

Accepted Manuscript

This is an Accepted Manuscript of the following article:

Saer Samanipour, Sarit Kaserzon, Soumini Vijayasathy, Hui Jiang, Phil Choi, Malcolm J. Reid, Jochen F. Mueller, Kevin V. Thomas. Machine learning combined with non-targeted LC-HRMS analysis for a risk warning system of chemical hazards in drinking water: A proof of concept. *Talanta*. Volume 195, 2019, pages 426-432, ISSN 0039-9140.

The article has been published in final form by Elsevier at

<http://dx.doi.org/10.1016/j.talanta.2018.11.039>

© 2019. This manuscript version is made available under the

CC-BY-NC-ND 4.0 license

<http://creativecommons.org/licenses/by-nc-nd/4.0/>

It is recommended to use the published version for citation.

Machine Learning Combined with Non-targeted LC-HRMS Analysis for a Risk Warning System of Chemical Hazards in Drinking Water: A Proof of Concept

Saer Samanipour^{a,b,*}, Sarit Kaserzon^b, Soumini Vijayasaraty^b, Hui Jiang^b, Phil Choi^b, Malcolm J. Reid^a, Jochen F. Mueller^b, Kevin V. Thomas^{a,b}

^a*Norwegian Institute for Water Research (NIVA), Gaustadalléen 21, 0349 Oslo, Norway*

^b*Queensland Alliance for Environmental Health Science (QAEHS), University of Queensland, 20 Cornwall Street, Woolloongabba, QLD, 4012, Australia.*

Abstract

Guaranteeing clean drinking water to the global population is becoming more challenging, because of the cases of water scarcity across the globe, growing population, and increased chemical footprint of this population. Existing targeted strategies for hazard monitoring in drinking water are not adequate to handle such diverse and multidimensional stressors. In the current study, we have developed, validated, and tested a machine learning algorithm based on the data produced via non-targeted liquid chromatography coupled with high resolution mass spectrometry (LC-HRMS) for the identification of potential chemical hazards in drinking water. The machine learning algorithm consisted of a composite statistical model including an unsupervised

*Saer Samanipour

Email address: saer.samanipour@niva.no (Saer Samanipour)

¹NIVA, Gaustadalléen 21, 0349 Oslo, Norway

Tel: +47 98222087

component (i.e. principal component analysis PCA) and a supervised one (i.e. partial least square discrimination analysis PLS-DA). This model was trained using a training set of 20 drinking water samples previously tested via conventional suspect screening. The developed model was validated using a validation set of 20 drinking water samples of which 4 were spiked with 15 labeled standards at four different concentration levels. The model successfully detected all of the added analytes in the four spiked samples without producing any cases of false detection. The same validation set was processed via conventional trend analysis in order to cross validate the composite model. The results of cross validation showed that even though the conventional trend analysis approach produced a false positive detection rate of $\leq 5\%$ the composite model outperformed that approach by producing zero cases of false detection. Additionally, the validated model went through an additional test with 42 extra drinking water samples from the same source for an unbiased examination of the model. Finally, the potentials and limitations of this approach were further discussed.

Keywords:

Machine learning; Non-target; LC-HRMS; Drinking water; Statistical modeling

1. Introduction

Providing clean drinking water is crucial for sustaining human health and it is therefore defined as one of the UN goals for sustainable development [1].

4 According to the World Health Organization, providing clean drinking water
5 to the global population may reduce the worldwide disease by $\sim 10\%$. To-
6 day, factors such as: urbanization, human chemical footprint (i.e. chemical
7 production, consumption, and release), global water scarcity due to climate
8 change, and population growth are making the production and distribution
9 of clean drinking water to the global population a challenging task [2–4].
10 The situation is far from static as the challenges grow and change, and this
11 is evident in the ever evolving water quality monitoring programs across the
12 globe. However, during the past two decades, it has become more and more
13 evident that existing water monitoring strategies are not adequate to address
14 these challenges [2–6].

15

16 Non-target analysis using liquid chromatography coupled with high reso-
17 lution mass spectrometry (LC-HRMS) has been the leading analytical strat-
18 egy to tackle the challenges faced by a diffuse and highly dynamic chemical
19 footprint [7–12]. This approach (i.e. non-targeted LC-HRMS), differently
20 from typically limited and targeted routine monitoring strategies, is not bi-
21 ased towards a small number of target analytes. However, it generates highly
22 complex datasets with thousands of features to be analyzed for each sample.
23 To deal with such large and complex datasets the analysts have to isolate the
24 environmentally relevant features from the generated features lists (i.e. pri-
25 oritization) [7, 9, 13, 14]. Prioritization may be performed using the intensity
26 of the features and/or based on the statistical significance of those features

27 when compared to the samples from different origins, for example [7, 13].
28 Intensity based prioritization is relatively fast, but it ignores the lower inten-
29 sity features, which may be relevant. Therefore, a statistical approach may
30 be more adequate for analysis of water samples, including drinking water.

31

32 Recently, advanced statistical tools such as machine learning algorithms
33 have been utilized for regression, dimension reduction, and sample classifi-
34 cation via simple or composite models [15]. This approach is a widely used
35 method for prediction of chemical and physical properties of compounds [16].
36 However, to our knowledge it has never been used in combination with non-
37 targeted LC-HRMS data for monitoring of water samples.

38

39 The aim of this study was to develop a risk warning system of potential
40 chemical hazards in drinking water by combining non-targeted LC-HRMS
41 and machine learning. The drinking water samples (i.e. 82 samples) were
42 divided in three groups: 20 samples for a training set, 20 samples for a
43 validation set, and 42 samples for a test set. The training set was used for
44 the model development whereas the test set and the validation sets were
45 utilized for the model validation. During the model validation and test, we
46 cross validated our model via trend analysis and suspect screening.

47 **2. Methods**

48 *2.1. Chemicals*

49 All chemical standards and solvents (ACS grade) were purchased from
50 Novachem Pty Ltd. (Victoria, Australia) whereas the technical grade filters
51 were obtained from Phenomenex. A complete list of the labeled internal stan-
52 dards, their measured retention time, and their measured masses is provided
53 in the Supporting Information, Section S1.

54 *2.2. Environmental Samples and Sample Processing*

55 In total 82 drinking water samples of 1 L each were received from South
56 East Queensland, Australia, during a 6 week time period between March and
57 April 2018. Each sampling day consisted of six water samples, except two
58 instances with five samples, taken during the day with intervals larger than
59 1 hr. The samples were treated drinking water directly from six treatment
60 plants with the same source water and treatment processes. The samples
61 were delivered to the lab at 4°C and were immediately processed and analyzed
62 (i.e filtered and spiked with internal standards). For the analysis, all 82
63 drinking water samples were filtered using 2 μm filters and an aliquot of each
64 was transferred into 1.5 mL vials having a final volume of 1 mL, without any
65 further processing. All the samples were spiked with 5 μL of a 1 ppm stock
66 solution of caffeine ^{13}C to obtain an injection standard (i.e. caffeine ^{13}C)
67 concentration of 5 ppb. The sample preparations were kept to minimum in
68 order to avoid any type of cross-contamination of the samples.

69 *2.3. Instrumental Conditions and Analysis*

70 All 82 drinking water samples were analyzed using a Sciex ExionLC chro-
71 matography system coupled to a Sciex X500R QTOF mass spectrometer (AB
72 SCIEX, USA). Ten μL of each sample was directly injected into the system
73 and separated with Kinetex Biphenyl column (50×2.1 mm, $2.6 \mu\text{m}$, Phe-
74 nomenex) at 50°C . The separations were carried out using 0.1% formic acid
75 in MilliQ water as mobile phase A and 0.1% formic acid in methanol as mo-
76 bile phase B at a flow rate of 0.4 mL/min. The gradient started at 0% B for
77 0.5 min, then ramped up to 100% B in 9.5 min with a non-linear Curve (con-
78 vex) and maintained at 100% B until 14.5 min before returning to 0% B for
79 equilibration. The mass spectrometer was equipped with a TurboIonSpray
80 ion source and operated employing Electron Spray Ionization (ESI) source in
81 positive mode with data-independent acquisition. During pseudo MS^2 scans,
82 the collision energy (CE) was set at 35 eV (more details are provided else-
83 where [12]). These instrumental conditions were previously optimized for
84 these type of analysis [12, 17, 18].

85

86 For quality control, all the glassware used during the analysis were baked
87 overnight at 450°C . We did not expect a large level of variability in the
88 samples due to the simplicity of the matrix (i.e. drinking water) as was
89 previously observed for similar matrices [12, 17, 18]. Moreover, each five
90 samples were followed by a blank injection, which consisted of a MilliQ water
91 spiked with the labeled internal standards (Section S1). All the analyzed

92 blanks were procedural blanks and were treated in the same way as the
93 samples. The samples were injected in a randomized order.

94 *2.4. Experimental Setup*

95 The 82 drinking water samples were divided into three categories: the
96 training set (20 samples), the validation set (20 samples), and the test set
97 (42 samples), Fig. 1. For the training set and validation set, we selected
98 a 50% division of the data in order to avoid any over-training of the model
99 [15, 19]. With regards to the test set, we used a large test set in order to as-
100 sess if a large enough training set was used for the model generation. In other
101 words, a small training set would result in a large number of false positive
102 detection during the model test. The training samples were employed dur-
103 ing the machine learning algorithm development (i.e. the composite model)
104 whereas the test set samples were used for an unbiased performance evalu-
105 ation of the model. Four out of 20 validation set samples were spiked with
106 a mixture of 15 labeled internal standard at 2.5 ppb, 5 ppb, 10 ppb, and 20
107 ppb of each internal standard in addition to caffeine ^{13}C . Prior to the model
108 development all these 40 water samples were subject to conventional suspect
109 screening in order to assess their quality (see Section 2.7 for more details).
110 Finally, we employed the developed and validated model to assess the quality
111 of the test set (i.e. 42 water samples). During the model validation and model
112 test steps we included two different cross validation steps, which consisted
113 of processing the same dataset with two conventional methods (i.e. trend

114 analysis and suspect screening). Further details regarding both the trend
115 analysis and the suspect screening are provided in Sections 2.5 and 2.7).
116 Moreover, the validated model was further examined via synthetic datasets
117 where 5 randomly selected samples from the validation set were added to the
118 test set. This process was repeated 50 times to further test the applicability
119 of the model for different drinking water samples. This implied during each
120 iteration a random combination of the spiked and unspiked samples were
121 added to the test set for further evaluation. Doing so enabled us to truly
122 evaluate the likelihood of false detection of the model. It should be noted
123 that the samples did not go through any sample pre-concentration and the
124 concentration of each spiked standard at the lowest concentration level (i.e.
125 2.5 ppb) was close to the measured limit of detection for the same standards
126 (i.e. ~ 1.5 ppb or 15 pg on column).

127

128 Using this experimental design, we were able to first build our model via
129 the training set, validate the model using the validation set, and test the
130 model through the synthetic test set.

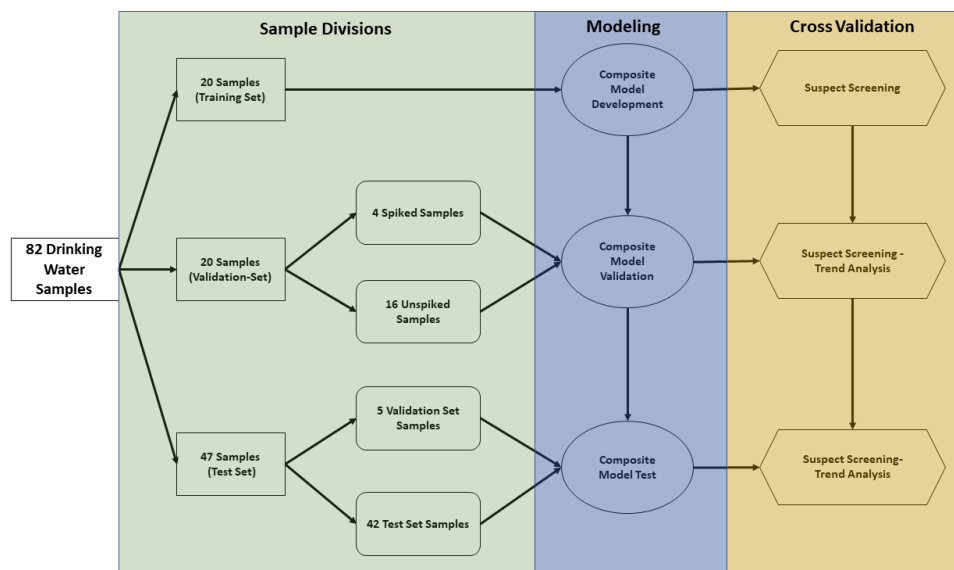


Figure 1: The schematic of the workflow employed in this study including sample division, modeling, and cross validation.

131 *2.5. Machine Learning Algorithm Workflow*

132 The acquired chromatograms for all the samples, including the training
133 set, validation set, and the test set went through the following steps sequen-
134 tially: 1) peak picking, 2) peak alignment, 3) correction for the background
135 variability, 4) standardization, and finally 5) modeling. The workflow was
136 divided in two parts pre-processing, which included steps 1 to 4 and the
137 modeling which was the fifth step in the complete workflow.

138 *2.5.1. Data Pre-processing*

139 All the chromatograms were peak picked using Sciex OS 1.4 (AB SCIEX,
140 USA) employing a minimum peak area of 1000 counts and a signal to noise
141 ratio of five. After the peak picking, we used Sciex OS for the alignment
142 of the chromatograms, which employed a maximum peak width of 6 s in
143 the time dimension whereas the mass window was set to 0.003 Da. Both
144 of these parameters were selected based on the reported peak boundaries in
145 the time dimension and the observed mass error in m/z values (i.e. ± 0.003
146 Da). The aligned peak list at this stage went through the correction for
147 the background variability. This step enabled us to correct for the variability
148 observed in the background signal caused by the instrument fluctuations [20].
149 We employed the C13 labeled caffeine signal for the background variability
150 correction of all the datasets, including the training set, validation set, and
151 the test set. For the background variability correction, the signal of all the
152 features in a sample was divided by the intensity recorded for the injection

153 standard (i.e. C13 caffeine) in that sample. The last step of the data pre-
154 processing consisted of standardization via Pareto method [20], which divides
155 the intensity of each feature (i.e. the tensor of m/z value, retention time, and
156 intensity) by the square root of the standard deviation of that feature across
157 all the chromatograms. The standardization reduces the variability range of
158 each feature, thus giving the same importance to each feature independently
159 from their intensity. Following the above-mentioned steps enabled us to
160 adequately prepare our data for the modeling steps.

161 *2.5.2. Composite Model Development*

162 The training set, consisting of the pre-processed peak-list (i.e. m/z, re-
163 tention time, and the relative intensities) of 20 drinking water samples, was
164 used for the composite model development. The purpose of this model was
165 to describe the chemistry of the unspiked water, through modeling the max-
166 imum variance in the training set for each feature. Therefore, an observed
167 larger variability for a certain feature during the validation step implied the
168 presence of an abnormality or a potential chemical hazard in that sample.
169 The validation set employed in this study included 20 drinking water samples
170 from which 4 were spiked with 15 labeled internal standards. The validation
171 set was generated in a double blind manner to comprehensively evaluate the
172 capability of the model in distinguishing the clean water samples from the
173 spiked ones. Both the training set and the validation set were also employed
174 for tuning the model parameters. For both the model building, model val-

175 idation, and model test we employed a non-targeted approach utilizing all
176 the features in the samples. Therefore, our model was based on the complete
177 chemical composition of the LC-HRMS analyzable fraction of the drinking
178 water samples. This implied that theoretically only one statistically mean-
179 ingful feature was enough for distinguishing the spiked samples from unspiked
180 ones.

181

182 Our model consisted of a linear combination of principal component anal-
183 ysis (PCA) [19] and partial least square discrimination analysis (PLS-DA)
184 [21, 22] modeling approaches, that enabled the confident separation of the
185 unspiked drinking water samples from the spiked ones. The PCA modeling
186 approach is an unsupervised method, which enables an unbiased evaluation
187 of the underlying trends in the data. However, given its nature [19], the
188 PCA is less sensitive towards small changes in the data. PLS-DA, on the
189 other hand, is a supervised approach, which takes advantage of the prior
190 knowledge of the data [21, 22]. In other words, this method utilizes the user
191 defined classification in the training set to create the model. This implies
192 that the model is forced to give a higher importance to certain variables,
193 that are causing the separation of the pre-defined groups from each other.
194 However, this method suffers from overfitting issues [21, 22]. We used a lin-
195 ear combination of the two modeling approaches in order to fully harvest the
196 higher sensitivity of PLS-DA and at the same time take full advantage of the
197 robustness of PCA.

199 For the PCA modeling, we used the singular value decomposition al-
200 gorithm [23] given the larger number of variables (i.e. features) than the
201 number of the measurements (i.e. the drinking water samples). We used
202 the sum of the absolute values of the scores for the first two PCs as the
203 output of the PCA model (S_{PCA}). The choice of using only the first two
204 PCs was based on the fact that these two PCs combined described $\geq 50\%$
205 of the observed variability in our dataset, which indicates the existence of
206 an underlying trend [19]. The same pre-processed training set was used for
207 PLS-DA model building. During the training step, the PLS-DA was trained
208 only using the unspiked samples, which enabled the generation of a highly
209 sensitive model. One of the crucial steps in the PLS-DA modeling is the se-
210 lection of the number of components to generate the model, in order to avoid
211 overfitting issues [21, 22]. This choice was carried out through an optimiza-
212 tion process employing the training set. We performed 100 simulations where
213 15 samples were randomly selected from the training set for each iteration.
214 A new PLS-DA model was generated during each simulation with new score
215 values and components. We also recorded the number of necessary compo-
216 nents to describe 95% (i.e. 95% confidence interval) of the variability in the
217 data for each simulation. The results of these simulations indicated that
218 four components were necessary to describe 95% of the variability in all the
219 simulated cases. Therefore, we limited the number of PLS-DA components
220 to three, in order to avoid overfitting issues [21, 22]. When calculating the

221 contribution of the score values of the components on the S_{PLS-DA} value, the
 222 first component contributed more than 50% of the S_{PLS-DA} . Consequently,
 223 for simplicity we only included the X-score (i.e. the score value associated to
 224 the predictor block) of the first component in the PLS-DA score calculations
 225 (i.e. S_{PLS-DA}). The selection of the number of components in the PLS-DA
 226 model is case dependent and must be evaluated during the model creation
 227 for each dataset. Finally, we generated a score value for the final model,
 228 hereafter referred to as final score (S_{final}), for each drinking water sample in
 229 the training set. The S_{final} was a weighted linear combination of the S_{PCA}
 230 and S_{PLS-DA} (Eq. 1). In Eq. 1 the S_{PCA} , S_{PLS-DA} , and S_{final} were the
 231 score values from PCA model, PLS-DA model, and the final model, respec-
 232 tively while the w_{PCA} and w_{PLS-DA} were the weight value associated with
 233 PCA and PLS-DA score values (Eq. 1). The training set was employed to
 234 optimize the weight values as such to produce S_{final} values ranging between
 235 -1 and 1. While performing the weight value optimization, we utilized the
 236 likelihood of false positive detection as the selection criteria for the tested
 237 weight values. The details of this process is described below, Section 3.1.

$$S_{final} = w_{PCA} \cdot S_{PCA} + w_{PLS-DA} \cdot S_{PLS-DA} \quad (1)$$

238 2.6. Trend Analysis

239 We performed trend analysis [5, 12, 24–26] on the validation set in order
 240 to compare the performance of the composite statistical model with a more

241 conventional approach. During these analysis, we produced the signal inten-
242 sity of each feature for all the samples including the pre-processed training set
243 and validation set. In this case we singled out the features that were enriched
244 at a statistically significant levels through the comparison of the median of
245 a feature across all the samples (i.e. background) to the intensity of that
246 feature in each sample (signal). For a feature to be considered statistically
247 significant, it had to produce a signal to background ratio of five in the vali-
248 dation set. Consequently, a feature that met all these criteria was considered
249 a statistically significant feature and was selected for post-processing (e.g.
250 identification). The signal to background ratio of five was selected based on
251 the observed variability of the features in the training set, which enabled us
252 to minimize the likelihood of false positive detection.

253 *2.7. Suspect Screening*

254 The samples for the case study were suspect screened using a local li-
255 brary of pesticides, pharmaceuticals, personal care products, illicit drugs,
256 and industrial chemicals (3000 chemicals), provided with the vendor soft-
257 ware package. We employed LibraryView package provided by Sciex OS for
258 these analyses. We utilized a mass accuracy of ± 0.003 Da and at least 3
259 matched fragments, in order to confidently identify a suspect analyte. These
260 criteria were previously shown to be effective in processing such datasets
261 [10, 12, 17, 18, 27, 28].

262 *2.8. Rate of False Detection*

263 We also evaluated the rate of false detection (i.e. false positive and/or
264 false negative) [29, 30] of the features that were isolated via the composite
265 model and/or trend analysis. A selected feature was considered a false posi-
266 tive when its accurate mass, retention time, or the sample order, during the
267 analysis, did not match the same parameters of the added internal standards.
268 On the other hand, a feature was assumed a false negative if it was added into
269 a sample as an added internal standard and it was not selected by either the
270 composite model or trend analysis as a statistically significant feature. This
271 appeared to be reasonable given that the thresholds for positive detection in
272 both the composite model and the trend analysis were set as such to produce
273 zero cases of false positive detection for the training set.

274

275 Using the rates of false detection, we were able to comprehensively com-
276 pare the performance of the composite statistical model to the more conven-
277 tional approach of trend analysis.

278 *2.8.1. Computations*

279 All the data manipulations and modeling were performed on a personal
280 computer with an i7 processor and 16 GB of memory using Matlab 2015b
281 [31].

282 3. Results and Discussion

283 A machine learning algorithm was developed and validated for a risk
284 warning system for chemical hazards in drinking water, using the data pro-
285 duced via non-targeted LC-HRMS. The machine learning algorithm took
286 advantage of a composite statistical model, which used a linear combination
287 of a supervised method (i.e. PLS-DA) and an unsupervised approach (i.e.
288 PCA). The composite statistical model utilized all the features present in
289 the sample, thus a non-targeted approach. This composite model utilizes the
290 training set to learn about the variability range of each feature in drinking
291 water samples. Consequently, if one or more of the features in the valida-
292 tion/test samples has a larger intensity compared to its observed variability
293 in the training set, the model will generate a large S_{final} value, which is
294 translated into a trigger for the risk warning system. We validated the de-
295 veloped model employing a validation set of 20 drinking water samples from
296 which 4 were spiked with 15 labeled internal standards at different concen-
297 tration levels, ranging from 2.5 ppb to 20 ppb. The spiked samples were
298 used for evaluation of false negative and false positive detection rates while
299 the unspiked samples were used for the assessment of false positive detec-
300 tion. We also compared the performance of the model with the conventional
301 trend analysis, typically used for processing this type of data. Finally, the
302 validated model was further tested in processing of 42 water samples along-
303 side with conventional suspect screening. This is, to our knowledge, the first
304 study using the combination of machine learning and non-target analysis for

305 a risk warning system of chemical hazards in water samples. Also it should
306 be noted that this is a proof of concept study and further implementation of
307 this approach on more complex samples are necessary and will be subject of
308 our future studies.

309 *3.1. Optimization of the Machine Learning Algorithm*

310 The training set was used to select the weight values as well as the thresh-
311 olds of false positive detection for the composite model. In order to select
312 these parameters, we ran 50,000 (400×125) simulations where for each it-
313 eration 18 randomly selected samples out of 20 samples in the training set
314 were used to generate the final composite model. In order to perform this
315 optimization, a squared matrix of weight values varying between 0 and 2
316 with steps of 0.1 was generated (i.e. a matrix of 20×20 , thus 400 members
317 in the matrix). At each point in this matrix 125 simulations took place for
318 false detection calculations. Employing this approach, we generated a dis-
319 tribution of S_{final} values for unspiked drinking water samples enabling us to
320 calculate the rate of false positive detection for different weight values. This
321 optimization process indicated that the best weight values were 0.1 and 1
322 for w_{PCA} and w_{PLS-DA} , respectively, producing the smallest cases of false
323 positive detections.

324

325 The S_{final} values of 1, 1.2, and 1.5 resulted in false positive detection
326 likelihoods of 5.0%, 1.0%, and 0.1%, respectively, employing the optimized

327 weight values (Fig. 2). In order to further evaluate the likelihood of false
328 positive detection, 5,000 simulations were performed using the pre-processed
329 training set and the optimized weight values, which resulted in a distribution
330 of the S_{final} values. These values in the generated distribution then were con-
331 verted into the likelihood of false positive detection [29, 30]. For this study, a
332 S_{final} value of 1.2 was selected as the threshold for a statistically significant
333 warning for a potential chemical hazard risk. The selected likelihood of false
334 positive detection enabled us to associate a high level of confidence to the
335 samples that produced an S_{final} value of ≥ 1.2 .

336

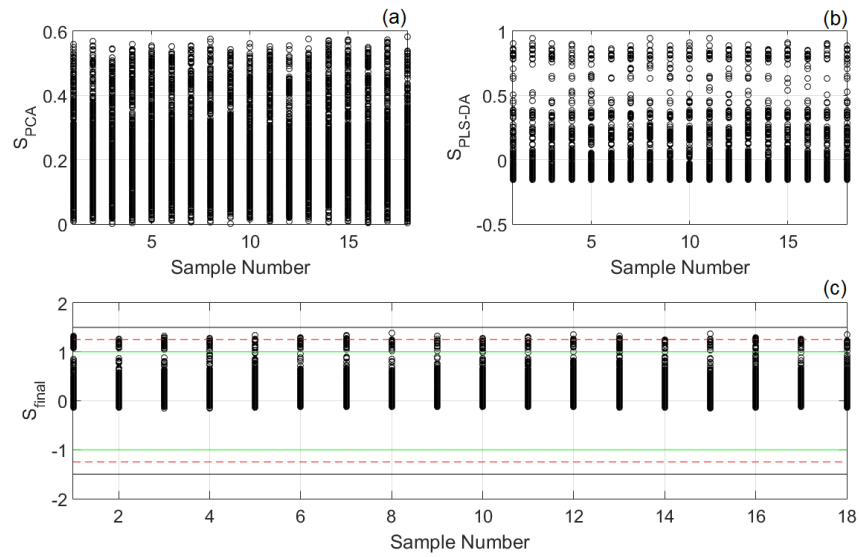


Figure 2: The (a) weighted S_{PCA} , (b) weighted S_{PLS-DA} , and the S_{final} values calculated via Eq. 1 for 5,000 simulations with weight values of 0.1 and 1 for PCA and PLS-DA models. The green line, red dotted line, and black line in panel (c) define the S_{final} values of 1, 1.2, and 1.5, respectively.

337 *3.2. Model Validation via Validation Set (i.e. Spiked Samples)*

338 The previously developed model was validated using the pre-processed
339 validation set. For the model validation, we replaced one of the samples in
340 the training set (randomly selected between sample 2 and sample 19 of the
341 training set) with one of the samples in the validation set. At this stage the
342 PLS-DA model was forced to consider the added sample as a spiked sam-
343 ple whether it was spiked or not. This implied that for a spiked validation
344 sample both models (i.e. PCA and PLS-DA) produced larger score values,
345 and consequently a large final score. On the other hand, for a non-spiked
346 validation sample, considered as spiked sample, the PLS-DA model produced
347 a large score value whereas the PCA model generated a small score, which
348 resulted in a small final score. The random selection of the location of the
349 added validation sample into the training set was due to the fact that we
350 wanted to be sure that the location of the sample addition did not affect the
351 outcome of the algorithm. We evaluated each sample in the validation set
352 using the above mentioned procedure in an iterative way.

353

354 The proposed machine learning algorithm (i.e. the composite statistical
355 model) detected all the 4 spiked samples without producing any cases of false
356 positive and/or false negative detections, Fig. 3. The S_{final} values ranged
357 from 1.27 (Fig. S1) for the sample in the validation set with the lowest spike
358 level (i.e. 2.5 ppb) to 2.5 (Fig. S2) for the sample spiked with 20 ppb of
359 the standard mixture. For all the samples that were not spiked with internal

standards the S_{final} was ≤ 1 , Fig. S3. By looking at the ratio of the loading values of PCA and PLS-DA model, we were able to identify the features that were the cause of the abnormality (Fig. S4). Based on the absolute intensity of the loading ratios, the top 95.0% of the features were selected for isolating those that were describing the large S_{final} values. This resulted in selection of 15 features, which belonged to the labeled standards. For example, a feature with loading value ratio of 8325 and 9330 for PCA and PLS-DA, respectively, was associated with the signal of carbamazepine D10, which was one of the added internal standards. Additionally, we evaluated the model limit of detection (LOD_{model}) for the tested 15 standards using the response factor calculated based on the slope of the standard addition calibration curve of the spiked samples. The composite model resulted in an averaged LOD_{model} of $\sim 1.8 \pm 0.3$ ppb for evaluated internal standards, which was comparable to the measured instrument LOD for these standards of ~ 1.5 ppb. This was performed by artificially reducing the signal of each internal standard in the validation set employing 0.01 ppb steps until the model was not able to distinguish the spiked samples from the unspiked training set. The last detectable signal for an internal standard was considered the LOD_{model} for that standard. Furthermore, we compared the LOD_{model} of the composite model (i.e. combined PCA and PLS-DA) to each of the models individually. The limit of detection of the PCA model (LOD_{PCA}) alone appeared to be $\sim 12.0 \pm 1$ ppb across all 15 labeled analyts, which was around 6 times larger than the LOD_{model} . When using the PCA model alone for analysis of the

383 validation set, this model produced four cases of false negative detections and
384 no cases of false positive detection. On the other hand, for the PLS-DA, this
385 model resulted in 6 cases of false positives and zero cases of false negative
386 detection for the processing of the validation set. This was due to the lower
387 LOD of PLS-DA model (LOD_{PLS-DA}) of $\sim 1.0 \pm 0.2$ ppb. These results
388 indicated the higher performance of the composite model compared to each
389 of the individual models suggesting high sensitivity and robustness of the
390 final composite model.

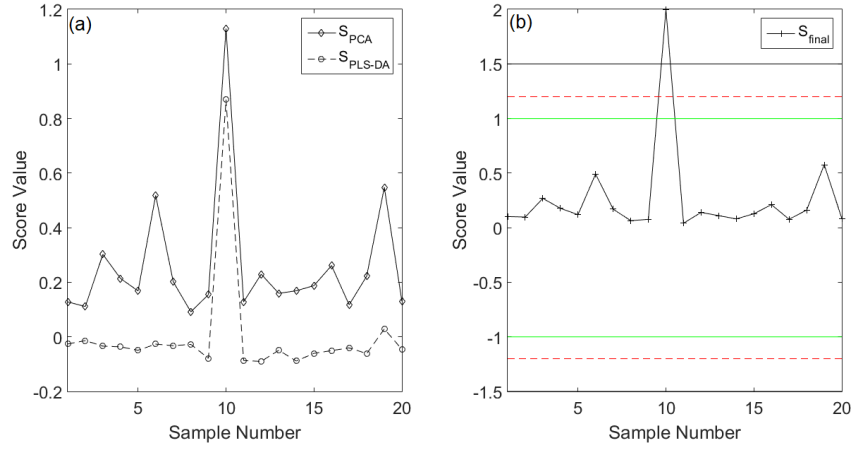


Figure 3: The score values for (a) PCA model (S_{PCA}), PLS-DA model (S_{PLS-DA}) and (b) the composite model score value calculated using Eq. 1. In this instance the sample number 10 was the spiked sample.

391 *3.3. Comparison Between the Composite Model and Conventional Trend Anal-*
392 *ysis*

393 We compared the performance of the composite model with the conven-
394 tional trend analysis, which is commonly used for detection of pulsed point
395 source into the water samples [24]. During the trend analysis, we selected a
396 signal to background ratio of five, for a feature to be considered statistically
397 significant. We used the pre-processed training set and the validation set for
398 the comparison between the two methods (Section 2.6).

399
400 The trend analysis approach produced 30 cases of false positive detections
401 (i.e. a false positive rate of $\leq 5\%$ [29, 30]) without producing any cases of
402 false negatives. On the other hand the composite model was able to detect
403 all 15 spiked analytes in all 4 samples without producing any cases of false
404 positive and/or false negative. For 27 out of 30 (i.e. 90%) of the false positive
405 cases manual inspection of the features caused their elimination from the list.
406 These features appeared to have low intensity and high level of variability
407 across all the samples, including the training set, Fig S5. The remaining
408 three features identified as false positives appeared to be the isotopes of a
409 real features. For example, a feature identified as a false positive with m/z
410 value of 181.083 and a retention time of 4.50 min was the M+2 isotope of
411 atrazine desisopropyl D5 with an accurate mass of 179.085 and the retention
412 time of 4.48 min. We further investigated the 3 meaningful features isolated
413 via trend analysis in the composite model. For those three features, their

414 loading values were smaller than 95%. Therefore, the composite model did
415 not consider these features as statistically significant, which further indicates
416 the robustness of this approach compared to the conventional method (i.e.
417 the trend analysis).

418

419 The composite model (i.e. the machine learning algorithm) performed
420 better than the conventional trend analysis when applied to the validation
421 set. This method was able to capture all the added analytes in the spiked
422 samples without producing any cases of false detections. This method showed
423 to be less sensitive to the high variability in the data compared to the con-
424 ventional trend analysis method. However, further tests are necessary to
425 comprehensively evaluate the effect of noise on such a model. Overall, this
426 showed to be a sensitive, accurate, and reliable tool for capturing contami-
427 nation in the drinking water.

428 *3.4. Further Testing via Test Set*

429 We further tested the capability of the composite model in distinguishing
430 a spiked water sample from an unspiked one. Additionally, this final test en-
431 abled us to evaluate the applicability of the same training set for a different
432 batch of water samples taken from the same source (Section 2.4).

433

434 The composite model produced 3 cases of false negative and zero cases
435 of false positive detection out of the total 2350 (i.e. 47 samples \times 50 sim-

436 ulations) evaluations during the test. All 3 cases belonged to the spiked
437 water samples at lowest concentration level (i.e. 2.5 ppb). Both the com-
438 posite model and the conventional suspect screening did not produce any
439 abnormality cases for the test set, which was expected considering that these
440 samples were treated drinking water.

441

442 The outcome of the composite model was in agreement with the conven-
443 tional suspect screening, which is indicative of its robustness. However, more
444 complex matrices should be tested in order to further evaluate the applicabil-
445 ity of this method. Analysis of more complex matrices such as ground water
446 and surface water will be subject of our future study. Finally, it should
447 be noted that this study is a proof of concept for applicability of such an
448 approach for water related matrices.

449 *3.5. Potential and Limitations*

450 The developed and validated composite model was shown to be a reliable,
451 robust and accurate method for detection of anomalies (i.e. potential con-
452 taminants) in drinking water samples. The thresholds for the risk warning
453 could be set by the acute and adverse toxicity of the drinking water samples,
454 which will expand the applicability of this method to monitoring of both
455 the produced drinking water as well as the source water used for producing
456 drinking water. At the current state, the samples were injected as is into
457 the instrument for analysis without any pre-concentration. However, addi-

458 tion of a pre-concentration step would drastically increase the sensitivity of
459 this method, based on the achieved model LODs that were similar to the
460 analytical LODs. In other words, the pre-concentration step may potentially
461 increase the sensitivity of the model by increasing the instrument sensitivity.
462 Moreover, this method is designed to screen the samples rapidly for anomalous
463 features. In addition to the triggered warning, the model will produce a list of
464 features that are causing anomalies, which should be evaluated by the analyst.
465 In practical terms, the analyst can focus only on the samples that
466 triggered a warning and the selected features rather than all the features and
467 samples, therefore simultaneous sample and feature prioritization. Finally,
468 this method could be employed for continuous monitoring of more complex
469 aqueous matrices as long as the observed variability in the training set is
470 representative of the normal state of that matrix.

471

472 It should be noted that this method was applied to the peak list in the
473 current study due to the cleanness of the drinking water matrix. However,
474 for more complex matrices, this method should be applied to the raw data in
475 order to be able to model the variability observed in the data. This implies
476 a drastic increase in its computational cost. The warning thresholds are
477 highly dependent on the observed within feature variability of the training
478 set. Consequently, the analyst must assure that the variability in the training
479 set is representative for the variability present in the test set in a normal state,
480 which is also necessary for the conventional trend analysis. In other words, if

481 the variability in the training set is too large, the model would lose sensitivity
482 (i.e. producing false negatives) whereas if the variability in the training
483 set is too small, then the model will become too sensitive (i.e. producing
484 false positives). Similarly to the trend analysis, given the dependency of
485 the explored chemical space on the analysis conditions[14, 32], the training
486 sets are specific to a sample set and analysis conditions. Therefore, a good
487 understanding of both the matrix and the analytical instrument is crucial to
488 the success of this approach.

489 **4. Acknowledgement**

490 The authors are thankful to the the Research Council of Norway for the
491 financial support of this project (RESOLVE, 243720). We are grateful to Dr.
492 Sharon Grant and Dr. Jake O'Brien for their comments during the project
493 development.

494 **5. Supporting Information**

495 Supporting Information containing the list of standards and complemen-
496 tary figures is available as stated in the text.

- 497 [1] UN, UN Goals for Sustainable Development,
498 <https://www.un.org/sustainabledevelopment/water-and-sanitation/>
499 (May 2010).
- 500 [2] S. R. Newton, R. L. McMahan, J. R. Sobus, K. Mansouri, A. J. Williams,
501 A. D. McEachran, M. J. Strynar, Suspect screening and non-targeted
502 analysis of drinking water using point-of-use filters, *Environmen. Pollu.*
503 234 (2018) 297–306.
- 504 [3] A. Pal, Y. He, M. Jekel, M. Reinhard, K. Y.-H. Gin, Emerging contami-
505 nants of public health significance as water quality indicator compounds
506 in the urban water cycle, *Environ. Int.* 71 (2014) 46–62.
- 507 [4] S. D. Richardson, T. A. Ternes, Water analysis: emerging contaminants
508 and current issues, *Anal Chem.* 90 (1) (2017) 398–428.
- 509 [5] S. D. Richardson, Water analysis: emerging contaminants and current
510 issues, *Anal. Chem.* 81 (12) (2009) 4645–4677.
- 511 [6] Y. Peng, S. Hall, L. Gautam, Drugs of abuse in drinking water—a review
512 of current detection methods, occurrence, elimination and health risks,
513 *TrAC Trends Anal. Chem.* 85 (2016) 232–240.
- 514 [7] T. Bader, W. Schulz, T. Lucke, W. Seitz, R. Winzenbacher, Applica-
515 tion of non-target analysis with lc-hrms for the monitoring of raw and
516 potable water: Strategy and results, in: *Assessing Transformation Prod-*

- 517 ucts of Chemicals by Non-Target and Suspect Screening- Strategies and
518 Workflows Volume 2, ACS Publications, 2016, pp. 49–70.
- 519 [8] N. A. Alygizakis, S. Samanipour, J. Hollender, M. Ibáñez, S. Kaserzon,
520 V. Kokkali, J. A. van Leerdam, J. F. Mueller, M. Pijnappels, M. J. Reid,
521 et al., Exploring the potential of a global emerging contaminant early
522 warning network through the use of retrospective suspect screening with
523 high-resolution mass spectrometry, *Environ. Sci. Technol.* 52 (9) (2018)
524 5135–5144.
- 525 [9] E. L. Schymanski, H. P. Singer, J. Slobodnik, I. M. Ipolyi, P. Oswald,
526 M. Krauss, T. Schulze, P. Haglund, T. Letzel, S. Grosse, et al, Non-
527 target screening with high-resolution mass spectrometry: critical re-
528 view using a collaborative trial on water analysis, *Anal. Bioanal. Chem.*
529 407 (21) (2015) 6237–6255.
- 530 [10] S. Samanipour, M. J. Reid, K. Bæk, K. V. Thomas, Combining a de-
531 convolution and a universal library search algorithm for the nontarget
532 analysis of data-independent acquisition mode liquid chromatography-
533 high-resolution mass spectrometry results, *Environ. Sci Technol.* 52 (8)
534 (2018) 4694–4701.
- 535 [11] E. L. Schymanski, H. P. Singer, P. Longrée, M. Loos, M. Ruff, M. A.
536 Stravs, C. Ripollés Vidal, J. Hollender, Strategies to characterize polar
537 organic contamination in wastewater: exploring the capability of high

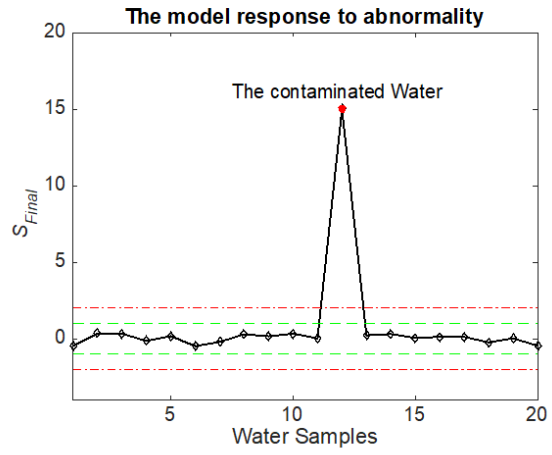
- 538 resolution mass spectrometry, *Environ. Sci. Technol.* 48 (3) (2014) 1811–
539 1818.
- 540 [12] S. L. Kaserzon, A. L. Heffernan, K. Thompson, J. F. Mueller, M. J. G.
541 Ramos, Rapid screening and identification of chemical hazards in surface
542 and drinking water using high resolution mass spectrometry and a case-
543 control filter, *Chemosphere* 182 (2017) 656–664.
- 544 [13] S. Samanipour, M. J. Reid, K. V. Thomas, Statistical variable selection:
545 An alternative prioritization strategy during the non-target analysis of
546 LC-HR-MS data, *Anal. Chem.* 89 (10) (2017) 5585–5591.
- 547 [14] S. Samanipour, J. A. Baz-Lomba, M. J. Reid, E. Ciceri, S. Rowland,
548 P. Nilsson, K. V. Thomas, Assessing sample extraction efficiencies for
549 the analysis of complex unresolved mixtures of organic pollutants: A
550 comprehensive non-target approach, *Anal. Chim. Acta* 1025 (2018) 92–
551 98.
- 552 [15] I. Goodfellow, Y. Bengio, A. Courville, Y. Bengio, *Deep learning*, Vol. 1,
553 MIT press Cambridge, 2016.
- 554 [16] T. H. Miller, J. A. Baz-Lomba, C. Harman, M. J. Reid, S. F. Owen,
555 N. R. Bury, K. V. Thomas, L. P. Barron, The first attempt at non-
556 linear in silico prediction of sampling rates for polar organic chemical
557 integrative samplers (pocis), *Environmen. Sci. Technol.* 50 (15) (2016)
558 7973–7981.

- 559 [17] S. Kaserzon, E. O'Malley, K. Thompson, C. Paxman, G. Elisei, G. Ea-
560 glesham, M. Gallen, J. Mueller, Catchment and drinking water quality
561 micro pollutant monitoring program–passive sampling. report 6–summer
562 2017 and summary report.
- 563 [18] S. Kaserzon, C. Gallen, K. Thompson, C. Paxman, J. O'Brien, G. Eagle-
564 sham, M. J. G. Ramos, M. Gallen, D. Drage, X. Wang, et al., Catchment
565 and drinking water quality micro pollutant monitoring program-passive
566 sampling. report 1 2014.
- 567 [19] R. G. Brereton, Applied chemometrics for scientists, John Wiley & Sons,
568 2007.
- 569 [20] R. A. van den Berg, H. C. Hoefsloot, J. A. Westerhuis, A. K. Smilde,
570 M. J. van der Werf, Centering, scaling, and transformations: improving
571 the biological information content of metabolomics data, BMC genomics
572 7 (1) (2006) 142.
- 573 [21] R. Kramer, Chemometric techniques for quantitative analysis, CRC
574 Press, 1998.
- 575 [22] R. G. Brereton, G. R. Lloyd, Partial least squares discriminant analysis:
576 taking the magic away, J. Chemometrics 28 (4) (2014) 213–225.
- 577 [23] G. H. Golub, C. Reinsch, Singular value decomposition and least squares
578 solutions, Numerische mathematik 14 (5) (1970) 403–420.

- 579 [24] M. Ruff, M. S. Mueller, M. Loos, H. P. Singer, Quantitative target and
580 systematic non-target analysis of polar organic micro-pollutants along
581 the river rhine using high-resolution mass-spectrometry–identification of
582 unknown sources and compounds, *Water Res.* 87 (2015) 145–154.
- 583 [25] Z. Li, S. L. Kaserzon, M. M. Plassmann, A. Sobek, M. J. G. Ramos,
584 M. Radke, A strategic screening approach to identify transformation
585 products of organic micropollutants formed in natural waters, *Environ.*
586 *Sci. Proc. & Imp.* 19 (4) (2017) 488–498.
- 587 [26] T. Bader, W. Schulz, K. Kuummerer, R. Winzenbacher, Lc-hrms data
588 processing strategy for reliable sample comparison exemplified by the
589 assessment of water treatment processes, *Anal. Chem.* 89 (24) (2017)
590 13219–13226.
- 591 [27] S. Samanipour, K. Langford, M. J. Reid, K. V. Thomas, A two stage
592 algorithm for target and suspect analysis of produced water via gas
593 chromatography coupled with high resolution time of flight mass spec-
594 trometry, *J. Chromatogra. A* 1463 (2016) 153–161.
- 595 [28] S. Samanipour, J. A. Baz-Lomba, N. A. Alygizakis, M. J. Reid, N. S.
596 Thomaidis, K. V. Thomas, Two stage algorithm vs commonly used ap-
597 proaches for the suspect screening of complex environmental samples
598 analyzed via liquid chromatography high resolution time of flight mass
599 spectroscopy: A test study, *J. Chromatogr. A* 1501 (2017) (2017) 68–78.

- 600 [29] D. S. Burke, J. F. Brundage, R. R. Redfield, J. J. Damato, C. A. Sch-
601 able, P. Putman, R. Visintine, H. I. Kim, Measurement of the false
602 positive rate in a screening program for human immunodeficiency virus
603 infections, *N. Engl. J. Med.* 319 (15) (1988) 961–964.
- 604 [30] Y. Benjamini, Y. Hochberg, Controlling the false discovery rate: a prac-
605 tical and powerful approach to multiple testing, *J. R. Stat. Soc.: Ser. B*
606 (Methodol.) (1995) 289–300.
- 607 [31] MATLAB version 9.1 Natick, Massachusetts: The MathWorks Inc.,
608 **2018**.
- 609 [32] S. Samanipour, M. Hooshyari, J. A. Baz-Lomba, M. J. Reid, M. Casale,
610 K. V. Thomas, The effect of extraction methodology on the recovery
611 and distribution of naphthenic acids of oilfield produced water, *Sci. Tot.*
612 *Environ.* 652 (2019) 1416–1423.

613 6. TOC



TOC for review only.