

1
2
3 1 **Classification of groundwater chemistry in Shimabara, using self -organizing maps**
4
5
6 2 (Short title: Classification of groundwater chemistry in Shimabara, using SOMs)
7
8
9 3 **Kei Nakagawa, Hiroki Amano, Akira Kawamura and Ronny Berndtsson**
10
11
12 4
13
14 5 Kei Nakagawa (corresponding author) and Hiroki Amano
15
16
17 6 Graduate School of Fisheries and Environmental Sciences, Nagasaki University, 1-14 Bunkyo-machi,
18
19
20 7 Nagasaki 852-8521, Japan
21
22
23 8 E-mail: kei-naka@nagasaki-u.ac.jp
24
25
26 9
27
28
29 10 Akira Kawamura
30
31
32 11 Graduate School of Urban Environmental Sciences, Tokyo Metropolitan University, 1-1 Minami-Oshawa,
33
34
35 12 Hachioji, Tokyo 192-0397, Japan
36
37
38 13
39
40
41 14 Ronny Berndtsson
42
43
44 15 Division of Water Resources Engineering & Center for Middle Eastern Studies, Lund University, Box 118
45
46
47 16 SE-221 00 Lund, Sweden
48
49
50 17
51
52
53 18
54
55
56 19
57
58 20
59
60
61
62
63
64
65

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

21 **Abstract**

22 Shimabara City in Nagasaki Prefecture, Japan, is located on a volcanic peninsula that has
23 abundant groundwater. Almost all public water supply use groundwater in this region. For this reason,
24 understanding groundwater characteristics is a pre-requisite for proper water supply management. Thus,
25 we investigated the groundwater chemistry characteristics in Shimabara by use of self-organizing maps
26 (SOM). The input to SOM was concentrations of eight major groundwater chemical components, namely
27 Cl⁻, NO₃⁻, SO₄²⁻, HCO₃⁻, Na⁺, K⁺, Mg²⁺, and Ca²⁺ collected at 36 sampling locations. The locations
28 constituted private and public water supply wells, springs, and a river sampled from April 2012 to May
29 2015. Results showed that depending on chemistry, surface and groundwater could be classified into five
30 main clusters displaying unique patterns. Further, the five clusters could be divided into two major water
31 types namely, nitrate- and non-polluted water. According to stiff and Piper trilinear diagrams the nitrate
32 polluted water represented Ca-(SO₄²⁻+NO₃) (calcium sulfate nitrate) type, while non-polluted water was
33 classified as Ca-HCO₃ (calcium bicarbonate) type. This indicates that recharging rain water in the upstream
34 areas is polluted by agricultural activities in the mid-slope areas of Shimabara.

35
36 **Keywords**

37 cluster analysis, groundwater, major ions, self-organizing map, water chemistry

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

41 **Introduction**

42 Groundwater is used for various purposes such as water supply, agriculture, and industry. During
43 recent decades, groundwater has been polluted by increasing fertilizer application to meet the demand of
44 food supply due to the population growth. Monitoring and protection of groundwater are essential to meet
45 the demand for safe groundwater. To understand effects of hydrogeological processes and anthropogenic
46 activities on regional groundwater it is important to study the chemical characteristics. The
47 hydrogeochemistry of groundwater is influenced by many factors such as climate, mineralogy of aquifers,
48 chemical composition of rainfall and surface water, topography, and anthropogenic activities. Thus,
49 hydrogeochemical interpretation of groundwater quality from representative water samples can provide
50 useful information on geochemical processes, hydrodynamics, origin, and interaction of the groundwater
51 with aquifer materials.

52 Shimabara City is known as a region that to a great extent relies on groundwater for public water
53 supply (Committee on nitrate reduction in Shimabara Peninsula 2011). However, the Shimabara
54 groundwater has been increasingly polluted by nitrate since 1988. We analyzed the present situation of
55 groundwater pollution by nitrate in Shimabara and showed that agricultural activities are the main polluter
56 of the groundwater (Nakagawa *et al.* 2016). To better understand characteristics of the water chemistry,
57 multivariate analysis such as principal component analysis (PCA), which can reduce data dimensionality
58 and extract synthetic indexes with minimum information loss, is often used (e.g., Aiuppa *et al.* 2003;
59 Cloutier *et al.* 2008; Banoeng-Yakubo *et al.* 2009; Sonkamble *et al.* 2012; Nadiri *et al.* 2013; Omonona *et*
60 *al.* 2014; Singaraja *et al.* 2014; Ghesquière *et al.* 2015; Marghade *et al.* 2015; Matiatos 2016). Using

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

61 groundwater chemistry, we classified Shimabara water by use of principal component and cluster analysis
62 (Nakagawa *et al.* 2016). The results showed that groundwater could be classified into 4 clusters where one
63 cluster expressed nitrate pollution and the other clusters showed ion dissolution from the aquifer matrix.
64 However, it is sometimes difficult to decipher PCA results due to bias resulting from complexity and
65 nonlinearity of large data (Choi *et al.* 2014). Recently, multivariate analysis using Self-Organized Maps
66 (SOM) has been applied to various research fields such as ecology (Céréghino *et al.* 2001; Bedoya *et al.*
67 2009), geomorphology (Hentati *et al.* 2010), hydrology (Kalteh & Berndtsson 2007), meteorology
68 (Nishiyama *et al.* 2007), and wastewater treatment (García & González 2004). SOM has also been used to
69 classify water chemistry of rivers and groundwater (Hong & Rosen 2001; Jin *et al.* 2011; Choi *et al.* 2014;
70 Nguyen *et al.* 2015). Thus, SOM is a powerful and effective tool for detection and interpretation of spatially
71 varying phenomena. Especially, SOM has a better ability to handle the nonlinearities, noisy or irregular
72 data, and multivariate data without mechanistic understanding of the system. SOM is also easily and quickly
73 updated when adding new data (Hong & Rosen 2001, Kalteh et al. 2008). The similarity of extracted pattern
74 classification can be visually compared using color gradients (Jin *et al.* 2011).

75 In the previous study (Nakagawa *et al.* 2016), we used field observed data from August 2011 to
76 November 2013. We continued to collect data and available data were extended to May 2015. Therefore,
77 in this study, we confirmed our previous results by using a more informative method, SOM, together with
78 an extended data base. Using SOM, visual representation of groundwater characteristics is easy and more
79 detailed clustering with better analyses results are possible as compared to conventional PCA. To improve
80 the understanding of groundwater characteristics in Shimabara we applied SOM combined with hierarchical

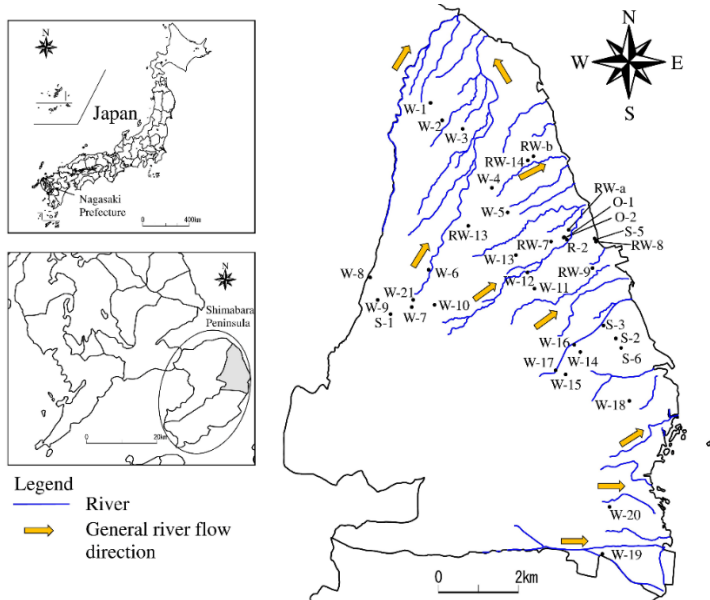
1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

81 cluster analysis using water chemistry as input. According to the results obtained by SOM analysis, we
82 discuss spatial trends of groundwater characteristics in Shimabara and the practical application of SOM for
83 future water use.

84

85 **Study area and data used**

86 Figure 1 shows the study area and sampling locations in Shimabara, Nagasaki Prefecture, Japan.
87 Shimabara has an area of 82.8 km² and is located in the northeastern part of Shimabara Peninsula. In the
88 center of the peninsula the active volcano Unzen (Mt. Fugendake) is located. The geology of Shimabara
89 area is thus formed by volcanic deposits composed of dacite, andecite, volcanic ash, and lapilli. Average
90 annual precipitation is about 2100 mm (1967-2013). The mean annual temperature is 16.9 °C, and average
91 monthly temperature ranges from 4.2 (January) to 29.0 °C (in August; Japan Meteorological Agency 2015).



92

93 Figure 1 Study area and sampling locations in Shimabara, Nagasaki Prefecture, Japan (RW: Residential
94 well, W: Public water supply well, O: Observation well, S: Spring, and R: River)

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

95 Figure 2 shows altitude and land use in Shimabara. According to the figure the land use can
96 generally be divided into forest, agriculture, and urban areas. Areas above an altitude of 200 m are generally
97 occupied by forest. According to the estimated regional groundwater flow, the forest areas, which compose
98 36.5 % of Shimabara, may be recognized as groundwater recharge zones. Upland and paddy fields are
99 concentrated to the northern parts of the area occupying 23.6 % and 7.5 % of Shimabara, respectively.
100 Buildings are usually located at altitudes below 100 m along the coast and represent 14.9 % of Shimabara.
101 Other land use is 17.5 %.

102 In total 353 water samples were collected from April 2012 to May 2015. Sampling was
103 performed at 7 resident wells (RW), 21 public water supply wells (W), 2 observation wells (O), 5 springs
104 (S), and 1 river (R) (Fig. 1). To ensure spatially representative groundwater conditions, sampling sites
105 covering the whole area of Shimabara except for forest and other land use (Figs. 1 and 2) were used.
106 Sampling was done four times annually with 2-4 months interval to ensure temporally varying groundwater
107 conditions. Sampling at specific locations (RW-14, b, W-21, O-2, S-2, 3, 5, and R-2) was done with less
108 frequency. The hydrogeochemical data used in this study consist of major dissolved ion concentration for
109 Cl^- , NO_3^- , SO_4^{2-} , HCO_3^- , Na^+ , K^+ , Mg^{2+} , and Ca^{2+} . Mean and standard deviation of 36 sampling sites using
110 averaged temporal ion concentrations for each sampling sites are summarized in Table 1. It is necessary to
111 normalize the data prior to application of SOM to ensure that all parameters are given the same importance.
112 SOM results are highly sensitive to data pre-processing method due to that the Euclidean distance between
113 input data is used (e.g., Jin *et al.* 2011). To solve this problem, the range between minimum and maximum
114 ion concentrations was standardized into [0, 1] (Nishiyama *et al.* 2007, Jin *et al.* 2011) as preprocessing in

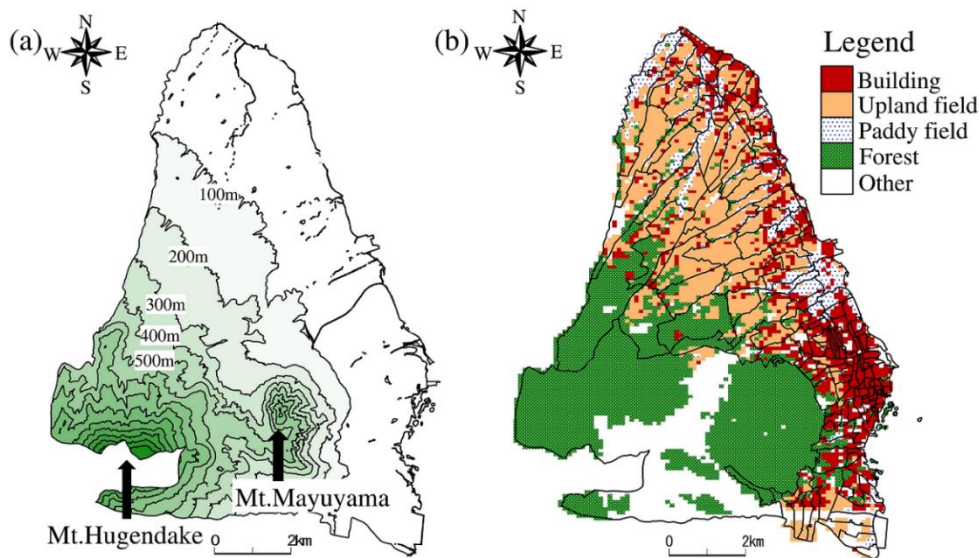


Figure 2 Altitude and land use map of Shimabara; (a) Altitude and (b) Land use

this study.

Methodology

The SOM is a modified artificial neural network characterized by unsupervised training that can project high-dimensional information onto a low-dimensional array (e.g., Vesanto *et al.* 2000). Many researchers have chosen a two-dimensional array (e.g., Jiang *et al.* 2014). The result is a readily understandable and visual pattern classification. The objective here of the SOM application was to obtain physically explainable reference vectors using input vectors. Thus, the input vectors were composed of in total 353 hydrogeochemical data points (approximately quarterly sampling at the 36 sampling locations) with 8 variables (major dissolved ion concentrations; Cl^- , NO_3^- , SO_4^{2-} , HCO_3^- , Na^+ , K^+ , Mg^{2+} , and Ca^{2+}). Reference vectors were obtained after iterative updates through a training phase that was composed by three

1
2
3 129 Table 1 Mean and standard deviations of 36 sampling sites using averaged temporal ion concentrations
4
5
6 130 for each sampling site used in the SOM
7

Major ion (mg L ⁻¹)	Mean	SD
Cl ⁻	12.4	1.4
NO ₃ ⁻	38.4	5.0
SO ₄ ²⁻	21.9	3.2
HCO ₃ ⁻	55.7	6.6
Na ⁺	12.1	2.4
K ⁺	6.4	1.2
Mg ²⁺	8.7	1.1
Ca ²⁺	22.4	2.9

8
9
10
11
12
13
14
15
16
17
18
19
20
21
22 131
23
24
25 132 main procedures: competition between nodes, selection of a winner node, and updating of the reference
26
27
28 133 vectors (e.g., Vesanto *et al.* 2000). Selection of proper initialization and data transformation methods is
29
30
31 134 important factors when designing a relevant SOM methodology. In SOM applications, in general a larger
32
33
34 135 map size gives a higher resolution for pattern recognition. The optimum number of SOM nodes is
35
36
37 136 determined applying the heuristic rule $m = 5\sqrt{n}$, where m denotes the number of SOM nodes and n
38
39
40 137 represents the number of input data (García & González 2004; Hentati *et al.* 2010; Jin *et al.* 2011). Herein,
41
42
43 138 this heuristic rule was used to determine the total number of nodes in the SOM. The ratio of the number of
44
45
46 139 rows and columns is determined by the square root of the ratio between the two largest eigenvalues of the
47
48
49 140 correlation matrix of input data. The eigenvalues are obtained from principal component analysis. In a
50
51
52 141 previous study fusing the sampled data from August 2011 to November 2013, two principal components
53
54 142 (Factor 1 and Factor 2) explained 86.5 % of the total variance (Nakagawa *et al.* 2016).
55

56
57 143 After organizing the SOM structure with the above rule, a linear initialization technique made
58
59
60
61
62
63
64
65

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

144 each node set with a reference vector. A linear initialization technique increases the speed of training phase
145 and proper abstracting pattern for limited data (Jeong *et al.* 2010). Further, when only limited data are
146 available, the linear initialization is more suitable for the pattern classification as compared to random
147 initialization because of small data sets and boundary effects (Nguyen *et al.* 2015). The linear initialization
148 used eigenvalues and eigenvectors of input data to set initial reference vectors on the structured SOM. This
149 means that the initial reference vectors already include prior information about the input data, resulting in
150 a quicker and more efficient training phase (Vesanto *et al.* 2000). In this study, each reference vector was
151 updated through the SOM training process using a batch mode with neighborhood function taking a
152 Gaussian form. Although, some issues on the implementation of the batch SOM are discussed at some detail
153 in Jiang *et al.* (2014), the results of the SOM analysis supported previous clustering results (Nakagawa *et*
154 *al.* 2016; shown below). The reference vectors obtained at the end of the training process were fine-tuned
155 using cluster analysis.

156 There are various clustering algorithms available in the literature (e.g., García & González 2004;
157 Jin *et al.* 2011). In this study, partitioned algorithms and hierarchical algorithms, which are k-means and
158 Ward's algorithms respectively, were applied for appropriate clustering of reference vectors. For partitioned
159 clustering methods, the k-means algorithm is most frequently used for SOM (e.g., Jin *et al.* 2011). The
160 Davies-Bouldin Index (DBI) applying k-means algorithm determines the optimal number of clusters
161 (García & González 2004; Jin *et al.* 2011). The DBI values, based on similarity within a cluster and
162 dissimilarity between clusters, were calculated from a minimum of two clusters to the total number of nodes.
163 Therefore, the smaller DBI value appears as the dissimilarity to each cluster becomes larger. In other words,

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

164 a minimum DBI represents the optimal number of clusters for the trained SOM. The Ward's linkage method,
165 which is the one of the hierarchical techniques, is the most commonly used clustering (Faggiano *et al.* 2010;
166 Hentai *et al.* 2010; Jin *et al.* 2011). In this study, the final fine-tuning cluster analysis was carried out using
167 Ward's method. The above calculation processes were carried out using a modified version of SOM Toolbox
168 2.0 (Vesanto *et al.* 2000). The output SOM clusters were plotted on Piper trilinear and stiff diagrams to
169 explain main features of each cluster. Furthermore, the SOM clusters were mapped spatially to clarify
170 influence from land use.

171

172 **Results and discussion**

173 Based on the methodology described above, the number of SOM nodes was determined equal to
174 91. The number of rows and columns was 7 and 13, respectively. Thus, this SOM design was used for the
175 cluster analysis of standardized water chemistry data from the 36 locations in Shimabara.

176 Figure 3 shows the obtained component planes for the 91 reference vectors (nodes) of the eight
177 ion component concentrations (standardized to a range between 0 and 1). Each component plane shows the
178 standardized value of each parameter (concentration) of the 91 reference vectors (nodes) using a color
179 gradient. Comparison between the component planes shows relationships (or correlation) among the
180 parameters. For example, a similar color gradient can be observed for Cl⁻ (Fig. 3 (a)) and NO₃⁻ (Fig. 3 (b)).
181 The same trend can be seen for Na⁺ (Fig. 3 (e)) and Mg²⁺ (Fig. 3 (g)) in their respective component planes.
182 This means that there is high positive correlation between these variables. A great advantage of SOM is
183 that relationships between nodes on the component plane are clearly visualized. For example, the node

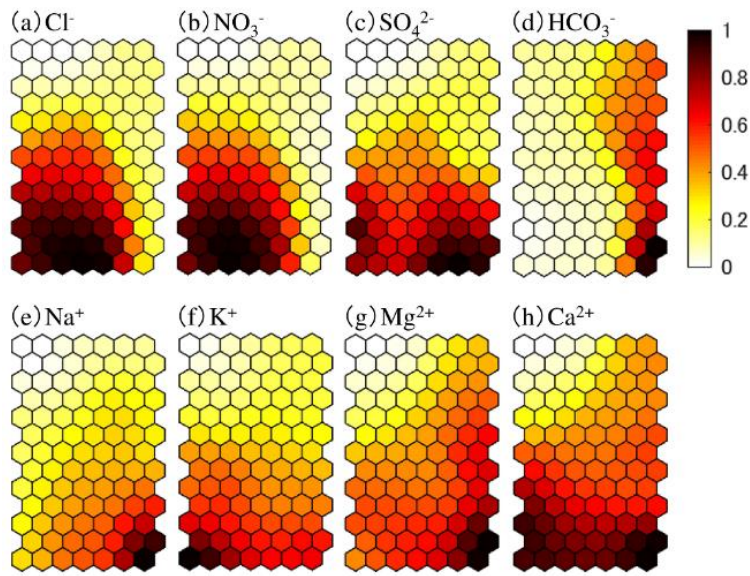


Figure 3 Component plane for (a) Cl^- , (b) NO_3^- , (c) SO_4^{2-} , (d) HCO_3^- , (e) Na^+ , (f) K^+ , (g) Mg^{2+} and (h) Ca^{2+}

located at the uppermost left end shows lower normalized concentrations for all ions ($\text{Cl}^-:0.00$, $\text{NO}_3^-:0.00$, $\text{SO}_4^{2-}:0.00$, $\text{HCO}_3^-:0.11$, $\text{Na}^+:0.00$, $\text{K}^+:0.00$, $\text{Mg}^{2+}:0.00$, and $\text{Ca}^{2+}:0.00$). The node located at the uppermost right end shows moderately higher normalized concentrations for HCO_3^- , Mg^{2+} , and Ca^{2+} ($\text{Cl}^-:0.13$, $\text{NO}_3^-:0.09$, $\text{SO}_4^{2-}:0.15$, $\text{HCO}_3^-:0.46$, $\text{Na}^+:0.09$, $\text{K}^+:0.18$, $\text{Mg}^{2+}:0.33$, $\text{Ca}^{2+}:0.40$). The node located at the lowermost left shows relatively higher normalized ion concentrations except for HCO_3^- ($\text{Cl}^-:0.85$, $\text{NO}_3^-:0.83$, $\text{SO}_4^{2-}:0.78$, $\text{HCO}_3^-:0.04$, $\text{Na}^+:0.30$, $\text{K}^+:1.00$, $\text{Mg}^{2+}:0.43$, $\text{Ca}^{2+}:0.90$). On the other hand, the node located at the lowermost right shows higher normalized ion concentrations except for Cl^- and NO_3^- ($\text{Cl}^-:0.29$, $\text{NO}_3^-:0.17$, $\text{SO}_4^{2-}:0.95$, $\text{HCO}_3^-:0.95$, $\text{Na}^+:1.00$, $\text{K}^+:0.65$, $\text{Mg}^{2+}:0.98$, $\text{Ca}^{2+}:1.00$).

To confirm quantitative relationships as mentioned above, correlation coefficients between reference vectors for each parameter were calculated (Table 2). There is a high correlation ($r = 0.99$) between Cl^- and NO_3^- . There is also a high correlation between Na^+ and Mg^{2+} ($r = 0.92$). Similarly, the color gradient for the

198 Table 2 Correlation between reference vectors for each parameter

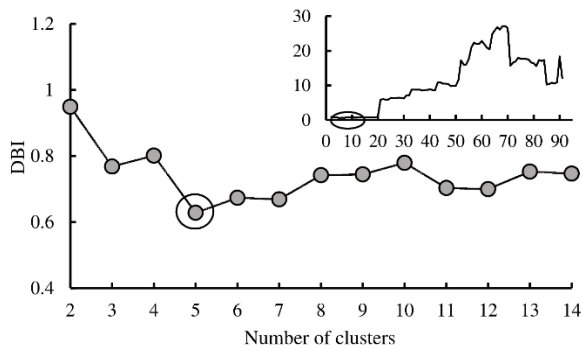
	NO ₃ ⁻	SO ₄ ²⁻	HCO ₃ ⁻	Na ⁺	K ⁺	Mg ²⁺	Ca ²⁺
Cl ⁻	0.99*	0.82*	-0.51*	0.47*	0.86*	0.46*	0.78*
NO ₃ ⁻		0.75*	-0.60*	0.38*	0.82*	0.36*	0.71*
SO ₄ ²⁻			-0.03	0.84*	0.92*	0.79*	0.94*
HCO ₃ ⁻				0.43*	-0.11	0.52*	0.11
Na ⁺					0.71*	0.92*	0.82*
K ⁺						0.72*	0.94*
Mg ²⁺							0.88*

199 * Correlations significant at $p = 0.01$

200

201 relationship between SO₄²⁻ and Ca²⁺ indicates a high correlation coefficient ($r = 0.94$). The relation between
 202 each ion indicates factors affecting groundwater chemistry. For example, a high co-variation ($R^2 = 0.72$)
 203 between higher concentrations of NO₃⁻ and Cl⁻ was observed, indicating that they originate from common
 204 sources such as human and animal waste (e.g., Diédhiou *et al.* 2012). Moreover, the same result can be
 205 observed between SO₄²⁻ and Ca²⁺ ($r = 0.79$). The high correlation implies that the dissolution of gypsum
 206 may be one of the key factors controlling the geochemical evolution of groundwater (Liu *et al.* 2015).

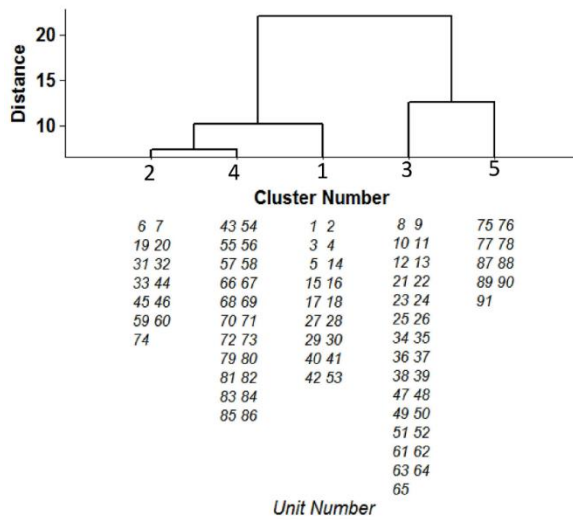
207 Figure 4 shows the variation of DBI with a magnified front between 2 and 14 clusters. The
 208 minimum DBI is shown for five clusters meaning that this number should be used as an optimal value.
 209 After determining the number of clusters, the hierarchical clustering algorithm by Ward was
 210 carried out for the five clusters to fine-tune pattern classification. Figure 5 shows the hierarchical cluster



211

212 Figure 4 Variation of DBI values with the optimal number of clusters marked by the circle on the figure

213



214

215 Figure 5 Dendrogram with node number classified into clusters

216

217 dendrogram. The 91 nodes of the SOM were classified into five different clusters. Figure 6 shows the

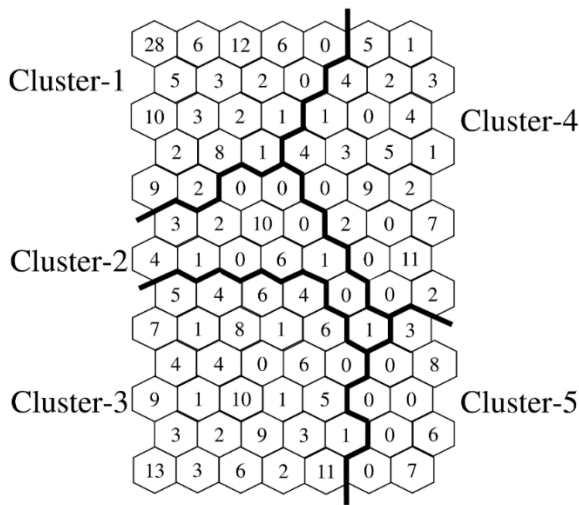
218 pattern classification map for these five clusters. The number for each node represents the raw data

219 classified into each node. Simultaneous analysis of the component planes (Fig. 3) and the pattern

220 classification result (Fig. 6) indicates what kind of data the respective clusters include. For example, cluster-

221 3 (lower left part of Fig. 6) is associated with high contents of Cl^- and NO_3^- . This pattern is observed in the

222 same part of the respective component planes for each parameter as shown in Fig. 3. On the other hand,



223

224 Figure 6 Pattern classification map of the five clusters by the SOM. The numbers on squares off the map
 225 represent the number of data classified into each node

226

227 groundwater samples in nodes with extremely low concentration of all ions, are located at the upper left
 228 part of each component plane (associated with cluster-1) as shown in Fig. 3.

229

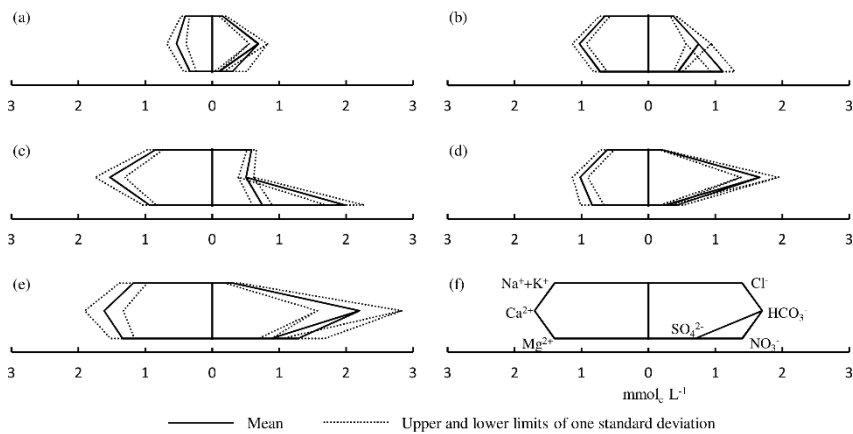
230 More quantitative information than the visualized pattern classification can be extracted and
 231 interpreted from the obtained reference vectors. Stiff diagrams for the respective clusters were represented
 232 by mean and upper and lower limits of one standard deviation using reference vectors of each cluster to

233

234 characterize the clustered data. For example, stiff diagram for cluster-1 is represented by reference vectors
 235 of 18 nodes classified into the cluster. Figure 7 shows stiff diagrams for the five clusters with eight
 236 parameters containing mean values and standard deviations. Cluster-1 (Fig. 7 (a)) shows low values for all

237

238 ions as compared to other clusters. The visible pattern of cluster-2 (Fig. 7 (b)) and cluster-3 (Fig. 7 (c)) is
 239 not similar as shown in the figure. However, they are characterized by high concentration of NO_3^- . Cluster-
 240 2 represents lower concentrations than that of cluster-3 for all ions except HCO_3^- . The pattern with the



238

239 Figure 7 Stiff diagrams for the respective clusters with mean value and upper and lower limits of one
 240 standard deviation by obtained reference vectors; (a) Cluster-1, (b) Cluster-2, (c) Cluster-3, (d)
 241 Cluster-4, (e) Cluster-5 and (f) Legend

242

243 highest Ca^{2+} in cations and HCO_3^- in anions is associated with cluster-4 (Fig. 7 (d)). In this cluster, the
 244 concentration of Na^+ , K^+ , and Mg^{2+} is slightly lower than that for Ca^{2+} . For anions, the concentration of
 245 HCO_3^- is significantly higher than other anions. This pattern is also shown in cluster-5 (Fig. 7 (e)). It is
 246 clear that all ion concentration except for Cl^- and NO_3^- of cluster-5 are higher than that of cluster-4.

247 The classified five clusters can generally be divided into two water quality types. Cluster-2 and
 248 -3 can be characterized as polluted water due to the high concentration of NO_3^- . The other group includes
 249 cluster-1, -4, and -5 representing non-polluted water (pristine water type).

250 Table 3 shows mean ion concentrations calculated from raw data and classified into the respective
 251 cluster. The NO_3^- for cluster-3 indicates higher mean value than 50 mg L^{-1} which is the maximum
 252 contamination level recommended by World Health Organization (WHO 2011) for drinking water. The
 253 NO_3^- for cluster-2 meets WHO standard. However, it is exceeding 13 mg L^{-1} which is the maximum nitrate

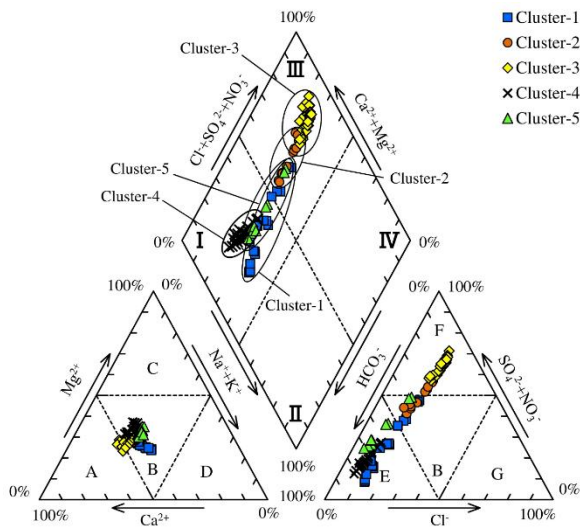
254 Table 3 Mean ion concentrations calculated from raw data and classified into clusters

	Cl ⁻	NO ₃ ⁻	SO ₄ ²⁻	HCO ₃ ⁻	Na ⁺	K ⁺	Mg ²⁺	Ca ²⁺
	(mg L ⁻¹)	(mg L ⁻¹)	(mg L ⁻¹)	(mg L ⁻¹)	(mg L ⁻¹)	(mg L ⁻¹)	(mg L ⁻¹)	(mg L ⁻¹)
Cluster-1	5.1	9.9	3.2	37.7	6.5	3.4	3.2	8.7
Cluster-2	14.3	42.1	22.5	39.0	11.2	6.2	8.1	20.5
Cluster-3	21.3	78.8	37.7	27.5	14.4	8.6	11.2	31.5
Cluster-4	6.4	9.9	10.5	108.5	11.1	4.9	10.6	21.0
Cluster-5	6.8	6.2	41.3	175.4	25.1	7.9	17.6	33.5

255

256 concentration unaffected by human activities (Eckhardt & Stackelberg 1995). It confirms that the two
 257 clusters include polluted water as mentioned above. Cluster-1, -4, and -5 display much lower mean NO₃⁻
 258 concentration. An NO₃⁻ concentration exceeding the maximum concentration level recommended by WHO
 259 has also been reported in other studies (e.g., Diédhiou *et al.* 2012; Hansen *et al.* 2012; Liu *et al.* 2015;
 260 Dragon *et al.* 2016; Matiatos 2016). In these investigations, maximum NO₃⁻ concentration ranged from 91
 261 to 855 mg L⁻¹.

262 Figure 8 shows Piper trilinear diagrams for all reference vectors (91) and respective cluster. With
 263 respect to cations, most vectors of all clusters are located in zone B in the lower left delta-shaped region,
 264 indicating a non-typical water. However, a part of the reference vectors for cluster-3 is located in zone A,
 265 indicating a calcium-type water. For anions, reference vectors are mostly located in zone B, E, and F in the
 266 lower right delta-shaped region, suggesting that the reference vectors of cluster-1, -4, and -5 are
 267 bicarbonate-type water and the reference vectors of cluster-2 and -3 are sulfate and nitrate-type water or
 268 non-typical water. Thus, in the Piper trilinear diagram, two main water types are revealed. These are
 269 calcium-magnesium bicarbonate type (zone I) including cluster-1, 4, and 5 (non-polluted water type) and
 270 calcium-magnesium chloride-sulfate-nitrate type (zone III) including cluster-2 and -3 (polluted water type).



271

272 Figure 8 Trilinear diagram for clusters obtained by reference vectors

273

274 Based on the stiff and Piper trilinear diagrams, the polluted water type is represented as Ca- ($\text{SO}_4^{2-} + \text{NO}_3^-$)
 275 (calcium sulfate nitrate type), while the non-polluted water type is classified as Ca- HCO_3^- (calcium
 276 bicarbonate type). Similar results were reported by Shin *et al.* (2013). According to the study, water samples
 277 collected from the upper reaches of Korean rivers were of Ca- HCO_3^- type, whereas water samples collected
 278 from lower reaches and with relative high nitrate concentration were classified as Na-Cl- NO_3^- type. This
 279 indicates that water samples are affected by anthropogenic factors such as fertilizer, manure, and septic
 280 waste.

281 Figure 9 shows the spatial distribution of the five clusters in Shimabara. All sampling locations
 282 belonging to cluster-2 and -3, representing the polluted water type, are located in the northern part of
 283 Shimabara encompassing a concentration of agricultural fields. In order to investigate the interaction
 284 between groundwater and river water, one sample was taken from the river (R-2) and included into the

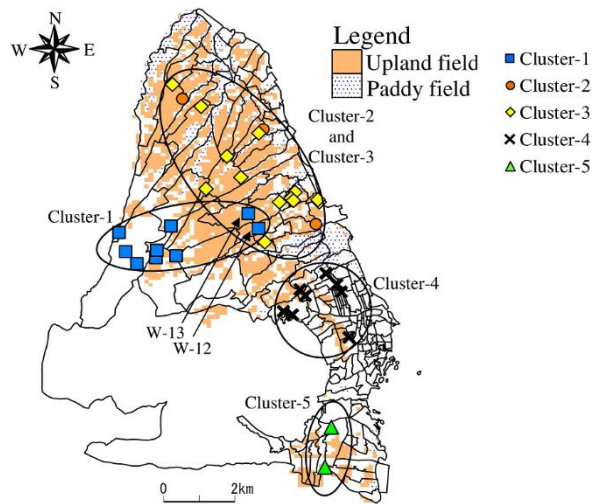


Figure 9 Spatial distribution of clusters

SOM analysis. The results showed that R-2 also is classified into cluster-3 as O-1 and 2. This revealed that they are connected and exchange water with each other. Samples with high nitrate concentrations often correspond agricultural land use (Babiker *et al.* 2004; Esmaeili *et al.* 2014). It confirms that agricultural activities are related to high nitrate concentrations in groundwater. Ishihara *et al.* (2002) reported that fecal coliforms were detected in the northern part of Shimabara. This means that the groundwater in this area is affected by livestock waste. It is observed that most sampling locations for cluster-1 are distributed in the mountainside forest area upstream the heavily polluted areas. This shows that groundwater is recharged in the area and typically is of pristine water type. The average NO_3^- concentration of cluster-1 is slightly lower than that of cluster-4 according to Table 3. Sampling points such as W-12 and 13 located in the agricultural area are thus affected by agricultural activities belonging to cluster-1. This suggests that cluster-1 shows a transition of water chemistry from pristine to polluted water type. The sampling locations for cluster-4 and -5, characterized by high ion concentrations, are located in the urban area at a lower altitude (below 100

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

300 m). This suggests that dissolution of ions from the aquifer matrix during groundwater flow from the
301 mountainside to the urbanized area may increase ion concentrations. Mayuyama avalanche debris deposits
302 are distributed in the east area of Mt. Mayuyama (Ozeki *et al.* 2005). This area corresponds to sampling
303 locations for cluster-5. The pattern of cluster-5 has high concentration for all ions as shown in Fig. 7. This
304 is due to the effect of volcanic deposit on the groundwater chemistry in the area.

305

306 **Summary and conclusion**

307 In this study, water chemistry data from 36 sampling locations, obtained from April 2012 to May
308 2015, were classified using SOM in combination with hierarchical cluster analysis to clarify groundwater
309 characteristics in Shimabara, Japan. The SOM provided readily understandable results for classifying the
310 water chemistry data into distinguishable hydrogeochemical types. The Piper trilinear and stiff diagrams
311 for the reference vectors were plotted to display fundamental characteristics of each cluster. In addition, the
312 spatial distribution of the respective clusters explained the spatial variability of the hydrogeochemical
313 characteristics determined by the SOM. Based on the SOM results, the water chemistry data could be
314 divided into five clusters that revealed two representative water types characterized by nitrate pollution
315 (cluster-2 and -3) and non-polluted (cluster-1, -4, and -5) water. The spatial distribution of cluster-2 and -3
316 shows that agricultural activities are causing groundwater pollution in the northern part of Shimabara. The
317 stiff and Piper trilinear diagrams based on the reference vectors for each cluster showed that non-polluted
318 water and polluted water are characterized by Ca-HCO₃ type and Ca-(SO₄²⁻+NO₃) type, respectively. This
319 indicates that nitrate pollution is a product from agricultural activities and classified into cluster-2 and -3.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

320 The SOM analysis showed that mountainside recharged pristine groundwater is classified into
321 cluster-1. Some groundwater of cluster-1 is also located close to the mid-slope hills. This means that non-
322 polluted water can be used from this agricultural area. For other purposes, water quality evaluation methods
323 such as the Wilcox classification diagram (Wilcox 1955), can be used to evaluate if water in cluster-2 or -3
324 can be used for, e.g., irrigation. The clusters from the SOM analysis are useful for further groundwater
325 remediation alternatives.

326 The application and results of the SOM support our previous conclusion (Nakagawa *et al.* 2016)
327 regarding the spatial distribution of nitrate pollution in the study area and its causes. Data that display a
328 scattered distribution in the piper trilinear diagram can be difficult to analyze by PCA. However, in this
329 case, SOM can be an alternative method (Choi *et al.* 2014). In this study, both PCA and SOM successfully
330 classified groundwater chemistry in the study area. However, SOM gives more robust and explainable
331 results that can be used to characterize groundwater chemistry. More detailed characteristics along this line
332 will be described in a new paper (Amano *et al.* submitted).

333

334 **Acknowledgements**

335 This work was supported by JSPS KAKENHI Grant Number 24360194 and 15KT0120.

336

337 **References**

338 Aiuppa A., Bellomo S., Brusca L., D’Alessandro W. & Federico C. 2003 Natural and anthropogenic factors
339 affecting groundwater quality of an active volcano (Mt. Etna, Italy). *Applied Geochemistry*, **18**(6),

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

340 863-882.

341 Amano H., Nakagawa K. & Kawamura A. Classification characteristics of multivariate analyses for
342 groundwater chemistry in the nitrate contaminated area. submitted. (in Japanese with English
343 abstract)

344 Babiker I. S., Mohamed M. A. A., Terao H., Kato K. & Ohta K. 2004 Assessment of groundwater
345 contamination by nitrate leaching from intensive vegetable cultivation using geographical
346 information system. *Environment International*, **29**(8), 1009-1017.

347 Banoeng-Yakubo B., Yidana S. M. & Nti E. 2009 Hydrochemical analysis of groundwater using
348 multivariate statistical methods - The Volta Region, Ghana. *KSCE Journal of Civil Engineering*,
349 **13**(1), 55-63.

350 Bedoya D., Novotny V. & Manolagos E.S. 2009 Instream and offstream environmental conditions and
351 stream biotic integrity importance of scale and site similarities for learning and prediction.
352 *Ecological Modelling*, **220**(19), 2393-2406.

353 Céréghino R., Giraudel J. L. & Compin A. 2001 Spatial analysis of stream invertebrates distribution in the
354 Adour-Garonne drainage basin (France), using Kohonen self organizing maps. *Ecological
355 Modelling*, **146**(1-3), 167-180.

356 Choi B. Y., Yun S. T., Kim K. H., Kim J. W., Kim H. M. & Koh Y. K. 2014 Hydrogeochemical interpretation
357 of South Korean groundwater monitoring data using Self-Organizing Maps. *Journal of
358 Geochemical Exploration*, **137**, 73-84.

359 Cloutier V., Lefebvre R., Therrien R. & Savard M. M. 2008 Multivariate statistical analysis of geochemical

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

360 data as indicative of the hydrogeochemical evolution of groundwater in a sedimentary rock
361 aquifer system. *Journal of Hydrology*, **353**(3-4), 294-313.

362 Committee on nitrate reduction in Shimabara Peninsula 2011 The second term of Shimabara peninsula
363 nitrate load reduction project. (in Japanese)

364 Diédhiou M., Cissé Faye S., Diouf O. C., Faye S., Faye A., Re V., Wohnlich S., Wisotzky F., Schulte, U. &
365 Maloszewski P. 2012 Tracing groundwater nitrate sources in the Dakar suburban area: an isotopic
366 multi-tracer approach. *Hydrological Processes*, **26**(5), 760–770.

367 Dragon K., Kasztelan D., Gorski J. & Najman J. 2016 Influence of subsurface drainage systems on nitrate
368 pollution of water supply aquifer (Tursko well-field, Poland). *Environmental Earth Sciences*, **75**,
369 100.

370 Eckhardt D. A. V. & Stackelberg P. E. 1995 Relation of ground-water quality to land use on Long Island,
371 New York. *Groundwater*, **33**(6), 1019–1033.

372 Esmaeili A., Moore F. & Keshavarzi B. 2014 Nitrate contamination in irrigation groundwater, Isfahan, Iran.
373 *Environmental Earth Sciences*, **72**(7), 2511-2522.

374 Faggiano L., Zwart D., García-Berthou E., Lek S., & Gevrey M. 2010 Patterning ecological risk of pesticide
375 contamination at the river basin scale. *Science of The Total Environment*, **408**(11), 2319-2326.

376 García H. L. & González I. M. 2004 Self-organizing map and clustering for wastewater treatment
377 monitoring. *Engineering Applications of Artificial Intelligence*, **17**(3), 215-225.

378 Ghesquière O., Walter J., Chesnaux R. & Rouleau A. 2015 Scenarios of groundwater chemical evolution
379 in a region of the Canadian Shield based on multivariate statistical analysis. *Journal of*

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

380 *Hydrology: Regional Studies*, **4(B)**, 246-266.

381 Hansen B., Dalgaard T., Thorling L., Sørensen B. & Erlandsen M. 2012 Regional analysis of groundwater
382 nitrate concentrations and trends in Denmark in regard to agricultural influence. *Biogeosciences*,
383 **9**, 3277–3286.

384 Hong Y. S. & Rosen M. R. 2001 Intelligent characterisation and diagnosis of the groundwater quality in an
385 urban fractured-rock aquifer using an artificial neural network. *Urban Water*, **3(3)**, 193-204.

386 Hentati A., Kawamura A., Amaguchi H. & Iseri Y. 2010 Evaluation of sedimentation vulnerability at small
387 hillside reservoirs in the semi-arid region of Tunisia using the Self-Organizing Map.
388 *Geomorphology*, **122(1-2)**, 56-64.

389 Ishihara T., Ura N. & Hamabe M. 2002 Investigation of ground water contaminated by nitrate-nitrogen.
390 *Annual Report of Nagasaki Prefectural Institute of Public Health and Environmental Sciences*,
391 **48**, 106-109 (in Japanese).

392 Jiang N., Luo K., Beggs P.J., Cheung K. & Scorgie Y. 2014 Insights into the implementation of synoptic
393 weather-type classification using self-organizing maps: an Australian case study. *International*
394 *Journal of Climatology*, **35(12)**, 3471-3485. doi: 10.1002/joc.4221.

395 Japan Meteorological Agency 2015 Weather observation data. *Japan Meteorological Agency Web*,
396 <http://www.jma.go.jp/jma/index.html> (accessed 28 January 2015)

397 Jeong K. S., Hong D. G., Byeon M. S., Jeong J. C., Kim H. G., Kim D. K. & Joo G. J. 2010 Stream
398 modification patterns in a river basin: Field survey and self-organizing map (SOM) application.
399 *Ecological Informatics*, **5(4)**, 293-303.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

400 Jin Y. H., Kawamura A., Park S. C., Nakagawa N., Amaguchi H. & Olsson J. 2011 Spatiotemporal
401 classification of environmental monitoring data in the Yeongsan River basin, Korea, using self-
402 organizing maps. *Journal of Environmental Monitoring*, **13**(10), 2886-2894.

403 Kalteh, A.M. & Berndtsson, R. 2007, Interpolating monthly precipitation by self-organizing map (SOM)
404 and multilayer perceptron (MLP), *Hydrological Sciences Journal*, **52**, 305-317.

405 Kalteh A. M., Hjorth P. & Berndtsson R. 2008 Review of the self-organizing map (SOM) approach in water
406 resources: Analysis, modeling and application. *Environmental Modelling & Software*, **23**(7), 835-
407 845.

408 Liu F., Song X., Yang L., Han D., Zhang Y., Ma Y. & Bu H. 2015 The role of anthropogenic and natural
409 factors in shaping the geochemical evolution of groundwater in the Subei Lake basin, Ordos
410 energy base, Northwestern China. *Science of the Total Environment*, **538**, 327-340.

411 Marghade D., Malpe D. B. & Subba Rao N. 2015 Identification of controlling processes of groundwater
412 quality in a developing urban area using principal component analysis. *Environmental Earth
413 Sciences*, **74**(7), 5919–5933.

414 Matiatos I. 2016 Nitrate source identification in groundwater of multiple land-use areas by combining
415 isotopes and multivariate statistical analysis: A case study of Asopos basin (Central Greece).
416 *Science of the Total Environment*, **541**, 802-814.

417 Nadiri A. A., Moghaddam A. A., Tsai F. T. C. & Fijani E. 2013 Hydrogeochemical analysis for Tasuj plain
418 aquifer, Iran. *Journal of Earth System Science*, **122**(4), 1091–1105.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

420 groundwater chemistry in Shimabara, Nagasaki, Japan. *Environmental Earth Sciences*, **75**, 234.

421 Nguyen T. T., Kawamura A., Tong T. N., Nakagawa N., Amaguchi H. & Gilbuena Jr. R. 2015 Clustering
422 spatio-seasonal hydrogeochemical data using self organizing maps for groundwater quality
423 assessment in the Red River Delta, Vietnam. *Journal of Hydrology*, **522**, 661-673.

424 Nishiyama K., Endo S., Jinno K., Uvo C. B., Olsson J. Berndtsson R. 2007 Identification of typical synoptic
425 patterns causing heavy rainfall in the rainy season in Japan by a Self-Organizing Map.
426 *Atmospheric Research*, **83**(2-4), 185-200.

427 Omonona O. V., Onwuka O. S. & Okogbue C. O. 2014 Characterization of groundwater quality in three
428 settlement areas of Enugu metropolis, southeastern Nigeria, using multivariate analysis.
429 *Environmental Monitoring and Assessment*, **186**(2), 651-664.

430 Ozeki N., Okuno M. & Kobayashi T. 2005 Growth history of Mayuyam, Unzen, Kyushu, Southwest Japan.
431 *Bulletin of the Volcanological Society of Japan*, **50**(6), 441-454 (in Japanese with English
432 Abstract).

433 Shin W. J., Ryu J. S., Lee K. S. & Chung G. S. 2013 Seasonal and spatial variations in water chemistry and
434 nitrate sources in six major Korean rivers. *Environmental Earth Sciences*, **68**(8), 2271–2279.

435 Singaraja C., Chidambaram S., Prasanna M. V., Thivya C. & Thilagavathi R. 2014 Statistical analysis of
436 the hydrogeochemical evolution of groundwater in hard rock coastal aquifers of Thoothukudi
437 district in Tamil Nadu, India. *Environmental Earth Sciences*, **71**(1), 451-464.

438 Sonkamble S., Sahya A., Mondal N. C. & Harikumar P. 2012 Appraisal and evolution of hydrochemical
439 processes from proximity basalt and granite areas of Deccan Volcanic Province (DVP) in India.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

440 *Journal of Hydrology*, **438-439**, 181-193.

441 Vesanto J., Himberg J., Alhoniemi E. & Parhankangas J. 2000 SOM Toolbox for Matlab 5. Helsinki
442 University of Technology Report A57.

443 WHO (World Health Organization) 2011 Guidelines for drinking-water quality-4th edn.

444 Wilcox L.V. 1955 Classification and use of irrigation water. United States Department of Agriculture,
445 Washington D.C., Circular No.969.

446

447 **Figure captions**

448 Figure 1 Study area and sampling locations in Shimabara, Nagasaki Prefecture, Japan (RW: Residential
449 well, W: Public water supply well, O: Observation well, S: Spring and R: River)

450 Figure 2 Altitude and land use map of Shimabara; (a) Altitude and (b) Land use

451 Figure 3 Component plane for (a) Cl⁻, (b) NO₃⁻, (c) SO₄²⁻, (d) HCO₃⁻, (e) Na⁺, (f) K⁺, (g)Mg²⁺ and (h) Ca²⁺

452 Figure 4 Variation of DBI values with the optimal number of clusters marked by the circle on the figure

453 Figure 5 Dendrogram with node number classified into clusters

454 Figure 6 Pattern classification map of the five clusters by the SOM. The numbers on squares off the map
455 represent the number of data classified into each node

456 Figure 7 Stiff diagrams for the respective clusters with mean value and upper and lower limits of one
457 standard deviation by obtained reference vectors; (a) Cluster-1, (b) Cluster-2, (c) Cluster-3, (d)
458 Cluster-4, (e) Cluster-5 and (f) Legend

459 Figure 8 Trilinear diagram for clusters obtained by reference vectors

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

460 Figure 9 Spatial distribution of clusters

461

462 **Table captions**

463 Table 1 Mean and standard deviations of 36 sampling sites using averaged ion concentrations for each

464 sampling sites used in SOM

465 Table 2 Correlation between reference vectors for each quality parameter

466 Table 3 Mean ion concentrations calculated from raw data and classified into clusters