1 **First, do no harm - missing data treatment to support lake ecological state**

2 **assessment**

3 **Grzegorz Chrobak [a,*], Tomasz Kowalczyk [b], Thomas B. Fischer [c], Szymon Szewrański [a],**

4 **Katarzyna Chrobak [d], Jan K. Kazak [a]**

5 [a] Institute of Spatial Management, Wroclaw University of Environmental and Life Sciences, Wroclaw, Poland,

6 [b] Institute of Environmental Development and Protection, Wroclaw University of Environmental and Life

7 Sciences, Wroclaw, Poland

8 [c] School of Environmental Sciences, Department of Geography and Planning, University of Liverpool, Liverpool,

9 United Kingdom,

10 [d] Department of Urban Design & Planning and Settlement Processes, Wroclaw University of Science and

11 Technology, Wrocław, Poland

12

13 [*] Corresponding author: Grunwaldzka St. 55, 50-357 Wroclaw, grzegorz.chrobak@upwr.edu.pl, +48 663 150 811

14

15 **Abstract:** Indicators of ecological potential of water bodies, that are associated with field measurements, are often

16 subject to data gaps. This is an obstacle to constructing reliable assessments of conditions of lakes, which can lead

17 to abandonment of assessment. Furthermore, it can lead to the use of methods, based merely on their availability.

18 In response to these problems, a systematic approach for expert-analyst interaction for missing data treatment is

19 proposed. In this context, a combination of algorithms with hierarchical clustering of results was used. A particular

20 emphasis is put on the stage of preparation and interpretation of input data and the role of an expert in the workflow

21 developed. The beneficiaries of this article are ecological data experts and analysts who work in teams to assess

22 and interpret the state of lake ecosystems, and who present the findings in reports that are used during public

23 consultations and discussions with key decision makers.

24 Keywords: ecological assessment, decision support system, missing data, lakes, water quality

25

26

27

28

29

## 1. Introduction

### 1.1 Data quality issues in ecological assessment

Since the publication of the Water Framework Directive in 2000, in te European Union (EU) management of water resources has become a priority, aiming to meet environmental objectives of water bodies (Di Quarto and Zinzani, 2021; Kallis and Butler, 2001). In this context, pro-ecosystem approaches require the use of methods that are based on a holistic understanding of dependencies in evaluation procedures, potentially leading to: 1) the emergence of innovative and genuine ecological approaches to water management practices (Gain et al., 2021; Giupponi, 2007; Poikane et al., 2015; Reis et al., 2017), and also to: (2) a rapid growth of methodologies, data and indicators produced by EU member states (Birk et al., 2013; Booty et al., 2001; Carey et al., 2021; Kelly et al., 2016; Zambelli et al., 2012). The number of approaches to assessing the ecological potential of water bodies is inextricably linked with issues of production, modeling and processing of observation and measurement data (Birk et al., 2012; Paruch et al., 2017; Posthuma et al., 2020). At each stage of the creation of environmental indicators, problems can arise related to the quality and availability of input values (Brito et al., 2020; Gobeyn et al., 2016; Lindholm et al., 2007; Matthies et al., 2007; Paruch et al., 2017). A key task and, at the same time, challenge are the intercalibration procedures that allow to obtain common reference levels for the classification of the ecological state of lakes. . Importantly, any unification of indicators requires a clear recognition of input data and the development of coherent methods for managing incomplete information (Gobeyn et al., 2016; Lahtinen et al., 2017). This can help to avoid undesirable consequences associated with ignoring unknowns.

### 1.2 Implications of missing information

The effects of a lack of data in the process of assessing ecological conditions of aquatic ecosystems can be seen at every level of data processing, including the ex-post evaluation of indicators (Yang et al., 2021; Zhang et al., 2019). Identification of the type of missing information is a critical element in the initial phase of dealing with measurement data (Little, 2021). The quantity, nature, and severity of data flaws have a direct impact on the methods that can be used to work with specific datasets. In the case of measurement sets used to assess the ecological status of lakes, deficiencies in observations often result from a type of defects, referred to as Missing At Random (MAR) (Seaman et al., 2013). In this context, there is a need to rely on substantive acceptability, as MAR is an assumption which is impossible to prove statistically (Little, 2021). Due to contingent emptiness in datasets, parameter bias can result in analyses (Schielzeth et al., 2020). How to best solve this problem depends on the assumptions made, as well as on the knowledge of the context (Koehler et al., 2017). In this context, the

60  most common consequences of mishandling gaps in data sets include; information loss, bias in statistical inference

61  or modeling, and results misinterpretation (Hossie et al., 2021; Noble and Nakagawa, 2021). Another problem

62  connected with an incomplete input dataset includes an inability to use certain data analysis methods / algorithms

63  (e.g. PCA, SVM, neural networks) (Ghannam and Techtmann, 2021). A consequence of these issues is that popular

64  methods ted to be used, such as partial deletion, interpolation, or imputation (Curley et al., 2019; Johnson et al.,

65  2021). Missing knowledge management requires informed decisions to be taken along the data analysis path

66  (Likmeta et al., 2021; Newman, 2014; Wang and Xue, 2020).

67

68  **1.3      Data imputation – ecological assessment perspective**

69  The assessment of the condition and potential of aquatic ecosystems is connected with the identification of

70  activities aimed at maintaining or improving the status of them, as required under Article 11 of the Water

71  Framework Directive. In practice, this is associated with a planning process that takes place in a 6-year cycle.

72  Responsible for them are water management boards together with the departments of boards of individual water

73  sub-regions (usually within river basins). Water administration are working together on: identifying anthropogenic

74  pressures; updating environmental objectives and protected water areas; restoring water bodies,; and setting

75  boundary values for heavily modified and artificial water bodies. An important stage is the preparation of strategic

76  environmental assessment (SEA; Mustow, 2021). At this key moment, assessors have the opportunity to influence

77  the shape of the analyzes, the interpretation of the results. Furthermore, they can apply for supplementing or

78  correcting the methodology. Comments are directed to the authors of the plan at the stage of public consultations.

79  Among other measures, indicators of the ecological status of lake ecosystems are used to obtain results that support

80  the definition of management practices. The evaluation of the structure and efficiency of surface water ecosystems

81  is known as ecological status. This demonstrates how stresses (such as pollution and habitat deterioration) have an

82  impact on specified quality components. Each surface water body has an ecological status that is assessed based

83  on biological quality components and supported by physico-chemical and hydromorphological quality elements.

84  According to the "one out, all out" approach, the element with the worst status out of all biological and supporting

85  quality factors determines the overall ecological status rating for a water body. Data used to evaluate the ecological

86  status of lakes are sets largely based on the results of field measurements. Observations are prone to errors that

87  can occur at the stage of collecting samples (Yanai et al., 2021). There is always uncertainty over results, even if

88  using different tools (Ejigu, 2021). Loss of data or a complete lack of it may result in abandoning the assessment,

89  which, in some cases, significantly reduces the pool of evaluated ecosystems. This often leads to gaps in data sets

90  that weaken results of individual measurement campaigns. Moreover, the same input data serve as components

91  necessary to construct different environmental indicators, placing additional emphasis on the validity of an

92  imputation attempt. In research on the ecological quality of ecosystems, various methods of supplementing missing

93  values are used (Muharemi et al., 2019; Said et al., 2019; Zhang and Thorburn, 2022).

94  The so-called hot deck imputation is used for handling missing data on large scale water quality indices (Ahmed

95  et al., 2021; Srebotnjak et al., 2012). Most extensively used are methods based on multiple imputation. These are

96  available for most data types (Ben Aissia et al., 2017; Betrie et al., 2016; Neri et al., 2018; Ngouna et al., 2020).

97  When faced with a high level of missingness data, machine learning techniques were adopted. These are able to

98  troubleshoot complex data issues (Irvin et al., 2021; Kim et al., 2020; Ngouna et al., 2020; Ratolojanahary et al.,

99  2019; Rodríguez et al., 2021). Furthermore, the spatial nature of the issue results in an introduction of time and

100 space variables (Koki et al., 2020; Labuzzetta et al., 2021; Liu et al., 2016; Lou and Obradovic, 2011; Sojka et al.,

101 2020; Yüksel, 2012; Zhang and Thorburn, 2021).  Research with ecological water quality indicators uses methods

102 based on a case study approach. This confirmins their effectiveness at the local scale (Bilgin and Bayraktar, 2021;

103 Liu et al., 2011; Ren et al., 2008; Sojka et al., 2019; Weerasinghe and Handapangoda, 2019). There is a noticeable

104 trend in the research indicating the need to develop methods that work well at the regional level, providing the

105 option of later intercalibration of the results (Akbar et al., 2011; Botha et al., 2020; Hu et al., 2018; Jiang et al.,

106 2017; Lepš and Šmilauer, 2006; Li et al., 2021; Luo et al., 2019). Holistic approaches facilitate macro-quality

107 management of water resources, which is important in the context of policy design and pan-regional impact

108 assessment. Moreover,  monitoring of ecological indicators and the impact of climate change on phenomena that

109 threaten the stability of ecosystems has lately been explored (Cheruvelil et al., 2017; Fazli et al., 2018; Hutjes,

110 2019; Krzeminski et al., 2019; Lizotte et al., 2014; Mankin et al., 1999; Mustajoki et al., 2004; Peters-Lidard et

111 al., 2021).

112

### 113  1.4     Research goals  structure of paper

114 The main goal of the research underpinning this paper was to present a workflow that can be used when an expert

115 group or an ecological assessor are faced with the problem of missing values in an input dataset. In this context, a

116 novel combined expert and analyst approach to ecological assessment is introduced. This approach gives experts

117 the opportunity to influence (and adjust) processes by making decisions in key nodes. A further goal is the

118 identification of possible techniques of data visualization, both with regards to raw data and analysis.

119

120   In the methodological approach, graphic representation of often complicated processes is crucial for effective

121   cooperation in an expert team. Featured data treatment schema takes the specificity of the work of experts into

122   account, dealing with various assessment objects with a different degree of data incompleteness in the assessment

123   process. Thus, there are certain cross-roads highlighted where a decision is necessary, made by a specialist or

124   requiring consultation before proceeding with the analysis (2.3. Proposed workflow). The data treatment

125   framework guides the user through the steps of pre-selecting data (3.1. Missing data identification and triage),

126   identifying and selecting imputation predictors (3.2. Predictor examination), the actual multiple data imputation

127   process using the random forest algorithm (3.3. Missing data imputation), and then introduces the step of clustering

128   similar complementary sets based on their characteristics in the context of the Ward criterion via hierarchical

129   clustering (3.4. Clustering imputation and 3.5. Data imputation results).
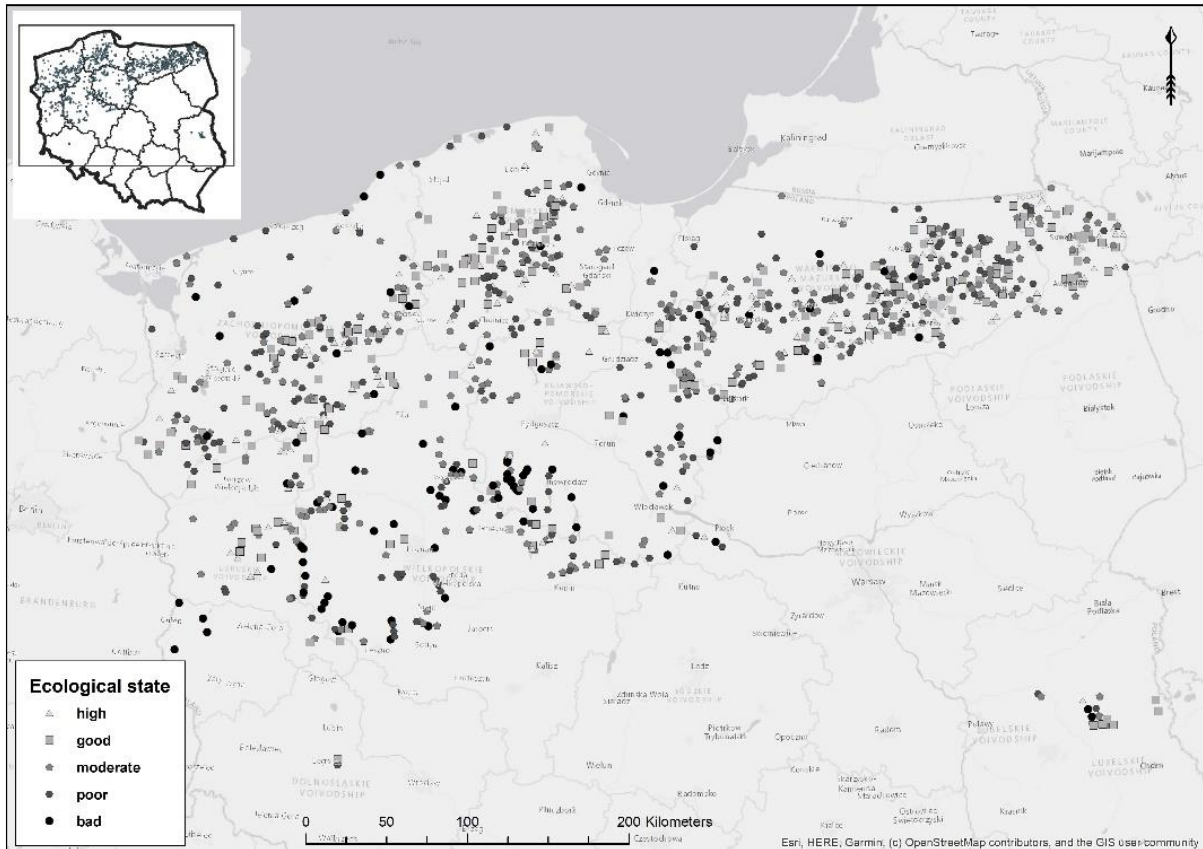
130

131   **2.       Materials and Methods**

132   **2.1      Data, software & previous research**

133   The input data used in this work come from the resources of the Chief Inspectorate of Environmental Protection

134   in Poland (Appendix A) (GIOŚ, 2015). These are measurements included in the data used to develop indicators of

135   the ecological condition of lake ecosystems. Results of the analyzes are reported to the European Commission data

136   repositories, including information on the state of water among the Member States of the European Union

137   (European Environment Agency, 2018) (Figure 1).

138   The analyzes concern a set of 499 objects for which measurements were made during the 2013-2015 measurement

139   campaign. Chlorophyll a, nitrogen, phosphorus, phytoplankton, Ecological State Macrophyte Index (ESMI),

140   Diatom Index for Lakes (IOJ), Phytoplankton Method for Polish Lakes (PMPL), visibility, and conductivity are

141   some of the measures used to determine a lake's ecological status. The basic information on data is provided in

142   Appendix B. Data were the subject of constructing a methodology aimed at improving effectiveness and

143   reproducibility of the procedure for determining ecological status indicators with the use of machine learning

144   algorithms (Chrobak et al., 2021b). In the next step, the set was used to extend the methodological approach to

145   include the use of an unsupervised tool, supporting the prioritization of lakes in the context of organizing remedial

146   measures necessary for the ecosystem to achieve environmental goals (Chrobak et al., 2021a). An important

147   element of working with data at each of these stages was the need to deal with the problem of missing information.

148   In this paper, the consequences of the lack of observations in the collection are addressed, and the missing data

149   imputation is performed and tested as a complementary solution working with workflow previously developed.

150

*Figure 1. The map presenting lakes with shapes representing resulting ecological state according to EU classes. Originally, missing data were not included in the calculations. Instead, they influenced the appropriate value of the uncertainty of the result in the tables attached to the report.*

## 2.2 Imputation and clustering techniques

In order to select the optimal technique for imputation of missing observations, the 'missingness' type of the dataset was identified (Zhou, 2020). The discovered systematic tendencies in the dataset show that missing observations can be predicted with use of other information present (see section 3.2 Predictor examination). It is due to existing correlations between fields and thanks to the knowledge of the data collection procedure that errors in measurements or deficiencies are not the result of a deliberate procedure. Thus, the missingness type was labelled as missing at random (MAR) (Bhaskaran and Smeeth, 2014). The following procedure of iterative imputation of missing values was preceded by stages (1); which involved applying the Pearson product-moment correlation method to analyze the degree and direction of data association (Russo, 2021) and (2); Principal Component Analysis (PCA) performed on the dataset with missing values to investigate uncertainty related absence of information (Husson et al., 2014).

167   Missing value imputation was done using methods of Multiple Imputation by Chained Equations (MICE) for

168   multivariate dataset cases (Zhang, 2016). The goal to replace missing values with plausible data to estimate a more

169   realistic layout dataset, which is affected only minimally by incomplete observations. Within the procedure, the

170   following steps were performed on the input dataset (Raghunathan et al., 2001):

171   Step 1. For every missing value in the dataset, random extraction is performed from non-missing data to provide

172   initial, basic imputation ($D$).

173   Step 2. The field with the least missing values ratio ($f$) is selected and transformed back to feature missing values.

174   Step 3. The $f$ is regressed as a dependent variable onto the initially imputed dataset as $f \sim D$.
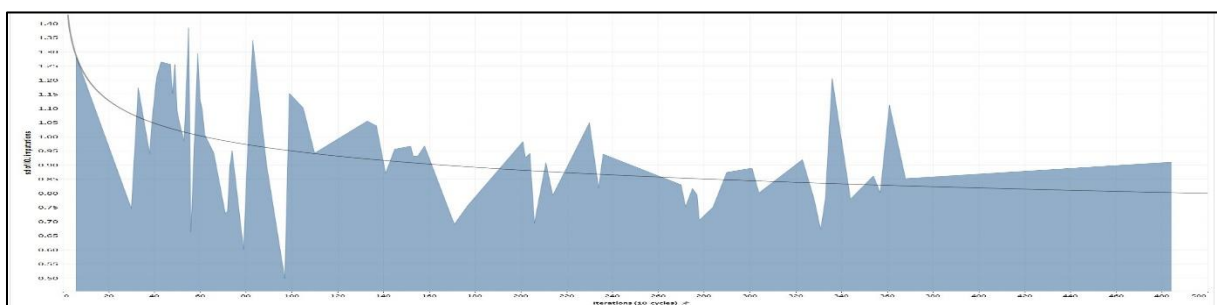
175   Step 4. The predicted values obtained from a regression model are used to fill missing data in $f$. Both, the non-

176   missing and imputed values are used once $f$ acts as an independent variable in regression modeling for the

177   following dependent variables.

178   Step 5. Steps 2-4 are repeated for each variable with missing data identified. One iteration is understood as an

179   operation of cycling through each of the variables. The cycle is finished once all missing values are replaced with

180   regression predictions that match the data relationships observed in the initial dataset.

181   The MICE model parameters selected for this research are:

182      a)   dataset: matrix (8 x 499) with missing values,

183      b)   data imputation method: random forest imputation (Shah et al., 2014),

184      c)   visit sequence: roman (left to right).

185   About 10 iteration cycles are performed in most research tasks (Gelman et al., 2011). However, at the conclusion

186   of iterative cycles, the distribution of the imputation parameters (for instance, the regression model coefficients)

187   should have converged and become stable. In order to eliminate the undesired dependency on the sequence in

188   which variables are imputed, the authors performed 50 iterations until reaching convergence (Figure 2). The

189   algorithm performance resulted in 30 imputed datasets, which were subject to a distribution-based clustering

190   process.



191
192 *Figure 2. The formation of standard deviation for successive imputation cycles led to selection of 50 initial iterations as a default parameter*
193 *for this analysis.*

194

195 For each of the fields with missing values, as a result of the data imputation method, 30 versions of the possible

196 information supplementation were obtained. The hierarchical clustering technique was used to select the

197 imputation sets that correspond to the formation of the original variable in the context of the parameters of the

198 similarity of the data distribution (Wu et al., 2009). Initially, each dataset was treated as a separate cluster in the

199 agglomerative version of the algorithm. Following that, similar clusters were merged to form larger units based on

200 predefined rules. When only one cluster emerged, the algorithm concluded that no further agglomeration is

201 possible (Murtagh and Legendre, 2014). The clustering procedure included the following steps: (Hartigan and

202 Wong, 1979):

203 Step 1. The distance matrix was computed between columns of versions of imputed columns (the original field is

204 a feature in proximity calculation, as well, with missing values allowed, but excluded from analysis) – resulting in

205 a cross-distance matrix.

206 Step 2. A cross-distance matrix was used as a dissimilarity structure for an agglomeration method to perform

207 proximity-based merging – every column was considered as an individual cluster.

208 Step 3. The clusters with similar characteristics (proximity) were merged.

209 Step 4. The cross-distance matrix was recalculated for each cluster.

210 Step 5. The steps 3-4 were repeated until a single cluster remained.

211 In the construction of the cross-distance matrix for each of the dataset fields, the form of squared Euclidean

212 distance matrix was used (Sarstedt and Mooi, 2014). The Ward's method, based on the optimal value of an

213 objective function – in this case – the minimum variance was used as a criterion for choosing a pair of clusters to

214 merge at each step (Ward, 1963). The overall within-cluster variance is reduced, using Ward's minimal variance

215 criterion (Kruskal and Black, 2012):

$$D_{1,2} = \sqrt{\frac{2 \cdot |k| \cdot |l|}{|k| + |l|}} \cdot \left\| \vec{k} - \vec{l} \right\|$$

216

217 where:

218 $D_{1,2}$ – dissimilarity between cluster 1 and cluster 2,

219 $k, l$ – observations from cluster 1 and cluster 2,

220 $\vec{k}, \vec{l}$ – centroids for clusters 1 and 2,

221 $\| . \|$ – Euclidean norm.

222 For using this approach, the pair of clusters was selected that, after merging, resulted in the least amount of total

223 within-cluster variance. A weighted squared distance between cluster centers was used to calculate this increase.

224   All clusters were singletons in the first stage (clusters containing a single point). The initial distance between

225   individual objects was proportional to the squared Euclidean distance in order to execute a recursive algorithm

226   under the objective function (Everitt, 1980) as:

227
$$D_{i,j} = \sum_{v}^{d} \left( x_{v_i} - x_{vj} \right)^2$$

228   where:

229   $D_{i,j}$ – distance between cells i and j,

230   $x_{v_i}$ – value of x variable at cell i,

231   $d$ – number of dataset dimensions.

232

233   Every feasible cluster pair is examined at each phase, and the two clusters whose merger results in the least amount

234   of information loss are combined. Ward defines information loss in terms of an error sum-of-squares criterion

235   (ESS) (Ward, 1963):

236
$$ESS = \sum_{i=1}^{n} x_i^2 - \frac{1}{n} \left( \sum_{i=1}^{n} x_i \right)^2$$

237   where:

238   $n$ – number of observations,

239   $x_i$ – the value of i-tj observation.

240   and $0$ being mean value of all the observations.
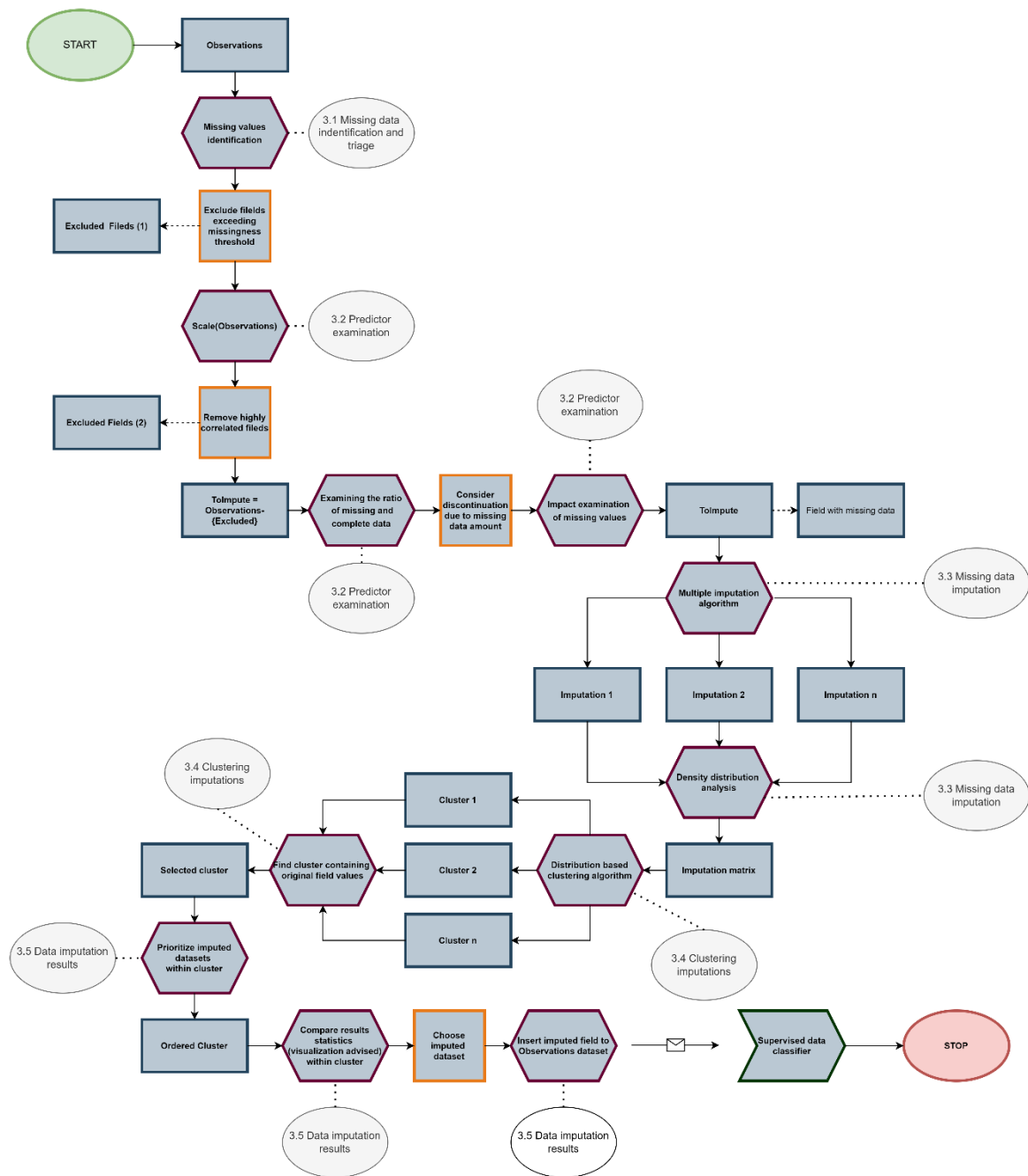
241

242   **2.3     Proposed workflow**

243   Within the block diagram of the suggested method, the proposed data analysis processes for the efficient

244   imputation of missing values have been systematized (Figure 3). The workflow was created to supplement the

245   methodology described in the authors' previous works on optimizing the assessment of the ecological state of lake

246   ecosystems (Chrobak et al., 2021a, 2021b). This enabled the evaluation solutions to be tailored to the framework

247   imposed by the Water Framework Directive, which indicates the need to conduct assessments involving expert

248   knowledge. From the technical point of view, the approach addresses cases where the analysis cannot be performed

249   effectively due to a significant number of missing observations. Thus, the decision whether to continue the analysis

250   with use of data imputation is made by the expert, who is guided by experience and aided with dataset recognition

251   led by skilled analyst. The aim is to obtain reliable premises for the implementation legitimacy of subsequent steps

252    of ecological assessment process. In the diagram of the analytical process shown below, the dataset objects (lakes

253    with measurements) appear as rectangles with blue border. Purple-outlined hexagonal blocks denote an analytical

254    or computational process that could produce new data objects or serve as the basis for decision-making. In some

255    places, these blocks are linked to orange-colored square blocks. In these cases, an expert decision is advised. Given

256    the number or severity of missing observations, the expert may decide to end the process. If the process is not

257    stopped during the data triage stage (section 3.1), the dataset is subjected to multivariate imputation, the results of

258    which are clustered. The sets of imputations proposed by the algorithm are reviewed again by an expert, who is

259    supported by the clustering results. Finally, the selected dataset with no missing values is submitted to further

260    analyses, serving as an input for the supervised classifier of the lake ecological state class. The operation of such

261    a classifier was described in the work that was published before this research (Chrobak et al., 2021b).

262

263    *Figure 3. The workflow of missing data curation and imputation. The main purpose of arranging the steps taken into a procedural form is to*
264    *systematize the methodology so that it is reproducible. Each of the process blocks enclosed by a purple frame symbolizes the action on the*
265    *data. The squares with an orange frame indicate the moment of the decision made by the analyst / expert. Each of the steps of the analysis*
266    *is discussed along with an example of implementation in the following subsections of this article.*
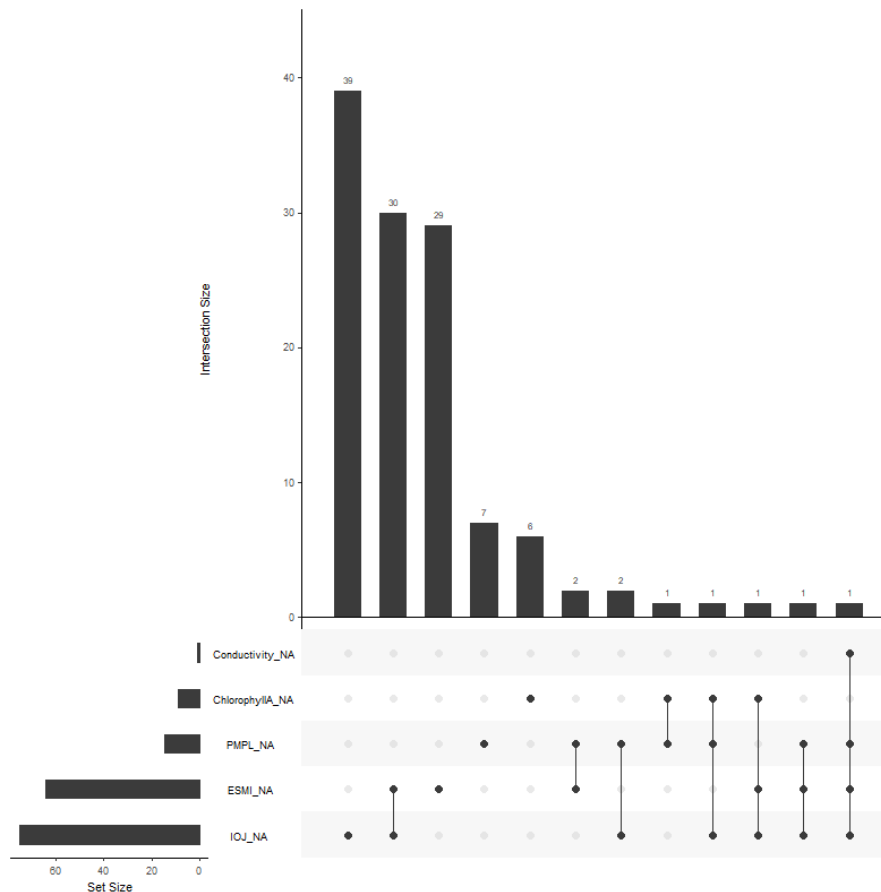
267

268    **3.    Results**

269    **3.1    Missing data identification and triage**

270    The input data of the analysis were characterized by a different number and structure of missing measurements.

271    The identification of the shortcomings started with the preparation of the chart showing the scale of the problem

272    (Figure 4). According to the adopted classification, the so-called "missing grade", deficiencies were identified in

273    5 out of 8 variables used in the process of assessing the ecological condition of lakes (Khorshidi et al., 2020). The

274    spread of NA's percentage ranged from 0.2% for the conductivity variable to 15% for the IOJ parameter.  It is

275    worth noting that the fields containing the measurement results for ESMI and IOJ together account for the existence

276    of approx. 80% of the deficiencies. Moreover, these deficiencies are characterized in the adopted methodology of

277    data triage as NotBad (missing <= 20% values), where the deficiencies in the field of PMPL, chlorophyll a, and

278    conductivity are labeled as Good (less than 5% missing). Despite the lack of fields with the Bad category, it is

279    important to remember that (1) the categories are arbitrary intervals that are largely dependent on the decision of

280    an expert who knows the data; and (2) it is possible that there are gaps in the intersection data that, when

281    accumulated at the intersections, will give a picture of real losses in the set of measurements' quality. IOJ and

282    ESMI parameters are components that strongly affect the results of ecological status classification, as indicated by

283    the PCA analysis by Chrobak et al., 2021. Leaving these fields out of the analysis may cause the final result to be

284    skewed.



285

*Figure 4. The visual representation of missing values across the dataset indicated deficiencies in five out of eight variables involved in the construction of the lake evaluation index. In addition, the number of objects (39) that have information gaps for more than one field is also indicated. The analysis did not reveal any cases where the object has gaps for each of the variables. The fields to note are IOJ and ESMI, together accounting for 80% of existing NA statements, which is a prerequisite for taking corrective action on the data.*

290

291

292

## 3.2    Predictor examination

One of the data preprocessing steps, crucial for later decisions made during data imputation, is the exploratory analysis of predictors (Braun and Oswald, 2011). The variables were subjected to the analysis of mutual linear dependencies, which allowed for an assumption of the situation earlier referred to as MAR in the context of missing observations. Strong correlations ($> = |0.5|$) were identified, e.g. for visibility-PMPL or nitrogen-chlorophyll pairs (Figure 5). Variables that are strongly associated with each other are not preferred candidates for following multiple data imputation (Ellington et al., 2015). In most situations, the selected imputation method should omit these variables during the algorithm implementation (Alice, 2015). For some instances, it is also possible for algorithms to fail or produce unreliable, overfitted results (Christie et al., 1984). Thus, highly correlated variables were excluded from the imputation process. For each of the fields with missing values a separate selection of predictors was performed, on the basis of which the calculations were continued.  As a result, in the case of the IOJ variable, each of the possible predictors was qualified (the weakest correlation concerned the relationship with nitrogen, the strongest with phosphorus). The following predictors were related to the ESMI index: phosphorus, IOJ, and conductivity. For the imputation of the field containing the PMPL measurement results, the variables: IOJ, conductivity, and phosphorus were specified. IOJ and conductivity variables were used to supplement deficiencies in the chlorophyll field. It can be seen that the IOJ variable, which is one of the imputation objects, has no correlations identified in the data set, which would rule out using any of the variables due to the concern about multicollinearity-induced bias. In that case, multiple imputation was performed using all of the available predictors. The PCA plot shows the effect of IOJ on data variability in the dataset (Figure 4).Furthermore, the plot includes variable-wise uncertainty due to the presence of empty observations (Husson et al., 2018). The analysis demonstrates that variability across different possible imputation scenarios is limited, implying that PCA results

314      may be perceived as plausible by a user (Benahmed and Houichi, 2018). It also shows the need to monitor the

315      impact of data imputation on the shaping of leading dimensions' explanatory skills (Chrobak et al., 2021b).
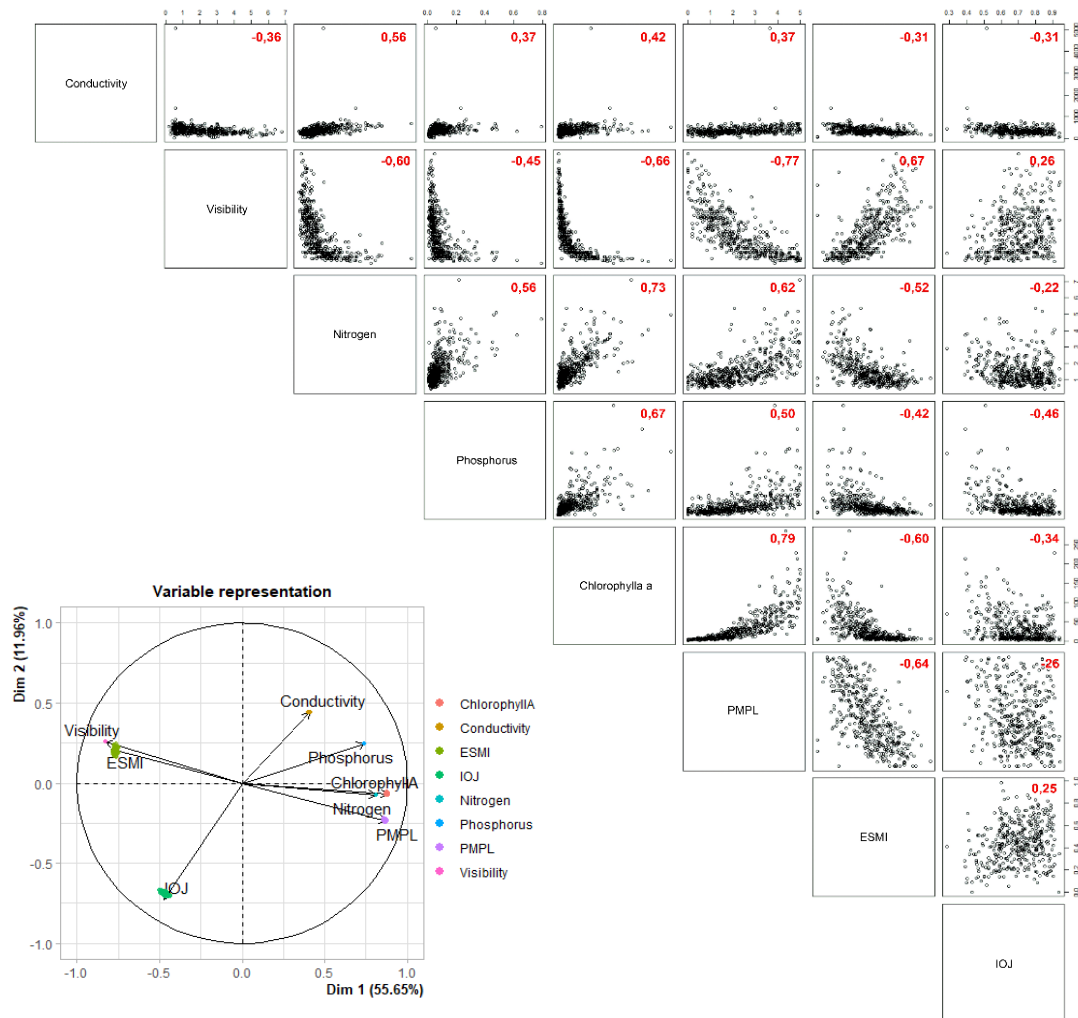
316



*Figure 5. Evaluation of predictors preceding the data imputation process. Correlation analysis using the Pearson product-moment correlation coefficient method indicated the existence of a linear relationship between some sets of observations. This information was used to select potential predictors of imputation of missing values. The results visible on the vector PCA indicated the importance of the IOJ parameter affecting the diversity of the data set, which affects the final ability of the variable to explain the differentiation in the shaping of the first coordinate variance in the reduced observation space.*
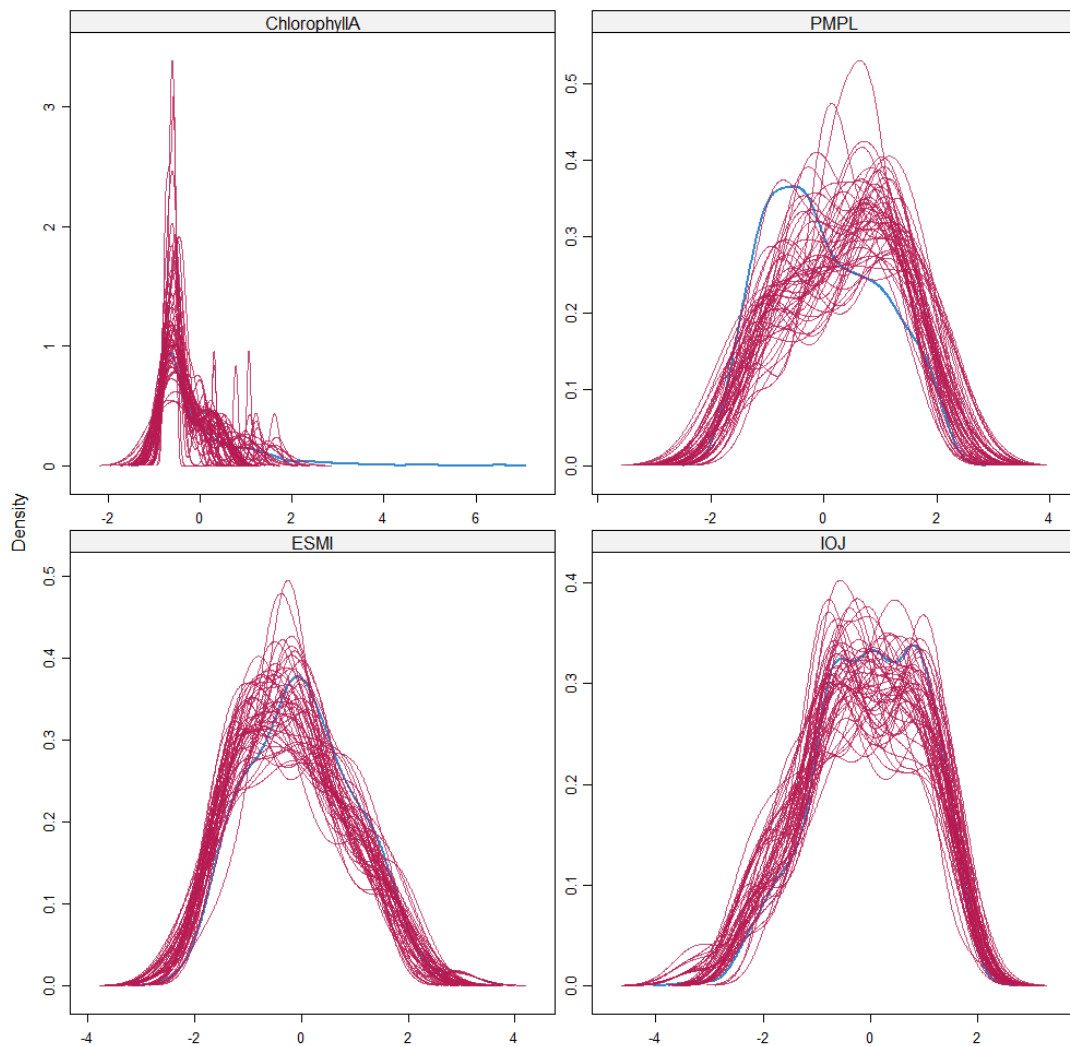
317

### 3.3     Missing data imputation

319      Missing data imputation concerned four variables (IOJ, ESMI, PMPL, chlorophyll a), for which individual sets of

320      predictors were selected in the previous stage of work. The applied method of multiple imputation is the MICE

321      approach, using the random forest algorithm (Xiao and Bulut, 2020). The method is effective when linear

322      relationships exist between variables and does not require the use of hyperparameter calibration practices. The

323      distributions were assumed for each variable and imputation was performed according to the distribution

324     characteristics obtained from the original, non-imputed dataset (Figure 6). It is not possible to know the true value

325     of intercept term due to missing data in the source field, thus introduction of a distribution assumptions was
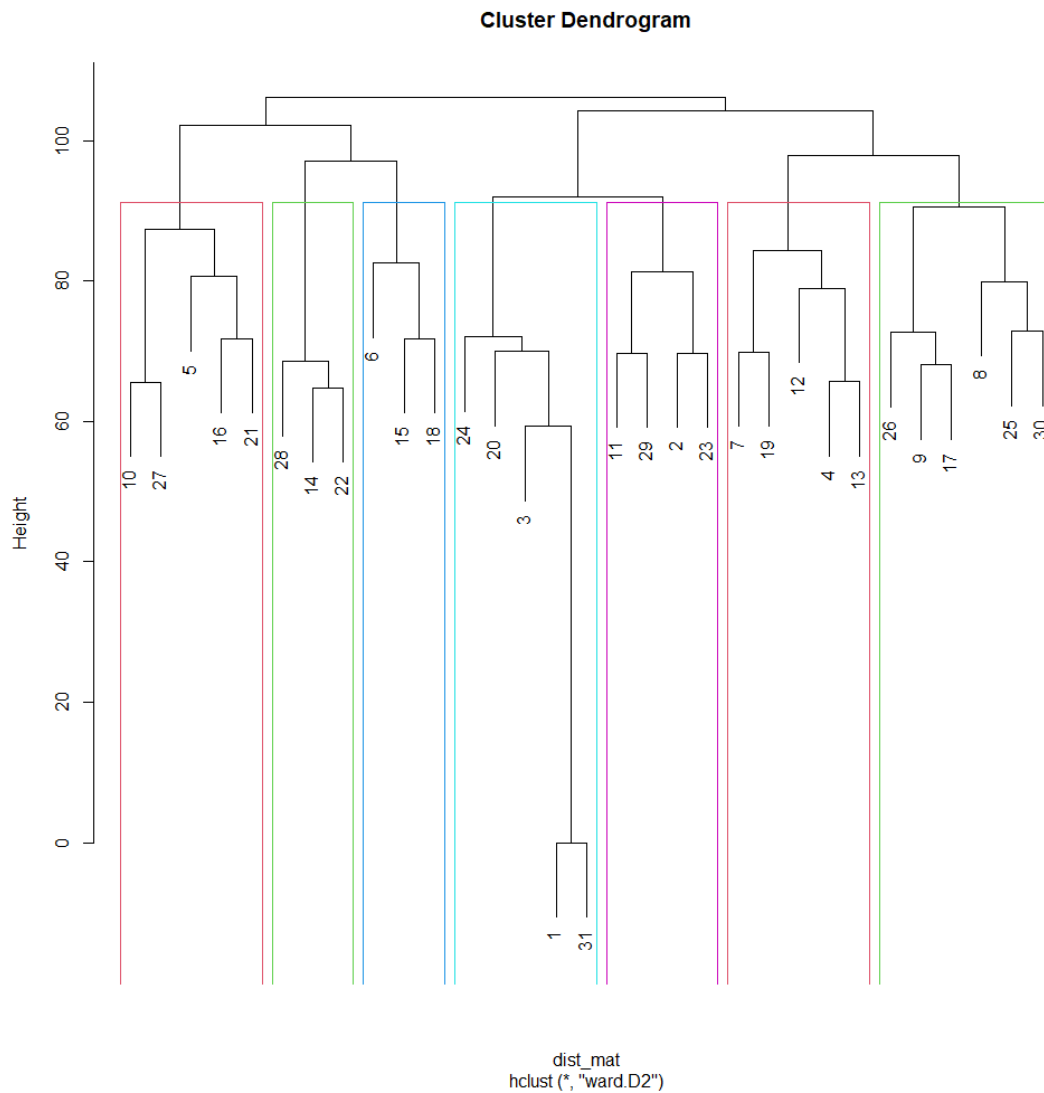
326     necessary.

327

328



329 *Figure 6. The density plots for each of imputation dataset are showed in red. The density of original field is displayed as blue line. The dataset*
330 *desired to be the best imputation option is expected to be similar in context of data density distribution. However, different results for*
331 *individual iterations of the 'rf' algorithm do not give an unambiguous fit of the optimal solution. The results also indicate the necessity of*
332 *continuous monitoring of the model results in order to avoid the use of distributions, the parameters of which (e.g. kurtosis) differ significantly*
333 *from the expected fit.*

334     The selection of the set of possible imputations was carried out for the IOJ variable as a presentation of the

335     functioning of the approach in practice. According to the results of the PCA analysis and the identification of

336     predictors, it is a variable that significantly influences the result of the final classification of the ecological state of

337     lakes in the adopted methodology. Gaps in observations of 15% make it an indicator that has the potential to be

338     the most difficult imputation, compared with e.g. chlorophyll (<5%). The plot in Figure 4 indicates the presence

339     of imputation sets that may result in an optimal but not overfitted match (Radosavljevic and Anderson, 2014).

### 3.4 Clustering imputations

According to the scheme of proceedings presented in the Materials & Methods section (Figure 3), the grouping of similar imputations was performed, using the hierarchical clustering method (Cohen-Addad et al., 2019). The aim of this part of the analysis was to use a tool that allows for fairly intuitive and quick interpretation of a given set of imputation sets, bearing in mind the possibility of carrying out more imputation iterations in specific cases or, if necessary, indicating many supplementary series (scenarios). In order to minimize the cluster-associated variance loss the Ward's method was applied, so that, at each algorithm performance step, the combination of every possible cluster pair was considered. It this case, the information loss was defined in terms of an error sum of squares criterion (ESS). Each of the leaves of the resulting dendrogram referred to the series obtained in the multiple imputation process. Sets of similar observations according to Ward's criterion were collected under the dendrogram branch (Figure 7). The height parameter of the combination displayed on the x axis indicated the similarity measure between two sets. Seven clusters within the data set were defined, using the so-called gap statistic method, which compared the total intra-cluster variation for different cluster quantities with their expected values under null reference distribution of the dataset generated with use of Monte Carlo simulations during the sampling procedure (Tibshirani et al., 2001). The original series of IOJ containing the missing observations (marked as 31) was introduced to the analysis, for which the distribution estimation was performed (Figure 7). The source set of observations was included within one cluster, marked as 4 with the sets: 1, 3, 20, and 24, which in the next steps will be considered as plausible and safe imputation options with regards to variance and distribution criteria. The distance obtained by the pair of objects 1 and 31 significantly differed from the other objects within cluster 4. Despite the fact that it indicates the best match according to the adopted criteria, it is advisable to perform a similarity test (e.g. z-statistic) in order to recognize the differences between the objects cluster (Ben-Zvi, 2004).
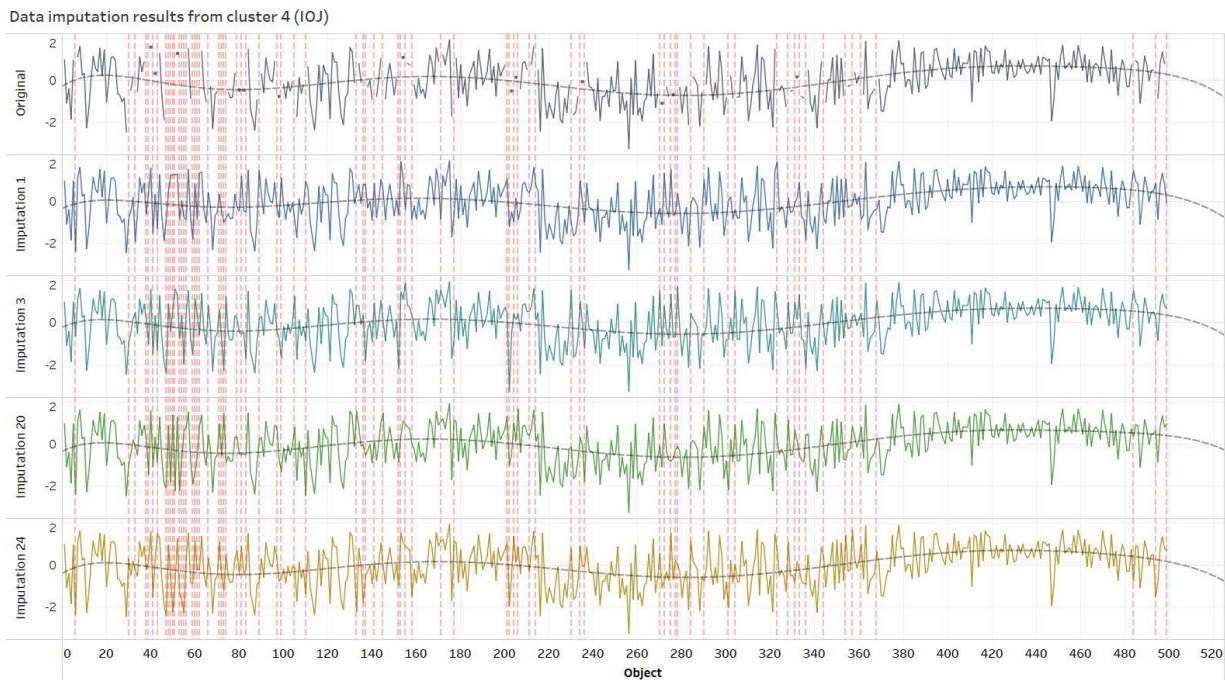
**Cluster Dendrogram**



dist_mat
hclust (*, "ward.D2")

*Figure 7. The dendrogram created for set of plausible imputation options for IOJ variable was based on bottom-up, distribution based hierarchical clustering algorithm. During consecutive model runs, seven separate clusters were distinguished. The 4th cluster (enclosed with blue frame) contains the original IOJ variable (marked with the number 31) entered for the analysis. High similarity in the context of distribution was recognized for imputation set no. 1. The next options of field completion with similar distribution are found in sets: 3, 20, and 24. The sets from the fourth cluster, in the given prioritization order, constituted a pool of plausible solutions to the problem of missing values.*

## 3.5    Data imputation results

The results of data imputation for the IOJ variable were presented in the form of sequences of corresponding series, arranged according to Lake ID in the original dataset (Figure 8). It allowed for the tracing of the imputation process within Cluster No. 4, as well as the final verification of the results, using polynomial regression on each of the retrieved series. Treating the process-aspect approach to data imputation is one of the most informative ways of presenting the process-aspect approach to data imputation. It proved to be highly informative to decision-makers and water-quality experts during the presentation of results and project-group meetings. The second way for visualizing the imputation process is to arrange lakes in order of catchment area, allowing for simultaneous
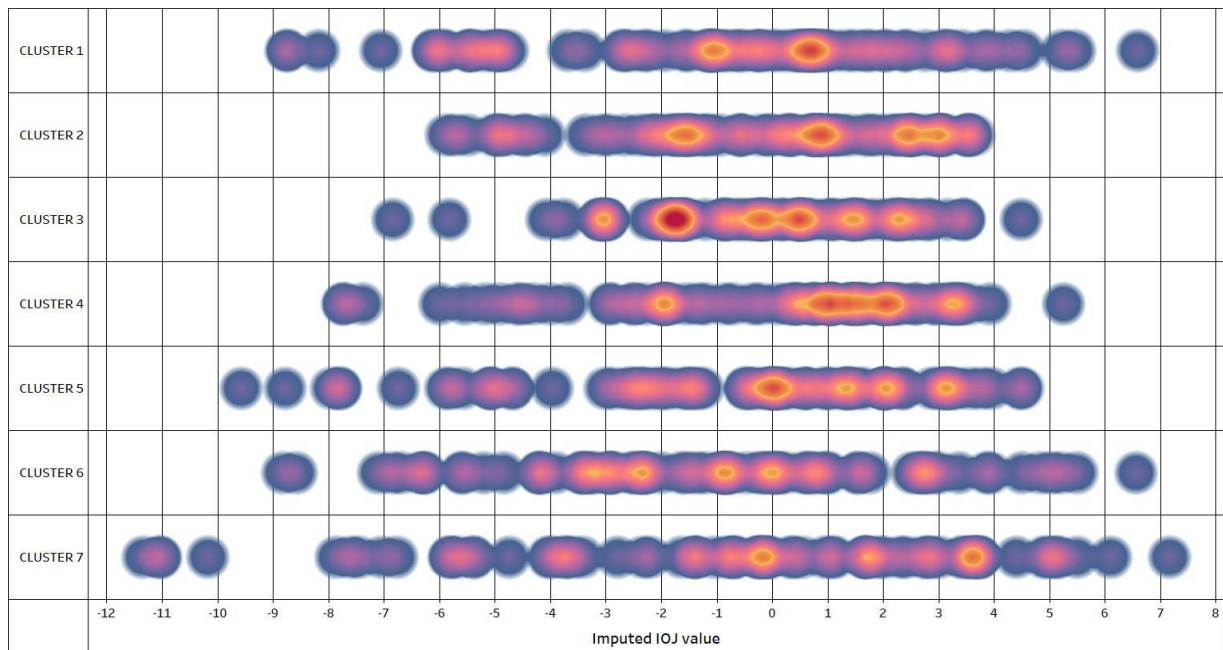
375      assessment of the degree of missing observations in spatial terms (Figure 9). The method also makes it simple to

376      partition the sets so that specialists working on specific catchment assessments can accurately evaluate the scope

377      of the problem in their work area and compare it to the situation in other task groups. Furthermore, the visualization

378      enables for cross-referencing of individual implementation outcomes across the cluster (red dashed lines) and

379      tracing of the data imputation process to identify undesired outliers generated by the method used.

380



*Data imputation results from cluster 4 (IOJ)*

381

382 *Figure 8. The chart shows a compilation of the four imputation sets (in order of priority) against the original IOJ value evolution of the*
383 *evaluation set. Dashed red lines indicate where data imputation has been performed. For each of the options within the cluster no. 4, the*
384 *statistics of the shaping of the variable allow for "safe" imputation of data and the use of the set in subsequent analyzes on the way to obtain*
385 *a reliable indicator of the ecological condition of lakes.*

386

*Figure 9. The distribution of clustered imputations shows point concentrations around values which triggered separation. The characteristics of each cluster can be distinguished during the reverse reasoning making it possible to determine entry requirements for next iterations of imputation algorithm when assessing dataset obtained in currently ongoing data collection campaign.*

## 4. Discussion

This research study underlying this paper focused on how to deal with missing-at-random data curation and imputation in the process of assessing the ecological status of lake ecosystems. The study was based on a collection of 499 lakes in Poland, with missing values detected to various degrees. A methodology was designed, based on the authors' knowledge and support in the field of expert evaluations, allowing for the imputation of data gaps to be implemented. The technique is demonstrated with an example from an authentic dataset used in the ecological status assessment with the goal of submitting the results to European Union bodies in relation to WFD obligations (Reyjol et al., 2014). The presented scheme of conduct is a complementary element to the previous works, where the stage of incomplete information management is part of an extensive algorithm of ecological assessment of lakes. The tools used in the study allowed for the selection of four ranked propositions of value imputation for the IOJ index, which was characterized by a 15% share of incomplete values. Data imputation, especially in the case of the identification of relatively large gaps in data sets (e.g.> 5%), is always associated with the risk of introducing bias into the process, which may negatively ('mis-informatively') affect the final results and their interpretation (Krueger, 2017). As a result, it's critical to understand the facts and intentionally employ the various strategies for addressing flaws. Testing the susceptibility of values to outliers is a useful practice which is part of the input data recognition stage (Jackson and Chen, 2004). Due to the emerging need to analyze lakes in a regional (or sub-basin) perspective, the future role of ecological status indicators, which will be used to make decisions at higher (supra-

409  local) levels of water resource quality management, should be taken into account (Mammides, 2020; Rivera-

410  Rondón and Catalan, 2020; Wu et al., 2021). It is connected with going beyond the locally understood and

411  evaluated indicators (Baldera et al., 2018; Kraemer et al., 2020). This is one of the challenges of the ecological

412  evaluation of aquatic ecosystems, as the management of gaps in large-scale data requires the development of

413  methods of analyzing the relationships between indicators and their components in the context of spatial and

414  temporal relationships between the objects of assessment (Kolada et al., 2014; Rossaro et al., 2012; Werner et al.,

415  2016). This may ultimately lead to the observation of a phenomenon referred to as data drift, defined as a difference

416  in variation of the data used to construct an initial assessment framework and the observations feeding the

417  assessment model in the next round of reporting (Brock and Carpenter, 2012; Koehnken et al., 2020). Taking the

418  changes in ecosystems and their internal relationships into account, especially in the era of the identified impact

419  of climate change effects, new factors may affect the variability of the ecological state of lakes over time. Thus, it

420  is critical to create a consistent procedure for detecting data drift, defining drift percentage criteria, and configuring

421  pro-active alerts so that the necessary action may be performed (Dong et al., 2018; Gupta et al., 2020). Shift may

422  manifest itself in the data at the level of their covariate shift, therefore steering with data imputation should

423  minimize the effect of completions on the distribution of the variable (Hilt et al., 2017; Martin et al., 2020).

424  The clustering approach used in this work to select plausible options is an alternative solution to the pooling stage

425  within the multiple imputation process. The classification algorithm used is, comparatively speaking, easy to

426  interpret (Cohen-Addad et al., 2019). The user also does not need to define the number of clusters a-priori.

427  However, during the process arbitrary decisions are made (distance metric, linkage criterion), which prompts the

428  expert to monitor the results in order to react quickly to noticeable errors, e.g. related to the use of mixed data

429  types (Karthikeyan et al., 2020; Zhang et al., 2013). In addition, the algorithm is sensitive to the increase in the

430  number of dimensions in the data, so an iterative analysis of successive variables requiring imputation is

431  recommended (Contreras and Murtagh, 2015). The Ward criterion used allowed for the creation of clusters based

432  on a minimal increase in degree in within cluster variance making the approach less susceptible to noise related to

433  multiple imputation results (McInnes et al., 2017).

434  Thus, the main limitations of the proposed approach are of two types. First, in terms of the algorithms used, the

435  method inherits some of their inherent limitations. In the case of the applied data imputation using the MICE

436  method with the use of random forest function, the limitations result from the need to control the results of

437  supplements. The expert should control the process so as not to allow indiscriminate acceptance of results

438  significantly deviating from the observed data. This may affect the second element of the process, which is

439     hierarchical clustering, which is sensitive to the presence of noise and outliers. This applies to both the original

440     input data and the imputation results. The second type of limitation is also related to noise, however, it concerns

441     noise generated on the side of expert judgment. The method does not allow for the complete elimination of

442     cognitive errors resulting from the participation of expert decisions characterized by their own systematic noise or

443     bias.

444     One of the indirect limitations of the whole assessment system, which this methodology also inherits, results from

445     the dependence on measurement timing and hydrological background for subsequent analyzes. As the analysts

446     work within a given time window, the measurement reports contain data that represent the ecological situation of

447     the reservoir considered to characterize it in terms of "typical state". In practice, this means that the samples of the

448     studied variables from the extremal hydrological periods (drought, flood) are included in the reports for separate

449     analyzes in the research dealing with extraordinary situations. Thus, the relationship between extraordinary

450     measures and "normal" periods is neglected. Undoubtedly, periods of ecological stress can affect the quality and

451     values of measurements, being for example a delayed ecosystem response to critical phenomena. Although striving

452     for normality of results through their early averaging and sampling in arbitrarily selected "typical" periods has a

453     mitigating effect on the variance of results, the noise generated at the early stage of the assessment is not measured

454     at present.

455     An important positive effect of the proposed imputation process is leading the data set to the smooth transition of

456     subsequent evaluation steps, where specialists often use tools that function only with non-missing input. Due to

457     the key nature of the input data management process, the transparency aspect of the analytical procedures used is

458     not without significance (Romañach et al., 2014; Zasada et al., 2017). Methods that include data visualizations as

459     inseparable elements of data processing are beneficial to supporting the ability to explain actions taken, especially

460     at the level of expert - decision makers interactions, which are critical for the often overlooked data-sense making

461     stage of ecological assessment (Arciniegas et al., 2013).

462

463     **5.      Conclusions**

464     The missing data treatment scheme presented in the paper is aimed at systematizing the value imputation stage so

465     that it is possible to perform an efficient, reproducible solution ready to implement within existing lake ecological

466     state assessment methods. The analyses included eight variables. There were gaps in the measurement data for five

467     of them. The number of missing items indicated the need to imputate data for four variables.  An approach was

468     used based on random forest multiple imputation with predictors examination. A hierarchical algorithm with a

469  Ward's variance minimization criterion was used to cluster plausible imputation solutions obtained in previous

470  step. There were seven clusters of similar additions found. Cluster 4. contained the original data set as well as four

471  completed sets that met the membership criteria. The results were presented as a dendrogram in the case of the

472  selection of clusters, as well as with the help of ordered trajectories of the shaping of the variable for the set

473  containing missing values in relation to the four possible supplementary series according to the adopted criteria.

474  The stage of missing data treatment was indicated as an integral part of the process of assessing the ecological

475  condition of lakes, influencing the selection of modeling and classification methods in subsequent stages of

476  analyzes related to the proper ecological assessment and prioritization of ecosystems in terms of the selection of

477  remedial solutions. The authors note the positive impact of methodological and visual communication on the

478  experts-analyst-decision maker line, which should be carried out with the transparency of the process (Moallemi

479  et al., 2020). This can be facilitated, for example, with the use of available data visualization techniques. This

480  research concludes the three-step approach to lake ecological assessment, which now consists of 1) data

481  preprocessing and missing values treatment, 2) model-based assessment, and 3) lake prioritization for remedial

482  purposes. Taking into account the holistic view of the research results, the proposed solutions are aimed at

483  systematization of the process of supplementing gaps in data on measurements, in contrast to the previous omission

484  of this issue in the reports on the assessment of the ecological state of lakes. The role of the expert limnologist was

485  also unclear in the course the analyzes. As a result, some lakes were only assessed by experts, while others using

486  analytical approaches. Some of the assessments were carried over from previous measurement campaigns. This

487  resulted in a conflict of results in the event that the lake apparently did not achieve environmental objectives,

488  despite the implemented remedial measures. Thus, a certain kind of data-result asymmetry occurred. The proposed

489  fragment of the methodology was therefore aimed at organizing the assessment process by: 1) defining the role of

490  an expert in the course of analyzes, 2) introducing a consistent methodology of data pre-processing, which will be

491  passed to expert judgment only in the next steps, 3) enabling the use of effective algorithms in the assessment,

492  which are sensitive to data deficiencies (e.g. kSVM or PCA ), and 4) enabling the preview of the entire assessment

493  process so that it can be corrected or further improved in the future. With reference to the results of the next

494  campaign to assess the ecological status of waters, future research should focus on assessing the scale of the

495  phenomenon of ecological data drift, which, based on the observed climate change, anthropological pressure and

496  loss of biodiversity, may have a significant impact on the broad concept indicator construction for lake water

497  ecological assessment.

498

## 6.  Software and data availability

The research was conducted with use of software providing: data visualization (*Tableau 2021.1.1, https://www.tableau.com/*), data modelling (*R 4.0.5 via RStudio 1.4.1106 „Tiger Daylily", https://www.r-project.org/, https://rstudio.com/*), and algorithm development (*draw.io 15.9.1*, https://www.diagrams.net/). Appendix B contains an R language script that converts all of the analysis procedures in this paper into an executable, reproducible workflow. The materials for this work are available from the HydroSource platform: https://www.hydroshare.org/resource/ebec024018be4c2ba04cbfa85bb14d8e/ in the repository titled "LakeEcoMissingData". Accessed as Resource: a) R-code for data preprocessing, imputation and clustering as "LakesMissingRcode.R", b) XML file of featured workflow schema as "LakeMissingWorkFlow", c) CSV file containing raw measurement results treated as input to this analysis, d) a set of results of the statistical analysis of the variables involved in the study.

## 7.  Literature

Ahmed, M., Mumtaz, R., Zaidi, S.M.H., 2021. Analysis of water quality indices and machine learning techniques for rating water pollution: A case study of Rawal Dam, Pakistan. Water Supply. https://doi.org/10.2166/ws.2021.082

Akbar, T.A., Hassan, Q.K., Achari, G., 2011. A Methodology for Clustering Lakes in Alberta on the basis of Water Quality Parameters. Clean - Soil, Air, Water. https://doi.org/10.1002/clen.201100050

Alice, M., 2015. Imputing missing data with R; MICE package | R-bloggers. R-bloggers.

Arciniegas, G., Janssen, R., Rietveld, P., 2013. Effectiveness of collaborative map-based decision support tools: Results of an experiment. Environ. Model. Softw. https://doi.org/10.1016/j.envsoft.2012.02.021

Baldera, A., Hanson, D.A., Kraft, B., 2018. Selecting indicators to monitor outcomes across projects and multiple restoration programs in the Gulf of Mexico. Ecol. Indic. https://doi.org/10.1016/j.ecolind.2018.01.025

Ben-Zvi, D., 2004. Reasoning about variability in comparing distributions. Stat. Educ. Res. J.

Ben Aissia, M.A., Chebana, F., Ouarda, T.B.M.J., 2017. Multivariate missing data in hydrology – Review and applications. Adv. Water Resour. https://doi.org/10.1016/j.advwatres.2017.10.002

Benahmed, L., Houichi, L., 2018. The effect of simple imputations based on four variants of PCA methods on the quantiles of annual rainfall data. Environ. Monit. Assess. https://doi.org/10.1007/s10661-018-6913-y

Betrie, G.D., Sadiq, R., Tesfamariam, S., Morin, K.A., 2016. On the Issue of Incomplete and Missing Water-

529      Quality Data in Mine Site Databases: Comparing Three Imputation Methods. Mine Water Environ.

530          https://doi.org/10.1007/s10230-014-0322-4

531      Bhaskaran, K., Smeeth, L., 2014. What is the difference between missing completely at random and missing at

532          random? Int. J. Epidemiol. https://doi.org/10.1093/ije/dyu080

533      Bilgin, A., Bayraktar, H.D., 2021. Assessment of lake water quality using multivariate statistical techniques and

534          chlorophyll-nutrient relationships: a case study of the Göksu Lake. Arab. J. Geosci.

535          https://doi.org/10.1007/s12517-021-06871-4

536      Birk, S., Bonne, W., Borja, A., Brucet, S., Courrat, A., Poikane, S., Solimini, A., Van De Bund, W., Zampoukas,

537          N., Hering, D., 2012. Three hundred ways to assess Europe's surface waters: An almost complete

538          overview of biological methods to implement the Water Framework Directive. Ecol. Indic.

539          https://doi.org/10.1016/j.ecolind.2011.10.009

540      Birk, S., Willby, N.J., Kelly, M.G., Bonne, W., Borja, A., Poikane, S., van de Bund, W., 2013. Intercalibrating

541          classifications of ecological status: Europe's quest for common management objectives for aquatic

542          ecosystems. Sci. Total Environ. https://doi.org/10.1016/j.scitotenv.2013.03.037

543      Booty, W.G., Lam, D.C.L., Wong, I.W.S., Siconolfi, P., 2001. Design and implementation of an environmental

544          decision support system. Environ. Model. Softw. https://doi.org/10.1016/S1364-8152(01)00016-0

545      Botha, E.J., Anstee, J.M., Sagar, S., Lehmann, E., Medeiros, T.A.G., 2020. Classification of Australian

546          waterbodies across a wide range of optical water types. Remote Sens. https://doi.org/10.3390/RS12183018

547      Braun, M.T., Oswald, F.L., 2011. Exploratory regression analysis: A tool for selecting models and determining

548          predictor importance. Behav. Res. Methods. https://doi.org/10.3758/s13428-010-0046-8

549      Brito, A.C., Garrido-Amador, P., Gameiro, C., Nogueira, M., Moita, M.T., Cabrita, M.T., 2020. Integrating in

550          situ and ocean color data to evaluate ecological quality under the water framework directive. Water

551          (Switzerland). https://doi.org/10.3390/w12123443

552      Brock, W.A., Carpenter, S.R., 2012. Early Warnings of Regime Shift When the Ecosystem Structure Is

553          Unknown. PLoS One. https://doi.org/10.1371/journal.pone.0045586

554      Carey, C.C., Woelmer, W.M., Lofton, M.E., Figueiredo, R.J., Bookout, B.J., Corrigan, R.S., Daneshmand, V.,

555          Hounshell, A.G., Howard, D.W., Lewis, A.S.L., McClure, R.P., Wander, H.L., Ward, N.K., Thomas, R.Q.,

556          2021. Advancing lake and reservoir water quality management with near-term, iterative ecological

557          forecasting. Inl. Waters. https://doi.org/10.1080/20442041.2020.1816421

558      Cheruvelil, K.S., Yuan, S., Webster, K.E., Tan, P.N., Lapierre, J.F., Collins, S.M., Fergus, C.E., Scott, C.E.,

559       Henry, E.N., Soranno, P.A., Filstrup, C.T., Wagner, T., 2017. Creating multithemed ecological regions for

560            macroscale ecology: Testing a flexible, repeatable, and accessible clustering method. Ecol. Evol.

561            https://doi.org/10.1002/ece3.2884

562    Christie, A.A., Kennelley, M.D., William King, J., Schaefer, T.F., 1984. Testing for incremental information

563            content in the presence of collinearity. J. Account. Econ. https://doi.org/10.1016/0165-4101(84)90025-9

564    Chrobak, G., Kowalczyk, T., Fischer, T.B., Chrobak, K., Szewrański, S., Kazak, J.K., 2021a. Combining

565            indicators for better decisions – Algorithms vs experts on lakes ecological status assessment. Ecol. Indic.

566            https://doi.org/10.1016/j.ecolind.2021.108318

567    Chrobak, G., Kowalczyk, T., Fischer, T.B., Szewrański, S., Chrobak, K., Kazak, J.K., 2021b. Ecological state

568            evaluation of lake ecosystems revisited: Latent variables with kSVM algorithm approach for assessment

569            automatization and data comprehension. Ecol. Indic. https://doi.org/10.1016/j.ecolind.2021.107567

570    Cohen-Addad, V., Kanade, V., Mallmann-Trenn, F., Mathieu, C., 2019. Hierarchical clustering: Objective

571            functions and algorithms. J. ACM. https://doi.org/10.1145/3321386

572    Contreras, P., Murtagh, F., 2015. Hierarchical clustering, in: Handbook of Cluster Analysis.

573            https://doi.org/10.1201/b19706

574    Curley, C., Krause, R.M., Feiock, R., Hawkins, C. V., 2019. Dealing with Missing Data: A Comparative

575            Exploration of Approaches Using the Integrated City Sustainability Database. Urban Aff. Rev.

576            https://doi.org/10.1177/1078087417726394

577    Di Quarto, F., Zinzani, A., 2021. European environmental governance and the post-ecology perspective: a

578            critical analysis of the Water Framework Directive. GeoJournal. https://doi.org/10.1007/s10708-021-

579            10402-9

580    Dong, F., Zhang, G., Lu, J., Li, K., 2018. Fuzzy competence model drift detection for data-driven decision

581            support systems. Knowledge-Based Syst. https://doi.org/10.1016/j.knosys.2017.08.018

582    Ejigu, M.T., 2021. Overview of water quality modeling. Cogent Eng.

583            https://doi.org/10.1080/23311916.2021.1891711

584    Ellington, E.H., Bastille-Rousseau, G., Austin, C., Landolt, K.N., Pond, B.A., Rees, E.E., Robar, N., Murray,

585            D.L., 2015. Using multiple imputation to estimate missing data in meta-regression. Methods Ecol. Evol.

586            https://doi.org/10.1111/2041-210X.12322

587    Europe Environment Agency, 2018. Ecological status of surface water bodies [WWW Document]. Eur. Environ.

588            Inf. Obs. Netw.

589    Everitt, B., 1980. Cluster analysis. Qual. Quant. https://doi.org/10.1007/BF00154794

590    Fazli, B., Shafie, A., Mohamed, A., Mohamad, M.F., Yahaya, N.K.E.M., Noordin, N., 2018. Development of

591        spatial similarity-based modelling to improve integrated lake water quality management in Malaysia.

592        Lakes Reserv. Res. Manag. https://doi.org/10.1111/lre.12204

593    G.-Tóth, L., Poikane, S., Penning, W.E., Free, G., Mäemets, H., Kolada, A., Hanganu, J., 2008. First steps in the

594        Central-Baltic intercalibration exercise on lake macrophytes: Where do we start? Aquat. Ecol.

595        https://doi.org/10.1007/s10452-008-9184-9

596    Gain, A.K., Hossain, S., Benson, D., Di Baldassarre, G., Giupponi, C., Huq, N., 2021. Social-ecological system

597        approaches for water resources management. Int. J. Sustain. Dev. World Ecol.

598        https://doi.org/10.1080/13504509.2020.1780647

599    Gelman, A., Levy, M.A., Abayomi, K., 2011. Diagnostics for Multivariate Imputations. SSRN Electron. J.

600        https://doi.org/10.2139/ssrn.1010415

601    Ghannam, R.B., Techtmann, S.M., 2021. Machine learning applications in microbial ecology, human

602        microbiome studies, and environmental monitoring. Comput. Struct. Biotechnol. J.

603        https://doi.org/10.1016/j.csbj.2021.01.028

604    GIOŚ, 2015. Bank danych pomiarowych [WWW Document]. URL https://powietrze.gios.gov.pl/pjp/archives

605    Giupponi, C., 2007. Decision Support Systems for implementing the European Water Framework Directive: The

606        MULINO approach. Environ. Model. Softw. https://doi.org/10.1016/j.envsoft.2005.07.024

607    Gobeyn, S., Bennetsen, E., Van Echelpoel, W., Everaert, G., Goethals, P.L.M., 2016. Impact of abundance data

608        errors on the uncertainty of an ecological water quality assessment index. Ecol. Indic.

609        https://doi.org/10.1016/j.ecolind.2015.07.031

610    Gupta, O., Goyal, N., Anand, D., Kadry, S., Nam, Y., Singh, A., 2020. Underwater Networked Wireless Sensor

611        Data Collection for Computational Intelligence Techniques: Issues, Challenges, and Approaches. IEEE

612        Access. https://doi.org/10.1109/ACCESS.2020.3007502

613    Hartigan, J.A., Wong, M.A., 1979. Algorithm AS 136: A K-Means Clustering Algorithm. Appl. Stat.

614        https://doi.org/10.2307/2346830

615    Hilt, S., Brothers, S., Jeppesen, E., Veraart, A.J., Kosten, S., 2017. Translating Regime Shifts in Shallow Lakes

616        into Changes in Ecosystem Functions and Services. Bioscience. https://doi.org/10.1093/biosci/bix106

617    Hossie, T.J., Gobin, J., Murray, D.L., 2021. Confronting Missing Ecological Data in the Age of Pandemic

618        Lockdown. Front. Ecol. Evol. https://doi.org/10.3389/fevo.2021.669477

619 Hu, Y., Peng, J., Liu, Y., Tian, L., 2018. Integrating ecosystem services trade-offs with paddy land-to-dry land

620     decisions: A scenario approach in Erhai Lake Basin, southwest China. Sci. Total Environ.

621     https://doi.org/10.1016/j.scitotenv.2017.12.340

622 Husson, F., Josse, J., Le, S., Mazet, J., 2018. FactoMineR: multivariate exploratory data analysis and data

623     mining. J. Stat. Softw.

624 Husson, F., Josse, J., Le, S., Mazet, J., 2014. Multivariate exploratory data analysis and data mining with R. R

625     Packag. version 1.26.

626 Hutjes, R., 2019. Service for Water Indicators in Climate Change Adaptation (SWICCA) [WWW Document].

627     Web page.

628 Irvin, J., Zhou, S., McNicol, G., Lu, F., Liu, V., Fluet-Chouinard, E., Ouyang, Z., Knox, S.H., Lucas-Moffat, A.,

629     Trotta, C., Papale, D., Vitale, D., Mammarella, I., Alekseychik, P., Aurela, M., Avati, A., Baldocchi, D.,

630     Bansal, S., Bohrer, G., Campbell, D.I., Chen, J., Chu, H., Dalmagro, H.J., Delwiche, K.B., Desai, A.R.,

631     Euskirchen, E., Feron, S., Goeckede, M., Heimann, M., Helbig, M., Helfter, C., Hemes, K.S., Hirano, T.,

632     Iwata, H., Jurasinski, G., Kalhori, A., Kondrich, A., Lai, D.Y., Lohila, A., Malhotra, A., Merbold, L.,

633     Mitra, B., Ng, A., Nilsson, M.B., Noormets, A., Peichl, M., Rey-Sanchez, A.C., Richardson, A.D., Runkle,

634     B.R., Schäfer, K.V., Sonnentag, O., Stuart-Haëntjens, E., Sturtevant, C., Ueyama, M., Valach, A.C.,

635     Vargas, R., Vourlitis, G.L., Ward, E.J., Wong, G.X., Zona, D., Alberto, M.C.R., Billesbach, D.P., Celis,

636     G., Dolman, H., Friborg, T., Fuchs, K., Gogo, S., Gondwe, M.J., Goodrich, J.P., Gottschalk, P., Hörtnagl,

637     L., Jacotot, A., Koebsch, F., Kasak, K., Maier, R., Morin, T.H., Nemitz, E., Oechel, W.C., Oikawa, P.Y.,

638     Ono, K., Sachs, T., Sakabe, A., Schuur, E.A., Shortt, R., Sullivan, R.C., Szutu, D.J., Tuittila, E.S.,

639     Varlagin, A., Verfaillie, J.G., Wille, C., Windham-Myers, L., Poulter, B., Jackson, R.B., 2021. Gap-filling

640     eddy covariance methane fluxes: Comparison of machine learning model predictions and uncertainties at

641     FLUXNET-CH4 wetlands. Agric. For. Meteorol. https://doi.org/10.1016/j.agrformet.2021.108528

642 Jackson, D.A., Chen, Y., 2004. Robust principal component analysis and outlier detection with ecological data.

643     Environmetrics. https://doi.org/10.1002/env.628

644 Jannicke Moe, S., Schartau, A.K., Bækken, T., Mcfarland, B., 2010. Assessing macroinvertebrate metrics for

645     classifying acidified rivers across northern Europe. Freshw. Biol. https://doi.org/10.1111/j.1365-

646     2427.2010.02413.x

647 Jiang, Q., Liang, Z., Zhao, L., Li, Y., Wu, S., Liu, Y., 2017. Integrated PCA-BN Approach for Identifying the

648     Water Quality Response Patterns for Lakes in Yunnan Plateau. Beijing Daxue Xuebao (Ziran Kexue

649      Ban)/Acta Sci. Nat. Univ. Pekin. https://doi.org/10.13209/j.0479-8023.2017.113

650    Johnson, T.F., Isaac, N.J.B., Paviolo, A., González-Suárez, M., 2021. Handling missing values in trait data.

651      Glob. Ecol. Biogeogr. https://doi.org/10.1111/geb.13185

652    Kallis, G., Butler, D., 2001. The EU water framework directive: Measures and implications. Water Policy.

653      https://doi.org/10.1016/S1366-7017(01)00007-1

654    Karthikeyan, B., George, D.J., Manikandan, G., Thomas, T., 2020. A comparative study on k-means clustering

655      and agglomerative hierarchical clustering. Int. J. Emerg. Trends Eng. Res.

656      https://doi.org/10.30534/ijeter/2020/20852020

657    Kelly, M.G., Birk, S., Willby, N.J., Denys, L., Drakare, S., Kahlert, M., Karjalainen, S.M., Marchetto, A., Pitt,

658      J.A., Urbanič, G., Poikane, S., 2016. Redundancy in the ecological assessment of lakes: Are

659      phytoplankton, macrophytes and phytobenthos all necessary? Sci. Total Environ.

660      https://doi.org/10.1016/j.scitotenv.2016.02.024

661    Khorshidi, H.A., Kirley, M., Aickelin, U., 2020. Machine learning with incomplete datasets using multi-

662      objective optimization models, in: Proceedings of the International Joint Conference on Neural Networks.

663      https://doi.org/10.1109/IJCNN48605.2020.9206742

664    Kim, Z., Jeong, H., Shin, S., Jung, J., Kim, J.H., Ki, S.J., 2020. Characterizing water quality and quantity profiles

665      with poor quality datin a machine learning algorithm. Desalin. Water Treat.

666      https://doi.org/10.5004/dwt.2020.25481

667    Kindsvater, H.K., Dulvy, N.K., Horswill, C., Juan-Jordá, M.J., Mangel, M., Matthiopoulos, J., 2018.

668      Overcoming the Data Crisis in Biodiversity Conservation. Trends Ecol. Evol.

669      https://doi.org/10.1016/j.tree.2018.06.004

670    Koehler, M., Bogatu, A., Civili, C., Konstantinou, N., Abel, E., Fernandes, A.A.A., Keane, J., Libkin, L., Paton,

671      N.W., 2017. Data context informed data wrangling, in: Proceedings - 2017 IEEE International Conference

672      on Big Data, Big Data 2017. https://doi.org/10.1109/BigData.2017.8258015

673    Koehnken, L., Rintoul, M.S., Goichot, M., Tickner, D., Loftus, A.C., Acreman, M.C., 2020. Impacts of riverine

674      sand mining on freshwater ecosystems: A review of the scientific evidence and guidance for future

675      research. River Res. Appl. https://doi.org/10.1002/rra.3586

676    Koki, I.B., Low, K.H., Zain, S.M., Juahir, H., Bayero, A.S., Azid, A., Zali, M.A., 2020. Spatial variability in

677      surface water quality of lakes and ex-mining ponds in malacca, malaysia: The geochemical influence.

678      Desalin. Water Treat. https://doi.org/10.5004/dwt.2020.25982

679    Kolada, A., Willby, N., Dudley, B., Nõges, P., Søndergaard, M., Hellsten, S., Mjelde, M., Penning, E., Van

680        Geest, G., Bertrin, V., Ecke, F., Mäemets, H., Karus, K., 2014. The applicability of macrophyte

681        compositional metrics for assessing eutrophication in European lakes. Ecol. Indic.

682        https://doi.org/10.1016/j.ecolind.2014.04.049

683    Kraemer, S.A., Barbosa da Costa, N., Shapiro, B.J., Fradette, M., Huot, Y., Walsh, D.A., 2020. A large-scale

684        assessment of lakes reveals a pervasive signal of land use on bacterial communities. ISME J.

685        https://doi.org/10.1038/s41396-020-0733-0

686    Krueger, T., 2017. Bayesian inference of uncertainty in freshwater quality caused by low-resolution monitoring.

687        Water Res. https://doi.org/10.1016/j.watres.2017.02.061

688    Kruskal, J.B., Black, P., 2012. Ward's hierarchical agglomerative clustering method: Which algorithms

689        implement Ward's criterion? J. Classif.

690    Krzeminski, P., Tomei, M.C., Karaolia, P., Langenhoff, A., Almeida, C.M.R., Felis, E., Gritten, F., Andersen,

691        H.R., Fernandes, T., Manaia, C.M., Rizzo, L., Fatta-Kassinos, D., 2019. Performance of secondary

692        wastewater treatment methods for the removal of contaminants of emerging concern implicated in crop

693        uptake and antibiotic resistance spread: A review. Sci. Total Environ.

694        https://doi.org/10.1016/j.scitotenv.2018.08.130

695    Labuzzetta, C., Zhu, Z., Chang, X., Zhou, Y., 2021. A submonthly surface water classification framework via

696        gap-fill imputation and random forest classifiers of landsat imagery. Remote Sens.

697        https://doi.org/10.3390/rs13091742

698    Lahtinen, T.J., Hämäläinen, R.P., Liesiö, J., 2017. Portfolio decision analysis methods in environmental decision

699        making. Environ. Model. Softw. https://doi.org/10.1016/j.envsoft.2017.04.001

700    Lazaridou, M., Ntislidou, C., Karaouzas, I., Skoulikidis, N., Birk, S., 2018. Harmonization of the assessment

701        method for classifying the ecological quality status of very large Greek rivers. Knowl. Manag. Aquat.

702        Ecosyst. https://doi.org/10.1051/kmae/2018038

703    Lepš, J., Šmilauer, P., 2006. Multivariate Analysis of Ecological Data. Bull. Ecol. Soc. Am.

704        https://doi.org/10.1890/0012-9623(2006)87[193:maoed]2.0.co;2

705    Li, J., Tian, L., Wang, Y., Jin, S., Li, T., Hou, X., 2021. Optimal sampling strategy of water quality monitoring at

706        high dynamic lakes: A remote sensing and spatial simulated annealing integrated approach. Sci. Total

707        Environ. https://doi.org/10.1016/j.scitotenv.2021.146113

708    Likmeta, A., Metelli, A.M., Ramponi, G., Tirinzoni, A., Giuliani, M., Restelli, M., 2021. Dealing with multiple

709    experts and non-stationarity in inverse reinforcement learning: an application to real-life problems. Mach.

710    Learn. https://doi.org/10.1007/s10994-020-05939-8

711    Lindholm, O., Greatorex, J.M., Paruch, A.M., 2007. Comparison of methods for calculation of sustainability

712    indices for alternative sewerage systems-Theoretical and practical considerations. Ecol. Indic.

713    https://doi.org/10.1016/j.ecolind.2005.10.002

714    Little, R.J., 2021. Missing data assumptions. Annu. Rev. Stat. Its Appl. https://doi.org/10.1146/annurev-

715    statistics-040720-031104

716    Liu, J., Liu, Q., Yang, H., 2016. Assessing water scarcity by simultaneously considering environmental flow

717    requirements, water quantity, and water quality. Ecol. Indic. https://doi.org/10.1016/j.ecolind.2015.07.019

718    Liu, W.C., Yu, H.L., Chung, C.E., 2011. Assessment of water quality in a subtropical alpine lake using

719    multivariate statistical techniques and geostatistical mapping: a case study. Int. J. Environ. Res. Public

720    Health. https://doi.org/10.3390/ijerph8041126

721    Lizotte, R.E., Knight, S.S., Locke, M.A., Bingner, R.L., 2014. Influence of integrated watershed-scale

722    agricultural conservation practices on lake water quality. J. Soil Water Conserv.

723    https://doi.org/10.2489/jswc.69.2.160

724    Lou, Q., Obradovic, Z., 2011. Modeling multivariate spatio-temporal remote sensing data with large gaps, in:

725    IJCAI International Joint Conference on Artificial Intelligence. https://doi.org/10.5591/978-1-57735-516-

726    8/IJCAI11-287

727    Luo, W., Zhu, S., Wu, S., Dai, J., 2019. Comparing artificial intelligence techniques for chlorophyll-a prediction

728    in US lakes. Environ. Sci. Pollut. Res. https://doi.org/10.1007/s11356-019-06360-y

729    Lyche Solheim, A., Rekolainen, S., Moe, S.J., Carvalho, L., Phillips, G., Ptacnik, R., Penning, W.E., Toth, L.G.,

730    O'Toole, C., Schartau, A.K.L., Hesthagen, T., 2008. Ecological threshold responses in European lakes and

731    their applicability for the Water Framework Directive (WFD) implementation: Synthesis of lakes results

732    from the REBECCA project. Aquat. Ecol. https://doi.org/10.1007/s10452-008-9188-5

733    Mammides, C., 2020. A global assessment of the human pressure on the world's lakes. Glob. Environ. Chang.

734    https://doi.org/10.1016/j.gloenvcha.2020.102084

735    Mankin, K.R., Koelliker, J.K., Kalita, P.K., 1999. Watershed and lake water quality assessment: An integrated

736    modeling approach. J. Am. Water Resour. Assoc. https://doi.org/10.1111/j.1752-1688.1999.tb04194.x

737    Martin, R., Radosavljevic, S., Schlüter, M., 2020. Short-term decisions in lake restoration have long-term

738    consequences for water quality. Reg. Environ. Chang. https://doi.org/10.1007/s10113-020-01643-4

739 Matthies, M., Giupponi, C., Ostendorf, B., 2007. Environmental decision support systems: Current issues,

740 methods and tools. Environ. Model. Softw. https://doi.org/10.1016/j.envsoft.2005.09.005

741 McInnes, L., Healy, J., Astels, S., 2017. hdbscan: Hierarchical density based clustering. J. Open Source Softw.

742 https://doi.org/10.21105/joss.00205

743 Moallemi, E.A., Zare, F., Reed, P.M., Elsawah, S., Ryan, M.J., Bryan, B.A., 2020. Structuring and evaluating

744 decision support processes to enhance the robustness of complex human–natural systems. Environ. Model.

745 Softw. https://doi.org/10.1016/j.envsoft.2019.104551

746 Muharemi, F., Logofătu, D., Leon, F., 2019. Machine learning approaches for anomaly detection of water quality

747 on a real-world data set. J. Inf. Telecommun. https://doi.org/10.1080/24751839.2019.1565653

748 Murtagh, F., Legendre, P., 2014. Ward's Hierarchical Agglomerative Clustering Method: Which Algorithms

749 Implement Ward's Criterion? J. Classif. https://doi.org/10.1007/s00357-014-9161-z

750 Mustajoki, J., Hämäläinen, R.P., Marttunen, M., 2004. Participatory multicriteria decision analysis with Web-

751 HIPRE: A case of lake regulation policy. Environ. Model. Softw.

752 https://doi.org/10.1016/j.envsoft.2003.07.002Mustow, S. E. 2021. Strategic environmental assessment in the

753 water sector, in: Fischer, T. B. and González, A. (eds.). *Handbook on Strategic Environmental Assessment*,

754 Cheltenham: Edward Elgar (chapter 13).

755 Neri, L., Coscieme, L., Giannetti, B.F., Pulselli, F.M., 2018. Imputing missing data in non-renewable empower

756 time series from night-time lights observations. Ecol. Indic. https://doi.org/10.1016/j.ecolind.2017.08.040

757 Newman, D.A., 2014. Missing Data: Five Practical Guidelines. Organ. Res. Methods.

758 https://doi.org/10.1177/1094428114548590

759 Ngouna, R.H., Ratolojanahary, R., Medjaher, K., Dauriac, F., Sebilo, M., Junca-Bourié, J., 2020. A data-driven

760 method for detecting and diagnosing causes of water quality contamination in a dataset with a high rate of

761 missing values. Eng. Appl. Artif. Intell. https://doi.org/10.1016/j.engappai.2020.103822

762 Noble, D.W.A., Nakagawa, S., 2021. Planned missing data designs and methods: Options for strengthening

763 inference, increasing research efficiency and improving animal welfare in ecological and evolutionary

764 research. Evol. Appl. https://doi.org/10.1111/eva.13273

765 Paruch, A.M., 2014. The impact of wastewater irrigation on the chemical quality of groundwater. Water

766 Environ. J. https://doi.org/10.1111/wej.12064

767 Paruch, L., Paruch, A.M., Blankenberg, A.G.B., Haarstad, K., Mæhlum, T., 2017. Norwegian study on microbial

768 source tracking for water quality control and pollution removal in constructed wetland treating catchment

769 run-off. Water Sci. Technol. https://doi.org/10.2166/wst.2017.303

770    Peters-Lidard, C.D., Rose, K.C., Kiang, J.E., Strobel, M.L., Anderson, M.L., Byrd, A.R., Kolian, M.J., Brekke,

771        L.D., Arndt, D.S., 2021. Indicators of climate change impacts on the water cycle and water management.

772        Clim. Change. https://doi.org/10.1007/s10584-021-03057-5

773    Poikane, S., Birk, S., Böhmer, J., Carvalho, L., De Hoyos, C., Gassner, H., Hellsten, S., Kelly, M., Lyche

774        Solheim, A., Olin, M., Pall, K., Phillips, G., Portielje, R., Ritterbusch, D., Sandin, L., Schartau, A.K.,

775        Solimini, A.G., Van Den Berg, M., Wolfram, G., Van De Bund, W., 2015. A hitchhiker's guide to

776        European lake ecological assessment and intercalibration. Ecol. Indic.

777        https://doi.org/10.1016/j.ecolind.2015.01.005

778    Posthuma, L., Zijp, M.C., De Zwart, D., Van de Meent, D., Globevnik, L., Koprivsek, M., Focks, A., Van Gils,

779        J., Birk, S., 2020. Chemical pollution imposes limitations to the ecological status of European surface

780        waters. Sci. Rep. https://doi.org/10.1038/s41598-020-71537-2

781    Radosavljevic, A., Anderson, R.P., 2014. Making better Maxent models of species distributions: Complexity,

782        overfitting and evaluation. J. Biogeogr. https://doi.org/10.1111/jbi.12227

783    Raghunathan, T., Lepkowski, J., Van Hoewyk, J., Solenberger, P., 2001. A multivariate technique for multiply

784        imputing missing values using a sequence of regression models. Surv. Methodol.

785    Ratolojanahary, R., Ngouna, R.H., Medjaher, K., Dauriac, F., Sebilo, M., 2019. Groundwater quality assessment

786        combining supervised and unsupervised methods, in: IFAC-PapersOnLine.

787        https://doi.org/10.1016/j.ifacol.2019.10.054

788    Reis, S., Voigt, K., Oxley, T., 2017. Thematic issue on modelling human and ecological health risks. Environ.

789        Model. Softw. https://doi.org/10.1016/j.envsoft.2017.02.029

790    Ren, C., Li, C., Jia, K., Zhang, S., Li, W., Cao, Y., 2008. Water quality assessment for Ulansuhai Lake using

791        fuzzy clustering and pattern recognition. Chinese J. Oceanol. Limnol. https://doi.org/10.1007/s00343-008-

792        0339-2

793    Reyjol, Y., Argillier, C., Bonne, W., Borja, A., Buijse, A.D., Cardoso, A.C., Daufresne, M., Kernan, M.,

794        Ferreira, M.T., Poikane, S., Prat, N., Solheim, A.L., Stroffek, S., Usseglio-Polatera, P., Villeneuve, B., van

795        de Bund, W., 2014. Assessing the ecological status in the context of the European Water Framework

796        Directive: Where do we go now? Sci. Total Environ. https://doi.org/10.1016/j.scitotenv.2014.07.119

797    Rivera-Rondón, C.A., Catalan, J., 2020. Diatoms as indicators of the multivariate environment of mountain

798        lakes. Sci. Total Environ. https://doi.org/10.1016/j.scitotenv.2019.135517

799    Rodríguez, R., Pastorini, M., Etcheverry, L., Chreties, C., Fossati, M., Castro, A., Gorgoglione, A., 2021. Water-

800    quality data imputation with a high percentage of missing values: A machine learning approach. Sustain.

801    https://doi.org/10.3390/su13116318

802    Romañach, S.S., McKelvy, M., Conzelmann, C., Suir, K., 2014. A visualization tool to support decision making

803    in environmental and biological planning. Environ. Model. Softw.

804    https://doi.org/10.1016/j.envsoft.2014.09.008

805    Rossaro, B., Boggero, A., Lods-Crozet, B., Free, G., Lencioni, V., Marziali, L., Wolfram, G., 2012. A benthic

806    quality index for European alpine lakes. Fauna Nor. https://doi.org/10.5324/fn.v31i0.1364

807    Rubin, D.B., 1976. Inference and missing data. Biometrika. https://doi.org/10.1093/biomet/63.3.581

808    Russo, R., 2021. The Pearson product-moment correlation coefficient r, in: Statistics for the Behavioural

809    Sciences. https://doi.org/10.4324/9780203641576-23

810    Said, N.M., Zin, Z.M., Ismail, M.N., Bakar, T.A., 2019. Comparative analysis of missing data imputation

811    methods for continuous variables in water consumption data. Int. J. Adv. Trends Comput. Sci. Eng.

812    https://doi.org/10.30534/ijatcse/2019/6981.62019

813    Sarstedt, M., Mooi, E., 2014. A Concise Guide to Market Research: Cluster analysis. Springer.

814    Schielzeth, H., Dingemanse, N.J., Nakagawa, S., Westneat, D.F., Allegue, H., Teplitsky, C., Réale, D.,

815    Dochtermann, N.A., Garamszegi, L.Z., Araya-Ajoy, Y.G., 2020. Robustness of linear mixed-effects

816    models to violations of distributional assumptions. Methods Ecol. Evol. https://doi.org/10.1111/2041-

817    210X.13434

818    Seaman, S., Galati, J., Jackson, D., Carlin, J., 2013. What is meant by "missing at random"? Stat. Sci.

819    https://doi.org/10.1214/13-STS415

820    Shah, A.D., Bartlett, J.W., Carpenter, J., Nicholas, O., Hemingway, H., 2014. Comparison of random forest and

821    parametric imputation models for imputing missing data using MICE: A CALIBER study. Am. J.

822    Epidemiol. https://doi.org/10.1093/aje/kwt312

823    Sojka, M., Choiński, A., Ptak, M., Siepak, M., 2020. The variability of lake water chemistry in the bory

824    tucholskie national park (Northern Poland). Water (Switzerland). https://doi.org/10.3390/w12020394

825    Sojka, M., Jaskuła, J., Wróżyński, R., 2019. ANALYSIS OF HEAVY METALS CONTAMINATION IN

826    BOTTOM SEDIMENTS OF LAKES LOCATED IN THE GNIEZNO LAKELAND. Acta Sci. Pol. Form.

827    Circumiectus. https://doi.org/10.15576/asp.fc/2019.18.4.137

828    Srebotnjak, T., Carr, G., De Sherbinin, A., Rickwood, C., 2012. A global Water Quality Index and hot-deck

829    imputation of missing data. Ecol. Indic. https://doi.org/10.1016/j.ecolind.2011.04.023

830  Tibshirani, R., Walther, G., Hastie, T., 2001. Estimating the number of clusters in a data set via the gap statistic.

831  J. R. Stat. Soc. Ser. B Stat. Methodol. https://doi.org/10.1111/1467-9868.00293

832  Wang, L., Xue, H., 2020. Group decision-making method based on expert classification consensus information

833  integration. Symmetry (Basel). https://doi.org/10.3390/sym12071180

834  Ward, J.H., 1963. Hierarchical Grouping to Optimize an Objective Function. J. Am. Stat. Assoc.

835  https://doi.org/10.1080/01621459.1963.10500845

836  Weerasinghe, V.P.A., Handapangoda, K., 2019. Surface water quality analysis of an urban lake; East Beira,

837  Colombo, Sri Lanka. Environ. Nanotechnology, Monit. Manag.

838  https://doi.org/10.1016/j.enmm.2019.100249

839  Werner, P., Adler, S., Dreßler, M., 2016. Effects of counting variances on water quality assessments:

840  Implications from four benthic diatom samples, each counted by 40 diatomists. J. Appl. Phycol.

841  https://doi.org/10.1007/s10811-015-0760-9

842  Wu, J., Xiong, H., Chen, J., 2009. Towards understanding hierarchical clustering: A data distribution

843  perspective. Neurocomputing. https://doi.org/10.1016/j.neucom.2008.12.011

844  Wu, Y., Duguay, C.R., Xu, L., 2021. Assessment of machine learning classifiers for global lake ice cover

845  mapping from MODIS TOA reflectance data. Remote Sens. Environ.

846  https://doi.org/10.1016/j.rse.2020.112206

847  Xiao, J., Bulut, O., 2020. Evaluating the Performances of Missing Data Handling Methods in Ability Estimation

848  From Sparse Data. Educ. Psychol. Meas. https://doi.org/10.1177/0013164420911136

849  Yanai, R.D., Mann, T.A., Hong, S.D., Pu, G., Zukswert, J.M., 2021. The current state of uncertainty reporting in

850  ecosystem studies: a systematic evaluation of peer-reviewed literature. Ecosphere.

851  https://doi.org/10.1002/ecs2.3535

852  Yang, Y., Xiong, Q., Wu, C., Zou, Q., Yu, Y., Yi, H., Gao, M., 2021. A study on water quality prediction by a

853  hybrid CNN-LSTM model with attention mechanism. Environ. Sci. Pollut. Res.

854  https://doi.org/10.1007/s11356-021-14687-8

855  Yüksel, I., 2012. Developing a Multi-Criteria Decision Making Model for PESTEL Analysis. Int. J. Bus. Manag.

856  https://doi.org/10.5539/ijbm.v7n24p52

857  Zambelli, P., Lora, C., Spinelli, R., Tattoni, C., Vitti, A., Zatelli, P., Ciolli, M., 2012. A GIS decision support

858  system for regional forest management to assess biomass availability for renewable energy production.

859  Environ. Model. Softw. https://doi.org/10.1016/j.envsoft.2012.05.016

860  Zasada, I., Piorr, A., Novo, P., Villanueva, A.J., Valánszki, I., 2017. What do we know about decision support

861  systems for landscape and environmental management? A review and expert survey within EU research

862  projects. Environ. Model. Softw. https://doi.org/10.1016/j.envsoft.2017.09.012

863  Zhang, S., Xia, Z., Wang, T., Williams, B.K., Szaro, R.C., Shapiro, C.D., Brown, E.D., Wieland, R., Gutzler, C.,

864  White, L.G., Welsh, W.D., Wassen, M.J., Runhaar, H., Barendregt, A., Okruszko, T., Walker, J.D.,

865  Chapra, S.C., Voinov, A.A., Seppelt, R., Reis, S., Nabel, J.E.M.S., Shokravi, S., Gaddis, E.J.B., Bousquet,

866  F., Costanza, R., Videira, N., Antunes, P., Santos, R., Lopes, R., Vayssières, J., Vigne, M., Alary, V.,

867  Lecomte, P., van der Zee, D.-J., Holkenborg, B., Robinson, S., Uusitalo, L., Lehikoinen, A., Helle, I.,

868  Myrberg, K., Tversky, A., Kahneman, D., Tsouvalis, J., Waterton, C., Tay, L., Diener, E., Swetnam, R.D.,

869  Fisher, B., Mbilinyi, B.P., Munishi, P.K.T., Willcock, S., Ricketts, T., Mwakalila, S., Balmford, A.,

870  Burgess, N.D., Marshall, A.R., Lewis, S.L., Surowiecki, J., Sun, A., Stanovich, K.E., Stanovich, P.J.,

871  Squires, H., Renn, O., Simon, H.A., Silvertown, J., Shirk, J.L., Ballard, H.L., Wilderman, C.C., Phillips,

872  T., Wiggins, A., Jordan, R., McCallie, E., Minarchek, M., Lewenstein, B. V., Krasny, M.E., Bonney, R.,

873  Sheppard, S.R.J., Cizek, P., Seidl, R., Scholten, L., Scheidegger, A., Reichert, P., Maurer, M., Sawyer, B.,

874  Sauvé, L., Renaud, L., Kaufman, D., Sanò, M., Richards, R., Medina, R., Sahin, O., Siems, R.S., Stewart,

875  R.A., Porter, M.G., Rockström, J., Steffen, W., Noone, K., Persson, A., Chapin, F.S., Lambin, E.F.,

876  Lenton, T.M., Scheffer, M., Folke, C., Schellnhuber, H.J., Nykvist, B., de Wit, C.A., Hughes, T., van der

877  Leeuw, S., Rodhe, H., Sörlin, S., Snyder, P.K., Svedin, U., Falkenmark, M., Karlberg, L., Corell, R.W.,

878  Fabry, V.J., Hansen, J., Walker, B., Liverman, D., Richardson, K., Crutzen, P., Foley, J.A., Röckmann, C.,

879  Ulrich, C., Dreyer, M., Bell, E., Borodzicz, E., Haapasaari, P., Hauge, K.H., Howell, D., Mäntyniemi, S.,

880  Miller, D.G.D., Tserpes, G., Pastoors, M., Robson, B.J., HAMILTON, D., WEBSTER, I, CHAN, T.,

881  Ritzer, G., Ritzema, H., Froebrich, J., Raju, R., Sreenivas, C., Kselik, R., Rinderknecht, S.L., Borsuk,

882  M.E., Schuwirth, N., Langhans, S., Reed, M.S., Raymond, C.M., Bryan, B.A., MacDonald, D.H., Cast, A.,

883  Strathearn, S., Grandgirard, A., Kalivas, T., Prestopnik, N.R., Crowston, K., Pettit, C.J., Lewis, H., Petsko,

884  G.A., Papathanasiou, J., Kenward, R., Page, T., Heathwaite, A.L.L., Thompson, L.J., Pope, L., Willows,

885  R., Oxley, T., Jeffrey, P., Lemon, M., Oliver, D.M., Fish, R.D., Winter, M., Hodgson, C.J., Chadwick,

886  D.R., O'Hagan, A., Nyaki, A., Gray, S.A.S., Lepczyk, C.A., Skibins, J.C., Rentsch, D., Nino-Ruiz, M.,

887  Bishop, I., Nicolson, C.R., Starfield, A.M., Kofinas, G.P., Kruse, J.A., Nettley, A., Desilvey, C., Anderson,

888  K., Wetherelt, A., Caseldine, C., Nativi, S., Mazzetti, P., Geller, G.N., Nash, U.W., Mustajoki, J.,

889  Hämäläinen, R.P., Marttunen, M., Murray-Rust, D., Rieser, V., Robinson, D.T., Miličič, V., Rounsevell,

890  M., Morris, D.E., Oakley, J.E., Crowe, J.A., McKinnon, J., McCall, M.K., Martinez, J., Verplanke, J.,

891  Dunn, C.E., Peters-Guarin., G., Matthews, K.B., Rivington, M., Blackstock, K., McCrum, G., Buchan, K.,

892  Maslow, A.H., Martin, G., Felten, B., Duru, M., Mackay, C., Lynam, T., Jong, W. de, Sheil, D.,

893  Kusumanto, T., Evans, K., Liu, S.B., Poore, B.S., Snell, R.J., Goodman, A., Plant, N.G., Stockdon, H.F.,

894  Morgan, K.L.M., Krohn, M.D., Lippe, M., Thai Minh, T., Neef, A., Hilger, T., Hoffmann, V., Lam, N.T.,

895  Cadisch, G., Li, Y., Zhu, Z., Leys, A.J., Vanclay, J.K., Latre, M.Á., Lopez-Pellicer, F.J., Nogueras-Iso, J.,

896  Béjar, R., Zarazaga-Soria, F.J., Muro-Medrano, P.R., Lange, E., Morgan, E., Romano, D., Lai, J.-S.,

897  Chang, W.-Y., Chan, Y.-C., Kang, S.-C., Tan, Y.-C., Labiosa, W.B., Forney, W.M., Esnard, A.-M.,

898  Mitsova-Boneva, D., Bernknopf, R., Hearn, P., Hogan, D., Pearlstine, L., Strong, D., Gladwin, H., Swain,

899  E., Kuhn, A., Britz, W., Willy, D.K., van Oel, P., Krueger, T., Hubacek, K., Smith, L., Hiscock, K.,

900  Kraker, J. de, Kroeze, C., Kirschner, P., Kragt, M.E., Macleod, C.J.A., Korff, Y. von, Daniell, K.A.,

901  Moellenkamp, S., Bots, P.W.G., Bijlsma, R.M., Koltko-Rivera, M.E., Knapp, C.N., Fernandez-Gimenez,

902  M., Kachergis, E., Rudeen, A., Kelly (Letcher), R.A., Jakeman, A.J.A.J., Barreteau, O., Elsawah, S.,

903  Hamilton, S.H., Henriksen, H.J., Kuikka, S., Maier, H.R., Rizzoli, A.E.A.E., van Delden, H., Kalaugher,

904  E., Bornman, J.F., Clark, A., Beukes, P., Kahan, D., Jones, N., Ross, H., Perez, P., Leitch, A., Jankowski,

905  P., Nyerges, T., Letcher, R.A., Norton, J.P., Irwin, A., IAP2, Howe, J., Hovmand, P.S., Andersen, D.F.,

906  Rouwette, E., Richardson, G.P., Rux, K., Calhoun, A., Hossard, L., Jeuffroy, M.H., Pelzer, E., Pinochet,

907  X., Souchere, V., Højberg, A.L., Troldborg, L., Stisen, S., Christensen, B.B.S., Hoffmann, M., Borenstein,

908  J., Heylighen, F., Chielens, K., Hewitt, R., Escobar, F., Hardcastle, A., Rambaldi, G., Long, B., Lanh, L.

909  Van, Son, D.Q., Hall, D.M., Lazarus, E.D., Swannack, T.M., Gilbertz, S.J., Horton, C.C., Peterson, T.R.,

910  Halbrendt, J., Crow, S., Radovich, T., Kimura, A.H., Tamang, B.B., Haklay, M., Groen, E.A., Heijungs,

911  R., Bokkers, E.A.M., de Boer, I.J.M., Greene, J.C., Greenblat, C.S., Chan, A., Clark, D., Cox, L.J., Henly-

912  Shepard, S., Graveline, N., Aunay, B., Fusillier, J.L., Rinaudo, J.D., Glynn, P.D., Giupponi, C., de Vries,

913  B.J.M., Hasselmann, K., Giordano, R., Liersch, S., Gerd Gigerenzer, Brighton, H., Gaillard, J.J.C.,

914  Monteil, C., Perrillat-Collomb, A., Chaudhary, S., Chaudhary, M., Chaudhary, O., Giazzi, F., Cadag,

915  J.R.D., Fung, A., Russon Gilman, H., Shkabatur, J., Fulton, E.A., Boschetti, F., Sporcic, M., Jones, T.,

916  Little, L.R., Dambacher, J.M., Gray, R., Scott, R., Gorton, R., Fritz, S., McCallum, I., Schill, C., Perger,

917  C., See, L., Schepaschenko, D., van der Velde, M., Kraxner, F., Obersteiner, M., Fraternali, P., Castelletti,

918  A., Soncini-Sessa, R., Vaca Ruiz, C., Foster, A., Dunham, I.M., Fisher, R., O'Leary, R.A., Low-Choy, S.,

919  Mengersen, K., Caley, M.J., Fischer, F., Fagin, R., Halpern, J.Y., Estelles-Arolas, E., Gonzalez-Ladron-de-

Guevara, F., Enserink, B., Patel, M., Kranz, N., Maestu, J., Guillaume, J.H.A., Filatova, T., Rook, J., Economist, Djaouti, D., Alvarez, J., Jessel, J.-P., Rampnoux, O., Delgado-Galván, X., Izquierdo, J., Benítez, J., Pérez-García, R., Debolini, M., Marraccini, E., Rizzo, D., Galli, M., Bonari, E., Dean, J., Ghemawat, S., d'Aquino, P., Bah, A., Creighton, J.L., Craig, R.K., Ruhl, J.B., Cohn, J.P., Cobb, A.N., Thompson, J.L., Chow, T.E., Sadler, R., Chingombe, W., Pedzisai, E., Manatsa, D., Mukwada, G., Taru, P., Chen, Y., Yu, J., Khan, S., Chen, S.H., Pollino, C.A., Chabris, C.F., Simons, D.J., Catenacci, M., Galelli, S., Ratto, M., Young, P.C., Carmona, G., Varela-Ortega, C., Bromley, J., Campo, P.C., Villanueva, T.R., Butler, M.P., Reed, P.M., Fisher-Vanden, K., Keller, K., Wagener, T., Buss, D., Brooking, C., Hunter, J., Bizikova, L., Burch, S., Robinson, J., Shaw, A., Wolters, H.A., Hoekstra, A.Y., BBC, Bastin, L., Cornford, D., Jones, R., Heuvelink, G.B.M., Pebesma, E., Stasch, C., Williams, M., Le Page, C., Barnaud, C., Page, C. Le, Dumrongrojwatthana, P., Trebuil, G., Aumann, C.A., Audubon, Arnstein, S.R., Arnold, T.R., Ariely, D., Argent, R.M., Arciniegas, G., Janssen, R., Rietveld, P., Anderson, C.A., Lepper, M.R., Ross, L., 2013. Serious games: Improving public policy through game-based learning and simulation. Environ. Model. Softw.

Zhang, Y., Thorburn, P.J., 2022. Handling missing data in near real-time environmental monitoring: A system and a review of selected methods. Futur. Gener. Comput. Syst. https://doi.org/10.1016/j.future.2021.09.033

Zhang, Y., Thorburn, P.J., 2021. A dual-head attention model for time series data imputation. Comput. Electron. Agric. https://doi.org/10.1016/j.compag.2021.106377

Zhang, Y.F., Thorburn, P.J., Vilas, M.P., Fitch, P., 2019. Machine learning approaches to improve and predict water quality data, in: 23rd International Congress on Modelling and Simulation - Supporting Evidence-Based Decision Making: The Role of Modelling and Simulation, MODSIM 2019. https://doi.org/10.36334/modsim.2019.d5.zhangyif

Zhang, Z., 2016. Multiple imputation with multivariate imputation by chained equation (MICE) package. Ann. Transl. Med. https://doi.org/10.3978/j.issn.2305-5839.2015.12.63

Zhou, X.H., 2020. Challenges and strategies in analysis of missing data. Biostat. Epidemiol. https://doi.org/10.1080/24709360.2018.1469810