

Document downloaded from:

<http://hdl.handle.net/10251/142526>

This paper must be cited as:

Zhang, M.; Muñoz Mas, R.; Martínez-Capel, F.; Qu, X.; Zhang, H.; Peng, W.; Liu, X. (09-2). Determining the macroinvertebrate community indicators and relevant environmental predictors of the Hun-Tai River Basin (Northeast China): A study based on community patterning. *The Science of The Total Environment*. 634:749-759.
<https://doi.org/10.1016/j.scitotenv.2018.04.021>



The final publication is available at

<https://doi.org/10.1016/j.scitotenv.2018.04.021>

Copyright Elsevier

Additional Information

1 **Determining the macroinvertebrate community indicators and relevant**
2 **environmental predictors of the Hun-Tai River Basin (Northeast**
3 **China): a study based on community patterning**

4
5 Min Zhang^{1,2}, Rafael Muñoz-Mas³, Francisco Martínez-Capel³, Xiaodong Qu^{1,2,*}, Haiping Zhang^{1,2},
6 Wenqi Peng^{1,2}, Xiaobo Liu^{1,2}

7 ¹ State Key Laboratory of Simulation and Regulation of Water Cycle in River Basin, China Institute
8 of Water Resources and Hydropower Research, Beijing 100038, China

9 ² Department of Water Environment, China Institute of Water Resources and Hydropower Research,
10 Beijing 100038, China

11 ³ Institut d'Investigació per a la Gestió Integrada de Zones Costaneres (IGIC), Universitat Politècnica
12 de València, C/ Paranimf 1, Grau de Gandia, València 46730, Spain.

13
14 *Corresponding author: Xiaodong Qu, e-mail: davidqu0108@126.com

15

16 **Abstract**

17 It is essential to understand the patterning of biota and environmental influencing factors for proper
18 rehabilitation and management at the river basin scale. The Hun-Tai River Basin was extensively
19 sampled four times for macroinvertebrate community and environmental variables during one year.
20 Self-Organizing Maps (SOMs) were used to reveal the aggregation patterns of the 355 samples.
21 Three community types (*i.e.*, clusters) were found (at the family level) based on the community
22 composition, which showed a clearly gradient by combining them with the representative
23 environmental variables: minimally impacted source area, intermediately anthropogenic impacted
24 sites, and highly anthropogenic impacted downstream area, respectively. This gradient was
25 corroborated by the decreasing trends in density and diversity of macroinvertebrates. Distance from
26 source, total phosphorus and water temperature were identified as the most important variables that
27 distinguished the delineated communities. In addition, the sampling season, substrate type, pH and
28 the percentage of grassland were also identified as relevant variables. These results demonstrated that
29 macroinvertebrates communities are structured in a hierarchical manner where geographic and water
30 quality prevail over temporal (season) and habitat (substrate type) features at the basin scale. In
31 addition, it implied that the local-scale environment variables affected macroinvertebrates under the
32 longitudinal gradient of the geographical and anthropogenic pressure. More than one families were
33 identified as the indicator for each type of community. Abundance contributed significantly for
34 distinguishing the indicators, while Baetidae with higher density indicated minimally and
35 intermediately impacted area and lower density indicated highly impacted area. Therefore, we
36 suggested the use of abundance data in community patterning and classification, especially in the
37 identification of the indicator taxa.

38 **Keywords**

39 Macroinvertebrate patterning, indicators, Self-Organizing Map, decision tree, site classification

40 **1 Introduction**

41 Currently, human activities are greatly influencing the flow rate, water yield, sediment transport and
42 nutrient releases in freshwater ecosystems at scales that far exceed those of natural phenomena
43 (Habersack et al., 2014). Accordingly, water resources are currently over-exploited in many regions,
44 which has resulted in 65% of rivers worldwide being under moderate-to-high threats in terms of
45 human water security and biodiversity loss (Vörösmarty et al., 2010).

46 Biotic assemblages in freshwater ecosystems integrate these impacts throughout the drainage basins;
47 thus, these assemblages can be considered as indicators of ecosystem health (Habersack et al., 2014).
48 Consequently, the classification and delineation of the ecological statuses of rivers based on the
49 biotic assemblages is an essential prerequisite for river ecosystem assessment, restoration and
50 management (Heino et al., 2002; Marchant et al., 2000; Siddig et al., 2016; Tsai et al., 2017).

51 Macroinvertebrate assemblages have been widely used as indicators of ecosystem changes because
52 macroinvertebrate communities encompass a diverse group with a wide range of life-history
53 requirements (O'Brien et al., 2016). Macroinvertebrates vary spatially and temporally and integrate
54 ecosystem changes as a result of their suite of feeding strategies and lifestyles and their different
55 sensitivities to changes in physical habitat and water quality (Milošević et al., 2016; Ogbeibu and
56 Oribhabor, 2002). According to a recent review on indicator species over the last 14 years, nearly 50%
57 of the taxa used as indicators were animals, and 70% of these were invertebrates (Siddig et al., 2016).
58 However, data on macroinvertebrate assemblages are highly complex and difficult to analyze
59 because macroinvertebrate assemblages consist of numerous species that respond in complex
60 manners to natural and anthropogenic pressures (Kim et al., 2013; Tsai et al., 2017). In this situation,
61 supervised machine learning approaches, which make use of techniques from mathematical
62 programming and statistics, have been used to scrutinize and model the environmental requirements

63 of relevant macroinvertebrate taxa, and these techniques include decision trees (C4.5 – D’heyere et
64 al., 2003) or multilayer perceptrons (Edia et al., 2010).

65 In addition, macroinvertebrate datasets include numerous taxa and a large number of samples, which
66 can also cause difficulties for community analysis and river regionalization (Kim et al., 2013). In
67 particular, ordination techniques and unsupervised machine learning approaches have been used to
68 explore patterns of occurrence and community shifts and their relationships with environmental
69 factors (Adriaenssens et al., 2007; Giraudel and Lek, 2001; Zhang et al., 2012a). Nevertheless, each
70 ordination technique may have important limitations and assumptions that are incompatible with the
71 over-dispersion and nonlinear nature of ecological data (Paliy and Shankar, 2016). Researchers have
72 advocated for the use of a type of unsupervised artificial neural network called Kohonen
73 Self-Organizing Maps (SOMs) (Kohonen, 1982), which have been demonstrated to be particularly
74 competent in analyses such as macroinvertebrate community delineations (Chon, 2011; Park et al.,
75 2007; Kim et al., 2013; Sroczyńska et al., 2017).

76 The freshwater ecosystems of China are a clear example of the abovementioned human-induced
77 impacts. For instance, more than 40% of the rivers in China are notably polluted, which has led to
78 poor drinking water quality for approximately 300 million rural residents (Liu and Yang, 2012). The
79 river ecosystems in the northeast have also degraded due to industrial and agricultural development;
80 thus, some river restoration work has been conducted in this area (Kong et al., 2013; Zhang et al.,
81 2011; Zhang et al., 2013). The Hun-Tai River Basin is a large river basin with a basin area of
82 2.73×10^4 km². It represents the overall status of the water in the Liaohe River Basin in northeast
83 China and is undergoing degradation. Many field surveys and studies using macroinvertebrates as
84 important indicators in river health assessments showed the ecosystem were not in good conditions
85 (Qu et al., 2016; Zhang et al., 2011; Zhang et al., 2013). However, most of these studies have mainly
86 focused on small rivers or tributaries, and cannot reflect the overall status of the whole watershed,
87 especially in such a large river basin. Up to now, little is known about the macroinvertebrate

88 community and related environmental variables in the entire basin. We hypothesized when data from
89 a large river basin and temporal span is merged, the geographical features could override the local
90 environmental variables (e.g. water quality) in structuring the macroinvertebrate community, because
91 the geographical features (e.g. elevation, distance from the river source) usually determine the
92 macroinvertebrate structures when communities are studied at broader scales (Dedieu et al., 2014;
93 Gaston, 2000).

94 This study analyzed the macroinvertebrate assemblages present in the large Hun-Tai River Basin to
95 elucidate the existence of different types of communities and determine the indicator families and
96 main environmental predictors for their occurrence. SOMs were used to reveal the existence of these
97 macroinvertebrate communities (*i.e.*, clusters) in different areas of the river basin. Then, a genetically
98 optimized C5.0 algorithm (Quinlan, 1992) (*i.e.*, a type of decision tree) was used to reveal the most
99 important set of environmental predictors and indicator taxa of each community in the Hun-Tai River
100 Basin.

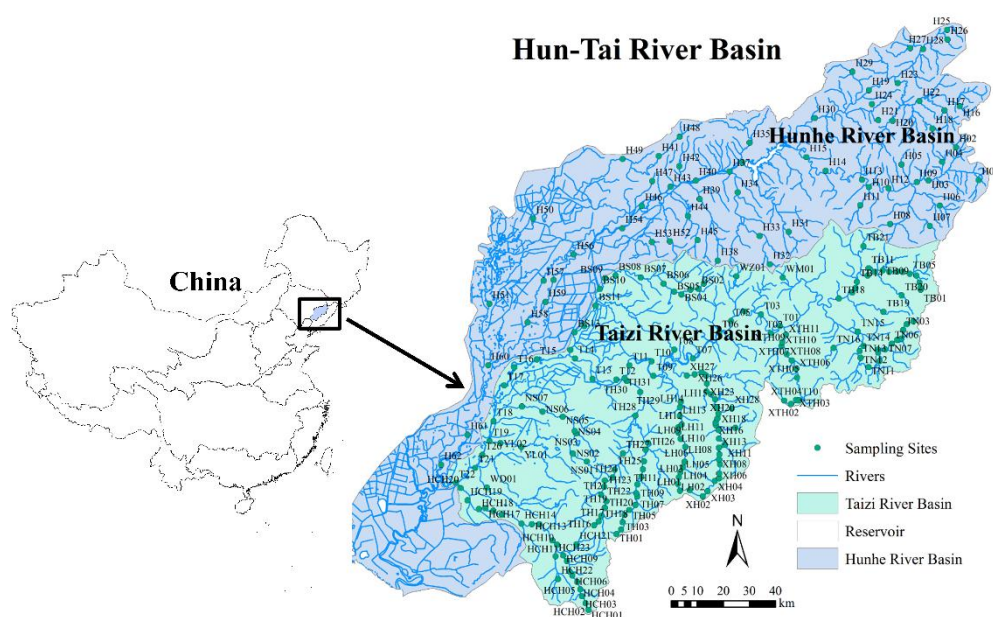
101

102 **2 Materials and methods**

103 **2.1 Study area**

104 The Hun-Tai River is located in Liaoning Province of Northeast China, and it has two main
105 tributaries, the Taizi and the Hunhe Rivers (Fig. 1). The lengths of the Hunhe and Taizi Rivers are
106 approximately 415 and 413 km, respectively. The climate in this area is typical continental monsoon,
107 with the highest temperature (34.3 °C) in the summer and lowest temperature (-25.2 °C) in the winter.
108 The precipitation follows the temperature pattern, with the annual average precipitation 778 mm, 63%
109 of which occurs in summer (Bu et al. 2014).

110 Field surveys were carried out in May 2009 (spring), August 2009 (summer), October 2009 (autumn)
 111 and May 2010 (spring). These surveys encompassed the entire river basin. In total, 287 sites from
 112 May 2009 to May 2010 in the Hun-Tai River Basin, in which 68 of Taizi River Basin were sampled
 113 twice (The number of sampling sites in each river and the codes of each river are located can be
 114 found in Appendix A). Consequently, 355 samples were ultimately collected where the selected
 115 environmental variables were measured *in situ* or obtained from reference databases.



116 Fig. 1. Location of the study area and sampling sites within the Hun-Tai River Basin. Further details
 117 about the codes depicted in this map can be found in Appendix A.
 118
 119

120 2.2 Data collection

121 2.2.1 Environmental variables

122 All sampling sites were characterized using variables determining the geography, hydrology, climate,
 123 landuse, water quality and habitat, and there were 30 variables and 1 binary control variable
 124 (wadeable or non-wadeable) in total (Table 1). For each site, a handheld global positioning system

125 (GPS, Trimble Juno SA) was used to obtain the latitude, altitude and elevation above sea level (m
126 a.s.l.). Distance from the source was extracted from the digital map of the river basin using ArcGIS
127 10.2. The river widths and water depths were directly measured with tape. The river width was
128 measured by randomly selecting three transects per sampling station While the water depth was
129 measured at the macroinvertebrate sampling sites. Air temperature, water temperature, pH, dissolved
130 oxygen (DO), electrical conductivity (EC), and total dissolved solids (TDS) were measured by a
131 multiparameter water quality probe (YSI-Pro Plus[®], YSI Inc., USA).

132 One liter of stream water was collected and transported in a portable fridge. Chemical oxygen
133 demand (COD), total nitrogen (TN), ammonia nitrogen (NH₃-N), nitrate-nitrogen (NO₃-N), and total
134 phosphorus (TP) were measured in the laboratory according to the “Environmental quality standards
135 for surface water” of China (GB3838-2002). COD and NO₃-N measurements were not collected in
136 BeiShahe, NanShahe, TangHe and XiaoTangHe Rivers in spring 2010 (Table 1). Suspended solids
137 (SS) were measured by filtration through pre-dried cellulose acetate membranes (0.45 μm) according
138 to the Chinese standard (GB11901-89). The particle sizes of the substrate were measured by using a
139 series of stainless steel mesh sizes according to the modified Wentworth classification of substrate
140 and then expressed using a percentage. The substrate sizes were classified as boulder (Ø > 256 mm),
141 large cobble (128 > Ø ≤ 256 mm), cobble (64 > Ø ≤ 128 mm), large pebble (32 > Ø ≤ 64 mm),
142 pebble (16 > Ø ≤ 32 mm), large gravel (8 > Ø ≤ 16 mm), gravel (4 > Ø ≤ 8 mm) and small gravel,
143 sand and silt (Ø < 4 mm) (Cummins, 1962). The substrates in the HaichengHe, LanHe and XiHe
144 Rivers were not monitored in summer 2009. Finally, the proportion of each landuse type was
145 extracted over 3 km upstream and 500 m wide buffer zones (Zhang et al., 2013) from a digital
146 landuse map using ArcGIS 10.2. The landuse map was interpreted from landsat TM data of the year
147 2010. These variables and the sampling season were used to infer the most determinant
148 environmental variables of the macroinvertebrate communities of the Hun-Tai River Basin. The
149 basic statistics and the number of unavailable data for the whole variables’ set are shown in Table 1.

150 Table 1. Summary and units of environmental variables collected in the Hun-Tai River Basin

Variable group	Variables	Min	Mean	Median	Max	Number of unavailable data
Geographic	Elevation (m a.s.l.)	4.53	205.38	188.00	663.00	0
	Distance from source (km)	0.12	65.55	34.52	380.68	0
Landuse	Grassland (%)	0.00	4.47	0.00	54.91	0
	Agricultural land (%)	0.00	48.59	48.11	97.04	0
	Residential land (%)	0.00	14.30	4.60	96.90	0
	Forest (%)	0.00	32.63	31.03	100.00	0
Hydrologic	Discharge (m ³ /s)	0.00	3.40	1.04	94.76	0
	River width (m)	0.30	40.41	15.00	420.00	0
	Water depth (cm)	0.80	24.60	22.33	130.00	0
	Velocity (m/s)	0.00	0.40	0.38	1.14	0
Climatic	Air temperature (°C)	5.20	21.22	21.80	36.80	0
	Water temperature (°C)	6.80	17.42	18.00	27.80	0
Water quality	pH	6.01	8.39	8.40	10.12	0
	DO (mg/l)	0.00	8.93	9.63	15.63	0
	EC (µs/cm)	3.90	317.51	272.00	1431.00	0
	SS (mg/l)	1.00	67.08	19.50	1110.00	0
	TDS (mg/l)	14.95	235.15	192.00	995.00	0
	COD (mg/l)	0.00	22.96	16.00	151.00	62
	TN (mg/l)	0.47	5.19	3.49	22.60	0
	NH ₃ -N (mg/l)	0.00	1.15	0.38	20.60	0
	NO ₃ -N (mg/l)	0.01	2.16	1.60	15.90	70
TP (mg/l)	0.00	0.19	0.07	3.03	0	
Habitat	Boulder (%)	0.00	12.74	0.00	90.56	119
	Large cobble (%)	0.00	19.91	18.64	79.70	119
	Cobble (%)	0.00	15.65	16.42	58.38	119
	Large pebble (%)	0.00	11.77	11.26	77.78	119
	Pebble (%)	0.00	8.82	8.18	47.13	119
	Large gravel (%)	0.00	3.64	3.01	20.00	119
	Gravel (%)	0.00	4.34	2.85	29.05	119
	Small gravel, sand and silt (%)	0.00	23.12	4.84	100.00	119

151

152 **2.2.2 Macroinvertebrate sampling**

153 Macroinvertebrates were collected using a Surber net (30×30 cm, 500 µm mesh). At the sites that
 154 could be waded, three replicates were obtained from two riffles and one shallow pool. For the sites
 155 that could not be waded, a Surber net was used to collect three replicates in shallow water along the
 156 riverside.

157 The samples were passed through a 500 μm mesh sieve, and organisms retained on the sieve were
158 fixed and preserved in 10% formaldehyde. Most taxa were identified based on the available
159 references (Brinkhurst, 1986; Merritt and Cummins, 1996; Morse et al., 1994; Wiggins, 1996).
160 Finally, 90 families were identified and they were used to cluster the sampling sites (community
161 delineation) and to obtain the community indicators. The density (individuals/ m^2 – ind./ m^2), richness
162 (number of species) and the Shannon–Wiener diversity and Pielou evenness indices were calculated
163 for each sample to characterize the macroinvertebrate communities of the Hun-Tai River Basin.
164 The wadeable vs non-wadeable nature of each site was used as a control variable during the process
165 to determine the most relevant environmental variables. Ruling out this variable indicated us that the
166 sampling protocol had no impact in the delineated communities (*i.e.*, clusters).

167

168 **2.3 Statistical analysis**

169 **2.3.1 Self-Organizing Maps (SOMs)**

170 Delineating macroinvertebrate communities with Self-Organizing Maps (SOMs) is done in two steps
171 (see e.g., Edia et al., 2010; Kim et al., 2013; Park et al., 2007) (Fig. 2). First, the SOM implements an
172 ordered dimensionality-reducing mapping of input variables (Kohonen, 1982). Therefore, SOM
173 provided a projection of the matrix of $n = 335$ rows and $p = 90$ families onto a topological structure
174 (*i.e.*, a XY 2D map of unit neurons nodes) of smaller dimensionality (*i.e.*, $X \ll n$ and $Y \ll p$). The
175 entire process is carried out preserving the original topology of the input data. In accordance, the
176 neurons that are located near to each other in the SOM had similar associated input samples (*i.e.*,
177 macroinvertebrate communities). Then, these neuron nodes are clustered, usually employing the
178 Ward's linkage method (Murtagh and Legendre, 2014), allowing exclusively the aggregation of
179 contiguous neuron nodes. Finally, the samples (*i.e.*, macroinvertebrate communities) assigned to the
180 neuron nodes aggregated in a given cluster are grouped (*i.e.* clustered) together, and upon these

181 communities or clusters, further analyses were performed (see e.g., Edia et al., 2010; Kim et al.,
182 2013; Park et al., 2007).

183 The training and visualization of the SOM was performed using the functionalities implemented
184 within the *R* package *kohonen* (Wehrens and Buydens, 2007). SOMs were trained for all possible
185 combinations from 2 to 20 neuron nodes for the *X* and *Y* axes and the quality of these alternative
186 topologies was evaluated using the quantization error (QE), which evaluates the resolution of the
187 map, and the topographic error (TE), which indicates the accuracy of the topology preservation of
188 the map (Kim et al., 2013; Tsai et al., 2017). The learning rate (α) varied linearly from 0.05 and 0.01,
189 and the neighborhood function was Gaussian. Conversely, the initial radius varied in accordance
190 with the SOM dimensions (*i.e.*, *X* and *Y*) after the $\operatorname{argmax}\left\{\left\lfloor\frac{X}{3}\right\rfloor, \left\lfloor\frac{Y}{3}\right\rfloor\right\}$. In terms of densities (ind./m²),
191 the macroinvertebrate data depicted large numerical differences. Therefore, data were transformed
192 (log+1) prior to training the SOMs (Adriaenssens et al., 2007; Kim et al., 2013; Tsai et al., 2017).
193 Finally, the input sample data were presented to each SOM 500 times and the selected distance
194 measure was the Euclidean distance (See Appendix B for further details about the optimization of
195 SOMs).

196 Once the optimal dimensions of the SOM were determined based on the QE and TE, the Ward's
197 linkage method was applied to the SOM to cluster the unit neurons and, consequently, the
198 macroinvertebrate samples (Chon, 2011; Park et al., 2007; Tsai et al., 2017). The function *NbClust*
199 included in the homonymous *R* package (Charrad et al., 2014) was used to determine the optimal
200 number of clusters. The latter function calculates 30 quality indices, from *ball* (Ball and Hall, 1965)
201 to *gap* (Tibshirani et al., 2001), and the optimal number of clusters is determined using the majority
202 rule. In this case, the optimal number of cluster between 2 and 15 was sought. Once the optimal
203 number of clusters was determined, the distribution of the categorical variables (season) was

204 visualized with the SOM, and the family density, richness and the diversity indices of the samples
205 assigned to each cluster were scrutinized with violin plots.

206

207 **2.3.2 C5.0 algorithm**

208 The most relevant indicator families and environmental predictors of the macroinvertebrate
209 communities that were delineated by the SOM (*i.e.*, clusters) were identified with a genetically
210 optimized C5.0 algorithm (Quinlan, 1992) (Fig. 2). The C5.0 algorithm was selected because it is a
211 kind of decision tree that is able to handle missing or unavailable data. In addition, C5.0 is able to
212 collapse the former tree-like structure into a compact list of IF-THEN rules. The missing data was
213 down-weighted when the entropy gain is calculated. The proportion of missing data do not
214 necessarily reduces the predictive capacity of a variable.

215 To prevent overfitting, a *wrapper* approach involving cross-validation and the Genetic Algorithm
216 (GA) (Holland, 1992) implemented within the *R* package *rgenoud* (Mebane Jr & Sekhon, 2011) was
217 used to find the optimal variable set and C5.0 hyperparameters (Muñoz-Mas et al., 2016). The
218 parameters of the GA were selected to avoid premature convergence (Muñoz-Mas et al., 2016)
219 whereas, compared to previous studies (D'heygere et al., 2003; Gobeyn et al., 2017), the population
220 size and number of generations were set very large (*i.e.* to 1000), and the optimization halted after
221 250 generations without improvement (See Appendix B for further detail about the optimization with
222 the GA). The optimization took place by maximizing the product of the individual sensitivities (S_n)
223 for each class (*i.e.*, cluster), as described in Equation 4 (Caballero et al., 2010; Pérez-Ortiz et al.,
224 2015). It was performed following a threefold cross-validation ($3 \times 3_{cross-validation}$) scheme
225 (Muñoz-Mas et al., 2016), with every fold presenting a similar proportion of samples per community
226 (*i.e.*, samples per cluster) and favoring the use of each variable in every of the nine decision trees
227 (Equation 4).

228

229
$$Fitness = \frac{1}{9} \sum_1^9 \left(\prod_i^k Sn_i \times \frac{\# variables\ used}{\# variables\ selected} \right) \text{ (Equation 4)}$$

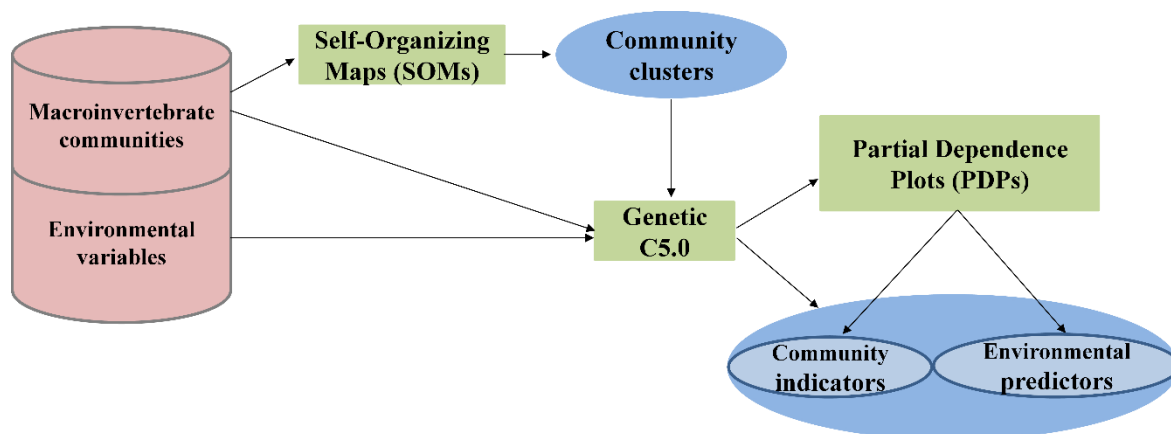
230

231 Where k corresponds to the ultimate number of communities delineated by the SOM (*i.e.*, the
232 ultimate number of different clusters). To avoid redundancy in the optimal variables' set (community
233 indicators or environmental predictors), the number of variables was restricted by preventing
234 correlated variable combinations (*i.e.*, $r^2 > 66\%$) (Additional information about data correlation can
235 be found in Appendix B).

236 To allow comparison with previous studies, other performance criteria were calculated, namely,
237 accuracy or correctly classified instances, Cohen's kappa, sensitivity, specificity, and balanced
238 accuracy (*i.e.* the number of correctly predicted cases weighted by the rarity of the community) (see
239 Mouton et al., 2010 for additional details about performance criteria).

240 Once the optimal hyperparameters and the most relevant environmental variables or
241 macroinvertebrate families were obtained, a single C5.0 decision tree was trained using the entire
242 dataset (*i.e.* without cross-validation) (Fukuda et al., 2013; Muñoz-Mas et al., 2016). The resulting
243 models were used to calculate the variable importance based on *usage*, which measures the
244 percentage of training set samples that fall into all the terminal nodes after the split, and *splits*, which
245 measures the percentage of splits associated with each variable (Fig. 2). Finally, the relationship
246 between the environmental variables or macroinvertebrate families and the probability of occurrence
247 of each community or cluster was scrutinized with partial dependence plots (Friedman, 2001) to
248 accommodate the tree-like or rule-based structure of the optimal C5.0s (Fig. 2). This was done
249 adapting the code implemented in the *randomForests* package (Liaw and Wiener, 2002).

250



251
252 Fig. 2. Flowchart depicting the process followed to delineate the macroinvertebrate communities and
253 to identify the family indicators and most relevant environmental drivers.

254

255 3 Results

256 3.1 Macroinvertebrate patterning with Self-Organizing Maps (SOMs)

257 The SOMs that simultaneously minimized the quantization and topographic errors (4.24 and 0.31,
258 respectively) had a lattice of 17 x 19 neurons. This SOM presented 32.8% of empty neurons, and
259 Ward's approach distinguished three clusters (hereafter clusters I, II and III) (Fig. 3).

260 Cluster I encompassed the largest number of samples ($n = 178$) and mainly included samples from
261 both spring and summer. Within this cluster, there were many samples in some neurons from the
262 LanHe River (LH), HunHe River (H) and Taizi River mainstream (T), and these neurons are located
263 in the upper part of the SOM. Some downstream sites in the XiHe (XH), HaiChengHe (HCH) and
264 TangHe (TH) Rivers were also located in this cluster.

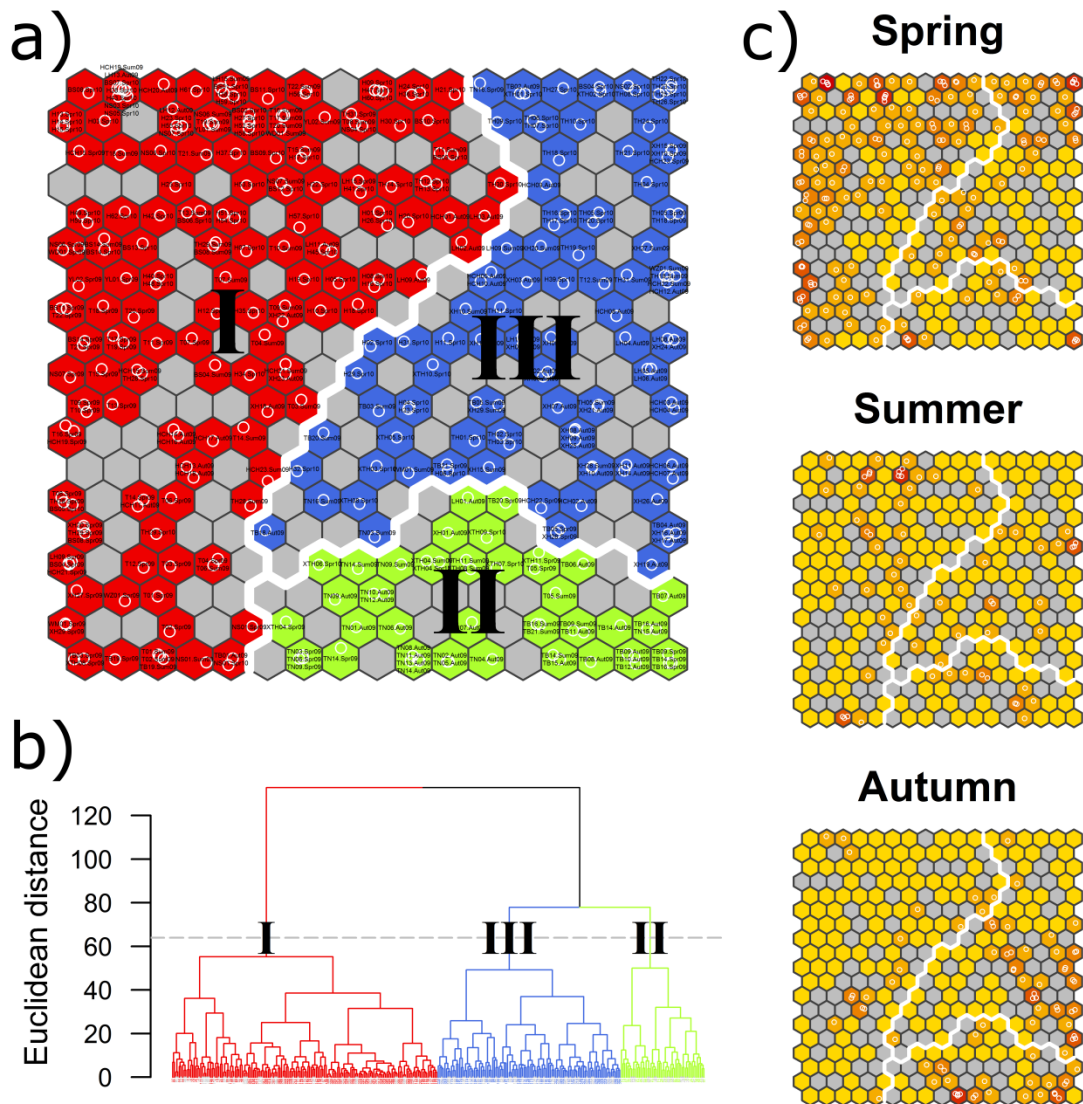
265 Cluster II mainly contained samples from all seasons ($n = 51$) from sites in the Taizi Nan River (TN)
266 and Taizi Bei River (TB), and it included most sites in the XiaoTangHe (XTH) river and some
267 upstream sites of LanHe, XiHe Rivers.

268 Cluster III included many sites in the autumn and spring and some sites in the summer ($n = 106$).

269 However, within this cluster, the spring sites exhibited higher similarity between each other, which

270 was determined based on the information shown by the neurons in the upper part of the SOM. The
 271 sites were relatively dissimilar in the summer and autumn and were further from the sites in the
 272 spring. Spatially, the sites in cluster III were mainly composed of sites in the XH, TH, XTH and
 273 HCH Rivers (especially the downstream sites), which were located in the lower reach of the Taizi
 274 River mainstream.

275



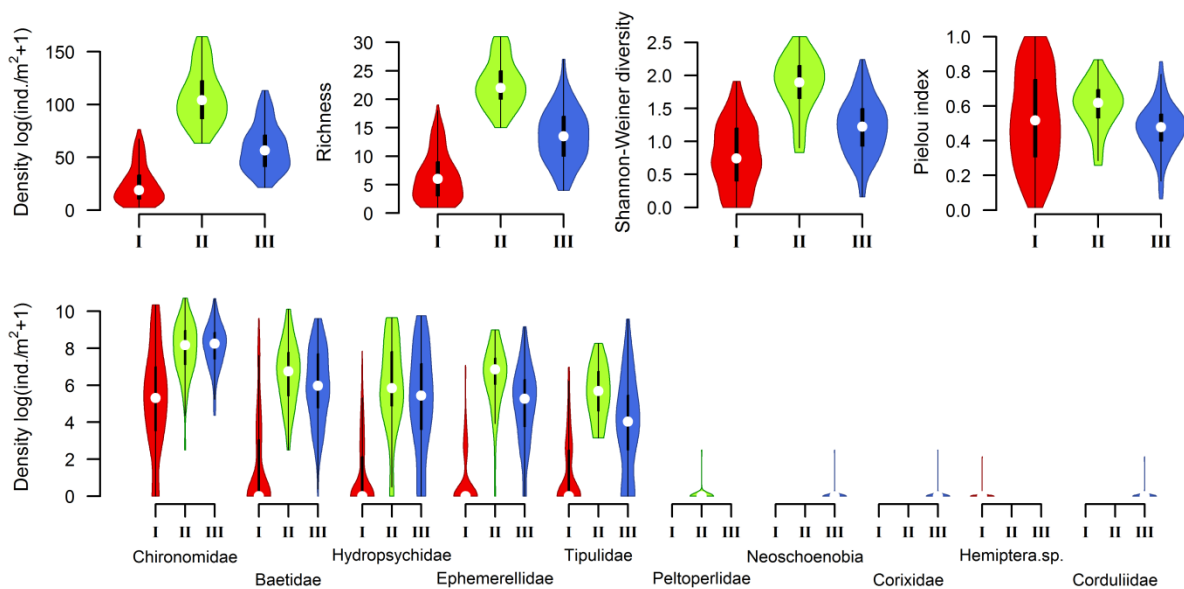
276 Fig. 3. a) Classification of the samples collected in Self-Organizing Map (SOM), b) dendrogram
 277 clustering of the SOM and c) mapping of the samples collected in each season (spring, summer and
 278 autumn). The darker the cell the larger the number of samples mapped in the corresponding neuron
 279 node.
 280

281

282 **3.2 Relevant Macroinvertebrate community characteristics**

283 In the three clusters, cluster I had the lowest macroinvertebrate density and family richness whereas
284 cluster II had the highest (Fig. 4). Shannon-Weiner diversity and Pielou evenness showed similar
285 patterns. The most abundant families were Chironomidae, Baetidae, Hydropsychidae,
286 Ephemerellidae and Tipulidae, while the least abundant ones were Peltoperlidae, Neoschoenobia,
287 Corixidae, Hemiptera sp., and Cordullidae. The most abundant families reproduced the general
288 patterns on abundance (Fig. 4). Therefore, cluster II encompassed the samples with higher densities
289 and cluster I those with the lower.

290



291
292 Fig. 4. Total density (log (ind./m²+1)) and richness (families) and diversity indices (Shannon-Wiener
293 diversity and Pielou evenness index) of the samples included in each cluster (upper sequence). The
294 density of the most and least abundant families per cluster is depicted in the lower sequence.

295

296 **3.3 Identification of the indicators and environmental predictors**

297 The genetically optimized C5.0 model to infer the most relevant macroinvertebrate families
298 outperformed the model to discover the main environmental predictors. However, both models

299 achieved high values for each performance criteria. In both cases, cluster II presented higher
 300 performance criteria, and cluster III presented lower performance (Table 2).

301
 302 Table 2. Summary of the accuracy or correctly classified instances Cohen's kappa, sensitivity,
 303 specificity and balanced accuracy calculated during the $3 \times 3_{cross-validation}$ (nine models).

Classification based on	Cluster	Accuracy/CCI	Cohen's kappa	Sensitivity	Specificity	Balanced accuracy
Environmental variables	I			0.82±0.04	0.85±0.11	0.83±0.06
	II	0.77±0.08	0.61±0.12	0.88±0.19	0.94±0.03	0.91±0.09
	III			0.70±0.16	0.85±0.05	0.78±0.10
Indicator families	I			0.81±0.07	0.88±0.1	0.85±0.07
	II	0.80±0.08	0.65±0.14	0.88±0.16	0.93±0.03	0.91±0.08
	III			0.74±0.12	0.87±0.06	0.81±0.09

304

305 3.3.1 Environmental predictors of the macroinvertebrate communities

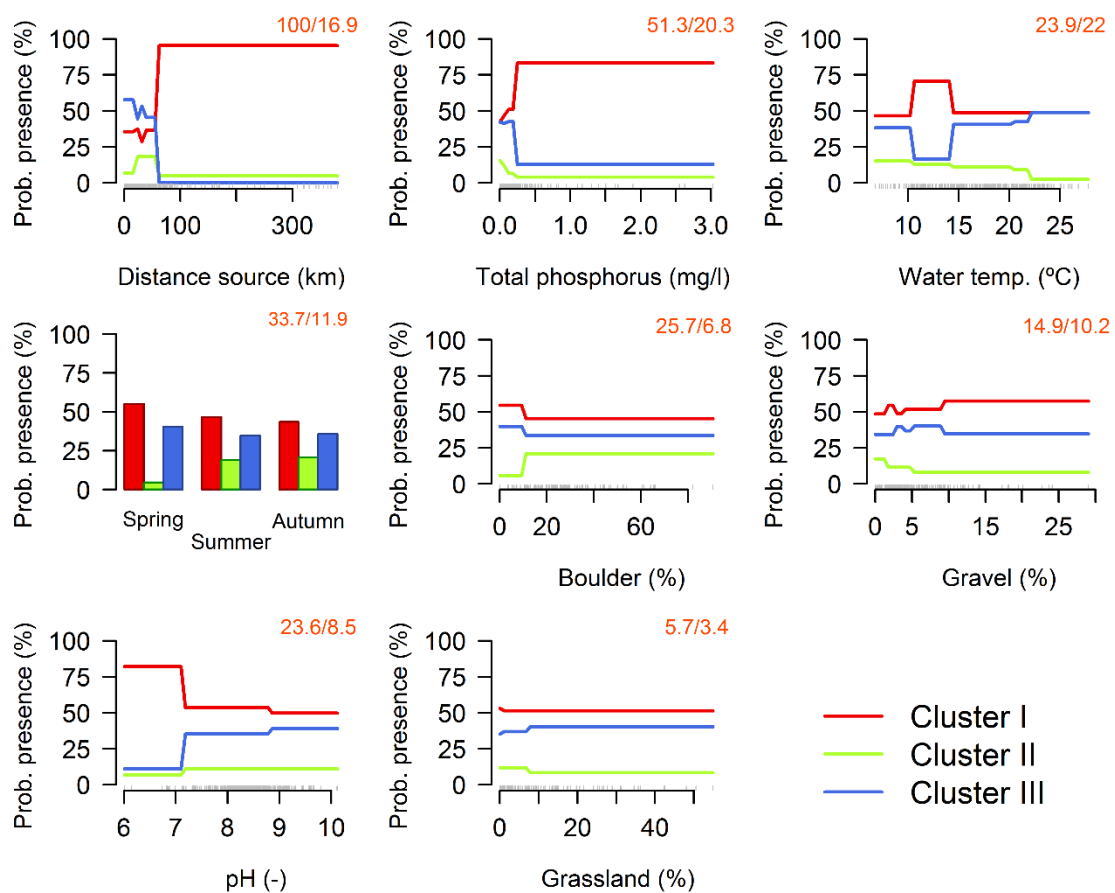
306 Based on the metrics *splits* and *usage*, the most relevant environmental predictors for the three
 307 macroinvertebrate communities (*i.e.*, clusters) were distance from source followed by total
 308 phosphorous (Fig. 5). The lesser important factors were pH and grassland whereas water
 309 temperature, season, boulder and gravel occupied intermediate positions. Nevertheless, there were
 310 differences in the ranges of the variables that characterized each cluster.

311 Cluster I, which principally encompassed samples that were collected in the spring and summer, was
 312 characterized by high values of total phosphorous, long distance from source and intermediate water
 313 temperature. The samples delineated in this community were collected in areas with low percentage
 314 of boulder, high percentage of gravel and lower pH whereas the percentage of grassland had almost
 315 no influence in its probability of presence.

316 Cluster II included samples collected in summer and especially in autumn. This community was
 317 characteristic of river segments with low total phosphorous, distance from source and water
 318 temperature. It occurred in segments with coarse substrate, high percentage of boulder and low of
 319 gravel and higher pH. Finally, low percentage of grassland had a positive effect on its presence.

320 Cluster III included samples collected in spring and autumn and the partial dependence plots for the
 321 most relevant variables resembled those for cluster II. Therefore, this community was characteristic
 322 of river segments with low total phosphorous and distance from source. However, its presence was
 323 favored by relatively higher water temperatures. It occurred in intermediate substrate granulometries
 324 (*i.e.*, low percentage of boulder and relatively low percentage of gravel) and higher pH whereas the
 325 percentage of grassland had a positive effect on its presence.

326



327

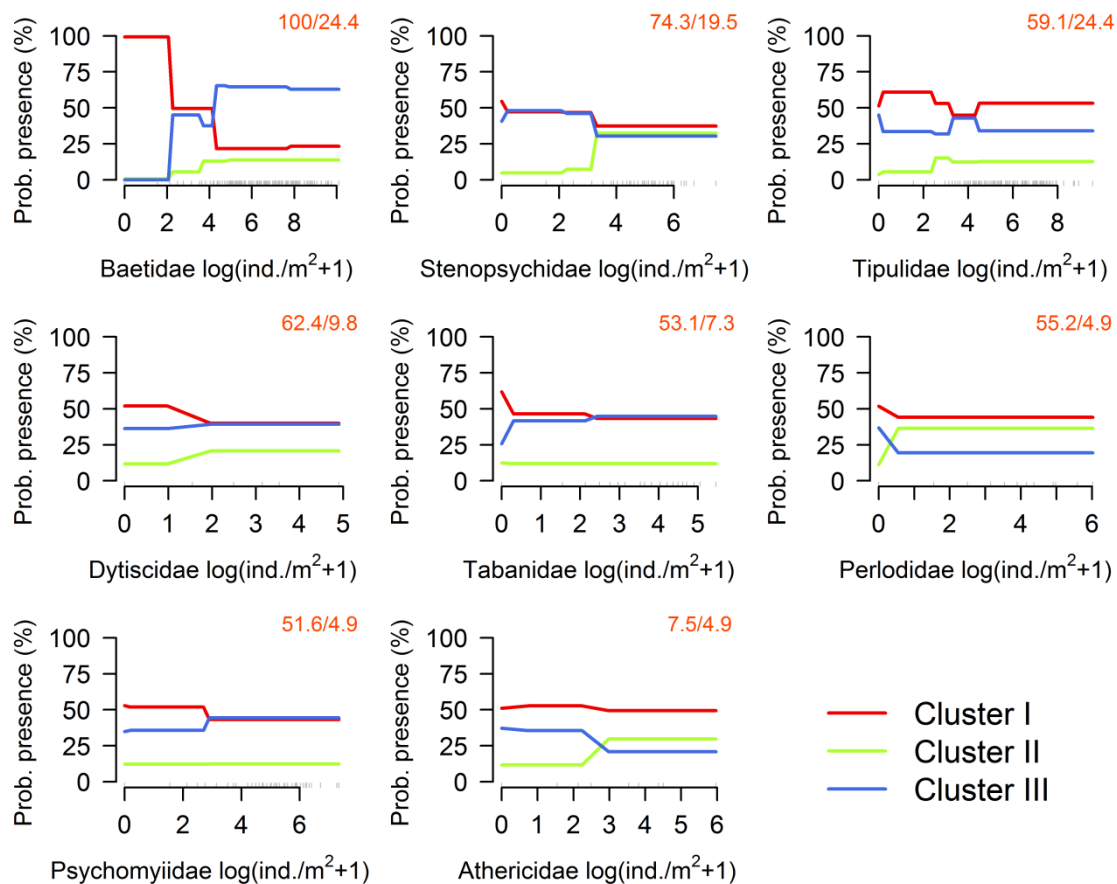
328 Fig. 5. Partial dependence plots relating the probabilities of the presence of each cluster and the
 329 selected environmental variables. The variable importance based on *usage/splits* is tagged above
 330 each panel. The tick marks near the x-axis depict the collected data.

331

332 **3.3.2 Relevant indicators of the macroinvertebrate communities**

333 The genetically optimized C5.0 algorithm found eight representative indicator families for this study
 334 area, namely, Baetidae, Stenopsychidae, Tipulidae, Dytiscidae, Tabanidae, Perlodidae
 335 Psychomyiidae, and Athericidae, (Fig. 6.). Among these families, Baetidae was the most
 336 characteristic family of the entire river basin. Cluster I was characterized by low densities of
 337 Baetidae, Stenopsychidae, Tipulidae, and Dytiscidae. For cluster II, higher densities of
 338 Stenopsychidae, Perlodidae and Athericidae were the most prominent mark. Cluster III was
 339 characterized by high densities of Baetidae, Tipulidae and Psychomyiidae, and low densities of
 340 Stenopsychidae, Perlodidae and Athericidae.

341



342 Fig. 6. Partial dependence plots relating the probabilities of the presence of each cluster and the
 343 selected community indicator families. The variable importance based on *usage/splits* is tagged
 344 above each panel. The tick marks near the x-axis depict the collected data.
 345

346

347 **4 Discussion**

348 **4.1 Community patterning and the environmental predictors**

349 All study sites were aggregated into three different clusters. This number of clusters is significantly
350 inferior when compared to previous studies (Kim et al., 2013; Park et al., 2007) but in line with the
351 spatial and temporal extension of the study (they do not encompass multiple river basins and long
352 periods). Given that our working dataset was fivefold larger than the datasets in previous studies, the
353 resulting number of clusters may highlight the robustness of the analysis.

354 Macroinvertebrate community patterns have been usually linked to anthropogenic stress gradients
355 (Álvarez-Cabria et al., 2011). In this study, the delineated communities (*i.e.*, clusters) portrayed a
356 very clear spatial gradient that can be observed by comparing the sites in the three clusters. By
357 combining the representative environmental variables, the three clusters could be defined as:
358 minimally impacted source area, intermediately anthropogenic impacted area, and highly
359 anthropogenic impacted downstream area, respectively. This classification was in concordance with
360 some previous assessments in this area (Li et al., 2013; Qu et al., 2016; Zhang et al., 2013). The sites
361 sampled in the TaiziNan and TaiziBei rivers, which were mostly encompassed in cluster II, are
362 located in a protected area, and, in these studies, they exhibited good ecological conditions. Most of
363 the sites in Hunhe River and the area in the tributaries and downstream of Taizi River, which were
364 encompassed in cluster I, were determined to be in a poor ecological status (Li et al., 2013), while
365 there were also some studies showing very low macroinvertebrate diversities in this area (Kong et al.,
366 2013).

367 This gradient in community structure likely reflected the longitudinal changes in natural (gradient,
368 temperature) and anthropogenic (water quality always deteriorated from source to downstream areas)
369 factors that influenced the study river basin, as previous studies have shown (Álvarez-Cabria et al.,

370 2011; Traversetti and Scalici, 2014). The macroinvertebrate density and diversity both showed a
371 decreasing trend with the pollution increase. This result is typical because reductions in
372 macroinvertebrate density and diversity have been observed in many studies as a response of the
373 benthic communities to pollution and habitat alterations (Boehme et al., 2016; Ogbeibu and
374 Oribhabor, 2002), which is in agreement with the predicted effects first of distance from source, total
375 phosphorous and water temperature and then on habitat quality (*i.e.*, substrate) (D'heyere et al.,
376 2003).

377 In our study, distance from source and total phosphorous had arguably almost equal importance to
378 the community classification, and to a lesser extent, water temperature, although the former one was
379 a little more important. In the original hypothesis, we expected the geographical variables to override
380 the local environment (e.g. water quality), because over larger scales, geographical gradient and
381 variability appear to have stronger influence (Allan, 2004; Mykrä et al., 2007), and spatial structuring
382 may mask the effect of the local environment on the macroinvertebrate community structures
383 (Tonkin et al., 2017). Our results support the hypothesis to some degree. The difference is the local
384 environment variables (e.g. water quality) were also a determinant factor for classifying the
385 macroinvertebrate communities. Interestingly, the landuse type (a larger scale variable) just showed,
386 apparently, minor influence to the community structure. In general terms, local habitat and biological
387 diversity of streams and rivers are strongly influenced by landuse type within the surrounding valley
388 at multiple scales (Allan, 2004). The percentage of the landuse type correlated well with distance
389 from source (See Appendix B). Human settlements, and hence anthropogenic impacts (*e.g.*, diffuse
390 pollution, landuse changes), are negatively correlated with elevation (Kummu et al., 2016), which is
391 a general pattern also observed in the Taizi River Basin. Therefore, it is plausible to consider
392 distance from source a proxy of these large scale processes. Then, we could conclude that the
393 influence of local scale variables on macroinvertebrate depended on larger-scale longitudinal

394 gradient under anthropogenic pressure, which is similar with what other studies have shown
395 (Manfrin et al., 2016).

396 In addition, within each cluster, the distribution patterns of some sites showed high relevance with
397 the seasonality. The sites subjected to seasonality had higher similarities and were grouped into
398 specific neurons (Fig. 3a). In the relatively clean areas (cluster II and III), sites in different seasons
399 were grouped into different neurons. This seasonal pattern overlapped with the spatial zonation
400 pattern, which was consistent with Sroczyńska et al. (2017) for a temporary Mediterranean river.
401 Conversely, in cluster I which occupied the upper part of the SOM, the sites in the spring and
402 summer and even some sites in the autumn were concentrated in nearby neurons, indicating the high
403 similarity of the community structure at different sites and in different seasons. Kim et al. (2013)
404 showed that the variability of the macroinvertebrate density (mainly tolerant species) was very small
405 in different seasons at severely polluted sites. This result may indicate that seasonality played an
406 important role in the community patterning or sites classification in clean areas with little or
407 intermediate anthropogenic influence, while this seasonality effect was minimized by other factors in
408 severely anthropogenic influenced areas, as the case in our study area.

409 Substrate is very important to macroinvertebrates (Connolly and Pearson, 2007; Sroczyńska et al.,
410 2017). However, in this case it showed intermediate influence for the community classification. This
411 might be caused by the taxonomic level employed to delineate the communities. Thus, low-level or
412 trait-based analyses may have rendered different variable rankings (Soininen et al., 2016; Wu et al.,
413 2014). Communities are product of the environmental variables that act at multiple spatial scales
414 (Boyero and Bailey, 2001; Liu et al., 2016). Therefore, the links between the environmental variables
415 and species may be masked at this level because different species adapt to changes in the
416 environment differently by a multitude of strategies (Lamouroux et al., 2004; Resh et al., 1994).
417 However, community incorporates all of the species information and reflects all environmental
418 changes at multiple scales (Heino et al., 2002; Marchant et al., 2000; Zhang et al., 2012b), and could

419 provide a comprehensive reflection of the ecosystem to indicate more relevant ecoregions for river
420 management and restoration plans.

421

422 **4.2 Family indicators and the relationships with the environmental variables**

423 As mentioned above, family level was used to identify the indicators. It was considered satisfactory
424 to characterize the ecological status because investigations at the species level at large spatial scales
425 are particularly costly (Heino et al., 2002; Rosenberg and Resh, 1993; Sroczyńska et al., 2017). In
426 this regard, some researcher compared the taxonomic resolution's influence to macroinvertebrate
427 community patterning and found that only little information (<6%) was lost using family level, as
428 opposed to species level. And they concluded that family level abundance was a better resolution for
429 patterning the macroinvertebrate community (Marshall et al., 2006). This also in agreement with the
430 standard taxonomic level employed in a number classic and new rapid bioassessment protocols (*e.g.*,
431 Kaaya et al., 2015).

432 Four out of the eight families that were considered as the most characteristic families (Fig. 9) in the
433 Hun-Tai River Basin belong to the orders Ephemeroptera, Plecoptera or Trichoptera (EPT), which
434 are widely accepted as sensitive indicators of habitat conditions (Boehme et al., 2016). Other
435 families had also lower tolerance except Tabanidae, which is considered a family slightly tolerant to
436 pollution (Mandaville, 2002; Qin et al., 2014).

437 Baetidae was the most characteristic family in our study area. A higher density of Baetidae
438 accompanied its higher probability of presence of the minimally impacted source area (cluster II) and
439 the intermediate impacted area (cluster III), which is in agree with other studies that performed
440 variable selection with GA which also found Baetidae intended to live in the conditions characterized
441 by no pollution (Gobeyn et al., 2017). In addition, this is a phenomenon observed in numerous
442 streams (Kasangaki et al., 2006; Törnblom et al., 2011).

443 The evenness in the three clusters were all very high, implying that there might be more than one
444 representative families as indicators. This has been proven by our indicator analysis results (Fig. 6).
445 The combination of three families with low tolerance values (Mandaville, 2002; Qin et al., 2014)
446 appearing at the same time with higher density was regarded as the indicators of the clean river
447 source areas (cluster II), while lower density of four sensitive families were identified as the
448 indicators of the most anthropogenic impacted area (cluster I).

449 Usually, the presence/absence of taxa is used to delineate the community patterning or for
450 identification of the indicator species (Paini et al., 2010; Tonkin et al., 2017). Some researchers also
451 demonstrated that species assemblage patterns were adequately reproduced at the resolution of
452 family using presence/absence data (Wright et al., 1989; Rutt et al., 1990). However, our indicator
453 analysis results showed that density (or abundance) could be better in community patterning For
454 instance, Baetidae appeared in many sites of both cluster I and III, but lower abundance indicated
455 cluster I while higher abundance indicated cluster III. Additionally, some species with broad niche
456 could live under different environmental condition, e.g. *Limnodrilus hoffmeisteri*, which could live in
457 many environment (Kim et al., 2013; Song et al., 2006; Zhang et al., 2012a), but only the very high
458 abundance could indicate an organic pollution (Chapman et al., 1982). In accordance, we concluded
459 that abundance could give more biological information than presence/absence data, which is in
460 agreement with other studies that highlighted the benefits of abundance data over presence/absence
461 data (Fukuda et al., 2011; Marshall et al., 2006). Therefore, we suggest the use of abundance rather
462 than presence/absence data in patterning the macroinvertebrate assemblages, especially to identify
463 the most indicative taxa.

464

465 **5 Conclusions**

466 In this study, the combination of SOMs and decision trees was demonstrated to be a proficient tool
467 for community patterning and identification of the most relevant environmental predictors and
468 indicator taxa (*i.e.*, families). A gradient of three types of sites (communities) were distinguished:
469 minimally impacted source area, intermediately anthropogenic impacted area, and highly
470 anthropogenic impacted downstream area. Distance from source and total phosphorus were
471 considered to be the most important environmental factors determining the presence of each
472 community (*i.e.*, cluster), which indicated that local environmental factors affect macroinvertebrate
473 community composition under the geographical gradient. At the same time, water temperature,
474 season, substrate, pH and the percentage of grassland were also identified as distinguishing factor.
475 These results support the hypothesis that the importance of environmental predictors are spatial-scale
476 dependent because macroinvertebrates' presence was primarily regulated by processes operating at
477 larger spatial scales (*i.e.*, summarized in the variable distance from source), while they responded,
478 only subsequently, to the water quality and habitat structure.

479 Eight families were identified as the indicators of the community in the study area with families of
480 the order Ephemeroptera, Plecoptera or Trichoptera (EPT) leading the ranking. Particularly, the
481 abundance of Baetidae, a sensitive family to pollution, was the most relevant indicator for the three
482 delineated communities. Baetidae was particularly abundant in less impacted sites, although
483 combinations of more than one family were identified as indicators for each community (*i.e.*, cluster)
484 in such an ecosystem with high species evenness. Abundance contributed a lot for distinguishing the
485 indicators in different areas, so that the indicators identified in our study were the families with their
486 density data. Therefore, we suggested to use the abundance rather than the presence/absence data in
487 community patterning or specific zonation of an ecosystem, especially in the identification of the
488 indicators. These results should help to improve Hun-Tai River management plans and to refine
489 future bio-monitoring programs in the region and the similar bioclimatic river regions.

490

491 **6 Acknowledgments**

492 This work was supported by the National Natural Science Foundation of China (51779275,
493 41501204, 51479219) and the IWHR Research & Development Support Program
494 (WE0145B532017).

495

496 **7 References**

- 497 Adriaenssens V., Verdonshot P.F.M., Goethals P.L.M., Pauw N.D., 2007. Application of clustering
498 techniques for the characterization of macroinvertebrate communities to support river restoration
499 management. *Aquat. Ecol.* 41, 387-398.
- 500 Allan, J. D., 2004. Landscapes and riverscapes: the influence of land use on stream ecosystems. *Annu. Rev.*
501 *Ecol. Evol. Syst.* 35, 257–284.
- 502 Álvarez-Cabria M., Barquín J., Juanes J.A., 2011. Macroinvertebrate community dynamics in a temperate
503 European Atlantic river. Do they conform to general ecological theory? *Hydrobiologia* 658, 277-291.
- 504 Ball G.H., Hall D.J., 1965. ISODATA, a novel method of data analysis and pattern classification. Stanford
505 research inst Menlo Park CA.
- 506 Boehme E.A., Zipper C.E., Schoenholtz S.H., Soucek D.J., Timpano A.J., 2016. Temporal dynamics of
507 benthic macroinvertebrate communities and their response to elevated specific conductance in
508 Appalachian coalfield headwater streams. *Ecol. Indic.* 64, 171-180.
- 509 Boyero L., Bailey R.C., 2001. Organization of macroinvertebrate communities at a hierarchy of spatial scales
510 in a tropical stream. *Hydrobiologia* 464, 219-225.
- 511 Brinkhurst R.O., 1986. Guide to the freshwater aquatic microdrile oligochaetes of North America. Ottawa,
512 Canada, Dept. of Fisheries and Oceans.
- 513 Brown B.L., Swan C.M., 2010. Dendritic network structure constrains metacommunity properties in riverine
514 ecosystems. *J. Anim. Ecol.* 79, 571-80.
- 515 Bu H., Meng W., Zhang Y., 2014. Spatial and seasonal characteristics of river water chemistry in the Taizi
516 River in Northeast China. *Environ. Monit. Assess.* 186, 3619-3632.
- 517 Caballero J.C.F., Martínez F.J., Hervás C., Gutiérrez P.A., 2010. Sensitivity versus accuracy in multiclass
518 problems using memetic pareto evolutionary neural networks. *IEEE Trans. Neural Networks* 21,
519 750-770.
- 520 Chapman, P.M., Farrell, M.A. and Brinkhurst, R.O., 1982. Effects of species interactions on the survival and
521 respiration of *Limnodrilus hoffmeisteri* and *Tubifex tubifex* (Oligochaeta Tubificidae) exposed to
522 various pollutants and environmental factors. *Water Res.* 16, 1405-1408.
- 523 Charrad M., Ghazzali N., Boiteau V., Niknafs A., 2014. An R Package for Determining the Relevant Number
524 of Clusters in a Data Set. *J. Stat. Softw.* 61, 1-36.
- 525 Chon T. S., 2011. Self-Organizing Maps applied to ecological sciences. *Ecol. Inform.* 6, 50-61.

- 526 Connolly N.M., Pearson R.G., 2007. The effect of fine sedimentation on tropical stream macroinvertebrate
527 assemblages: a comparison using flow-through artificial stream channels and recirculating
528 mesocosms. *Hydrobiologia* 592, 423-438.
- 529 Cummins K.W., 1962. An evaluation of some techniques for the collection and analysis of benthic samples
530 with special emphasis on lotic waters. *Am. Midl. Nat.* 67, 477-504.
- 531 Dedieu N., Vigouroux R., Cerdan P., Céréghino R., 2014. Invertebrate communities delineate
532 hydro-ecoregions and respond to anthropogenic disturbance in East-Amazonian streams.
533 *Hydrobiologia* 742, 95-105.
- 534 D'heygere, T., Goethals P. L. M., De Pauw N., 2003. Use of genetic algorithms to select input variables in
535 decision tree models for the prediction of benthic macroinvertebrates. *Eco. Model.* 160: 291–300.
- 536 Edia E.O., Gevrey M., Ouattara A., Brosse S., Gourène G., Lek S., 2010. Patterning and predicting aquatic
537 insect richness in four West-African coastal rivers using artificial neural networks. *Knowl. Manag.*
538 *Aquat. Ecosyst.* 63, 78-83.
- 539 Friedman J.H., 2001. Greedy function approximation: a gradient boosting machine. *Ann. Stat.*, 1189-1232.
- 540 Fukuda S., De Baets B., Waegeman W., Verwaeren J., Mouton A.M., 2013. Habitat prediction and knowledge
541 extraction for spawning European grayling (*Thymallus thymallus* L.) using a broad range of species
542 distribution models. *Environ. Model. Software* 47, 1-6.
- 543 Fukuda, S., Mouton, A.M., and De Baets, B., 2011. Abundance versus presence/absence data for modelling
544 fish habitat preference with a genetic Takagi-Sugeno fuzzy system. *Environ. Monit. Assess.* 184,
545 6159–6171.
- 546 Gaston K.J., 2000. Global patterns in biodiversity. *Nature* 405, 220-227.
- 547 Giraudel J.L., Lek S., 2001. A comparison of self-organizing map algorithm and some conventional statistical
548 methods for ecological community ordination. *Ecol. Model.* 146, 329-339.
- 549 Gobeyn, S., Volk M., Dominguez-Granda L., Goethals P. L. M., 2017. Input variable selection with a simple
550 genetic algorithm for conceptual species distribution models: A case study of river pollution in
551 Ecuador. *Environ. Model. Software* 92: 269–316.
- 552 Habersack H., Haspe D., Kondolf M., 2014. Large Rivers in the Anthropocene- Insights and tools for
553 understanding climatic, land use, and reservoir influences. *Water Resour. Res.* 50, 3641-3646.
- 554 Heino J., Muotka T., Paavola R., Hämäläinen H., Koskenniemi E., 2002. Correspondence between regional
555 delineations and spatial patterns in macroinvertebrate assemblages of boreal headwater streams. *J. N.*
556 *Am. Benthol. Soc.* 21, 397-413.
- 557 Kaaya, L.T., Day, J.A., and Dallas, H.F. 2015. Tanzania River Scoring System (TARISS): a
558 macroinvertebrate-based biotic index for rapid bioassessment of rivers. *African J. Aquat. Sci.* 40,
559 109–117.
- 560 Kasangaki, A., Babaasa, D., Efitre, J., McNeilage, A., and Bitariho, R. (2006) Links Between Anthropogenic
561 Perturbations and Benthic Macroinvertebrate Assemblages in Afromontane Forest Streams in Uganda.
562 *Hydrobiologia* 563, 231–245.
- 563 Kim D.-H., Cho W.-S., Chon T.-S., 2013. Self-organizing map and species abundance distribution of stream
564 benthic macroinvertebrates in revealing community patterns in different seasons. *Ecol. Inform.* 17,
565 14-29.
- 566 Kohonen T., 1982. Self-organized formation of topologically correct feature maps. *Biol. Cybern.* 43, 59-69.
- 567 Kong W., Meng W., Zhang Y., Gippel C., Qu X., 2013. A freshwater ecoregion delineation approach based on
568 freshwater macroinvertebrate community features and spatial environmental data in Taizi River
569 Basin, northeastern China. *Ecol. Res.* 28, 581-592.

- 570 Kummu, M., de Moel, H., Salvucci, G., Viviroli, D., Ward, P.J., and Varis, O. (2016) Over the hills and
571 further away from coast: global geospatial patterns of human and environment over the 20th–21st
572 centuries. *Environ. Res. Lett.* 11 (3), 34010.
- 573 Lamouroux N., Dolédec S., Gayraud S., 2004. Biological traits of stream macroinvertebrate communities:
574 effects of microhabitat, reach, and basin filters. *J. N. Am. Benthol. Soc.* 23, 449-466.
- 575 Li Y., Xu Z., Yang X., 2013. Health assessment by using a benthic-index of biotic integrity in the Huntai River
576 basin. *J Beijing Normal Univ (Nat Sci)* 49, 297-303.
- 577 Liaw A., Wiener M., 2002. Classification and regression by randomForest. *R news* 2, 18-22.
- 578 Liu J., Yang W., 2012. Water sustainability for China and beyond. *Science* 337, 649-650.
- 579 Liu S., Xie G., Wang L., Cottenie K., Liu D., Wang B., 2016. Different roles of environmental variables and
580 spatial factors in structuring stream benthic diatom and macroinvertebrate in Yangtze River Delta,
581 China. *Ecological Indicators* 61, 602-611. Mandaville S.M., 2002. Benthic macroinvertebrates in
582 freshwaters_ Taxa tolerance values, metrics, and protocols, Soil & Water Conservation Society of
583 Metro Halifax.
- 584 Manfrin, A., Traversetti, L., Pilotto, F., Larsen, S. and Scalici, M., 2015. Effect of spatial scale on
585 macroinvertebrate assemblages along a Mediterranean river. *Hydrobiologia*, 765(1): 185-196.
- 586 Marchant R., Wells F., Newall P., 2000. Assessment of an ecoregion approach for classifying
587 macroinvertebrate assemblages from streams in Victoria, Australia. *J. N. Am. Benthol. Soc.* 19,
588 497-500.
- 589 Marshall, J.C., Steward, A.L. and Harch, B.D., 2006. Taxonomic Resolution and Quantification of Freshwater
590 Macroinvertebrate Samples from an Australian Dryland River: The Benefits and Costs of Using
591 Species Abundance Data. *Hydrobiologia* 572, 171-194.
- 592 Mebane Jr W.R., Sekhon J.S., 2011. Genetic optimization using derivatives: the rgenoud package for R. *J. Stat.*
593 *Softw.* 42, 1-26.
- 594 Merritt R.W., Cummins K.W., 1996. An introduction to the aquatic insects of North America, Kendall Hunt.
- 595 Milošević D., Čerba D., Szekeres J., Csányi B., Tubić B., Simić V., Paunović M., 2016. Artificial neural
596 networks as an indicator search engine: The visualization of natural and man-caused taxa variability.
597 *Ecol. Indic.* 61, 777-789.
- 598 Morse J.C., Yang L., Tian L., 1994. Aquatic insects of China useful for monitoring water quality, Hohai
599 University Press.
- 600 Mouton, A.M., De Baets, B., Goethals, P.L.M., 2010. Ecological relevance of performance criteria for species
601 distribution models. *Ecol. Model.* 221 (16), 1995–2002.
- 602 Muñoz-Mas R., Fukuda S., Vezza P., Martínez-Capel F., 2016. Comparing four methods for decision-tree
603 induction: A case study on the invasive Iberian gudgeon (*Gobio lozanoi*; Doadrio and Madeira, 2004).
604 *Ecol. Inform.* 34, 22-34.
- 605 Murtagh F., Legendre P., 2014. Ward's hierarchical agglomerative clustering method: which algorithms
606 implement Ward's criterion? *J. Classif.* 31, 274-295.
- 607 Mykrä, H., Heino, J. and Muotka, T., 2007. Scale-related patterns in the spatial and environmental
608 components of stream macroinvertebrate assemblage variation. *Global Ecol. Biogeogr.* 16, 149–159.
- 609 O'Brien A., Townsend K., Hale R., Sharley D., Pettigrove V., 2016. How is ecosystem health defined and
610 measured? A critical review of freshwater and estuarine studies. *Ecol. Indic.* 69, 722-729.
- 611 Ogbeibu A.E., Oribhabor B.J., 2002. Ecological impact of river impoundment using benthic
612 macro-invertebrates as indicators. *Water Res.* 36, 2427–2436.
- 613 Pérez-Ortiz M., Gutiérrez P.A., Hervás-Martínez C., Yao X., 2015. Graph-based approaches for over-sampling
614 in the context of ordinal regression. *IEEE Tran. Knowledge and Data Engineering* 27, 1233-1245.

- 615 Paini, D.R., Worner, S.P., Cook, D.C., De Barro, P.J., Thomas, M.B., 2010. Using a self-organizing map to
616 predict invasive species: sensitivity to data errors and a comparison with expert opinion. *J. Appl. Ecol.*
617 47, 290–298.
- 618 Paliy O., Shankar V., 2016. Application of multivariate statistical techniques in microbial ecology. *Mol. Ecol.*
619 25, 1032-1057.
- 620 Park Y.S., Céréghino R., Compin A., Lek S., 2003. Applications of artificial neural networks for patterning
621 and predicting aquatic insect species richness in running waters. *Ecol. Model.* 160, 265-280.
- 622 Park Y.-S., Song M.-Y., Park Y.-C., Oh K.-H., Cho E., Chon T.-S., 2007. Community patterns of benthic
623 macroinvertebrates collected on the national scale in Korea. *Ecol. Model.* 203, 26-33.
- 624 Qin C.Y., Zhou J., Cao Y., Zhang Y., Hughes R.M., Wang B.X., 2014. Quantitative tolerance values for
625 common stream benthic macroinvertebrates in the Yangtze River Delta, Eastern China. *Environ.*
626 *Monit. Assess.* 186, 5883-5895.
- 627 Qu X., Zhang H., Zhang M., Liu M., Yu Y., Xie Y., Peng W., 2016. Application of multiple biological indices
628 for river health assessment in northeastern China. *Annales de Limnologie - International Journal of*
629 *Limnology* 52, 75-89.
- 630 Quinlan, J.R., 1992. C4.5: programs for machine learning. Morgan Kaufmann Publishers, Inc., San Mateo,
631 CA (USA). Resh V.H., Hildrew A.G., Statzner B., Townsend C.R., 1994. Theoretical habitat templates,
632 species traits, and species richness: a synthesis of long-term ecological research on the Upper Rhone
633 River in the context of concurrently developed ecological theory. *Freshwat. Biol.* 31, 539-554.
- 634 Rosenberg D.M., Resh V.H., 1993. Freshwater biomonitoring and benthic macroinvertebrates. New York,
635 Chapman and Hall.
- 636 Rutt, G.P., Weatherley, N.S. and Ormerod, S.J., 1990. Relationships between the physicochemistry and
637 macroinvertebrates of British upland streams: the development of modelling and indicator systems for
638 predicting fauna and detecting acidity. *Freshwat. Biol.* 24(3), 463-480.
- 639 Siddig A.A.H., Ellison A.M., Ochs A., Villar-Leeman C., Lau M.K., 2016. How do ecologists select and use
640 indicator species to monitor ecological change? Insights from 14 years of publication in *Ecological*
641 *Indicators*. *Ecol. Indic.* 60, 223-230.
- 642 Song, M.-Y., Park, Y.-S., Kwak, I.-S., Woo, H. and Chon, T.-S., 2006. Characterization of benthic
643 macroinvertebrate communities in a restored stream by using self-organizing map. *Ecol. Inform.* 1(3):
644 295-305.
- 645 Soininen, J., Jamoneau, A., Rosebery, J. and Passy, S.I., 2016. Global patterns and drivers of species and trait
646 composition in diatoms. *Global Ecol. Biogeogr.* 25, 940-950.
- 647 Sroczyńska K., Claro M., Kruk A., Wojtal-Frankiewicz A., Range P., Chícharo L., 2017. Indicator
648 macroinvertebrate species in a temporary Mediterranean river: Recognition of patterns in binary
649 assemblage data with a Kohonen artificial neural network. *Ecol. Indic.* 73, 319-330.
- 650 Traversetti, L. and Scalici, M., 2014. Assessing the influence of source distance and hydroecoregion on the
651 invertebrate assemblage similarity in central Italy streams. *Knowl. Manag. Aquat. Ecosyst.* 414, 02.
- 652 Tibshirani R., Walther G., Hastie T., 2001. Estimating the number of clusters in a data set via the gap statistic.
653 *J. Roy. Stat. Soc. Ser. B. (Stat. Method.)* 63, 411-423.
- 654 Tonkin J.D., Tachamo Shah R.D., Shah D.N., Hoppeler F., Jähniq S.C., Pauls S.U., 2017. Metacommunity
655 structuring in Himalayan streams over large elevational gradients: the role of dispersal routes and
656 niche characteristics. *J. Biogeogr.* 44, 62-74.
- 657 Törnblom, J., Roberge, J.-M., and Angelstam, P., 2011. Rapid assessment of headwater stream
658 macroinvertebrate diversity: an evaluation of surrogates across a land-use gradient. *Fundam. Appl.*
659 *Limnol.* 178, 287–300.

- 660 Tsai W.P., Huang S.P., Cheng S.T., Shao K.T., Chang F.J., 2017. A data-mining framework for exploring the
661 multi-relation between fish species and water quality through self-organizing map. *Sci. Total Environ.*
662 579, 474-483. Vörösmarty C.J., McIntyre P.B., Gessner M.O., Dudgeon D., Prusevich A., Green
663 P., Glidden S., Bunn S. E., Sullivan C.A., Reidy Liermann C. & Davies P.M., 2010. Global threats to
664 human water security and river biodiversity. *Nature* 467, 555–561.
- 665 Wehrens R., Buydens L.M., 2007. Self- and super-organizing maps in R: the Kohonen package. *J. Stat. Softw.*
666 21, 1-19.
- 667 Wright, J.F., Armitage, P.D., Furse, M.T. and Moss, D., 1989. Prediction of invertebrate communities using
668 stream measurements. *River. Res. App.* 4(2), 147-155.
- 669 Wiggins G.B., 1996. Trichoptera families. An introduction to the aquatic insects of North America, 309-349.
- 670 Wu, N., Qu, Y., Guse, B., Makarevičiūtė, K., To, S., Riis, T. and Fohrer, N., 2018. Hydrological and
671 environmental variables outperform spatial factors in structuring species, trait composition, and beta
672 diversity of pelagic algae. *Ecol. Evol.* 8, 2947-2961.
- 673 Zhang M., Cai Q., Xu Y., Kong L., Tan L., Wang L., 2012a. Spatial Distribution of Macroinvertebrate
674 Community along a Longitudinal Gradient in a Eutrophic Reservoir- Bay during Different
675 Impoundment Stages, China. *Int. Rev. Hydrobiol.* 97, 169-183.
- 676 Zhang M., Cai Q., Xu Y., Wang L., Kong L., 2012b. Zonation and its influencing factors of a large subtropical
677 reservoir (Danjiangkou Reservoir in central Chian), based on macroinvertebrates. *Fresenius Environ.*
678 *Bull.* 21, 2095-2104.
- 679 Zhang Y., Guan D., Jin C., Wang A., Wu J., Yuan F., 2011. Analysis of impacts of climate variability and
680 human activity on streamflow for a river basin in northeast China. *J. Hyd.* 410, 239-247.
- 681 Zhang Y., Zhao R., Kong W., Geng S., Bentsen C.N., Qu X., 2013. Relationships between macroinvertebrate
682 communities and land use types within different riparian widths in three headwater streams of Taizi
683 River, China. *J. Freshwat. Ecol.* 28, 307-328.

Supplementary Appendix A

[Click here to download Supplementary material for on-line publication only: Appendix A.docx](#)

