

On the prediction of the ecological status of human-altered streams and its rule-based interpretation

Terence A. Etchells^a, Alfredo Vellido^b, Eugenia Martí^c, Paulo J.G. Lisboa^a, Joaquim Comas^d

^aCMS. Liverpool John Moores University, Liverpool, UK.

^bLSI. Universitat Politècnica de Catalunya. Barcelona, Spain

^cCEAB-CSIC. Blanes, Barcelona, Spain

^dLEQUIA. Universitat de Girona. Girona, Spain. quim@lequia1.udg.es

Abstract: The recent Water Framework Directive of the European Union set year 2015 as their target for freshwater and coastal ecosystems all across Europe to achieve good ecological status. This study concerns the analysis of the empirical data from the knowledge base of an environmental decision support system developed within the European project STREAMES. These data, which come from several low-order streams located mostly on the Mediterranean region, consist of measurements of several physical, chemical and biological variables. We aim to classify these data according to the ecological status of the streams they correspond to, where stream nutrient retention efficiency (a functional ecosystem attribute) is used as an indicator of ecological status. This classification task is performed using supervised Artificial Neural Networks. The interpretability of the obtained classification results can be improved by their description in terms of simple, actionable rules. This is accomplished through the application of Orthogonal Search-based Rule Extraction, a novel overlapping rule extraction method. All the newly acquired knowledge should help water managers to focus their efforts on strategies that minimize the negative human impacts on vulnerable low-order streams.

Keywords: Human-altered streams; ecological status; nutrient retention; Artificial Neural Networks; Orthogonal Search-based Rule Extraction.

1. INTRODUCTION

The Water Framework Directive (WFD) of the European Union (Council of the European Communities, 2000) set year 2015 as their target to achieve a category of “good” for the ecological status of freshwater and coastal ecosystems all across Europe. The WFD considers five categories of ecological status according to an ordinal scale: bad / poor / moderate / good / high.

The data analysed in this study are part of the knowledge base of the environmental decision support system (EDSS) that was the object of the STREAMES (STream REAch Management, an Expert System) European project. According to WFD requirements, one of the main goals of this EDSS was the evaluation of water quality and, to a larger extent, the ecological status of fluvial ecosystems.

The data were obtained from streams affected to different degrees by inputs of nutrients from diffuse and/or point sources. Several streams were selected throughout Europe and Israel, with emphasis on streams located in the Mediterranean

region, for which the effects of nutrient inputs are amplified by their usually irregular and relatively low flows. In the *Water Resources* section of the world-wide Pilot 2006 Environmental Performance Index report by Esty et al. [2006], which includes a 100-point ranking of 133 countries, the European countries under study in STREAMES rank from position 29 for Austria (99.4 points) to position 108 for Spain (62.4 points). Physical, chemical (nutrient and major ions concentrations), and biological parameters, including nutrient retention metrics, were measured at these streams. Stream nutrient retention efficiency (a functional ecosystem attribute) is actually considered here as a descriptor and surrogate measure of ecological status. This is in contrast with the most commonly used approaches reported in the existing literature, where ecological status is evaluated by means of communities of organisms or habitat descriptors (i.e., structural ecosystem attributes).

Artificial Neural Networks (ANN) are gaining ground as environmental modelling tools and have been used for water research analysis,

including the prediction of water quality (See, for instance, recent work by Kralisch et al. [2003], Almasri and Kaluarachchi [2005], and Gatts et al. [2005]). In this study, we use ANNs beyond water quality modelling, and aim to classify the data according to ecological status, where this is measured by a discretised version of a nutrient retention metric, used as the target for the ANN.

One of the potential drawbacks affecting the application of ANN models to classification problems is that of the limited interpretability of their results. One way to overcome this limitation is by pairing the ANN model with a rule extraction method. The interpretability of the ecological status classification results can be greatly improved by their description in terms of reasonably simple and actionable rules. This is accomplished in this study through the application of Orthogonal Search-based Rule Extraction (OSRE), a novel overlapping rule extraction method by Etchells and Lisboa [2006].

The paper is structured as follows. First, the ANN model used for classification and the OSRE method are, in turn, introduced. The analysed stream sites and data are then summarily described. This is followed by the presentation of the experimental results and their discussion. Some brief conclusions are finally provided.

2. NEURAL NETWORKS FOR ECOLOGICAL STATUS PREDICTION

In keeping with good practice for ANN design, the model used for classification in this application is regularised to ensure that it generalises well to unseen data. It consists of a Multi-Layer Perceptron (MLP) with a single hidden layer and with a weight decay term adjusted by cross-validation to optimise the classification rate in an out-of-sample test data set. The number of hidden nodes used was also set to ensure sufficient modelling capacity to obtain good classification results on the test data.

Conditional class probabilities were obtained using sigmoidal transfer functions in the hidden and output layers, with an objective function given by the log-likelihood and weight decay terms. This model therefore treats all classes as independent. An automatic adjustment was made for the prevalence of data records within- and out-of-class, by re-weighting the log-likelihood function and adjusting the predicted class-conditional probability, as in Lisboa et al. [2000].

The ANN model evaluation was performed within the Receiver Operating Characteristic (ROC) framework. This measures the success rate for

accurate detection of data records in-class (sensitivity) and out-of-class (specificity), as well as the proportion of in-class predictions that are accurate (positive predictive value). Sensitivity and specificity define the ROC curve. These measures of performance are used to rank OSRE rules and to describe their predictive accuracy.

3. ORTHOGONAL SEARCH-BASED RULE EXTRACTION

The OSRE algorithm, by Etchells and Lisboa [2006] is a method to efficiently extract comprehensible rules from smooth models, such as those created, for instance, by supervised ANNs for data classification. Rule extraction can help to extract actionable knowledge from trained ANNs that would otherwise be difficult to interpret. In this study, this process should ease water managers' decision making tasks. OSRE is a principled approach underpinned by a theoretical framework of continuous valued logic developed by Tsukimoto [2000]. In essence, the algorithm extracts rules by taking each data record that the model predicts to be within a particular class, and searching in the direction of each data feature to find the limits of the space regions for which the model prediction is in that class (Figure 1, top).

These regions form hyper-boxes that capture in-class data and they are converted to conjunctive rules in terms of the data features and their values (Figure 1, bottom). The obtained set of rules is subjected to a number of refinement steps: removing repetitions; filtering rules of poor specificity and sensitivity (see definitions in section 2.1); and removing rules that are subsets of other rules. The rules are then ranked in terms of their sensitivity values to form a hierarchy describing the in-class data. OSRE has been assessed and compared with alternative methods in Etchells and Lisboa [2006].

4. STREAM DATA

The STREAMES project focussed on the effects of high nutrient loads on low-order streams. Eleven third-order streams were selected across seven European countries plus Israel. Two of them were discarded for this study due to extreme data incompleteness for the data features selected in this study. Sites were selected to cover a broad range of climate, geomorphology and environmental conditions. Scenarios were differentiated according to hydrologic conditions

and the dominant land-use within the selected stream catchment. In addition, and in order to estimate the effect of nutrient inputs from point sources on the structure and function of the streams, two reaches located upstream and downstream of a wastewater treatment plant effluent input were selected for each stream. Details can be found at www.streames.org.

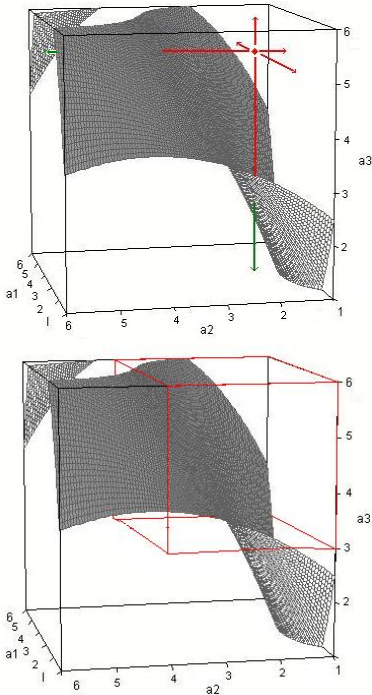


Figure 1. Illustration of the basic OSRE procedure. (Top): orthogonal searching to find decision boundaries on the available data features; (bottom): overlapping hyper-boxes constructed from the search results.

In every reach, six (on average) experimental campaigns were conducted over a year to cover a wide range of environmental conditions. In each reach and on each date, physical (hydrology, hydraulics, morphology), chemical (nutrient and major ions concentrations), and biological (both structural: biofilm biomass and chlorophyll; and functional: nutrient retention and ecosystem metabolism) parameters were measured.

A nutrient retention metric, *Nutrient Uptake Length* (S_w), is used in this study as a surrogate indicator of the streams' ecological status. S_w is the average distance that a nutrient molecule travels before it is removed from the water column and, therefore, it gauges the nutrient retention efficiency of the stream. For space limitations, and although this metric is available for ammonium, nitrates and phosphates, only the S_w for ammonium was used. This continuous

metric was initially discretised to correspond with the WFD (Council of the European Communities, 2000) categorization of ecological status into five classes: bad / poor / moderate / good / high. Nevertheless, and given that there were only a few data records available for some of those categories, we decided to simplify the procedure and consider only three categories: good (30.9% of the data) / moderate (40% of the data) / bad (29.1% of the data). The discretised S_w for ammonium was therefore used as the target for the ANN classifier.

The original data records-to-features ratio was far too low to implement any reliable analytical model. Therefore, experts in the areas of chemistry, biogeochemistry and stream ecology set up to agree on a more parsimonious dataset, consisting on 110 records and 22 descriptive features, which are detailed in Table 1.

Table 1. List of the 22 features selected for this study, grouped by their typology.

TYPE	FEATURE
Ion Concentrations (chemical)	Cations ($\text{Na}^+ + \text{K}^+ + \text{Mg}^{2+} + \text{Ca}^{2+} + \text{NH}_4^+$)
	Anions ($\text{Cl}^- + \text{SO}_4^{2-} + \text{NO}_3^-$)
Nutrient Concentrations (chemical)	Alkalinity
	$\text{NH}_4^+\text{-N}$
	$\text{NO}_3^-\text{-N}$
	$\text{PO}_4^{3-}\text{-P}$
	Dissolved Organic Carbon (DOC)
Hydrological, Hydraulic & Morphologic (physical)	Conductivity
	Dissolved Inorganic Nitrogen (DIN)
	Depth (Wet channel average depth)
	Wet Perimeter
	Substrate Ratio (Percentage of {Cobbles + Pebbles} substrata, divided by percentage of {Gravel + Sand + Silt} substrata)
	Wet Perimeter / Depth Ratio
	K1 (Water transient storage exchange coefficient: from water column to transient storage zone)
	K2 (Water transient storage exchange coefficient: from transient storage zone to water column)
	Respiration (Daily rate of ecosystem respiration)
	GPP (Daily rate of gross primary production)
Stream Metabolism & Biofilm (biological)	GPP:R (GPP to Respiration ratio per day)
	Daily Light (PAR)
	Temperature
	Chlorophylla
	Biomass

5. EXPERIMENTAL RESULTS AND DISCUSSION

5.1 ANN classification results

Due to the paucity of available data records and the high data dimensionality, the classification task is considerably difficult. For this purpose, a two layer MLP consisting of 22 input nodes (one for each data feature), 8 hidden nodes and 3 output nodes (1 for each class), was implemented as described in section 2. The network weights were randomly initialised. The 110 data records were split into halves for training and testing. The training parameters and network architecture were selected by cross-validation as those maximizing the accuracy in the classification of test data; the parameters were set at: learning rate = 0.01; momentum = 0.9; weight decay = 0.05; number of epochs = 60.

The final overall classification accuracy for the training data was 95% and 71% for the test data. One of the reasons for the relatively low test accuracy might be low records-to-features ratio mentioned at the beginning of this section.

5.2 OSRE results

The classification results, by themselves, would tell us little about what is behind the decisions made by the MLP for each of the three classes. This is why the rule extraction performed by OSRE is convenient if we want the results to become both interpretable and actionable.

Tables 2, 3, and 4 provide the collections of rules defining each class (good / moderate / bad ecological status). Each rule is a conjunction of the data features and their values. In all tables, the rules are arranged in decreasing order of importance. In order to make these rules easier to understand and discuss, the value ranges of the 22 features in Table 1 are plotted in Figure 2.

According to tables 2 to 4, and with the help of Figure 2, the three classes of ecological status can be interpreted as follows:

- **Class 1 (Good ecological status):** According to Table 2, good ecological status is mainly characterised by a combination of very low levels of $\text{PO}_4^{3-}\text{-P}$ (lower tenth of its range); most of the *Substrate Ratio* range, excluding extreme values (extremely large or small substrata); the two lowest fifths of *Gross Primary Production to Respiration Ratio*;

and, finally, the three highest fourths of the *Temperature* range. The second rule is far simpler and consists of a combination of the lowest tenth of the *K2 water transient storage exchange coefficient* range, and the upper three fifths of the *Respiration* range. For brevity, only these rules are described.

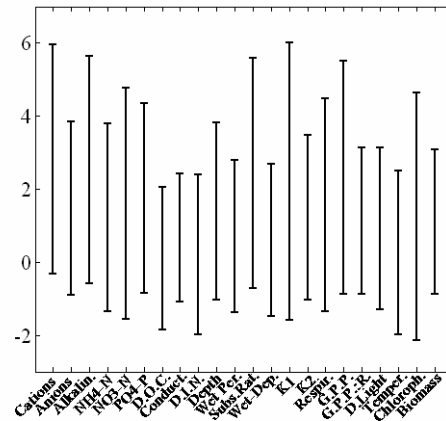


Figure 2. Plot of the standardized values range for each of the 22 data features.

- **Class 2 (Moderate ecological status):** According to Table 3, this class is the most difficult to discriminate, which is reflected in the proliferation and relative complexity of the rules. This could also be understood as moderate ecological status being caused by multiple factors. For the sake of brevity, only the two main rules are described. The first rule characterises moderate ecological status as a combination of most of the *Dissolved Inorganic Nitrogen* range but the highest and lowest values; a combination of three physical parameters (most of the two upper thirds of *Wet Perimeter*, excluding the highest values; the lowest tenth of the *Substrate Ratio* range -predominance of gravel, sand and silt-, and most the upper two thirds of the *Wet Perimeter / Depth ratio* range, excluding the highest values); and most of the two lowest thirds of the *Biomass* range, excluding the lowest values. The second rule is far simpler and consists of a combination of the lowest two fifths of the $\text{PO}_4^{3-}\text{-P}$ range, and most of the three upper fourths of the *Gross Primary Production* parameter, excluding the most extreme values.
- **Class 3 (Bad ecological status):** According to Table 4, stream records corresponding to bad ecological status show the most parsimonious and easiest rule description, suggesting that bad ecological status is easier

to discriminate than good or moderate. Again, only the two main rules are described. Bad ecological status seems to be mainly characterised by a combination of the lower three fourths of the NO_3^- -N range; most of the PO_4^{3-} -P range but the lowest values; the lowest half of the *Wet Perimeter* range; the lowest fourth of the *Respiration* range; and the lowest two thirds of the *Temperature* range. The second rule is a combination of physical (lowest two thirds of *Wet Perimeter* and most of the *K2* range but the lowest values) and biological (lowest two thirds of *Respiration* and lowest two fifths of *Gross Primary Production*) parameters.

A full interpretation of the rules is beyond the scope of this paper. It is worth noting, though, that, overall, most of the relevant rules include biological features, in particular rates of metabolism. This suggests a certain coupling between stream nutrient/ammonium retention and in-stream biological activity. Good ecological status seems to be linked to those environmental conditions characterized by low eutrophication (or oligotrophic conditions: low PO_4^{3-}), relatively high water transient storage (low *K2*, higher contact between the sediments and available nutrients) and high metabolic activity. Moderate ecological status appears to be a transition in between the two other classes. Finally, bad ecological status is conditioned by low metabolic activity coupled with low habitat availability and eutrophic conditions (as expressed by high PO_4^{3-} availability). Overall, these results support our simplification of the problem for the available data, using 3 instead of 5 classes of ecological status.

6. CONCLUSION

The recent Water Framework Directive of the European Union set year 2015 as their target for freshwater and coastal ecosystems in Europe to achieve good ecological status. In this study we have used data from several low-order streams and attempted to classify them according to ecological status, where this status is described by a nutrient retention metric. The results are encouraging and reveal that the selected data features provide reasonably good ecological status class discrimination.

The interpretation of the ANN classification results has been assisted by the application of a novel rule extraction algorithm: OSRE. The description of the ecological status classes provided by the extracted rules suggests that the categorization of the ecological status of a stream

in terms of nutrient retention could be indirectly inferred using a subset of variables (those defining the rules) that are easier to obtain and monitor by water managers. Interestingly, OSRE results indicate that bad ecological status can be described by more concise rules than better scenarios.

ACKNOWLEDGEMENTS

This work is partially supported by the European Union through the EVK1-CT-2000-00081 project. Alfredo Vellido and Eugenia Martí are research fellows within the Ramón y Cajal program of the Spanish Ministry of Education and Science.

REFERENCES

- Almasri, M.N., and J.J. Kaluarachchi, Modular neural networks to predict the nitrate distribution in ground water using the on-ground nitrogen loading and recharge data, *Environmental Modelling & Software*, 20 (7), 851-871, 2005.
- Council of the European Communities, Directive 2000/60/EC of the European Parliament and of the Council of 23 October 2000 establishing a framework for Community action in the field of water policy, 2000.
- Esty, D.C., Levy, M.A., Srebotnjak, T., de Sherbinin, A., Kim, C.H., and B. Anderson, *Pilot 2006 Environmental Performance Index*, Yale Center for Environmental Law & Policy, New Haven, 2006.
- Etchells, T.A., and P.J.G. Lisboa, Orthogonal Search-based Rule Extraction (OSRE) method for trained neural networks: A practical and efficient approach, *IEEE Transactions on Neural Networks*, 17(2), 2006.
- Gatts, C.E.N., Ovalle, A.R.C., and C.F. Silva, Neural pattern recognition and multivariate data: water typology of the Paraíba do Sul River, Brazil, *Environmental Modelling & Software*, 20(7), 883-889, 2005.
- Kralisch, S., Fink, M., Flügel, W.-A., and C. Beckstein, A neural network approach for the optimisation of watershed management, *Environmental Modelling & Software*, 18(8-9), 815-823, 2003.
- Lisboa, P.J.G., Vellido, A., and H. Wong, Bias reduction in skewed binary classification with Bayesian neural networks, *Neural Networks*. 13(4-5), 407-410, 2000.

Table 2. OSRE rules for Class 1 (Good ecological status). *Spec* stands for Specificity; *Sens* for Sensitivity; *PPV* is the Positive Predictive Value: the ratio of the number of in-class data that the rule predicts to the total number of data the rule predicts. See definitions in section 2.1.

Rule	For this Rule Only			For all rules up to this one		
	Sens	Spec	PPV	Sens	Spec	PPV
-0.85 ≤ PO ₄ ³⁻ -P ≤ -0.23	0.71	0.91	0.77	0.71	0.91	0.77
-0.41 ≤ Substrate Ratio ≤ 5.4						
-0.87 ≤ GPP:R ≤ 0.89						
-0.89 ≤ Temperature ≤ 2.46						
-1.04 ≤ K2 ≤ -0.50	0.09	0.99	0.75	0.76	0.89	0.76
1.03 ≤ Respiration ≤ 4.44						
-0.56 ≤ Alkalinity ≤ 3.18	0.32	0.97	0.74	0.82	0.87	0.74
-1.35 ≤ NH ₄ ⁺ -N ≤ 0.21						
-0.85 ≤ PO ₄ ³⁻ -P ≤ -0.12						
-1.38 ≤ Wet Perimeter ≤ 1.26						
-1.49 ≤ Wet Perimeter / Depth Ratio ≤ 0.84						
-1.04 ≤ K2 ≤ -0.25						
-0.83 ≤ GPP:R ≤ -0.03						
-1.29 ≤ Temperature ≤ 2.41						
-1.04 ≤ K2 ≤ -0.61	0.09	1.00	1.00	0.85	0.87	0.74
-0.49 ≤ Temperature ≤ 0.75						

Table 3. OSRE rules for Class 2 (Moderate ecological status). *Spec*, *Sens* and *PPV* as in Table 2.

Rule	For this Rule Only			For all rules up to this one		
	Sens	Spec	PPV	Sens	Spec	PPV
-1.77 ≤ DIN ≤ 2.22	0.23	0.97	0.83	0.23	0.97	0.83
0.08 ≤ Wet Perimeter ≤ 2.63						
-0.73 ≤ Substrate Ratio ≤ -0.10						
-0.86 ≤ Wet Perimeter / Depth Ratio ≤ 2.58						
-0.73 ≤ Biomass ≤ 2.05	0.19	1.00	1.00	0.42	0.97	0.9
-0.85 ≤ PO ₄ ³⁻ -P ≤ 1.38						
0.92 ≤ GPP ≤ 5.46						
0.29 ≤ Wet Perimeter ≤ 2.66	0.21	0.99	0.9	0.53	0.96	0.88
0.33 ≤ GPP:R ≤ 3.1						
1.43 ≤ Wet Perimeter / Depth Ratio ≤ 2.55	0.07	1	1	0.61	0.96	0.9
-0.55 ≤ K2 ≤ 0.83						
-1.29 ≤ Daily Light ≤ 1.08						
0.81 ≤ Temperature ≤ 2.48 OR -1.99 ≤ Temperature ≤ -0.09						
-1.15 ≤ NH ₄ ⁺ -N ≤ 3.77	0.07	0.96	0.5	0.67	0.91	0.82
1.05 ≤ PO ₄ ³⁻ -P ≤ 4.26						
0.03 ≤ Respiration ≤ 4.47						
-0.87 ≤ GPP ≤ 2.46						
-1.1 ≤ Conductivity ≤ 0.32	0.14	0.97	0.75	0.7	0.9	0.75
-1.77 ≤ DIN ≤ 2.36						
-0.13 ≤ Wet Perimeter ≤ 2.74						
-0.72 ≤ Substrate Ratio ≤ -0.08						
-0.24 ≤ Wet Perimeter / Depth Ratio ≤ 2.59						

Table 4. OSRE rules for Class 3 (Bad ecological status). *Spec*, *Sens* and *PPV* as in Table 2.

Rule	For this Rule Only			For all rules up to this one		
	Sens	Spec	PPV	Sens	Spec	PPV
-1.55 ≤ NO ₃ ⁻ -N ≤ 3.23	0.61	0.96	0.87	0.61	0.96	0.87
-0.44 ≤ PO ₄ ³⁻ -P ≤ 4.32						
-1.38 ≤ Wet Perimeter ≤ 0.92						
-1.35 ≤ Respiration ≤ 0.13						
-1.99 ≤ Temperature ≤ 1.07						
-1.33 ≤ Wet Perimeter ≤ 1.60	0.57	0.95	0.83	0.91	0.92	0.83
-0.35 ≤ K2 ≤ 3.38						
-1.35 ≤ Respiration ≤ 2.52						
-0.87 ≤ GPP ≤ 1.65						
-1.03 ≤ Depth ≤ 1.17	0.52	0.99	0.94	0.94	0.92	0.84
-1.38 ≤ Wet Perimeter ≤ 1.05						

$$-0.06 \leq \mathbf{K2} \leq 3.41$$
