

Environmental Science Processes & Impacts

Accepted Manuscript



This is an *Accepted Manuscript*, which has been through the Royal Society of Chemistry peer review process and has been accepted for publication.

Accepted Manuscripts are published online shortly after acceptance, before technical editing, formatting and proof reading. Using this free service, authors can make their results available to the community, in citable form, before we publish the edited article. We will replace this *Accepted Manuscript* with the edited and formatted *Advance Article* as soon as it is available.

You can find more information about *Accepted Manuscripts* in the [Information for Authors](#).

Please note that technical editing may introduce minor changes to the text and/or graphics, which may alter content. The journal's standard [Terms & Conditions](#) and the [Ethical guidelines](#) still apply. In no event shall the Royal Society of Chemistry be held responsible for any errors or omissions in this *Accepted Manuscript* or any consequences arising from the use of any information it contains.



rsc.li/process-impacts

1 **Contaminant classification using cosine distance based on multiple conventional**
2 **sensors**

3
4 Shuming Liu^{1*}, Han Che¹, Kate Smith¹, Tian Chang¹

5 ¹School of Environment, Tsinghua University, Beijing, 100084, China

6 *corresponding author: shumingliu@tsinghua.edu.cn, +86 1062787964

7 **Abstract**

8 Emergent contamination events have a significant impact on water systems. After
9 contamination detection, it is important to classify the type of contaminant quickly to
10 provide support for remediation attempts. Conventional methods generally either rely
11 on laboratory-based analysis, which requires long analysis time, or on
12 multivariable-based geometry analysis and sequence analysis, which is prone to being
13 affected by contaminant concentration. This paper proposes a new contaminant
14 classification method, which discriminates contaminants in a real time manner
15 independent of contaminant concentration. The proposed method quantifies the
16 similarities or dissimilarities between sensors' responses to different types of
17 contaminants. The performance of the proposed method was evaluated using data
18 from contaminant injection experiments in a laboratory and compared with a
19 Euclidean distance-based method. The robustness of the proposed method was
20 evaluated using an uncertainty analysis. Results show that the proposed method
21 performed better in identifying the type of contaminant than the Euclidean distance

22 based method and that it could classify the type of contaminant in minutes without
23 significantly compromising the correct classification rate (*CCR*).

24

25 **Keywords**

26 contaminant classification, conventional sensor, cosine distance, early warning system,
27 water quality

28

29 **Introduction**

30 Water systems are vulnerable to contamination accidents¹⁻². For example, in April
31 2014, crude oil leaked from a petrochemical pipeline in Lanzhou, China,
32 contaminating the water source of a local water plant and introducing hazardous
33 levels of benzene into the city's tap water. Water supply to Lanzhou city was
34 suspended as a result. An intense effort is currently underway to improve analytical
35 monitoring and detection of biological, chemical, and radiological contaminants in
36 water systems. One approach for avoiding or mitigating the impact of contamination
37 is to establish an Early Warning System (EWS). EWS should provide a fast and
38 accurate means of distinguishing between normal variations and contamination events,
39 and should be able to classify the type of contaminant³.

40

41 After an EWS detects the presence of contamination, the next important issue is to
42 classify the type of contaminant. The most commonly used method for contaminant

43 classification is laboratory-based analysis, e.g. ICP-MS. The advantage of this type of
44 analysis is that it can accurately qualify and quantify the contaminant. The
45 disadvantage is that it is time-consuming. In the event of an emergent contamination
46 event, the key to all remediation attempts is time. Therefore, methods of fast
47 classification of contaminants are in great demand. One possible solution is online
48 compound-specific sensors, which need less time than laboratory-based methods⁴⁻⁸.
49 However, compound specific sensors can normally only identify one type or a small
50 group of contaminants. In this case, low efficiency or failure in contaminant
51 classification can be expected.

52

53 To overcome this drawback, several researchers have attempted to develop real-time
54 contaminant classification methods. Kröll⁹ reported the Hach HST approach using
55 multiple types of sensors for event detection and contaminant classification. In the
56 Hach HST approach, signals from 5 separate orthogonal measurements of water
57 quality (pH, conductivity, turbidity, chlorine residual, TOC) were processed from a
58 5-parameter measure into a single scalar trigger signal. The deviation signal was
59 compared to a preset threshold level. If the signal exceeded the threshold, the trigger
60 was activated⁹. The deviation vector was then used for further classification of the
61 cause of the contamination. The direction of the deviation vector relates to the agent's
62 characteristics. Seeing that this is the case, laboratory agent data can be used to build
63 a threat agent library of deviation vectors. A deviation vector from the monitor can be

64 compared to agent vectors in the threat agent library to see if there is a match within a
65 given tolerance level. This system can be used to classify what caused the trigger
66 event. Yang et al.¹⁰ reported a real-time event adaptive detection, classification and
67 warning (READiw) method for event detection and contamination classification. In
68 this method, four discrimination systems were developed to differentiate the 11 tested
69 contaminants according to the various responses of sensors. The classification process
70 was more based on geometry analysis. The similarity or dissimilarity between
71 examples and classes were not quantitatively evaluated. Olikar and Osfield¹¹
72 developed a contamination event detection method for water distribution systems,
73 which comprised a weighted support vector machine for the detection of outliers, and
74 subsequent sequence analysis for the classification of contamination events. It was
75 noticed that either geometry analysis or sequence analysis was prone to being affected
76 by the magnitude of sensor responses, which were normally related to contaminant
77 concentrations. This could then lead to misclassification.

78

79 Although effort has been put into developing methods for contaminant classification
80 in recent years, more attention is necessary. Therefore, the objectives of this study are
81 1) to develop a classification method which is independent of contaminant
82 concentration; 2) to compare the performance of the proposed method with a
83 Euclidean distance-based method.

84

85 **Materials and methods**

86 **Data collection**

87 In order to collect contamination data, a pilot-scale contaminant injection experiment
88 (CIE) platform was developed. A process flow schematic of the CIE platform is
89 shown in Figure 1. The water tank is approximately 85 cm high with a diameter of 70
90 cm, and has a total capacity of 300 L. The tank is linked with online water quality
91 sensors via a peristaltic pump at 0.5 L per minute. Eight types of sensors developed
92 by Hach Homeland Security Technologies were utilized in this study. They can
93 measure the following 8 parameters simultaneously and continuously: temperature,
94 pH, turbidity, conductivity, oxidation reduction potential (ORP), UV-254,
95 nitrate-nitrogen and phosphate. The CIE platform was operated in recirculation mode
96 for baseline establishment. Generally, the process of establishing baseline takes 4-6
97 hours before any contaminant experiments can be carried out. When operating in
98 single-pass contaminant mode, the contaminant is injected into the pipe connecting
99 the tank and sensors via another peristaltic pump. It is injected at a rate of 2-20 mL
100 per minute depending on concentration requirement. For more information about the
101 CIE platform and the injection experiment, the readers could refer to Liu et al.¹².

102 (Figure 1)

103 **Contaminants investigated**

104 Specific quantities of various contaminants were injected into the system simulator.
105 The contaminants investigated were determined according to statistical reports on

106 water pollution incidents in urban water supply systems in China over the past 20
107 years. Three groups of the most common six pollutants were selected: atrazine,
108 glyphosate, cadmium nitrate, nickel nitrate, sodium fluoride and sodium nitrate. They
109 were also selected based on China's national standards regarding source water quality
110 GB3838-2002 and drinking water quality GB5749-2006. The concentration ranges of
111 tested contaminants are provided in the supplementary material (Table T1) and were
112 decided using the concentration limit given in the above national standards.

113

114 **Classification method**

115 Clustering or cluster analysis is the process of grouping a set of objects into classes of
116 similar objects. Objects in any one cluster share similar features. Although
117 definitions of similarity vary from one clustering model to another, in most of these
118 models the concept of similarity is based on distances, e.g., Euclidean distance and
119 cosine distance¹³⁻¹⁵.

120 (Figure 2)

121 In cluster analysis, similar objects are assumed to have close values. If the distance of
122 an *object* to a particular *class* is shorter than the distances to other classes, the *object*
123 is deemed as belonging to that *class* (Figure 2). In this way, cluster analysis can be
124 used to identify the type of contaminant. An *object* can be an *example* or *instance* of
125 the *class*. In this study, the term *instance* refers to the object in a pre-defined *class*,
126 while *example* refers to the object to be classified. Both *instances* and *examples* are

127 vectors consisting of *features*. The *features* are extracted and derived from the sensor
 128 responses for contaminants.

129

130 Figure 3 shows the responses to cadmium nitrate and atrazine at time $t1$ and $t2$ for 8
 131 types of sensors. If the sensor reading is taken as the *feature*, p^{t1} , p^{t2} , q^{t1} and q^{t2} are
 132 8-dimensional vectors. As shown in Figure 3, the graphs for p^{t1} and p^{t2} are clearly
 133 similar to each other, while the graph for q^{t1} is closer to the graph for q^{t2} . An essential
 134 task of this study is to quantify the similarity or dissimilarity between two vectors,
 135 which is then used for contaminant identification.

136 (Figure 3)

137 **Similarity measure**

138 There are several methods of measuring the similarity between two objects (i.e. two
 139 l -dimensional vectors). In this study, cosine similarity was adopted. Cosine similarity
 140 is a measure of similarity between two vectors of an inner product space that
 141 measures the cosine of the angle between them¹⁶⁻¹⁷. The cosine of two vectors can be
 142 derived by using the Euclidean dot product formula.

$$143 \quad p \cdot q = \|p\| \|q\| \cos \theta \quad (1)$$

144 Given two vectors of attributes, p and q , the cosine similarity, $\cos(\theta)$, is represented
 145 using

$$146 \quad \text{similarity}(p, q) = \cos(\theta) = \frac{p \cdot q}{\|p\| \|q\|} = \frac{\sum_{i=1}^n p_i q_i}{\sqrt{\sum_{i=1}^n p_i^2} \sqrt{\sum_{i=1}^n q_i^2}} \quad (2)$$

147 in which n is the dimension of vector p and q .

148

149 This function gives a similarity measure in the sense that the cosine value gets larger

150 as the two vectors become more parallel to each other in the l -dimensional space. Or,

151 in other words, as the two data segments become more similar, their cosine similarity

152 approaches 1.0 and their distance approaches 0.0. Therefore, cosine similarity can be

153 used as a distance metric in the following way:

$$154 \quad D(p, q) = 1 - \text{similarity}(p, q) \quad (3)$$

155 Since the cosine similarity reflects the magnitude of the angle between two vectors in

156 the l -dimensional space, it is a many-to-one function. Compared with the other

157 distance measures, like Euclidean distance, the cosine similarity ignores the

158 magnitude difference between the two vectors, i.e.

$$159 \quad \text{similarity}(Ap, q) = \frac{\sum_{i=1}^n Ap_i q_i}{\sqrt{\sum_{i=1}^n (Ap_i)^2} \sqrt{\sum_{i=1}^n q_i^2}} = \text{similarity}(p, q) \quad (4)$$

160 Therefore, when the cosine distance is used for contaminant identification, the

161 variation range of sensor data need not be predetermined.

162

163 **Contaminant Classification**

164 The distance from a point p to a *class* c is given by:

$$165 \quad D(p, C) = 1 - \text{similarity}(p, \mu_c) \quad (5)$$

166 in which, $D(p, C)$ is the distance from a point to a *class* and μ_c is the mean of all

167 *instances in class C.*

168

169 The type of contaminant is identified by comparing the distances from examples to

170 *classes*. Assuming there are n types of contaminants, C_1, C_2, \dots, C_n , (or n *classes*),

171 each *class* contains many vectors (i.e. *instance of class*). For any example p to be

172 identified, if there exists

173
$$D(p, C_i) < D(p, C_j), j = 1, 2, \dots, n, i \neq j \quad (6)$$

174 then it is deemed that $p \in C_i$.

175

176 **Evaluation of classification performance**

177 The performance of the classification method is evaluated using the correct

178 classification rate (*CCR*). *CCR* can be calculated by

179
$$CCR = \frac{CC}{CC+IC} \times 100\% \quad (7)$$

180 where *CC* refers to the correct classification of a contaminant, *IC* is the incorrect

181 classification of a contaminant as another type of contaminant. A greater *CCR* means

182 the method is more capable of contamination identification.

183

184 **Robustness of the proposed method**

185 The proposed method relies on the readings of online water quality sensors. Inevitably,

186 fluctuations exist in online readings, which might come from equipment noise or

187 ambient variation. An important issue for a contaminant classification method is how

188 robust it is when dealing with fluctuations in readings. To evaluate the robustness of
 189 the proposed method, artificial uncertainties were added to the raw readings. It is
 190 assumed that the uncertainty obeys Gaussian distributions. The uncertainty
 191 quantification is achieved through a sampling-based method, Latin hypercube
 192 sampling (LHS) technique. In LHS¹⁸, values of stochastic tested vectors are generated
 193 in a random, yet constrained way. First, the values of variables in the original tested
 194 vectors are taken as means (i.e. raw readings of sensors) and the standard deviation is
 195 equal to 1% of the mean value (i.e., coefficient of variation $C_v=0.01$, for example).
 196 The range of each vector variable can be calculated using a Gaussian distribution
 197 equation, which is then divided into N_s non-overlapping intervals on the basis of
 198 equal probability. After that, a single random value is selected from each interval. This
 199 process is repeated for all variables in a *feature* vector. Once that is done, the N_s
 200 values obtained for the first vector variable are paired in a random manner with N_s
 201 values obtained for the second vector variable and so on. N_s *feature* vectors are
 202 generated from the original *feature* vector. By repeating the same process, *feature*
 203 vectors and the associated uncertainty can be obtained for all time steps. The *CCRs*
 204 for *feature* vectors with uncertainty can then be obtained. Finally, the robustness is
 205 evaluated using equation 8.

$$206 \quad \text{robustness} = \frac{CCR_{\text{confidence}}}{CCR_o} \quad (8)$$

207 in which, $CCR_{\text{confidence}}$ is the 95% confidence limit of the *CCRs* with
 208 uncertainty and CCR_o is the original *CCR*. For example, if the 95% confidence limit

209 of the *CCRs* with uncertainty is 0.8 and the original *CCR* is 1, then the robustness
210 value is $0.8/1=0.8$. A higher robustness value means that the method is more robust.

211

212 **Experiments and Results**

213 **Formation of *classes* of contaminants**

214 In this study, *features* were extracted to facilitate the quantitative evaluation of
215 similarity or dissimilarity between different types of contaminants. For all sensors in
216 this study, the sensor responses obtained at each time step were adopted to form a
217 *feature* vector (8 dimensions). For instance, the vector at the 1st minute for glyphosate
218 was [1.32 7.06 757.67 10.76 276.96 3.16 9.42 0.08] with the vector sequence being
219 turbidity, pH, conductivity, temperature, ORP, nitrate, UV and phosphate. Figure 4
220 shows the corresponding *feature* vectors at different concentrations. As shown in
221 Figure 4, the extracted *features* share some similarity, but dissimilarity also exists. By
222 extracting such data from all time steps, the *class* for glyphosate was established. The
223 same procedure was repeated for the other contaminants examined in this study and a
224 library containing 6 *classes* was obtained.

225 (Figure 4)

226 **Contaminant classification**

227 Glyphosate and cadmium nitrate were chosen to demonstrate the performance of the
228 contaminant classification method. The concentrations for glyphosate were 1.4mg/l,
229 2.8mg/l, 7.0mg/l and 14.0mg/l. For cadmium nitrate, the concentrations were

230 0.004mg/l, 0.008mg/l, 0.016mg/l and 0.032mg/l. A new group of contaminant
231 injection experiments were conducted to produce data for contaminant classification.
232 The raw experimental data contained sensor responses for baseline and presence of
233 contaminant at 4 concentrations. As reported in Liu et al.¹², the contamination events
234 were detected 1 minute after introduction of contaminants. The sensor response data
235 after detection were separated from the raw data and used in the classification. They
236 were treated using the procedure above to obtain *example feature* vectors. In total,
237 there were 110 glyphosate and 200 cadmium nitrate *example* data to be tested.

238 (Figure 5)

239 For glyphosate, the cosine distances to all *classes* for each *example* (or 1 minute time
240 step) was calculated using equation 4 and are shown in Figure 5. The green dots show
241 the distance between *examples* and the glyphosate *class*. For all time steps (from 1 to
242 110), it can be noted that, although the concentration varies, the cosine distances from
243 the examples to glyphosate *class* are rather stable and small. They are mostly in the
244 range of [0 0.02]. The distances to the other *classes* are much greater. For example,
245 the distance to chromium nitrate is around 0.16 (the blue circles in Figure 5). This is
246 shown in Table 1, along with mean and standard deviation values of the distances.
247 The proposed method classified the type of contaminant by comparing the cosine
248 distance. The one with the closest distance is deemed to be the correct class. Table 1
249 and Figure 5 reveal that the examples are closer to the glyphosate *class*. On the basis
250 of equation 6, the *feature* vectors of the example are more similar to the ones for

251 glyphosate. Therefore, it can be concluded that the contaminant is glyphosate. Using
252 equation 7, the *CCR* of the classification was calculated to be 0.918, which suggests
253 that the tested contaminant is correctly classified in 91.8% of situations in this study.

254 (Figure 6)

255 For cadmium nitrate, Figure 6 shows the cosine distances to different classes, in
256 which the red dots indicate the distances to the cadmium nitrate *class*. For all time
257 steps, the distances to the cadmium nitrate class were in the range of 0.01 to 0.04 with
258 the mean of 0.0277 (Table 1). It is obvious that the distances to cadmium nitrate are
259 smaller than the ones to other classes in most cases in this study. The *CCR* was
260 calculated to be 0.975.

261 (Table 1)

262 In terms of the time needed for classification, once a contamination event is detected
263 by an EWS, the contaminant classification module will be activated. Theoretically, the
264 type of contaminant can be classified within 1 minute (i.e. the sensor reporting step).
265 However, in practice, the time might be a bit longer since the sensor responses to
266 presence of contaminant might sometimes need to stabilize. As shown in Figure 5, the
267 contaminant was classified correctly to be glyphosate 1 minute after the
268 contamination event alarm. This means that the distance to the correct class was the
269 smallest from the 1st minute onwards. For the case of cadmium nitrate, the proposed
270 method can classify correctly 6 minutes after activation. In the first 5 minutes, the
271 tested examples were incorrectly classified. The key strength of the proposed method

272 is that it classifies the type of contaminant in a real time manner. Compared to
273 laboratory-based methods, classification in 6 minutes with no significant compromise
274 of *CCR* is an advantage.

275

276 **Discussion**

277 **Comparison to Euclidean distance based method**

278 In previous studies, Liu et al.¹² reported that the magnitudes of sensor responses vary
279 with the concentration of contaminant (or see Figure F1, F2, F3, F4 and F5 in
280 supplement documents). This is typically obvious for pH, nitrate, phosphate and ORP.
281 For example, the pH and ORP values for the glyphosate concentration of 1.4, 2.8, 7.0,
282 14.0mg/l are 6.89, 6.71, 6.41, 6.10 and 277.66, 283.33, 291.67, 299.29 mV
283 respectively. The aim of this study is to establish a method to classify the type of
284 contaminant by evaluating the similarity between examples and classes. The
285 classification method should be independent of or less related to the concentration of
286 the contaminants since this is not known in advance in a real event. In other words,
287 the distance evaluation method should not be too dependent of contaminant
288 concentration. If the distance evaluation is closely related to magnitude of sensor
289 response, the classification method might fail to differentiate events caused by the
290 same type of contaminant with different concentrations.

291

292 There are several types of evaluation methods for the distance of vectors. The most

293 commonly used one is the Euclidean distance, which is the "ordinary" distance
294 between two points¹⁹⁻²⁰. The Euclidean distance between points p and q is the length
295 of the line segment connecting them, which can be calculated using

$$296 \quad E(q, p) = \|q - p\| = \sqrt{(q - p) \cdot (q - p)} \quad (9)$$

297 in which $E(q, p)$ is the Euclidean distance between points p and q .

298

299 Figure 7 schematically shows the Euclidean distances and cosine distance between
300 points p_1 , p_2 , q_1 and q_2 . Points p_1 and p_2 are the sensor response vectors of
301 contaminant 1 at concentrations 1 and 2. Points q_1 and q_2 are the vectors for
302 contaminant 2 at concentrations 1 and 2. As shown in Figure 7, the Euclidean distance
303 between p_1 and p_2 is $\|p_1 - p_2\|$ and the cosine distance is 0. For p_2 and q_2 , the
304 Euclidean distance is $\|p_2 - q_2\|$ and the cosine distance is $1 - \cos(\theta)$. Therefore, by
305 using the cosine distance method, p_1 and p_2 (also q_1 and p_2) can be classified to the
306 correct class. However, if the Euclidean distance were used, it might group p_2 and q_2
307 into the same class because $\|p_2 - q_2\| < \|p_1 - p_2\|$. To further explain this, the
308 vectors associated with glyphosate and cadmium nitrate at different concentrations
309 were taken as examples to calculate the Euclidean and cosine distances.

310 (Figure 7)

311

312 Table 2 shows the cosine and Euclidean distances between points a , b , c and d , in
313 which a is the vector of sensor responses to glyphosate at concentration of 1.4mg/l, b

314 is the vector for glyphosate at 14.0 mg/l, c is the vector for cadmium nitrate at
315 concentration of 0.008 mg/l and d is the vector for cadmium nitrate at concentration
316 of 0.032 mg/l. In Table 2, the numbers above the diagonal are cosine distances and the
317 ones below are Euclidean distances. As shown in Table 2, the cosine distances for
318 points from the same class are smaller than those for points from different classes. For
319 example, $D(a,b)=0.0027$, while $D(a,c)=0.1091$. This explains the correctness of the
320 assumption in Figure 2, i.e. similar objects have shorter distance.

321

322 It is also observed that the cosine distance is not 'sensitive' to the magnitude the
323 vector (in other words, the concentration of the contaminants). As shown in Figure 4,
324 the magnitude of sensor response vectors at 1.4 mg/l and 14mg/l is obviously
325 different. However, their cosine distances to other points are close. For example,
326 $D(a,c)=0.1091$, $D(b,c)=0.1081$. Euclidean distance, on the other hand, is related to the
327 magnitude of the vector. For example, $E(a,c)=92.3888$, while $E(b,c)= 158.4424$.
328 Furthermore, the case may arise where the Euclidean distance between points from
329 the same class might be greater than that between points from two different classes.
330 This is shown in Table 2. For example, $E(a,b) = 95.5981$ and $E(a,c) = 92.3888$. In this
331 case, incorrect classification would occur if Euclidean distance were used for
332 contaminant classification. Point c would be wrongly classified as being in the same
333 class as point a if Euclidean distance was adopted. Therefore, it was concluded that
334 cosine distance is more suitable than Euclidean distance for classifying the type of

335 contaminant since the evaluation for similarity is more related to the contaminant's
336 characteristics rather the magnitude of sensor responses.

337 (Table 2)

338 **Robustness**

339 The level of uncertainty is given by the value of C_v . In this study, four values of C_v
340 (0.005, 0.01, 0.02 and 0.03) were used. The value of N_s is determined according to
341 the literature. For a given C_v , by setting $N_s=2000$, 220000 *feature* vectors with
342 uncertainty were finally generated for glyphosate. These *feature* vectors were divided
343 into 2000 groups. Each contains 110 *feature* vectors. By feeding the 2000 groups of
344 *feature* vectors into the proposed contaminant classification method, the *CCRs* for
345 every group were obtained. The histograms of these *CCRs* are displayed in Figure 8,
346 which shows that the proposed method has robustness of over 0.82 for uncertainty
347 $C_v=0.005$, $C_v=0.01$ and $C_v=0.02$. This suggests that the performance of the proposed
348 contaminant classification method is steady and reliable and can cope well with the
349 uncertainty from the online sensors. For the case of $C_v=0.03$, the performance of the
350 method is less satisfactory. The *CCR* for this case is 0.75, which is much lower than
351 the original *CCR* (0.92). It should be noted that the uncertainty examined in this study
352 is assumed to be from equipment noise or ambient variation. A change of sensor
353 reading due to sudden sensor failure or presence of contaminant is not treated as noise,
354 but instead as an event, which normally means a 1-20% change of sensor reading.
355 Therefore, it is deemed that the uncertainty levels adopted in this study are significant

356 enough.

357

358 It is worth noting that previous studies about online sensors in water supply systems
359 generally assumed perfect sensors, which means that sensors worked in good
360 condition²¹⁻²³. Although this assumption has allowed researchers to make significant
361 progress in early warning system design, sensor failures can significantly impact the
362 reliability of an early warning system design. It is commonly known that sensor
363 readings do contain uncertainties as they are easily affected by ambient variation and
364 equipment condition. From an implementation perspective, it is essential that an early
365 warning system using online sensor data is robust enough and can cope with
366 uncertainty from sensors. The analysis here shows that the proposed method has good
367 capacity of tolerate uncertainty in sensor readings.

368 (Figure 8)

369 (Table 3)

370 **Future studies**

371 This study proposed a concentration-independent contaminant classification method
372 based on conventional water quality sensors. The basis of this method is that points in
373 one class stay close and are separate from other classes. In spite of great improvement
374 in recent years, readings from online sensors are still affected by noise and ambient
375 variation. For a method based on online sensor readings, it is important to understand
376 the impact of uncertainty from sensor readings on the model output. Although this

377 study demonstrated the robustness of the proposed method in the event of sensor
378 uncertainty or ambient variation through an initial uncertainty analysis, a global
379 sensitivity analysis would be more helpful to understand the extent of uncertainty
380 from each sensor. This should be conducted in future study.

381

382 Meanwhile, since the proposed method classifies by comparing the distances to
383 predefined classes, incorrect classification error would occur if two (or more) classes
384 overlap each other. This study involved a limited number of contaminants and no
385 overlaps were noticed, but the possibility does exist. In a future study, this has to be
386 addressed. A possible solution to this is that the classification decision could be made
387 based on distances from more than one type of features. For example, if the features
388 using original sensor responses from two types of contaminants overlap, another type
389 of feature (e.g. the deviation between real readings and baseline) can be employed to
390 differentiate these two classes.

391

392 **Conclusion**

393 By using data from online water quality sensors, this study proposed a real time and
394 concentration-independent contaminant classification method. From the analysis, the
395 following conclusions were drawn.

- 396 1) The proposed method classifies the type of the contaminant by comparing their
397 cosine distances to predefined classes. Results from the analysis show that the

398 proposed method can identify glyphosate and cadmium nitrate 1 and 6 minutes
399 after detection with the *CCR* of 91.8% and 97.5%. Compared to laboratory-based
400 methods, classification in minutes without significant compromising the *CCR* is
401 an advantage.

402 2) Results show that the performance of the proposed method was not related to the
403 contaminant concentration. This implies that the proposed method is more
404 suitable than the Euclidean distance method for contaminant classification since
405 the concentration of contaminant is not known a priori.

406

407 **Acknowledgements**

408 This work is jointly supported by Tsinghua Independent Research Program
409 (2011Z01002) and Water Major Program (2012ZX07408-002).

410

411 **Reference**

412 1. USEPA, *Baseline threat information for vulnerability assessments of community*
413 *water systems*, Washington, DC, 2002.

414 2. USEPA, *Planning for and responding to drinking water contamination threats and*
415 *incidents*, Washington, DC, 2003.

416 3. M. V. Storey, B. van der Gaag and B. P. Burns, *Water Res*, 2011, 45, 741-747.

417 4. C. J. de Hoogh, A. J. Wagenvoort, F. Jonker, J. A. Van Leerdam and A. C.

418 Hogenboom, *Environ. Sci. Technol.*, 2006, 40, 2678-2685.

- 419 5. J. Jeon, Kim, J. H. Lee, B.C., S. D. Kim, *Sci. Total Environ.*, 2008, 389, 545-556.
- 420 6. R. K. Henderson, A. Baker, K. R. Murphy, A. KHambly, R. M. Stuetz and S. J.
- 421 Khan, *Water Res.*, 2009, 43 (4), 863-881.
- 422 7. P. R. Hawkins, S. Novic, P. Cox, Neilan, B. A. Burns, B. P. Shaw, G., W.
- 423 Wickramasinghe, Y. Peerapornpisal, W. Ruangyuttikarn, T. Itayama, T. Saitou, M.
- 424 Mizuochi and Y. Inamori, *J. Water Supply: Res. Technol.-AQUA*, 2005, 54, 509-518.
- 425 8. C. P. Marshall, S. Leuko, C. M. Coyle, M. R. Walter, B. P. Burns and B. A. Neilan,
- 426 *Astrobiology*, 2007, 7, 631-643.
- 427 9. D. Kroll, *Securing our water supply: Protecting a vulnerable resource*, Pennwell,
- 428 2006.
- 429 10. J. Y. Yang, R. C. Haught and J. A. Goodrich, *J. Environ. Manage.*, 2009, 90,
- 430 2494-2506.
- 431 11. N. Alikar and A. Ostfeld, *Water Res.*, 2014, 51, 234-245.
- 432 12. S. M. Liu, H. Che, K. Smith and L. Chen, *Environmental Science Processes &*
- 433 *Impacts*, 2014, 16, 2028-2038.
- 434 13. M. Tabacchi, C. Asensio, I. Pavon, M. Recuero, J. Mir, and M.C. Artal, *Applied*
- 435 *Acoustics*, 2013, 74, 1022-1032.
- 436 14. Z.P. Zhao, P. Li, and X.Z. Xu, *Applied Mathematics & Information Sciences*, 2013,
- 437 7, 1243-1250.
- 438 15. K.A. Nguyen, R.A. Stewart, and H. Zhang, *Environmental Modeling & Software*
- 439 47, 108-127.

- 440 16. A.A. Pascasio, *Linear Algebra Appl*, 2001, 325, 147-159.
- 441 17. M.W. Sohn, *Soc Networks*, 2001, 23, 141-165.
- 442 18. G. Manache, C.S. Melching, *J Water Res Plan Man*, 2004, 130, 232-242.
- 443 19. L. Liberti, C. Lavor, N. Maculan and A. Mucherino, *Siam Rev*, 2014, 56, 63-69.
- 444 20. C. Stoecker, S. Welter, J.H. Moltz, B. Lassen, J.M. Kuhnigk, S. Krass and H.O.
445 Peitgen, *Med Phys*, 2013, 40, 091912.
- 446 21. M. Weickgenannt, Z. Kapelan, M. Blokker and D.A. Savic, *J Water Res Plan Man*,
447 2010, 136, 629-636.
- 448 22. A. Kessler, A. Ostfeld, G. Sinai, *J Water Res Plan Man*, 1998, 124, 192-198.
- 449 23. A. Ostfeld, E. Salomons, *J. Water. Res. Pl.-Asce*, 2004, 130(5), 377-385.
- 450

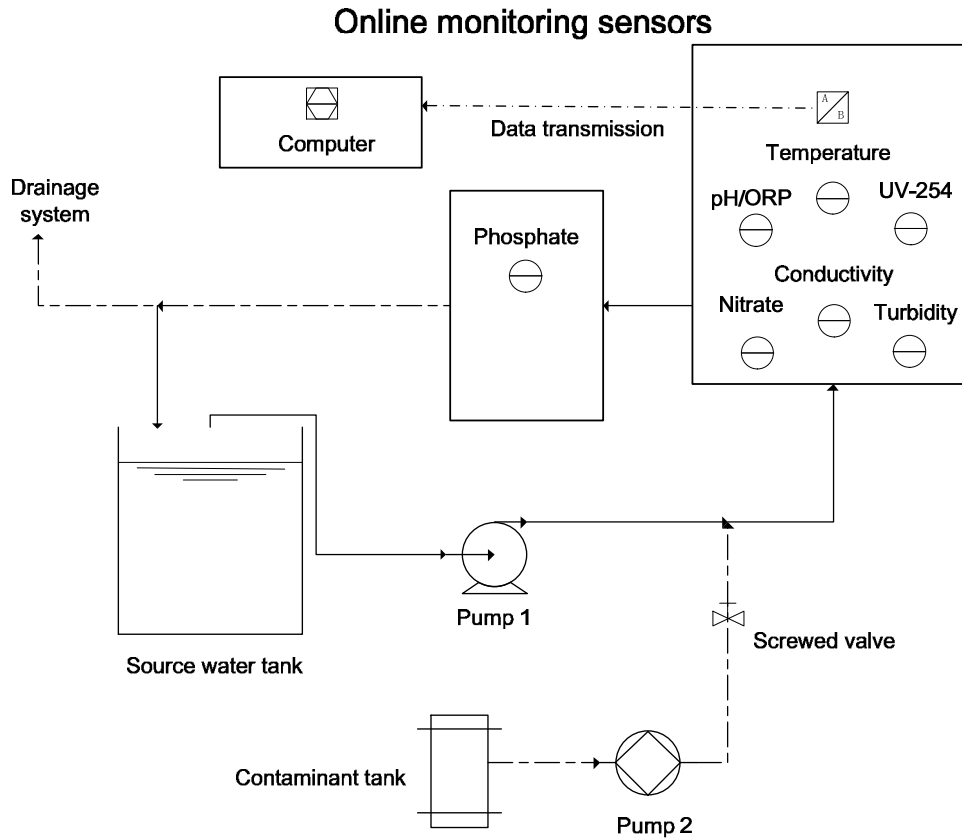


Figure 1 A process flow schematic of the pilot-scale system

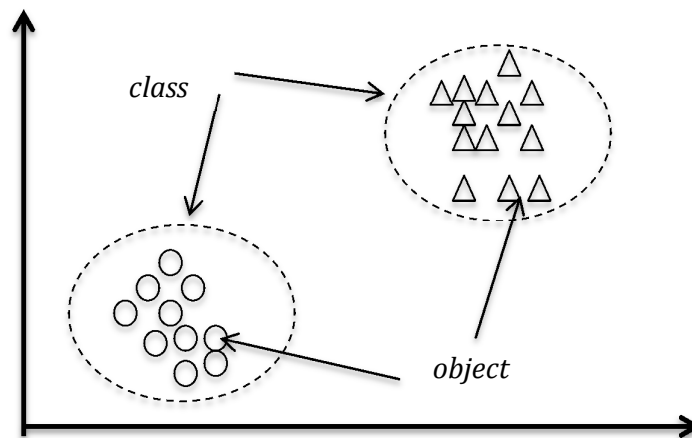


Figure 2 Schematic graphs of *class* and *instance*

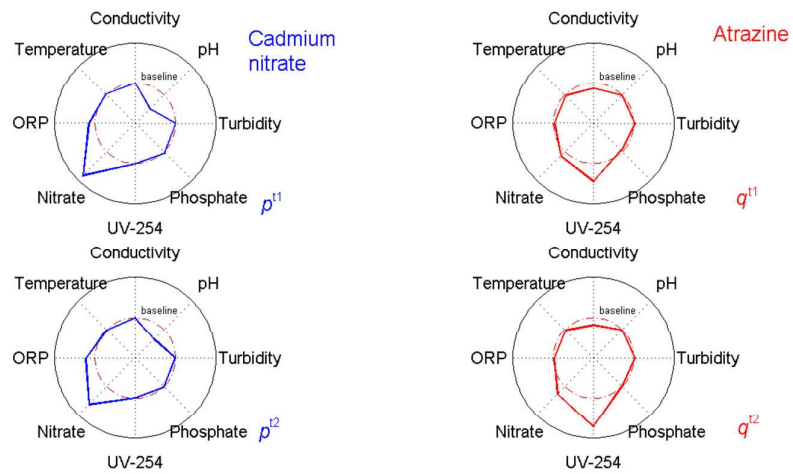


Figure 3 Four instances of *features* of cadmium nitrate and atrazine

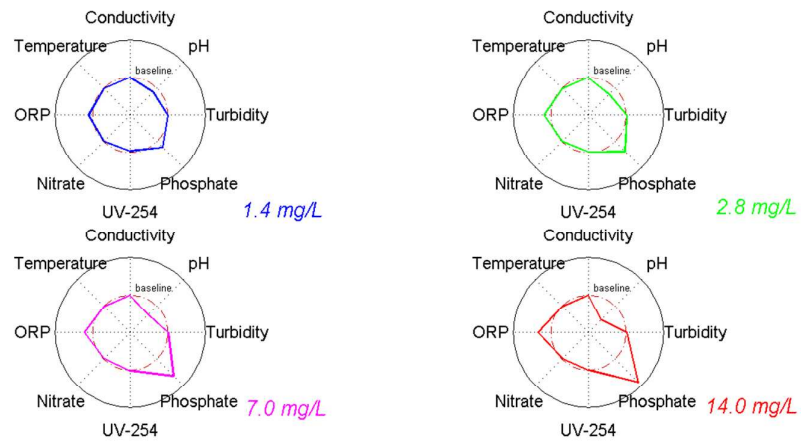


Figure 4 The demonstration of *feature* vectors at glyphosate 1.4, 2.8, 7, 14mg/l

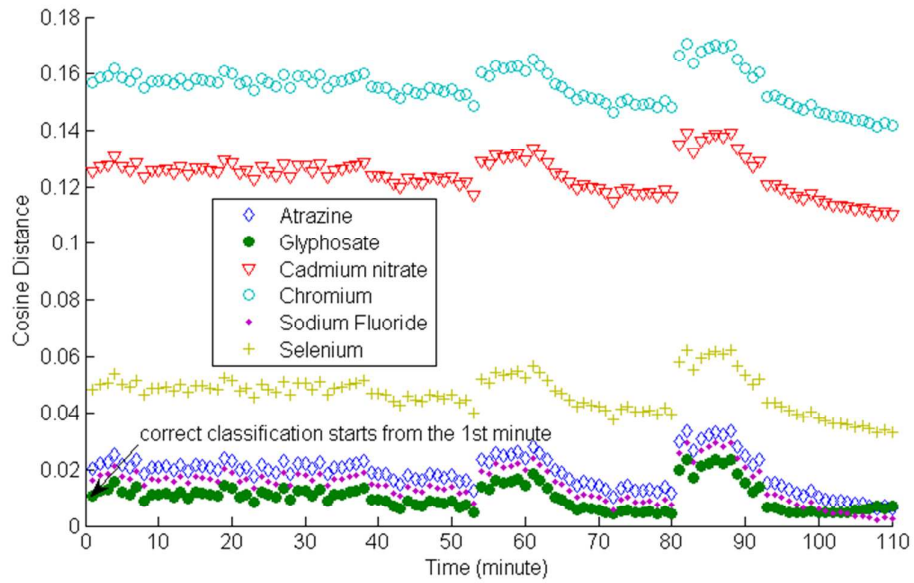


Figure 5 The cosine distance of glyphosate to 6 classes

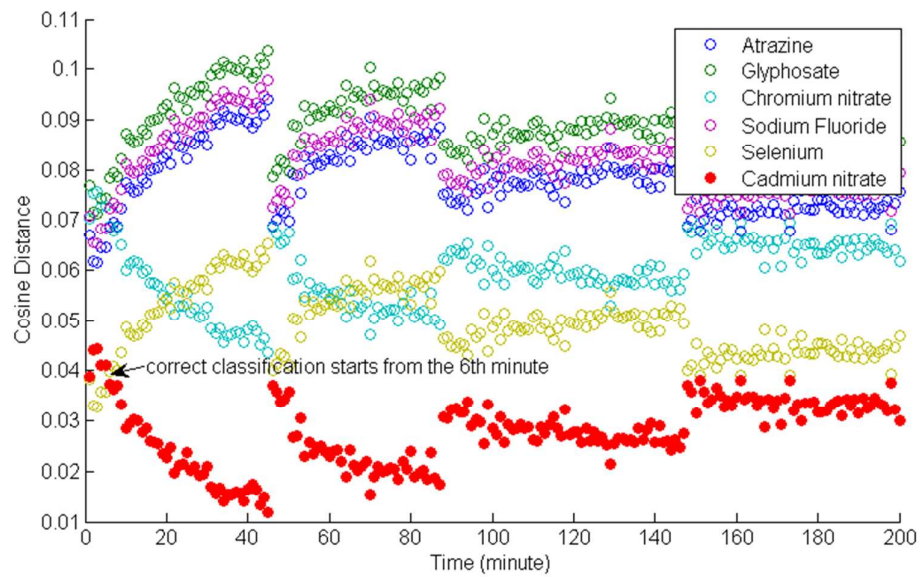


Figure 6 The cosine distance of cadmium nitrate to 6 *classes*

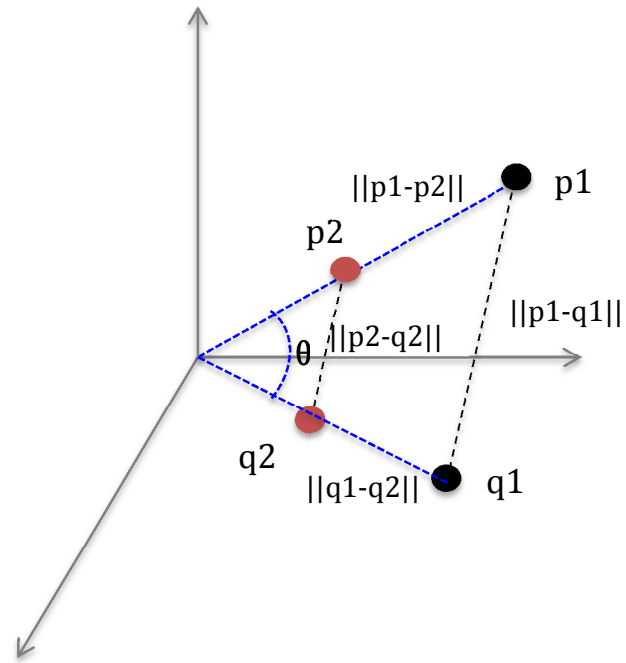


Figure 7 Schematic drawing of cosine and Euclidean distances

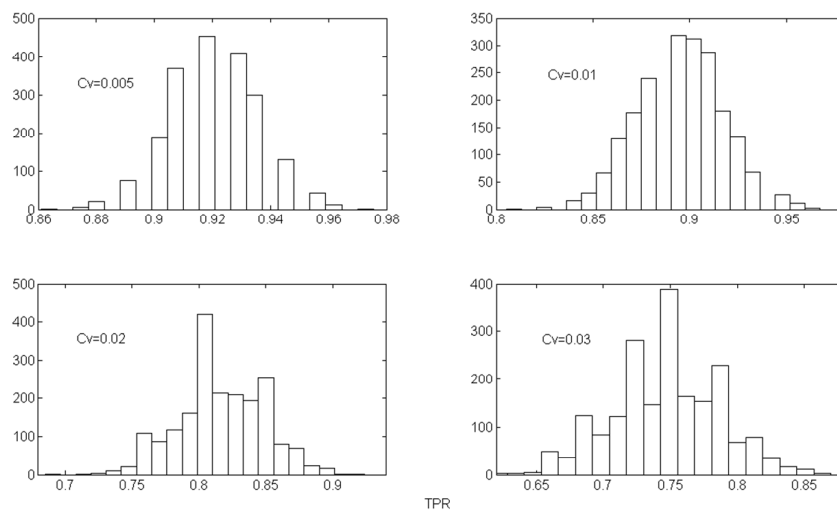


Figure 8 The histogram of *CCRs* with uncertainty and robustness

Table 1 Averaged cosine distances from examples to *classes*

| Tested example | | Cosine distances from examples to classes | | | | | |
|-----------------|---------------------------|---|---------------|-----------------|----------------|-----------------|----------------|
| | | Atrazine | Glyphosate | Cadmium nitrate | Nickel nitrate | Sodium fluoride | Sodium nitrate |
| Glyphosate | <i>mean</i> | 0.0190 | 0.0105 | 0.1240 | 0.1555 | 0.0148 | 0.0470 |
| | <i>Standard deviation</i> | 0.0063 | 0.0048 | 0.0065 | 0.0065 | 0.0063 | 0.0065 |
| Cadmium nitrate | <i>mean</i> | 0.0781 | 0.0880 | 0.0277 | 0.0593 | 0.0819 | 0.0494 |
| | <i>Standard deviation</i> | 0.0065 | 0.0066 | 0.0065 | 0.0065 | 0.0065 | 0.0066 |

Table 2 The cosine and Euclidean distances

| Euclidean \ Cosine | | Glyphosate | | Cadmium nitrate | |
|--------------------|-------------------------------|--------------------|-------------------|----------------------|----------------------|
| | | <i>a</i> - 1.4mg/l | <i>b</i> - 14mg/l | <i>c</i> - 0.008mg/l | <i>d</i> - 0.032mg/l |
| Glyphosate | <i>a</i> - 1.4mg/l | 0 | 0.0027 | 0.1091 | 0.1391 |
| | <i>b</i> - 14mg/l | 95.5981 | 0 | 0.1081 | 0.1381 |
| Cadmium nitrate | <i>c</i> - 0.008mg/l | 92.3888 | 158.4424 | 0 | 0.0302 |
| | <i>d</i> - 0.032mg/l | 113.1857 | 166.8406 | 25.9895 | 0 |

Note: The numbers above the diagonal are cosine distances, while the ones below are

Euclidean distances.

Table 3 Statistics of *CCR* under uncertainty (original *CCR*: 0.92)

| <i>CCR</i> | $C_v = 0.005$ | $C_v = 0.01$ | $C_v = 0.02$ | $C_v = 0.03$ |
|--------------------|---------------|--------------|--------------|--------------|
| Mean | 0.92 | 0.90 | 0.82 | 0.75 |
| Standard deviation | 0.01 | 0.02 | 0.03 | 0.04 |
| Robustness | 0.97 | 0.93 | 0.82 | 0.73 |