



HAL
open science

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, Kerrie L. Mengersen

► **To cite this version:**

Javier Rodriguez Perez, Catherine Leigh, Benoit Liquet, Claire Kermorvant, Erin Peterson, et al..
Detecting technical anomalies in high-frequency water-quality data using Artificial Neural Networks.
Environmental Science and Technology, 2020, 54 (21), pp.13719-13730. 10.1021/acs.est.0c04069 .
hal-02929461

HAL Id: hal-02929461

<https://hal.science/hal-02929461v1>

Submitted on 27 Feb 2021

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

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^{*,†,§,#}

[†]*Université de Pau et des Pays de l'Adour, E2S UPPA, CNRS, Pau, France*

[‡]*Institute for Multidisciplinary Research in Applied Biology, Departamento Ciencias del
Medio Natural, UPNA Pamplona, Spain*

[¶]*Biosciences and Food Technology Discipline, School of Science, RMIT University,
Bundoora, Australia*

[§]*Australian Research Council Centre of Excellence for Mathematical and Statistical
Frontiers (ACEMS), Australia*

^{||}*Institute for Future Environments, Queensland University of Technology, Brisbane,
Australia*

[⊥]*Department of Mathematics and Statistics, Macquarie University, Sydney, Australia*

[#]*School of Mathematical Sciences, Queensland University of Technology, Brisbane,
Australia*

[@]*Université de Toulon, Aix Marseille Université, CNRS, IRD, Mediterranean Institute of
Oceanography (MIO), La Garde, France*

E-mail: k.mengersen@qut.edu.au

Abstract

3 Anomaly detection (AD) in high-volume environmental data requires one to tackle
4 a series of challenges associated with the typical low frequency of anomalous events,
5 the broad-range of possible anomaly types and local non-stationary environmental con-
6 ditions, suggesting the need for flexible statistical methods that are able to cope with
7 unbalanced high-volume data problems. Here, we aimed to detect anomalies caused by
8 technical errors in water-quality (turbidity and conductivity) data collected by auto-
9 mated *in-situ* sensors deployed in contrasting riverine and estuarine environments. We
10 first applied a range of Artificial Neural Networks (ANN) that differed in both learning
11 method and hyper-parameter values, then calibrated models using a Bayesian multi-
12 objective optimisation procedure, and selected and evaluated the "best" model for each
13 water-quality variable, environment and anomaly type. We found that semi-supervised
14 classification was better able to detect sudden spikes, sudden shifts and small sudden
15 spikes whereas supervised classification had higher accuracy for predicting long-term
16 anomalies associated with drifts and periods of otherwise unexplained high-variability.

17 Introduction

18 Monitoring the water quality of rivers is becoming increasingly relevant in the Anthropocene
19 as the need for management strategies to safeguard water resources in human dominated
20 landscapes accelerates. With the advent of new technologies, data on water properties can
21 be obtained with high frequency in near-real time, which can facilitate fast and adaptive
22 management strategies.^{1,2} However, the generation of high-volume, high-velocity data can
23 create problems of quality control due to a combination of (1) technical issues with the
24 water-quality sensors that result in technical anomalies, and (2) operational problems with
25 traditional, often manual and periodic, methods of quality control that are no longer feasible.
26 In order to improve the quality and utility of high-volumes of data, we therefore need to
27 explore efficient methods to complement or replace traditional methods of quality control,
28 including manual anomaly detection and correction and quality coding.^{3,4} Improving water-

29 quality monitoring will also require robust and flexible statistical methods that can handle
30 data streaming, while also providing big data tools for data-driven decision making.⁵⁻⁷

31 The accurate detection of anomalies in high-frequency, high-volume water-quality data
32 is additionally challenged by several factors. For example, the common low frequency of
33 anomalies in water-quality time series constrains the ability to not only detect anomalies
34 but also to evaluate the performance of detection methods.^{8,9} There is also a broad range of
35 possible anomaly types, ranging from sudden changes in water-quality values to anomalously
36 constant values maintained for long periods.^{10,11} In addition, the differentiation of technical
37 anomalies, e.g. related to sensor malfunction or database failure, from the real deviation of
38 the natural system from its typical behaviour is not always straightforward.¹ This is because
39 natural fluctuations in environmental conditions at a given site can influence the level of
40 detectability of sensor-related anomalies from the normal signal: the greater the fluctuation,
41 the harder the detection of technical errors. All of these issues can constrain the generality
42 and robustness of detection methods under a broad set of environmental conditions.^{10,12}

43 A combination of methods may be required to address the challenges raised above. For
44 instance, sudden spikes and shifts in values and impossible and/or out-of-sensor-range val-
45 ues on single or just a few observation points can be easily differentiated from adjacent,
46 non-anomalous values. Despite the extremely low frequency of such point-based anoma-
47 lies in time-series data, rule-based, regression-based and feature-based methods may provide
48 optimal performance in terms of their detection.^{1,11} Common anomalies in chemical and bio-
49 chemical data from *in-situ* sensors also comprise multiple observations, such as sensor drift
50 and periods of unusually high or low variability, which may indicate the need for sensor cal-
51 ibration or maintenance.¹³ The detection of these multiple-points anomalies can be tackled
52 by various methods,^{7,11} but remains challenging and often still requires user intervention.
53 Several methods can be combined to extend the range of detected anomalies. Recent work
54 by Leigh et al.¹¹ showed that a combination of rule-, feature and regression-based methods
55 facilitated the correct classification of impossible values, sudden isolated spikes and level

56 shifts, although drift and periods of high variability still tended to be associated with high
57 rates of false positives.

58 Considering the current global effort to produce real-time, high-frequency and long-
59 term monitoring of aquatic ecosystems, and the limitations of existing quality control ap-
60 proaches,^{1,4,13} there is a pressing need to explore new methods. Artificial Neural Networks
61 (ANNs) are a promising alternative.^{6,7,14} One basic strength of ANNs in analysing high-
62 frequency time-series data is that they do not require *a priori* knowledge of the underly-
63 ing physical and environmental processes. ANNs have therefore been used tackle multiple
64 problems of data streaming including predicting temporal patterns in both marine^{15,16} and
65 freshwater systems.^{6,17,18} They have potential to detect anomalies and sudden changes that
66 occur over short-time intervals in water-quality data,^{6,18,19} suggesting that they can be read-
67 ily applied to optimise anomaly detection in complex and fluctuating environmental systems.
68 Previous work on AD applied to water-quality variables compared ANNs to other machine
69 learning methods and found that ANNs have comparable performance to other methods in
70 terms of AD.^{12,18,20,21} However, authors also acknowledge that remains understudied the op-
71 timisation techniques (especially on those machine learning methods whose learning process
72 depends on a large number of hyper-parameters) that aim to improve detection capacity
73 under real-world conditions.^{7,20}

74 The aim of this study is to test the ability of ANNs to detect technical anomalies in
75 water-quality time-series data in contrasting riverine and estuarine environments. We were
76 especially interested in testing the ability of ANNs to detect multiple-point anomalies, for
77 example those that may result from sensor drift, given that methods previously examined
78 were limited in their ability to accurately detect such context-dependent events. The studied
79 variables are turbidity and conductivity collected at high-frequency by autonomous *in-situ*
80 sensors. The monitored sites in subtropical Australia and temperate France provide con-
81 trasting ranges of environmental conditions, particularly in terms of the studied variables
82 and their magnitudes and dynamics. The comparison and evaluation of detection methods

83 in regions that vary in environmental conditions and across a common range of variables
84 is essential to determine the methods' broad suitability for anomaly detection.⁷ Given that
85 the environmental fluctuations in each river (freshwater and estuarine sites) may influence
86 the ability to detect the different types of anomalies, we calibrated ANNs using models
87 that differed in learning method and hyper-parameter values. To do so, we implemented a
88 Bayesian optimisation method using a multi-objective approach (i.e. based on scores spe-
89 cially suited to unbalanced classification problems), in order to find the combination of ANN
90 hyper-parameters best suited for anomaly detection at each location. The performance of the
91 calibrated ANNs in detecting different types of anomalies was also evaluated and compared
92 to regression-based methods developed for and applied to similar data.¹¹

93 **Materials and Methods**

94 **Study sites**

95 The studied water-quality data derive from two freshwater rivers in north eastern Australia,
96 the Pioneer River and Sandy Creek, and an estuarine river in south western France, the
97 Ardour River estuary. Both study areas have seasonality in climate and therefore provide
98 non-stationary environmental conditions of water properties throughout the year.

99 The Australian study area is characterised by humid subtropical climate and strong
100 seasonality: the wet season (typically occurring between December and April) has higher
101 rainfall and air temperatures than the dry season which has lower rainfall and is associated
102 with low to zero river flows.²² Pioneer River (PR) is in the Mackay Whitsunday region of
103 northeast Australia, with a length of 120 km and a monitored catchment area of 1,466 km²,
104 with the upper reaches flowing predominantly through National or State Parks and its middle
105 and lower reaches flowing through land dominated by sugarcane farming. Sandy Creek (SC)
106 is a low-lying coastal-plain stream, 72 km long with a monitored catchment area of 326 km²
107 and with a similar land-use and land-cover profile to that of the lower Pioneer River. Both

108 study sites are in the freshwater reaches of these rivers.

109 In contrast, south-western France is characterised by an oceanic temperate climate, with
110 drier and warmer conditions from April to October. The Adour River (AR) in the Aquitaine
111 region of southwestern France is 330 km in length and has a catchment area of 16,880
112 km², with reaches flowing throughout a mosaic of agricultural and forested lands. The
113 study site is in the lower estuary between Bayonne city and the river mouth about 3 km
114 downstream. The tidal range varies from 1 to 4.5 m. Under the combined influence of rainfall
115 and snowmelt during late spring, river flow is highly variable, with minima and maxima of
116 about 80 and 3000 m³/s, respectively.²³ Driven by such strong fluctuations in tidal range and
117 river discharge, the estuary is classified as a highly variable, time-dependent salt-wedge.²³

118 **Water-quality data**

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

129 At each site, the *in-situ* sensors record turbidity (NTU) and electrical conductivity (con-
130 ductivity; $\mu\text{S}/\text{cm}$). The temporal resolution of data differs among rivers according to local
131 variability in environmental conditions; measurements are taken every 60 and 90 minutes in
132 PR and SC, respectively, and every 10 minutes in the AR. At PR and SC, measurements
133 are sometimes also taken at more frequent intervals during high flow events. The data se-

134 quences are one year (12 March 2017 to 12 March 2018) and five months (11 March to
135 13 August 2019) for the Australian and French sites, respectively, totalling 6280, 5402 and
136 14,185 observations at PR, SC and AR, respectively.

137 Water-quality parameters were strongly affected by the non-stationary flow regimes of
138 each site (see Fig. 1). Turbidity is a visual property of water clarity with higher values
139 indicating lower clarity. High turbidity often occurs during flood events, when rivers become
140 loaded with sediment eroded from the catchment, and during local re-suspension events that
141 can result from salt-wedge arrival, flow instability, or navigational and dredging activities.
142 In contrast, low values tend to occur during low discharge periods, e.g. during the dry
143 season, when the slow currents promote sediment deposition or, in estuarine environments,
144 when low-turbidity marine waters have entered the system. These overall trends can become
145 much more complex, and sometimes reversed, under the influence of specific local features.
146 Conductivity reflects the concentration of ions in the water. The dominant effect is generally
147 due to salts, such as sodium chloride, with marine waters typically having much higher
148 conductivity than inland waters. Other dissolved ionic compounds such as nutrients also
149 contribute to conductivity.

150 In the Australian sites, the water regime is driven by rainfall patterns, generating strong
151 positive or negative phases for either turbidity and conductivity values, coinciding with new
152 and sudden inputs of fresh water or periods of low and zero flow (see upper and mid panels in
153 Fig. 1). By contrast, AR is strongly conditioned by the interplay between river discharge and
154 tide forcing. The basic trend is that during high-discharge events, the turbidity reaches high
155 values corresponding to massive seaward fluxes of suspended sediments, and the conductivity
156 is minimal because the marine waters are flushed out from the estuary (see lower panels in
157 Fig. 1). During moderate discharge events, the tidal influence becomes dominant and drives
158 tidal fluctuations of conductivity and turbidity inside the estuary: at high tide, the estuary
159 is filled by salty (high conductivity) and clearer (low turbidity) marine waters while the
160 reverse behaviour is observed at low tide. This pattern is further affected by strong vertical

161 gradients, the salty marine waters being heavier than fresh riverine waters, which induces
162 complex time-dependent dynamics at the tidal scale.²³

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

170 **Definitions and types of anomalies**

171 Environmental sensors generating high-frequency data sets are usually subject to technical
172 anomalies that can derive from fouling of sensors, sensor calibration shifts, power supply
173 problems or unforeseen environmental conditions that adversely affect the sensor equip-
174 ment.^{1,26} In general, such technical anomalies strongly depart from the expected pattern of
175 the non-anomalous observations.¹⁰ In addition, strong inter- and intra-seasonal variations
176 in environmental conditions from site to site may create local differences in the occurrence
177 and types of anomalies, which can constrain their transferability to other sites.²⁷ We thus
178 need to develop methods for detecting these technical anomalies, that can overcome these
179 challenges,^{1,11,27} particularly in the context of data streaming.

180 The present analysis focuses on technical anomaly detection (AD) in water-quality time
181 series collected by *in-situ* sensors. Types of such anomalies have been described^{1,11} and
182 previous work has indicated the need to explore new statistical methods to better detect the
183 full suite of these types, notably those occurring over multiple, continuous observations (for
184 example as a result of sensor drift), given that such anomalies commonly occur in water-
185 quality time series yet are not as successfully detected as other types of anomalies (e.g.
186 single-point anomalies, out-of-range values).¹¹

187 Following the framework provided by,¹¹ a range of technical anomalies in the time-series
188 was first identified by local water-quality experts (i.e. for SC and PR by¹¹ and for AR by
189 DS). For each site and variable (i.e. turbidity and conductivity), the identified anomalies
190 were labelled along with their types based on anomaly classes defined by.¹¹ In the present
191 work, we first focused on anomaly types of Class 1 defined as a single observation generating
192 sudden changes in value from the previous observation. Class 1 anomalies were also consid-
193 ered a high priority in terms of detection, given that (true) sudden changes in turbidity and
194 conductivity may be used as early warning signals of water quality by local environmental
195 agencies and that they can strongly influence water-quality assessments and consequent man-
196 agement decisions. We additionally focused on Class 3 anomalies, which include technical
197 anomalies such as long-term calibration offsets and changes comprised of multiple dependent
198 observations. Class 3 was considered lower priority than Class 1 and may require *a posteriori*
199 user intervention (i.e. after data collection rather than in real time) to confirm observations
200 as anomalous. Specifically, using the nomenclature introduced by,¹¹ we based our analyses of
201 AD on those types of Class 1 anomalies defined as (a) large sudden spikes (type A), sudden
202 shifts (type D), small sudden spikes (type J), and Class 3 anomalies defined as drift (type
203 H), high variability (type E) and untrustworthy data not defined by other types (type L).
204 For more information about definitions of each anomaly type and class see.¹¹

205 The present analysis does not focus on AD of those types associated with impossible,
206 out-of-sensor-range and missing values (Class 2) given that they can be easily detected by
207 automated, hard-coded classification rules.¹¹ For each site, we therefore removed all Class
208 2 anomalies from the data prior to analysis. The filtered data were then log-transformed
209 to remove exponential data variance (see Fig. 1) to produce the time-series used for the
210 presented analysis.

211 **Learning methods and data processing**

212 One of the most challenging issues for AD in time-series data is that most data are usually
213 "normal" or non-anomalous, while anomalous values are rare.^{8,9} One of the most common
214 discriminating approaches to solve such a problem is to compare the similarity between
215 two sequences of time-dependent variables, specifically between the sequence of observed
216 and predicted values.^{11,28} When the anomalous values are pre-labelled, we can apply semi-
217 supervised classification based on training the learning process with the non-anomalous data,
218 fitting models with prediction errors and predicting anomalous events as those observed
219 values falling outside prediction intervals.^{11,28} Supervised classification feeds the learning
220 process with both a sequence of labelled values for AD (including both anomalous and non-
221 anomalous values, tagged accordingly). The learning process can then generate a sequence
222 of probabilities which can be binary-classified as anomalous or non-anomalous according to
223 a predefined threshold parameter.

224 For each water-quality variable and site, time-series data were partitioned according to
225 each learning process (i.e. semi-supervised or supervised classification). For semi-supervised
226 classification, we retained the "normal" values (i.e. we discarded the anomalous values) and
227 divided the time-series into four contiguous sequences of equal size: namely, two adjacent
228 sequences of values for training and two other adjacent sequences for validation (for details,
229 see Figure S1 in SI). Each pair of contiguous sequences for training or validation comprised
230 one sequence for prediction (predictor variable) followed by an adjacent sequence for the out-
231 come (outcome variable), respectively. For supervised classification, the sequence of values
232 (including both anomalous and "normal" values) for prediction and the sequence of labelled
233 values for AD were each divided into two adjacent sequences of equal length, for training
234 and validation (for details, see Figure S1 in SI). Thus, predictor and outcome variables were
235 composed of one sequence for training followed by an adjacent sequence for validation, re-
236 spectively. Given the proportionally low number of anomalies in our data set, we did not
237 use a "test" data set during the optimisation process.

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

252 Metrics for model performance and evaluation

253 To evaluate and compare the classification performance of AD, we calculated the four cate-
254 gories of the confusion matrix (i.e. true and false positives and true and false negatives; TP ,
255 FP , TN , FN , respectively) based on the discrimination threshold value t (typically 0.5; see
256 below for detail on the threshold values we used). From these, we calculated accuracy (Acc),
257 sensitivity (sn) and specificity (sp), and positive and negative predictive values (PPV and
258 NPV , respectively), which allowed us to compare results directly with those of Leigh et al.¹¹

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

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

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

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

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

259 *Acc*, *sn*, *sp*, *NPV* and *PPV* range between 0 and 1. *sn* assesses the probability of
260 determining *TP* correctly, whereas *sp* of determining *TN* correctly. *Acc* assesses the ability
261 to differentiate *TP* and *TN* correctly. Finally, *PPV* and *NPV* define the proportion of
262 anomalous versus "normal" observations, specifically the negative and positive predictive
263 values, respectively. When dealing with unbalanced classification problems (i.e. when there
264 are far fewer anomalous than non-anomalous observations), many evaluation metrics are
265 biased towards the majority class, maximising TN classification while minimising TP classi-
266 fication.^{29,30} We therefore also calculated evaluation metrics specially formulated to provide
267 and optimal classification for both positive and negative values in unbalanced data sets.³⁰
268 The first was balanced accuracy (*b.Acc*), which is defined as the arithmetic mean between
269 the *sn* and *sp* values:

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

270 The second and third were the F_1 and the Matthews Correlation Coefficient (MCC),
 271 defined as:

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

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

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

276 Artificial Neural Networks and their implementation in time-series

277 AD

278 Artificial Neural Networks (ANN) are machine learning methods, in which neurons (the basic
 279 unit of the learning process) are typically aggregated in layers that generate connections
 280 between input and output data. The typical structure of an ANN comprises input, hidden
 281 and output layers, usually connected sequentially. The input layer contains the observed
 282 data as variables (input neurons) and the output layer the predicted values (output or target
 283 neurons).³¹ The input and output layers are connected by hidden layers, and the relationships
 284 between the neurons are set by non-linear activation functions. For building and checking

285 the performance of different types of model structure (i.e. the network structure typically
286 associated with the number of layers and neurons), the ANN must be trained, by which the
287 weights associated with connections between neurons are optimised using various methods
288 and training algorithms.

289 A special case of ANNs adapted to time series applications are Recurrent Neural Net-
290 works (RNN), which can be considered a special case of Auto-Regressive Integrated Mov-
291 ing Average (ARIMA) and non-linear autoregressive moving average (NARMA) models.³²
292 RNN architecture mimics the cyclical connectivity of neurons, making them well suited for
293 the analysis and prediction of non-stationary time series.³³⁻³⁵ Long Short-Term Memory
294 (LSTM) networks are a re-design of the RNNs capable of learning long-term correlations in
295 a sequence.³⁶ LSTMs are structured in network units known as memory blocks composed
296 of self-connected memory cells and three multiplicative units, namely the "input", "output"
297 and "forget gates" connected to all cells within each memory block.³⁴ Unlike RNNs, LSTM
298 avoids the "vanishing gradients problem", that is when the error signal is used to train,
299 meaning that the network exponentially decreases the further one goes backwards in the
300 network, by means of the gates of the network units.³⁷ As a result, LSTMs allow the model
301 to be trained successfully using backpropagation through time, which is key for accounting
302 for long-term dependent time-series sequences. Although we applied LSTM in subsequent
303 analyses, we use the more generic term ANN from this point onwards.

304 In this study, ANNs were computed and fitted with "Keras", a model-level library which
305 provides high-level building for programming for developing deep-learning models. "Keras"
306 allows the implementation of a wide variety of neural-network building blocks (e.g. layers,
307 activation functions, optimisers) and supports the latest and most effective advances in deep
308 network training (including recurrent neural networks).³⁸ "Keras" is a high-level wrapper of
309 "TensorFlow" and helps to provide a simplified way of building neural networks from stan-
310 dard types of layers while facilitating a reproducible platform for developing deep learning
311 approaches in computational environmental sciences.³⁹ "Keras" is written in "Python" and

312 the "Keras" package provides an "R" interface to the deep-learning native functions.⁴⁰

313 In each model run, we compiled a model with a pre-defined set of hyper-parameters (see
314 below), in which the internal model parameters were iteratively updated throughout the
315 training steps (epochs; for definitions of ANNs hyper-parameters, see Text S1 in SI). During
316 each model run, iterations were computed until the error from the model was minimised or
317 reached a pre-defined value. In order to assess the performance of the learning process, we
318 used a pre-defined set of standard metrics of "Keras" commonly used for classification and
319 regression problems. Specifically, for both training and validation data sets, we calculated the
320 loss function for semi-supervised classification and the mean-square error and the accuracy for
321 supervised classification. For each model run, we accounted for over-fitting using learning
322 curves which showed how error changes as the training set size increases.³⁷ To do so, we
323 compared the shape and dynamics of the learning curves of the training and validation data,
324 which can be used to diagnose the bias and variance of the learning process. For instance,
325 the "training learning curve" is calculated from the training data and provided information
326 on the learning process, whereas the "validation learning curve" is calculated from the hold-
327 out data and gives information on the model generality. For details of the dynamics of
328 the learning curves for either semi-supervised and supervised classification with our data
329 sets, see Supporting Information (SI). SI contains the code and R functions to allow the
330 implementation of ANNs using "Keras" in the "R" statistical language.⁴¹

331 **Optimisation of hyper-parameters and their influence on AD perfor-** 332 **mance**

333 ANNs include a broad suite of hyper-parameters that affect the ability to learn patterns
334 from the data and the performance of model predictions. Before training, we selected ten
335 hyper-parameters that affect (i) the network structure and (ii) the training algorithm; see SI
336 for details on hyper-parameter definitions. Furthermore, we defined two additional hyper-
337 parameters related to (iii) the "sliding window" defined by $n \times p$ matrices delimited at

338 regular temporal intervals (see above) and (iv) the "threshold classification" as a measure
339 of the discrimination threshold value to compute the four categories of the confusion matrix
340 (see above). For each hyper-parameter we then defined a range of values (see details in SI)
341 that affect the performance of ANNs. The range of values depended on the type of the
342 hyper-parameter and varied from a continuous searching space (for those hyper-parameters
343 characterised by double precision, such as the learning rate) and discrete values (for those
344 defined by integer values, such as number of layers or units, and functions or algorithms,
345 such as the optimisation algorithm). In the case of the "sliding window" we decided to use
346 discrete values, defining the temporal resolution of the recorded time-series sequence and
347 covering meaningful and regular environmental processes occurring in rivers (e.g. 24 h, 12
348 h, 6 h). For more information about this issue, see SI.

349 Given the large number of possible models to be tested with all the combinations of
350 hyper-parameter values, we tuned the ANN models using a Bayesian optimisation method,
351 which is a common class of optimisation methods, especially in deep-learning networks.⁴² The
352 Bayesian optimisation method works by constructing a posterior distribution of functions
353 (assuming a Gaussian process) associated with the variability of each hyper-parameter, which
354 best describes the "objective function", defined as the "cost" associated with the optimisation
355 problem. With each iteration of the algorithm, the posterior distribution improves and
356 the algorithm becomes more accurate in those regions of the parameter space with higher
357 likelihood to maximise or minimise the objective function.⁴² The Bayesian statistical model
358 comprises two components: (i) a Bayesian statistical model for modelling the objective
359 function, and (ii) an acquisition function for deciding where to sample next. In our case, we
360 applied the *mlrMBO* toolbox implemented in the R statistical language. Compared with
361 other black-box benchmark optimisers, the *mlrMBO* toolbox performs well for expensive
362 optimisation scenarios for single- and multi-objective optimisation tasks, with continuous or
363 mixed parameter spaces.⁴³

364 Our aim for the optimisation procedure was to find combinations of hyper-parameters

365 that resulted in the best performance for AD classification (see above). The optimisation
366 procedure followed two steps, firstly generating a random search space and secondly focusing
367 on search shrinks, based on the results generated during the random search. First, we started
368 the algorithm by generating a random design ($n = 250$), which included a varying number of
369 hyperparameters with the aim to generate enough variability to detect the most promising
370 values. In our case we generated iterations by varying randomly the parameter combinations
371 of the 12 hyper-parameters. After computing the model and the optimisation scores of each
372 iteration, we secondly focused on search shrinks of the search space following a Bayesian
373 optimisation method with $n = 250$ iterations based on maximising performance metrics or
374 objective functions. In our case, we followed a multi-objective Bayesian optimisation method
375 based on maximising, in each k iteration run, the value of $b.Acc$, f_1 and MCC scores,
376 defined above. Such an approach allowed us to maximise the different and complementary
377 properties summarised by each score for unbalanced classification. Following the multi-
378 objective Bayesian method, the optimisation resulted in a total of $n = 500$ iterations (i.e.
379 $n = 250$ for random search and $n = 250$ for Bayesian optimisation) for each combination of
380 water-quality variables, sites and learning procedures.

381 To examine the effects of hyper-parameters predictor variables (hyper-parameters) on
382 the dependent variables (performance metrics), we applied methods for causal inference
383 using random forests.⁴⁴ Specifically, we determined the statistical importance of each hyper-
384 parameter as a predictor on the dependent variables or optimisation scores (e.g. $b.Acc$, f_1 or
385 MCC) by calculating the "Variable Importance" (VI). VI reflects the model performance
386 across the entire range of predictor and response variables, converted into a set of ordinal
387 ranks. VI can be further used to test how the model response changes as the value of
388 any of the predictor variables is changed, meaning that VI is similar to a standard "one-
389 parameter-at-a-time" sensitivity analysis.⁴⁵ For each hyper-parameter, we computed VI to
390 measure the relative importance (or dependence, or contribution) of such hyper-parameter
391 predictor variables in terms of their effect on optimisation scores. In our case, we computed

392 the out-of-bag error, which is an error estimation technique used to evaluate the accuracy of
393 a random forest after permuting each predictor variable. We used the R statistical language
394 and the "randomForest" library.⁴⁶

395 For each water-quality variable and site, we retained the "best" model as that which
396 maximised the optimisation scores. Specifically, for the complete set of candidate models,
397 we averaged the value of $b.Acc$, f_1 or MCC and we retained the "best" model as that
398 maximising the averaged optimisation scores. We compared the shape and dynamics of the
399 learning curves from the "best" models (a) to diagnose whether or not the training and
400 validation data sets were sufficiently representative (i.e. one data set could capture the
401 statistical characteristics relative to other data sets) and (b) to test the behaviour of the
402 learning process (i.e. underfit, overfit, good fit).

403 Results

404 Hyper-parameters optimisation for AD

405 The learning rate for AD stabilised early in the optimisation of ANN models, with few
406 improvements on the performance beyond $n = 200$ model iterations (although see the ANN
407 model for conductivity at SC). For semi-supervised classification, the costs of computing
408 ranged from 0.312 h for turbidity in SC and 3.53 h for turbidity in AR, whereas for supervised
409 classification the cost ranged from 1.11 h for conductivity in PR and 99.5 h turbidity in
410 AR after $n = 500$ model iterations (see Table S2 in SI). Comparing learning methods,
411 supervised classification had better performance and generated consistently higher values for
412 $b.Acc$, f_1 or MCC , showing a similar and consistent pattern of the accumulative curve along
413 the optimisation process (Fig. 2). Compared to supervised classification, semi-supervised
414 classification required a larger number of model iterations for maximising any of the three
415 performance metrics, notably for turbidity in SC and AR.

416 Considering the whole suite of model iterations, we found a consistent pattern with re-

417 spect to VI of those hyper-parameters affecting optimisation scores. Overall, the "Learning
418 hyper-parameters" had higher VI values than "Model hyper-parameters" for model perfor-
419 mance (see Tables S2 to S4 in SI). In addition, *th.class* had higher VI in semi-supervised
420 classification, whereas *s.win* had minor VU for both supervised and semi-supervised clas-
421 sification. Specifically, the hyper-parameters with higher VI for $b.Acc$, f_1 and MCC , re-
422 spectively, were *th.class* (78.2, 78.4 and 76.5, respectively; VI values averaged across sites,
423 learning methods and water-quality variables), *dropout* (46.2, 47.8 and 49.0), *b.size* (43.5,
424 48.4 and 51.8), *momen* (31.8, 34.3 and 36.2), *l.rate* (31.3, 33.2 and 35.6). Comparing learn-
425 ing methods, we found that VI of hyper-parameters differed depending on the combination
426 of types of anomalies and water-quality variables. For instance, semi-supervised classifi-
427 cation had higher VI values for *th.class*, *dropout*, *momen* and *l.rate*, whereas supervised
428 classification for *b.size* and *activ*. Comparing sites, *th.class* had the highest VI values in
429 semi-supervised classification for conductivity in PR, *dropout* and *l.rate* in semi-supervised
430 classification for turbidity in SC, *b.size* for supervised classification for turbidity in PR,
431 *momen* for semi-supervised classification for turbidity in AR and *activ* for supervised classi-
432 fication for turbidity in AR. For details of the VI of those hyper-parameters independently
433 affecting $b.Acc$, f_1 and MCC , see Tables S3 to S5 in SI.

434 After hyper-parameter optimisation, the "best" models also had performed well in terms
435 of $b.Acc$, f_1 and MCC scores, for all water-quality variables and learning methods; Table 1);
436 for details of hyper-parameters and values of the best models, see SI. However semi-supervised
437 classification had less balanced predictions (i.e. lower values for $b.Acc$, f_1 score and MCC)
438 than supervised classification, meaning that anomaly cases (positives) were proportionally
439 less-correctly predicted than "normal" cases (negatives). Specifically, semi-supervised classi-
440 fication had lower rates of TP that we classified as true ($sn = 0.652$; values averaged across
441 sites and water-quality variables) and lower proportions of positives and negatives that were
442 true ($PPV = 0.512$), and that generated a moderately balanced detection rate for either
443 TP and TN results ($b.Acc = 0.757$, $f_1 = 0.490$ and $MCC = 0.665$). By contrast, supervised

444 classification provided a higher and balanced detection rate for both TP and TN ($b.Acc =$
445 0.822 , $f_1 = 0.622$ and $MCC = 0.762$), and thus anomaly cases (positives) were as predicted
446 correctly as "normal" cases (negatives). In contrast, supervised classification had higher
447 detection rates of TP ($sn = 0.704$) and higher rates of correct classification of TP and TN
448 ($PPV = 0.643$). For detailed information about the performance of each model, see Table
449 1.

450 For "best" models, we additionally checked that our training and validation data sets were
451 sufficiently representative, based on the learning process of either training or validation data.
452 Overall, we found strong variation for each water-quality variable, site and learning method,
453 suggesting that the presence of certain types of anomalies were not always consistent between
454 training and validation data sets. Both for semi-supervised and supervised classification, we
455 found that the validation data were easier to predict than training data (i.e. the "training
456 learning curve" had poorer performance than the "validation learning curve"), suggesting
457 that the validation data had lower complexity of anomaly types; for details of learning curves,
458 see Fig. S2 to S5 in SI. For supervised classification, we additionally found that the validation
459 data were unable to produce a good fit (i.e. both training and validation learning curves were
460 almost flat), probably as a combination of the low number of anomaly cases (conductivity
461 in PR; see Fig. 1 and Fig. S8 in SI) and the presence of a long-term anomaly event at the
462 end of the data (turbidity in AR; see Fig. 1 and Fig. S9 in SI); the latter pattern did not
463 happen when fitting semi-supervised classification in both data sets, for which the learning
464 process produced a good fit for both training and validation processes. Overall, we did not
465 detect over-fitting of the training data, a result confirmed by the relatively medium-to-low
466 values of *dropout* (i.e. <0.5) in most of the "best" models (see Tables S2 to S4 in SI for
467 details of the values of hyper-parameters of the "best" models).

468 AD among types, learning methods and sites

469 Anomaly types of Class 1 were present, but at low abundances, in all water-quality data sets,
470 comprising, on average, 0.137% of cases (Table 2). By contrast, the majority of anomalies
471 were classified as Class 3, which provided 11.6% of cases (Table 2). Anomalies of Class 3
472 were context-dependent with respect to each site and water-quality variable. Specifically,
473 SC had two anomalous periods of high-variability (type E) occurring during the first four
474 months of monitoring (see Fig. 1). PR had two drift sequences (type H) for turbidity and
475 one period of untrustworthy data (type L) and one drift sequence (type H) for conductivity
476 during the first period of monitoring. Finally, AR had a long-drift sequence (type H) at the
477 end of the monitoring period.

478 Comparing the performance of each "best" model with respect to detecting of the differ-
479 ent types of anomalies, we found, on average, that semi-supervised classification had higher
480 capacity for detecting Class 1 anomalies (45.7% vs 23.6% for semi-supervised and supervised
481 classification, respectively; averaged values across sites and water-quality variables), whereas
482 supervised classification had proportionally higher capacity for detecting Class 3 anomalies
483 (72.6% vs 68.8% respectively) (Table 3). However, such differences between learning meth-
484 ods were particularly context-dependent and based on the different combinations of types
485 of anomalies at each site. Semi-supervised classification better detected large-sudden spikes
486 (type A) for turbidity in AR, and small sudden spikes (J) for turbidity in PR and AR.
487 Supervised classification, by contrast, had better performance for detecting drift (type H).
488 Both semi-supervised and supervised classification performed well for detecting high vari-
489 ability (E) and untrustworthy anomalies for turbidity in SC and PR and conductivity in PR.
490 Sudden shifts (type D) for turbidity in SC and drift (H) for conductivity in PR were not
491 detected by any learning method.

492 Discussion

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

507 Previous work has shown that regression time-series methods are useful for AD in sta-
508 tionary and non-stationary time-series data sets,,^{47,48} but the detection of certain anomaly
509 types, such as sensor drift and periods of anomalously high variability, remain challenging.
510 ANNs are a versatile method that can train models using different learning methods, and as
511 shown by our study, can provide the flexibility required to detect a broad suite of anomaly
512 types, including improved performance for detecting extended periods of untrustworthy data.
513 In our case, we applied regression-based and semi-supervised ANNs similarly for AD; that
514 is, models first predict data sequences based on "normal" cases and then classify as anoma-
515 lous cases those departing from a given threshold value. Despite the good performance of
516 ANNs for AD, our findings demonstrate that ANNs could have limitations regardless of
517 the underlying statistical method used and their applicability in near-real time AD. In our
518 study, c. 100 model iterations, which were necessary to obtain acceptable model perfor-

519 mances, required three orders of magnitude of computing time longer than regression-based
520 ARIMA. Compared to these regression-based time-series methods, the Bayesian optimisation
521 of hyper-parameter values allowed us to optimise classification of anomalies, at the expense
522 of increased computing costs for hyper-parameter optimisation.

523 For the Australian sites, we found that the performance of semi-supervised ANNs (pro-
524 posed here) and regression-based ARIMA (proposed by¹¹) were equivalent in terms of correct
525 classification, notably by providing high false-detection rates (falses) of both anomalies (pos-
526 itives) and "normal" cases (negatives): semi-supervised classification had higher rates of FP
527 and FN which resulted in low values of sn and PPV . We also found that supervised ANNs
528 considerably minimised false detection rates (i.e. low values of FP and FN and maximised
529 sn and PPV values), when compared with the regression-based models and semi-supervised
530 ANNs. Semi-supervised ANNs and regression-based ARIMA generated higher rates of false
531 alarms (i.e. both FP and FN), which overestimate anomalous events, than supervised ANNs.

532 We did not detect over-training in the "best" models, but after checking the shape and
533 dynamics of the learning curves we found that the training and validation data sets had
534 inconsistent numbers and presence of anomaly types. For instance, the validation data set
535 in turbidity in AR had a long-term anomaly event which is not present in the training data
536 (see Fig. 1). As a result, there is inconsistency between training and validation data sets
537 in the performance, typically the validation data had higher performance than the training
538 data (see Figs S2 to S5 in SI). The difference in complexity was mainly a consequence of
539 the presence of Class 3 anomalies in the data sets, usually occurring as single long-term
540 anomalous events in each time series. Cross-validation could solve this problem, but it
541 is also true that multiple partitioning of the data sets could have divided the time-series
542 into multiple smaller independent time-series sequences with inconsistent representation of
543 anomaly types among data sets. Multiple partitioning of the time-series sequence could be
544 especially critical during the optimisation process, notably on those hyper-parameters tuning
545 the size of the sub-samples processed during the learning process (i.e. $b.size$ and $momen$,

546 and also for *s.win*). All the above issues were a consequence of the proportionally lower
547 abundance of anomalous events compared to "normal" events in the data sets, a phenomenon
548 that makes the detection of anomalies challenging.

549 We found that semi-supervised ANNs were especially suited to AD of short-term anoma-
550 lous events (i.e. sudden spikes, sudden shifts and small sudden spikes defined as Class 1
551 anomalies, following the terminology of¹¹). Although semi-supervised and supervised ANNs
552 were comparable regarding their ability to detect long-term anomalous events (i.e. drift
553 and periods of high-variability or Class 3 anomalies), the latter had better performance in
554 terms of correct AD (see above). At least for AD of long-term events, our results with
555 ANNs thus outperformed those of regression-based ARIMA,¹¹ which is an important step
556 forward for AD in water-quality data given such anomalies have consistently proved chal-
557 lenging to detect by other automated AD methods. This is because such anomalous, often
558 very context-dependent events can behave similarly to natural, non-stationary water-quality
559 events such that they are only detected manually by a trained eye very familiar with the
560 local conditions. Our novel use of ANNs and Bayesian optimisation for hyper-parameter
561 selection therefore holds much promise for AD in high-frequency water-quality data from a
562 broad range of environments and ecosystems, including rivers, estuaries and marine waters.

563 Although our study is based on the analysis of data from three sites only, it provides
564 a methodological framework for AD in high-frequency water-quality data collected from
565 sites with contrasting non-stationary environmental processes.^{49,50} We found that the non-
566 stationary environmental conditions of each site played a substantial role in both (i) deter-
567 mining the types of anomalies present and (ii) increasing the uncertainty around the accurate
568 detection of anomalies. The water-quality variables analysed here are affected by long-term
569 environmental processes occurring at each site (i.e. seasonal precipitation patterns in both
570 Australian and French sites), as well as regular and short-term events (e.g. cyclone floods in
571 Australian sites and tidal regimes in the French site). Both short- and long-term environmen-
572 tal processes may interact and this could affected the AD performance. For Australian sites,

573 the detection of sudden spikes and shifts (i.e. short-term anomalies) and of drifts or periods
574 with high variability (i.e. long-term anomalies) were rarely masked by natural environmental
575 processes (i.e. seasonal rainfall patterns or cyclone floods), and that resulted in the optimal
576 classification of short- and long-term anomalous events by using either semi-supervised and
577 supervised classification, respectively (Table 3). At the French site, by contrast, strong tidal
578 regimes occurred alongside medium-to-high seasonal discharge events, and that conditioned
579 and limited the accuracy of detecting both short- and long-term anomalies (Table 3).

580 In this work, we calibrated ANNs models using Bayesian optimisation of hyper-parameter
581 values, by means of multi-objective optimisation. Although optimisation methodology is
582 well-established for the detection of anomalies in water-quality data,⁷ multi-objective op-
583 timisation procedures have rarely been applied for AD in time-series analysis for either
584 prediction or detection anomalies using ANNs methods; but see Perelman et al.¹⁹ for an
585 application of detecting anomalies using a single-objective optimisation. The application of
586 multi-objective methodology allow us to get a balanced performance of models for detect-
587 ing anomalies in a broad set of environmental conditions. Notwithstanding the potential
588 utility of ANNs for AD demonstrated above, there is substantial room for improvement in
589 the performance under real-world conditions.¹² First, we found that the ability for AD is
590 context-dependent, meaning that accuracy is conditioned on the spatio-temporal environ-
591 mental variability of the data set available. Anomalies are rare events, meaning that we
592 need large data sets spanning the entire environmental variability of each locality to ensure
593 their inclusion in a data set.¹¹ Second, we performed our detections based on the analysis of
594 independent environmental variables (i.e. turbidity or conductivity), but the application of
595 ANNs with multivariate time-series could enable us to account for the temporal correlation of
596 multiple variables monitored at the same time.^{19,24,51} Finally, methodological improvements
597 provide additional avenues to increase the performance of ANNs in AD, such as Bayesian
598 RNNs,^{52,53} which allow quantification of the uncertainty and the use of ensemble averag-
599 ing for ANNs, combining unsupervised and supervised classification⁵⁴ or the combination of

600 ANNs with other methodologies.^{6,55}

601 The study of AD in high-frequency data has numerous applications in environmental
602 and health monitoring, and fault and fraud detection, where there is a need to provide
603 near-real time solutions and optimal performance for monitoring and to ensure the quality
604 of data streaming. We have demonstrated that ANNs are a flexible method for providing
605 optimal performance for AD, given they are able to cope with both long- and short-term
606 non-stationary processes that condition high-frequency data. However, ANNs and their
607 hyper-parameter optimisation have been understudied in the context of AD in water-quality.
608 Given the promising results from our study, we therefore recommend further investigation
609 and development of such methods to improve the accurate detection of a broad suite using
610 machine learning or deep-learning methods of anomalies detection under a wide range of
611 environmental conditions.⁷ Environmental data are intrinsically variable though time and
612 space and thus our approach is transferable for AD in complex spatio-temporal applications
613 and ecosystems. Our findings will therefore be of relevance to water scientists and man-
614 agers throughout the world in order to broaden the applicability of ANNs for efficient water
615 monitoring.

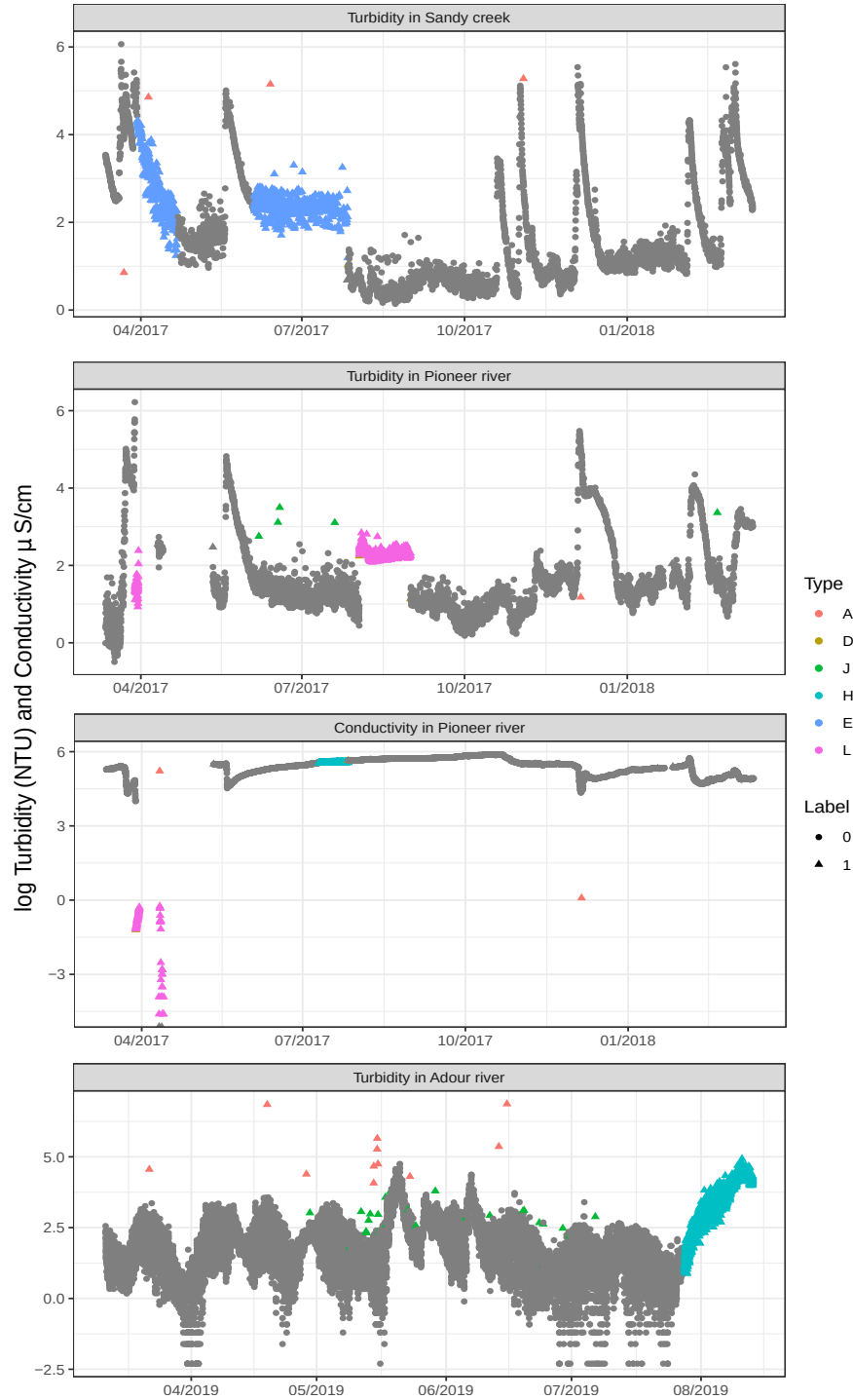


Figure 1: Observed trends for each water-quality variable and site, including the type of anomaly. Shapes correspond to the different data values (i.e. circles for "normal" values and triangles for anomalous values) and colours to the different anomaly types. Anomaly types were classified by local water-quality 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).

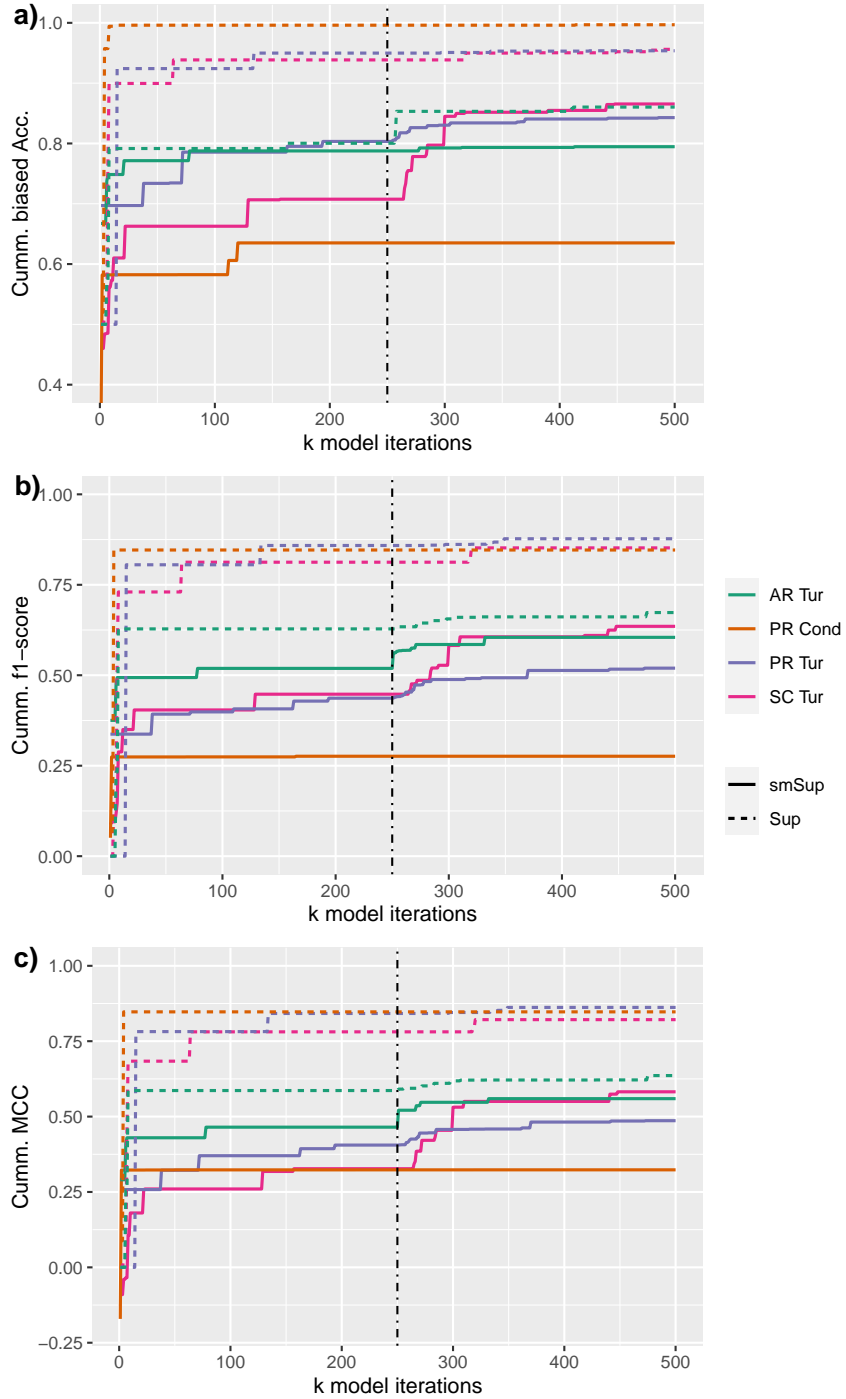


Figure 2: Cumulative optimisation scores for hyper-parameter optimisation along k model iterations. Each panel showed the procedure of multi-objective optimisation procedure of (a) balanced Accuracy, (b) f1-score and (c) Matthew's Correlation Coefficient occurring in each k iteration run. The optimisation begins searching with 250 random iterations, and then with 250 iterations of the Bayesian optimization procedure. Colours define each water-quality variable and site, whereas line patterns indicate the learning method. Abbreviations: Sandy creek (SC), Pioneer river (PR) and AR Adour river (AR); Turbidity (Tur) and Conductivity (Cond); semi-supervised classification (smSup) and supervised classification (Sup).

616 Acknowledgement

617 Funding was provided by the Energy Environment Solutions (E2S-UPPA) consortium and
618 the BIGCEES project from E2S-UPPA ("Big model and Big data in Computational Ecol-
619 ogy and Environmental Sciences"), the Queensland Department of Environment and Sci-
620 ence (DES) and the ARC Centre of Excellence for Mathematical and Statistical Frontiers
621 (ACEMS). A repository of the water-quality data from the in situ sensors used herein and
622 the code used to implement methods of Artificial Neural Networks for anomaly detection are
623 provided in the Supporting Information.

624 Supporting Information Available

625 A listing of the contents of each file supplied as Supporting Information are included. The
626 following files are available free of charge.

- 627 • Figure S1: Scheme of the data processing for time-series AD for each learning method.
- 628 • Text S1 and Table S1: Brief definition of the hyper-parameters and range of hyper-
629 parameter values for fitting and optimising ANN models.
- 630 • Table S2: Costs of computing time during learning process.
- 631 • Tables S3 to S5: Variable importance (VI) of values of hyper-parameters to maximise
632 the *biased Accuracy*, the f_1 -score and the *Matthews Correlation Coefficient*.
- 633 • Figures S2 to S5: Learning curves for semi-supervised and supervised learning processes
634 for the "best" ANN models after Bayesian optimisation.
- 635 • Figures S6 to S9: Data of observed trend and the probability of anomaly detection for
636 the "best" ANN models after Bayesian optimisation.

²Abbreviations: 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)

Table 1: Performance of the best model for AD. For each site and water-quality variable, details of the hyper-parameters and their values of the best model are shown in SI. We calculated performance scores for ARIMA with Anomaly Detection (ArAD), semi-supervised ANNs (smSup) and supervised ANNs classification (Sup). For abbreviations of sites and variables see Figure 2.

Site	Var	Train	TN	FN	FP	TP	acc	sn	sp	PPV	NPV	b_acc	f1	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	12744	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	13017	925	179	495	0.924	0.349	0.986	0.734	0.934	0.668	0.47	0.55

¹ 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*), *f1* score (*f1*) and Matthew's Correlation Coefficient (*MCC*).

Table 2: Number of anomalous events according to each type (columns), site and water-quality variable (rows). 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.¹¹ For abbreviations of sites, variables and training see Figure 2.

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

2

Table 3: Performance of the "best" model for AD by anomaly type. For each site and water-quality variable, values represent the percentage of data values detected relative to the total number of respective anomalies labelled 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.

	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	-	-

References

- (1) Horsburgh, J. S.; Jones, A. S.; Stevens, D. K.; Tarboton, D. G.; Mesner, N. O. *Environmental Modelling & 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. Sensors in the stream: the high-frequency wave of the present. 2016.
- (3) Hill, D. J.; Minsker, B. S. *Environmental Modelling & Software* **2010**, *25*, 1014–1022.
- (4) Horsburgh, J. S.; Reeder, S. L.; Jones, A. S.; Meline, J. *Environmental Modelling & Software* **2015**, *70*, 32–44.
- (5) Jiang, J.; Wang, P.; Lung, W.-s.; Guo, L.; Li, M. *Journal of hazardous materials* **2012**, *227*, 280–291.
- (6) Shi, B.; Wang, P.; Jiang, J.; Liu, R. *Science of the Total Environment* **2018**, *610*, 1390–1399.
- (7) Dogo, E. M.; Nwulu, N. I.; Twala, B.; Aigbavboa, C. *Urban Water Journal* **2019**, *16*, 235–248.
- (8) Chandola, V.; Banerjee, A.; Kumar, V. *ACM computing surveys (CSUR)* **2009**, *41*, 15.
- (9) Gupta, M.; Gao, J.; Aggarwal, C.; Han, J. *Synthesis Lectures on Data Mining and Knowledge Discovery* **2014**, *5*, 1–129.
- (10) Goldstein, M.; Uchida, S. *PloS one* **2016**, *11*, 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. *Science of The Total Environment* **2019**, *664*, 885–898.

- 659 (12) Muharemi, F.; Logofătu, D.; Leon, F. *Journal of Information and Telecommunication*
660 **2019**, *3*, 294–307.
- 661 (13) Bourgeois, W.; Romain, A.-C.; Nicolas, J.; Stuetz, R. M. *Journal of Environmental*
662 *Monitoring* **2003**, *5*, 852–860.
- 663 (14) Shipmon, D. T.; Gurevitch, J. M.; Piselli, P. M.; Edwards, S. T. *arXiv preprint*
664 *arXiv:1708.03665* **2017**,
- 665 (15) Makarynskyy, O.; Makarynska, D.; Kuhn, M.; Featherstone, W. *Estuarine, Coastal and*
666 *Shelf Science* **2004**, *61*, 351–360.
- 667 (16) Makarynska, D.; Makarynskyy, O. *Computers & Geosciences* **2008**, *34*, 1910–1917.
- 668 (17) Wu, W.; Dandy, G. C.; Maier, H. R. *Environmental Modelling & Software* **2014**, *54*,
669 108–127.
- 670 (18) Tinelli, S.; Juran, I. *Water Supply* **2019**, *19*, 1785–1792.
- 671 (19) Perelman, L.; Arad, J.; Housh, M.; Ostfeld, A. *Environmental science & technology*
672 **2012**, *46*, 8212–8219.
- 673 (20) Muharemi, F.; Logofătu, D.; Andersson, C.; Leon, F. *Modern Approaches for Intelligent*
674 *Information and Database Systems*; Springer, 2018; pp 173–183.
- 675 (21) Fehst, V.; La, H. C.; Nghiem, T.-D.; Mayer, B. E.; Englert, P.; Fiebig, K.-H. Automatic
676 vs. manual feature engineering for anomaly detection of drinking-water quality. *Pro-*
677 *ceedings of the Genetic and Evolutionary Computation Conference Companion*. 2018;
678 pp 5–6.
- 679 (22) Brodie, J. *ACTFR Technical Report No. 02/03*; Australian Centre for Tropical Fresh-
680 water Research, James Cook University, 2004.

- 681 (23) Defontaine, S.; Sous, D.; Morichon, D.; Verney, R.; Monperrus, M. *Estuarine, Coastal*
682 *and Shelf Science* **2019**, 106445.
- 683 (24) Leigh, C.; Kandanaarachchi, S.; McGree, J. M.; Hyndman, R. J.; Alsibai, O.;
684 Mengersen, K.; Peterson, E. E. *PloS one* **2019**, *14*.
- 685 (25) Leigh, C.; Burford, M. A.; Connolly, R. M.; Olley, J. M.; Saeck, E.; Sheldon, F.;
686 Smart, J. C.; Bunn, S. E. *Water* **2013**, *5*, 780–797.
- 687 (26) Wagner, R. J.; Boulger Jr, R. W.; Oblinger, C. J.; Smith, B. A. *Guidelines and standard*
688 *procedures for continuous water-quality monitors: station operation, record computa-*
689 *tion, and data reporting*; 2006.
- 690 (27) Cox, B. *Science of the total environment* **2003**, *314*, 335–377.
- 691 (28) Malhotra, P.; Vig, L.; Shroff, G.; Agarwal, P. Long short term memory networks for
692 anomaly detection in time series. *Proceedings*. 2015; p 89.
- 693 (29) Kelleher, J. D.; Mac Namee, B.; D’arcy, A. *Fundamentals of machine learning for*
694 *predictive data analytics: algorithms, worked examples, and case studies*; MIT Press,
695 2015.
- 696 (30) Boughorbel, S.; Jarray, F.; El-Anbari, M. *PloS one* **2017**, *12*, e0177678.
- 697 (31) Haykin, S. *Neural Networks and Learning Machines, 3/E*; Pearson Education India,
698 2010.
- 699 (32) Siami-Namini, S.; Tavakoli, N.; Namin, A. S. A Comparison of ARIMA and LSTM in
700 Forecasting Time Series. 2018 17th IEEE International Conference on Machine Learning
701 and Applications (ICMLA). 2018; pp 1394–1401.
- 702 (33) Hilas, C. S.; Mastorocostas, P. A. *Knowledge-Based Systems* **2008**, *21*, 721–726.

- 703 (34) Graves, A. *Supervised sequence labelling with recurrent neural networks*; Springer, 2012;
704 pp 37–45.
- 705 (35) Maier, H. R.; Dandy, G. C. *Environmental modelling & software* **2000**, *15*, 101–124.
- 706 (36) Hochreiter, S.; Schmidhuber, J. *Neural computation* **1997**, *9*, 1735–1780.
- 707 (37) Goodfellow, I.; Bengio, Y.; Courville, A. *Deep learning*; MIT press, 2016.
- 708 (38) Gulli, A.; Pal, S. *Deep Learning with Keras*; Packt Publishing Ltd, 2017.
- 709 (39) Rampasek, L.; Goldenberg, A. *Cell systems* **2016**, *2*, 12–14.
- 710 (40) Allaire, J.; Chollet, F. *R package version* **2017**,
- 711 (41) Team, R. C. **2017**,
- 712 (42) Snoek, J.; Larochelle, H.; Adams, R. P. Practical bayesian optimization of machine
713 learning algorithms. *Advances in neural information processing systems*. 2012; pp 2951–
714 2959.
- 715 (43) Bischl, B.; Richter, J.; Bossek, J.; Horn, D.; Thomas, J.; Lang, M. *arXiv preprint*
716 *arXiv:1703.03373* **2017**,
- 717 (44) Wager, S.; Athey, S. *Journal of the American Statistical Association* **2018**, *113*, 1228–
718 1242.
- 719 (45) Mishra, S.; Datta-Gupta, A. *Applied statistical modeling and data analytics: A practical*
720 *guide for the petroleum geosciences*; Elsevier, 2017.
- 721 (46) Breiman, L. *Machine learning* **2001**, *45*, 5–32.
- 722 (47) Hyndman, R. J.; Athanasopoulos, G. *Forecasting: principles and practice*; OTexts,
723 2018.

- 724 (48) Box, G. E.; Jenkins, G. M.; Reinsel, G. C.; Ljung, G. M. *Time series analysis: fore-*
725 *casting and control*; John Wiley & Sons, 2015.
- 726 (49) Clarke, R. T. **2007**,
- 727 (50) Sivakumar, B. *Chaos in Hydrology*; Springer, 2017; pp 29–62.
- 728 (51) Sánchez-Fernández, A.; Baldán, F.; Sainz-Palmero, G.; Benítez, J.; Fuente, M. *Chemo-*
729 *metrics and Intelligent Laboratory Systems* **2018**, *182*, 57–69.
- 730 (52) Mirikitani, D. T.; Nikolaev, N. *IEEE Transactions on Neural Networks* **2009**, *21*, 262–
731 274.
- 732 (53) Sun, W.; Paiva, A. R.; Xu, P.; Sundaram, A.; Braatz, R. D. *arXiv preprint*
733 *arXiv:1911.04386* **2019**,
- 734 (54) Comar, P. M.; Liu, L.; Saha, S.; Tan, P.-N.; Nucci, A. Combining supervised and unsu-
735 pervised learning for zero-day malware detection. 2013 Proceedings IEEE INFOCOM.
736 2013; pp 2022–2030.
- 737 (55) Dairi, A.; Cheng, T.; Harrou, F.; Sun, Y.; Leiknes, T. *Sustainable Cities and Society*
738 **2019**, *50*, 101670.

739 Graphical TOC Entry

740

