fédération IPRA, EA4581, 64600 Anglet, France

Complete contact information is available at: <https://pubs.acs.org/doi/10.1021/acs.est.0c04069>

### Notes

The authors declare no competing financial interest.

## ■ ACKNOWLEDGMENTS

Funding was provided by the Energy Environment Solutions (E2S-UPPA) consortium and the BIGCEES project from E2S-UPPA ("Big model and Big data in Computational Ecology and Environmental Sciences"), the Queensland Department of Environment and Science (DES) and the ARC Centre of Excellence for Mathematical and Statistical Frontiers (ACEMS). A repository of the water-quality data from the in situ sensors used herein and the code used to implement methods of ANNs for AD are provided in the [Supporting Information](#).

## ■ REFERENCES

- (1) Horsburgh, J. S.; Spackman Jones, A.; Stevens, D. K.; Tarboton, D. G.; Mesner, N. O. A sensor network for high frequency estimation of water quality constituent fluxes using surrogates. *Environ. Model. Software* **2010**, *25*, 1031–1044.
- (2) Rode, M.; Wade, A. J.; Cohen, M. J.; Hensley, R. T.; Bowes, M. J.; Kirchner, J. W.; Arhonditsis, G. B.; Jordan, P.; Kronvang, B.; Halliday, S. J.; Skeffington, R. A.; Rozemeijer, J. C.; Aubert, A. H.; Rinke, K.; Jomaa, S. Sensors in the Stream: The High-Frequency Wave of the Present. *Environ. Sci. Technol.* **2016**, *50*, 10297–10307.
- (3) Hill, D. J.; Minsker, B. S. Anomaly detection in streaming environmental sensor data: A data-driven modeling approach. *Environ. Model. Software* **2010**, *25*, 1014–1022.
- (4) Horsburgh, J. S.; Reeder, S. L.; Jones, A. S.; Meline, J. Open source software for visualization and quality control of continuous hydrologic and water quality sensor data. *Environ. Model. Software* **2015**, *70*, 32–44.
- (5) Jiang, J.; Wang, P.; Lung, W.-s.; Guo, L.; Li, M. A GIS-based generic real-time risk assessment framework and decision tools for chemical spills in the river basin. *J. Hazard. Mater.* **2012**, *227*–228, 280–291.
- (6) Shi, B.; Wang, P.; Jiang, J.; Liu, R. Applying high-frequency surrogate measurements and a wavelet-ANN model to provide early warnings of rapid surface water quality anomalies. *Sci. Total Environ.* **2018**, *610*–611, 1390–1399.
- (7) Dogo, E. M.; Nwulu, N. I.; Twala, B.; Aigbavboa, C. A survey of machine learning methods applied to anomaly detection on drinking-water quality data. *Urban Water J.* **2019**, *16*, 235–248.
- (8) Chandola, V.; Banerjee, A.; Kumar, V. *ACM Comput. Surv.* **2009**, *41*, 15.
- (9) Gupta, M.; Gao, J.; Aggarwal, C.; Han, J. Outlier Detection for Temporal Data. *Synth. Lect. Data Min. Knowl. Discov.* **2014**, *5*, 1–129.
- (10) Goldstein, M.; Uchida, S. A Comparative Evaluation of Unsupervised Anomaly Detection Algorithms for Multivariate Data. *PLoS One* **2016**, *11*, No. e0152173.
- (11) Leigh, C.; Alsibai, O.; Hyndman, R. J.; Kandanaarachchi, S.; King, O. C.; McGree, J. M.; Neelamraju, C.; Strauss, J.; Talagala, P. D.; Turner, R. D. R.; Mengersen, K.; Peterson, E. E. A framework for automated anomaly detection in high frequency water-quality data from in situ sensors. *Sci. Total Environ.* **2019**, *664*, 885–898.
- (12) Muharemi, F.; Logofătu, D.; Leon, F. Machine learning approaches for anomaly detection of water quality on a real-world data set. *J. Inf. Telecommun.* **2019**, *3*, 294–307.
- (13) Bourgeois, W.; Romain, A.-C.; Nicolas, J.; Stuetz, R. M. The use of sensor arrays for environmental monitoring: interests and limitations. *J. Environ. Monit.* **2003**, *5*, 852–860.
- (14) Shipmon, D. T.; Gurevitch, J. M.; Piselli, P. M.; Edwards, S. T. Time Series Anomaly Detection; Detection of anomalous drops with limited features and sparse examples in noisy highly periodic data. **2017**, arXiv:1708.03665. arXiv preprint.
- (15) Makarynsky, O.; Makarynska, D.; Kuhn, M.; Featherstone, W. E. Predicting sea level variations with artificial neural networks at Hillarys Boat Harbour, Western Australia. *Estuarine, Coastal Shelf Sci.* **2004**, *61*, 351–360.
- (16) Makarynska, D.; Makarynsky, O. Predicting sea-level variations at the Cocos (Keeling) Islands with artificial neural networks. *Comput. Geosci.* **2008**, *34*, 1910–1917.
- (17) Wu, W.; Dandy, G. C.; Maier, H. R. Protocol for developing ANN models and its application to the assessment of the quality of the ANN model development process in drinking water quality modelling. *Environ. Model. Software* **2014**, *54*, 108–127.
- (18) Tinelli, S.; Juran, I. Artificial intelligence-based monitoring system of water quality parameters for early detection of non-specific bio-contamination in water distribution systems. *Water Supply* **2019**, *19*, 1785–1792.
- (19) Perelman, L.; Arad, J.; Housh, M.; Ostfeld, A. Event Detection in Water Distribution Systems from Multivariate Water Quality Time Series. *Environ. Sci. Technol.* **2012**, *46*, 8212–8219.
- (20) Muharemi, F.; Logofătu, D.; Andersson, C.; Leon, F. *Modern Approaches for Intelligent Information and Database Systems*; Springer, 2018; pp 173–183.
- (21) Fehst, V.; La, H. C.; Nghiem, T.-D.; Mayer, B. E.; Englert, P.; Fiebig, K.-H. Automatic vs. manual feature engineering for anomaly detection of drinking-water quality. *Proceedings of the Genetic and Evolutionary Computation Conference Companion*, 2018; pp 5–6.
- (22) Brodie, J. *ACTFR Technical Report No. 02/03*; Australian Centre for Tropical Freshwater Research, James Cook University, 2004.
- (23) Defontaine, S.; Sous, D.; Morichon, D.; Verney, R.; Monperrus, M. Hydrodynamics and SPM transport in an engineered tidal estuary: The Adour river (France). *Estuarine, Coastal Shelf Sci.* **2019**, *231*, 106445.
- (24) Leigh, C.; Kandanaarachchi, S.; McGree, J. M.; Hyndman, R. J.; Alsibai, O.; Mengersen, K.; Peterson, E. E. *PLoS One* **2019**, *14*, No. e0215503.
- (25) Leigh, C.; Burford, M.; Connolly, R.; Olley, J.; Saeck, E.; Sheldon, F.; Smart, J.; Bunn, S. Science to Support Management of Receiving Waters in an Event-Driven Ecosystem: From Land to River to Sea. *Water* **2013**, *5*, 780–797.
- (26) Wagner, R. J.; Boulger, R. W., Jr.; Oblinger, C. J.; Smith, B. A. *Guidelines and Standard Procedures for Continuous Water-Quality Monitors: Station Operation, Record Computation, and Data Reporting*, 2006.
- (27) Cox, B. A review of currently available in-stream water-quality models and their applicability for simulating dissolved oxygen in lowland rivers. *Sci. Total Environ.* **2003**, *314*–316, 335–377.
- (28) Malhotra, P.; Vig, L.; Shroff, G.; Agarwal, P. Long short term memory networks for anomaly detection in time series. *ESANN 2015 Proceedings, European Symposium on Artificial Neural Networks, Computational Intelligence and Machine Learning*, 2015; p 89.
- (29) Kelleher, J. D.; Mac Namee, B.; D'arcy, A. *Fundamentals of Machine Learning for Predictive Data Analytics: Algorithms, Worked Examples, and Case Studies*; MIT Press, 2015.
- (30) Boughorbel, S.; Jarray, F.; El-Anbari, M. Optimal classifier for imbalanced data using Matthews Correlation Coefficient metric. *PLoS One* **2017**, *12*, No. e0177678.
- (31) Haykin, S. *Neural Networks and Learning Machines*, 3rd ed.; Pearson Education India, 2010.
- (32) Siami-Namini, S.; Tavakoli, N.; Namin, A. S. A Comparison of ARIMA and LSTM in Forecasting Time Series. *2018 17th IEEE International Conference on Machine Learning and Applications (ICMLA)*, 2018; pp 1394–1401.
- (33) Hilas, C. S.; Mastorocostas, P. A. An application of supervised and unsupervised learning approaches to telecommunications fraud detection. *Knowl. Base Syst.* **2008**, *21*, 721–726.
- (34) Graves, A. *Supervised Sequence Labelling with Recurrent Neural Networks*; Springer, 2012; pp 37–45.
- (35) Maier, H. R.; Dandy, G. C. Neural networks for the prediction and forecasting of water resources variables: a review of modelling issues and applications. *Environ. Model. Software* **2000**, *15*, 101–124.
- (36) Hochreiter, S.; Schmidhuber, J. Long Short-Term Memory. *Neural Comput.* **1997**, *9*, 1735–1780.



- (37) Goodfellow, I.; Bengio, Y.; Courville, A. *Deep Learning*; MIT press, 2016.
- (38) Gulli, A.; Pal, S. *Deep Learning with Keras*; Packt Publishing Ltd., 2017.
- (39) Rampasek, L.; Goldenberg, A. TensorFlow: Biology's Gateway to Deep Learning? *Cell Syst.* **2016**, *2*, 12–14.
- (40) Allaire, J.; Chollet, F. *R Package*, version, 2017.
- (41) *R Core Team*, 2017.
- (42) Snoek, J.; Larochelle, H.; Adams, R. P. Practical bayesian optimization of machine learning algorithms. *Adv. Neural Inf. Process. Syst.* **2012**, *4*, 2951–2959.
- (43) Bischl, B.; Richter, J.; Bossek, J.; Horn, D.; Thomas, J.; Lang, M. mlrMBO: A Modular Framework for Model-Based Optimization of Expensive Black-Box Functions. **2017**, arXiv:1703.03373. arXiv preprint.
- (44) Wager, S.; Athey, S. Estimation and Inference of Heterogeneous Treatment Effects using Random Forests. *J. Am. Stat. Assoc.* **2018**, *113*, 1228–1242.
- (45) Mishra, S.; Datta-Gupta, A. *Applied Statistical Modeling and Data Analytics: A Practical Guide for the Petroleum Geosciences*; Elsevier, 2017.
- (46) Breiman, L. *Mach. Learn.* **2001**, *45*, 5–32.
- (47) Hyndman, R. J.; Athanasopoulos, G. *Forecasting: Principles and Practice*; OTexts, 2018.
- (48) Box, G. E.; Jenkins, G. M.; Reinsel, G. C.; Ljung, G. M. *Time Series Analysis: Forecasting and Control*; John Wiley & Sons, 2015.
- (49) Clarke, R. T. Hydrological prediction in a non-stationary world. *Hydrology and Earth System Sciences Discussions*, European Geosciences Union: 2007, *11*, pp. 408–414.
- (50) Sivakumar, B. *Chaos in Hydrology*; Springer, 2017; pp 29–62.
- (51) Sánchez-Fernández, A.; Baldán, F. J.; Sainz-Palmero, G. I.; Benítez, J. M.; Fuente, M. J. Fault detection based on time series modeling and multivariate statistical process control. *Chemom. Intell. Lab. Syst.* **2018**, *182*, 57–69.
- (52) Mirikitani, D. T.; Nikolaev, N. *IEEE Trans. Neural Network.* **2009**, *21*, 262–274.
- (53) Sun, W.; Paiva, A. R.; Xu, P.; Sundaram, A.; Braatz, R. D. Fault Detection and Identification using Bayesian Recurrent Neural Networks. **2019**, arXiv:1911.04386. arXiv preprint.
- (54) Comar, P. M.; Liu, L.; Saha, S.; Tan, P.-N.; Nucci, A. Combining supervised and unsupervised learning for zero-day malware detection. *2013 Proceedings IEEE INFOCOM*, 2013; pp 2022–2030.
- (55) Dairi, A.; Cheng, T.; Harrou, F.; Sun, Y.; Leiknes, T. Deep learning approach for sustainable WWTP operation: A case study on data-driven influent conditions monitoring. *Sustain. Cities Soc.* **2019**, *50*, 101670.