



ELSEVIER

Ecological Modelling 146 (2001) 207–217

ECOLOGICAL
MODELLING

www.elsevier.com/locate/ecolmodel

Modelling population dynamics of aquatic insects with artificial neural networks

Michael Obach ^{a,*}, Rüdiger Wagner ^a, Heinrich Werner ^b,
Hans-Heinrich Schmidt ^a

^a *Limnologische Fluss-Station Schlitz der Max-Planck-Gesellschaft, Damenweg 1, D-36110 Schlitz, Germany*

^b *Department of Mathematics and Computer Sciences, Research Group Neural Networks, Heinrich-Plett-Str. 40, D-34132 Kassel, Germany*

Abstract

We modelled the total number of individuals of selected water insects based on a 30-year data set of population dynamics and environmental variables (discharge, temperature, precipitation, abundance of parental generation) in a small stream in central Germany. For data exploration, visualisation of data, outlier detection, hypothesis generation, and to detect basic patterns in the data, we used Kohonen's self organizing maps (SOM). They are comparable to statistical cluster analysis by ordinating data into groups. Based on annual abundance patterns of Ephemeroptera, Plecoptera and Trichoptera (EPT), species groups with similar ecological requirements were distinguished. Furthermore, we applied linear neural networks, general regression neural networks, modified multi-layer perceptrons, and radial basis function networks combined with a SOM (RBF-SOM) and successfully predicted the annual abundance of selected species from environmental variables. Results were visualised in three-dimensional plots. Relevance detection methods were sensitivity analysis, stepwise method and Genetic Algorithms. Instead of a sliding windows approach we computed the in- and output data of fixed periods for two caddis flies. In order to assess the quality of the models we applied several reliability measures and compared the generalisation error with the long-term mean of the target variable. RBF-SOMs were used to denominate and visualise local and general model accuracy. Results were interpreted on the basis of known species traits. We conclude that it is possible to predict the abundance of aquatic insects based on relevant environmental factors using artificial neural networks. © 2001 Elsevier Science B.V. All rights reserved.

Keywords: Radial basis function neural network; Self-organizing maps; Hybrid training; General regression neural network; Visualisation of multidimensional data; Reliability measure

1. Introduction

An ecosystem is a local biological community and its pattern of interaction with its environment. The community is the set of biological species interacting with each other (e.g. predator-

* Corresponding author.

E-mail address: mobach@mpil-schlitz.mpg.de (M. Obach).

prey, intra- and interspecific competition) in this environment. Ecosystems or parts of ecosystems have been measured, described, analysed and modelled by scientists with a number of different approaches. However, many methods are restricted in use because of the requirements to the data (e.g. normal distribution). Particularly, the amount of data often is too low with biological systems. In these cases artificial neuronal networks (ANN) are a promising tool to investigate and model ecosystems, or populations because there are only limited requirements to the data.

In previous papers (Borchardt et al., 1997; Dapper, 1998; Schleiter et al., 1999; Wagner et al., 2000) the number of aquatic insects in emergence traps was predicted with ANNs using the 'sliding window' technique. Species chosen were part of the EPT (Ephemeroptera, Plecoptera, Trichoptera) community of the Breitenbach, a small stream in central Germany. The abundance (specimen numbers per month) was predicted based on the information of environmental variables (precipitation, discharge and water temperature) of the last 13 months as well as on the abundance of the parent generation.

Because computation was time consuming and the models were complex, various methods to reduce the amount of data necessary for adequate modelling by pre-processing were tested and compared. Models with the five most relevant input variables chosen either by correlation, regression or by sensitivity analysis resulted in good predictions for most species. However, the predictions for taxa with larvae living on sandy substratum (unstable against variation in discharge) needed the total amount of information available for adequate prediction. There was not one single optimal method to be applied to every species with a similar quality of prediction. Further, no model predicted the abundance of an individual species at four sites along the same stream with similar accuracy (Borchardt et al., 2001). Models for most species were sufficient. However, the high proportion of months with no specimens (trivial predictions) was the main reason to try to relate information on the environment exclusively for months with abundance values greater than 0.

Spatial and temporal changes in habitats provide a shifting mosaic of environmental and biotic conditions that play a major role in organizing stream communities (Townsend, 1989; Palmer and Poff, 1997). Physical disturbance by stream flow structures stream communities by frequency, harshness, or seasonal predictability (e.g. Reice, 1985; Southwood, 1988; Poff, 1992), the availability of refugia in streams buffers populations against disturbances (e.g. Townsend and Hildrew, 1994). We test whether abundance patterns of aquatic insects based on knowledge of their life history and biological traits are related to the patterns of environmental variables. Complete long-term information is available for more than 100 EPT taxa for 1969–1998. Test organisms were three out of the 40 most common EPT species of the Breitenbach (Table 1).

In this study, we select input variables (predictors) relevant to the respective output. Output is the monthly or yearly insect abundance influenced by environmental factors and the density of the parental generation. Furthermore, visualisation is improved to assess the quality of the data in early stages of processing and to compare the prediction with the variability of the training data. Firstly, high variability of training data causes low confidence in the network output in concerning areas of the input space. Secondly, the assessment of the results was restricted to simple error measures and diagrams comparing predictions and observations. In areas with low information provided by the training data the quality of the network output may also be limited. Adequate quality measures proposed by Werner and Obach (2001) and applied by Schleiter et al., (2001a,b) are used also in this study.

2. Material and methods

The Breitenbach is a thoroughly studied first-order stream near Schlitz in central Germany (50°40'N, 9°45'E). It flows through an area underlain by Bunter Sandstone. Its most important source is a spring in the middle course at an altitude of about 310 m, below that the stream runs perennial. It flows into the river Fulda at

about 220 m. The total stream length is about 4 km, but only the lower 2 km were investigated. The stream bottom consists of varying grain sizes. There are slabs and large stones as well as sand, silt, and clay. In its upper course, the stream runs close to the forest edge (mainly *Pinus* sp., *Fagus sylvatica*, *Quercus robur*, *Carpinus betulus*), but further downstream in the study area it is bordered almost exclusively by meadows. Detailed descriptions of the stream, including botanical aspects of the riparian vegetation, were given by Cox (1990) and Ringe (1974). A map and information on physical and chemical characteristics of the stream were provided by Wagner et al. (1994). At the study site the catchment area upstream is about 10 km², mean monthly precipitation was 55 l/m² (min = 4.9, max = 181), mean monthly water temperature 11.6°C (min = 3, max = 25), and mean monthly maximum discharge 40 l/s (min = 0.8, max = 616).

Community data consist of monthly or yearly specimen numbers (Fig. 1) of taxa collected in an emergence trap per 6 m stream segment

length, representing a surface area of ≈ 5 m², depending on discharge. Environmental variables were monthly maxima of water temperature, discharge, and the monthly sum of precipitation.

All (input and output) variables were mainly linearly rescaled into the interval [0,1]. In order to reduce the influence of few extreme events, discharge data were logarithmically rescaled (upper limit 250 l/s).

Data were divided into a training and a test set. Test data to assess the generalisation performance of ANN models consisted of measurements and observations of the arbitrarily chosen years 1974, 1982, 1984, 1988, 1992, 1996; the remaining 24 years were training data. Crossvalidation on the training data selected relevant input variables and optimised network parameters.

For data exploration and visualisation we used self-organising maps (SOM) (Kohonen, 1995; Chon et al., 1996; Foody, 1999). General regression neural networks (GRNN) (Specht, 1991; Obach, 1998), that are non-parametric Nadaraya–Watson kernel regression estimators, were applied to select relevant predictors with a

Table 1

EPT species used for modelling; short code; Ephemeroptera (E), Plecoptera (P), Trichoptera (T)

Short code	Name	Short code	Name
A.fim	<i>Apatania fimbriata</i> (Pictet) (T)	M.lon	<i>Micrasema longulum</i> MacLachlan (T)
A.fus	<i>Agapetus fuscipes</i> Curtis (T)	N.cam	<i>Nemoura cambrica</i> Stephens (P)
A.red	<i>Adicella reducta</i> MacLachlan (T)	N.cin	<i>Nemoura cinerea</i> Retzius (P)
A.sta	<i>Amphinemura standfussi</i> Ris (P)	N.fle	<i>Nemoura flexuosa</i> Aubert (P)
B.rho	<i>Baetis rhodani</i> Pictet (E)	N.mar	<i>Nemoura marginata</i> Pictet (P)
B.ver	<i>Baetis vernus</i> Curtis (E)	N.pic	<i>Nemurella pictetii</i> Klapálek (P)
C.lut	<i>Centroptilum luteolum</i> Müller (E)	O.alb	<i>Odontocerum albicorne</i> Scopoli (T)
C.vil	<i>Chaetopteryx villosa</i> Fabricius (T)	P.aub	<i>Protonemura auberti</i> Illies (P)
D.ann	<i>Drusus annulatus</i> Stephens (T)	P.cin	<i>Potamophylax cingulatus</i> Stephens (T)
E.dan	<i>Ephemerella danica</i> Müller (E)	P.con	<i>Plectrocnemia conspersa</i> Curtis (T)
E.ign	<i>Ephemerella ignita</i> Poda (E)	P.int	<i>Protonemura intricata</i> Ris (P)
E.muc	<i>Ephemerella mucronata</i> Bengtsson (E)	P.luc	<i>Potamophylax luctuosus</i> Pill & Mitt. (T)
E.ven	<i>Ecdyonurus venosus</i> Fabricius (E)	P.nit	<i>Protonemura nitida</i> Pictet (P)
H.dig	<i>Halesus digitatus</i> Schrank (T)	P.sub	<i>Paraleptophlebia submarginata</i> Stephens (E)
H.ins	<i>Hydropsyche instabilis</i> Curtis (T)	R.fas	<i>Rhyacophila fasciata</i> Hagen (T)
H.sax	<i>Hydropsyche saxonica</i> MacLachlan (T)	S.pal	<i>Silo pallipes</i> Fabricius (T)
I.goe	<i>Isoperla goertzi</i> Illies (P)	S.per	<i>Sericostoma personatum</i> Kirby & Spence (T)
L.dig	<i>Leuctra digitata</i> Kempný (P)	S.tor	<i>Siphonoperla torrentium</i> (P)
L.nig	<i>Leuctra nigra</i> Oldenberg (P)	T.ros	<i>Tinodes rostocki</i> MacLachlan (T)
L.pri	<i>Leuctra prima</i> Kempný (P)	W.occ	<i>Wormaldia occipitalis</i> Pictet (T)

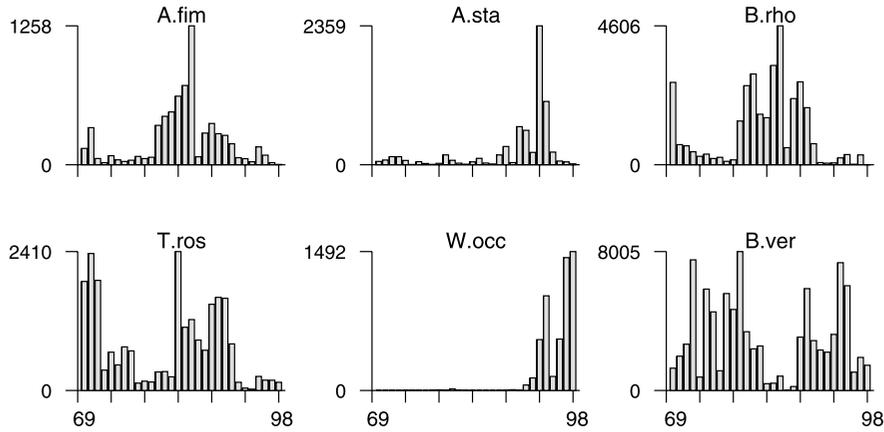


Fig. 1. Annual abundance of selected EPT species, 1969–1998; for short codes compare Table 1.

stepwise method, and genetic algorithms (Goldberg, 1989). GRNNs and linear neural networks (LNN), comparable to multiple linear regression, were used as fast-learning and easy to interpret benchmarks. Special multi-layer-perceptron (MLP), called Senso (Dapper, 1998), were used for comparison. Kohonen's SOMs served to position the centres of radial basis function networks (Bishop, 1995; Poggio and Girosi, 1990). This hybrid network was used to provide quality measures and respective visualisation and is described in Werner and Obach (2001).

For the prediction of monthly sums of individuals (y_i) found in the trap we introduce the ratio of the root mean squared errors of the model (\hat{y}_i) and the long term mean (\bar{y}_i) of the respective variable:

$$E_{\text{ltm}} = \frac{\sqrt{\sum_{i=1}^n (y_i - \hat{y}_i)^2}}{\sqrt{\sum_{i=1}^n (y_i - \bar{y}_i)^2}} \quad (1)$$

This measure is adequate for the assessment of predictions of seasonal time series as given by the emergence data of aquatic insects. Appropriate models developed with ANNs should be better than the phase average, i.e. $E_{\text{ltm}} < 1$. The substitution of \hat{y}_i by the global mean \bar{y} in Eq. (1), equal to the ratio of the RMSE and the standard deviation, describes the relative RMSE (rRMSE). It is similar to the determination coefficient R^2 .

3. Results

3.1. Visualisation of annual abundance patterns

The annual abundance patterns of 40 EPT species were compared and visualised using a SOM with a 12×10 topology. Results are displayed as a Sammon-map (Sammon, 1969) in Fig. 2. The relative positions of the nodes denominate the similarity of the abundance patterns of the species (Fig. 1). In the upper left corner, a group of dominant grazers is mapped (T.ros, B.rho, A.fim, A.fus). Along the upper right margin lie rare taxa

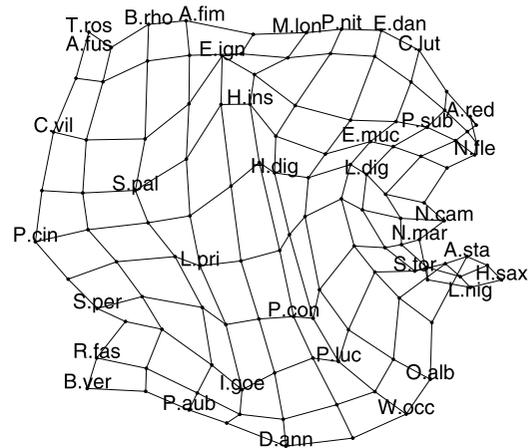


Fig. 2. Sammon-map of EPT species (annual abundance) on a SOM with a 12×10 topology; for short codes compare Table 1.

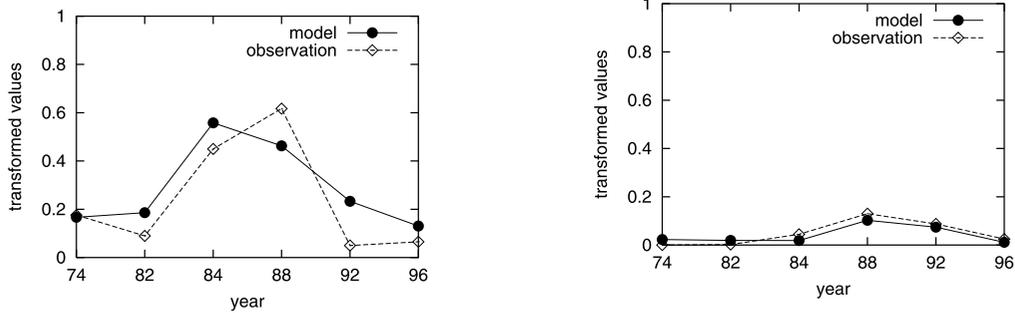


Fig. 3. Prediction of the annual abundance of *Tinodes rostocki* (left) and *Amphine-mura standfussi* (right) for six test years.

(M.lon, P.nit, E.dan) and along the right margin there is a number of fine particle feeding (C.lut to L.nig) Ephemeroptera and Plecoptera. In the lower left corner is a group of taxa with more or less stable populations over 30 years. Rare taxa with periodically high abundances are found in a small sector bordered by P.con, W.occ, P.luc, and H.dig. The following species pairs are found in identical positions as T.ros/A.fus: L.nig/N.pic, P.aub/P.int, and N.fle/P.mey; the first species superimposes the second in Fig. 2. We assume that adjacent taxa have correlated ecological requirements that lead to similar or identical predictors. However, this is not necessarily an indication of inter-specific relations.

3.2. Prediction of annual abundance with feed-forward nets

The prediction of the annual abundance of *T. rostocki* was based on the abundance of the parental generation, precipitation in April, July and December of the previous year, maximum temperature in January, and maximum discharge in June (Fig. 3).

The best linear model (RMSE = 0.1363) described by a linear neural network (LNN) for the prediction of the annual abundance of *T. rostocki* was

$$\hat{Y} = 0.73X_1 - 0.61X_2 + 0.55X_3 + 0.39X_4.$$

The most relevant predictors in the equation were the maximum monthly abundance of the parent generation (X_1), the maximum water temperature in November (X_2), precipitation in April (X_3), and

precipitation in September (X_4). Compared with others, this linear model appeared transparent.

Prediction of annual abundance of *Amphine-mura standfussi* using a GRNN resulted in an RMSE = 0.0207. This measure was low compared with *T. rostocki*, RMSE = 0.1172. In the case of *A. standfussi* only a low amount of variability is to explain in the test data, leading to a low error. From the ecological point of view, the model of *T. rostocki* is of similar quality, because data with high variability were accurately predicted. An appropriate measure for comparison is the relative RMSE, i.e. rRMSE = 0.44 for *A. standfussi* and rRMSE = 0.54 for *T. rostocki*. Both models were clearly superior than the mean of the long-term data, our minimum requirement.

Within these multidimensional models we visualised a smaller amount of information in surface graphics (Fig. 4). Two of the six input variables were altered, and the remaining four were kept constant at their mean values. The smooth surfaces of interrelations describe almost linear dependencies. It is evident that precipitation in December does not cause any significant variation in the model in contrast to precipitation in April. However, December precipitation is not necessarily unimportant if other predictors vary, that so far have remained constant (Fig. 4, left). Further, the model indicates that low precipitation in July and the abundance of the parental generation result in increased specimen numbers. The importance of precipitation in July and the abundance of the parental generation is similar, but the effect is reversed (Fig. 4, right).

To find more relevant predictors for *T. rostocki*, genetic algorithms were applied. However, the model was not improved after 8400 generations, with ten specimens per generation, a mutation rate of 0.08, a crossover rate of 0.8, and a penalty for each predictor of 0.015.

3.3. Discharge patterns and abundance

Prototypes (codebook vectors) of discharge patterns that fit the lifecycles of the respective species are visualised in Fig. 5. They are typical patterns of the respective variables modelled and distinguished by SOMs. Test data were related to these prototypes according to their smallest Euclidean distance by a 6×1 SOM.

The abundance of *Apatania fimbriata* based only on the discharge patterns during 13 months is illustrated in Fig. 6 (left). The clear differences

between codebook vectors CBV_5 and CBV_1 are evident (right) when compared with Fig. 5 (left). The discharge pattern of the years 1974 and 1992 were similar to prototype CBV_2, 1982 to CBV_5, 1984 to CBV_3, 1988 to CBV_6 and 1996 to CBV_1. No test pattern corresponded to CBV_4.

These prototype vectors of a 6×1 SOM were used as centres of an RBF network to predict the yearly abundance of *A. fimbriata*. Model and observation of *A. fimbriata* fit on the left and the right side of the diagram, but differences are distinct in the centre (Fig. 6). The prediction of test data for the years 1974 and 1992 is close to the class mean (CBV) of the training data, i.e. no improvement of the model compared to the class mean. The prediction of test data for almost every year by the class mean was slightly better than the RBF SOM model (Table 2). This depends on the

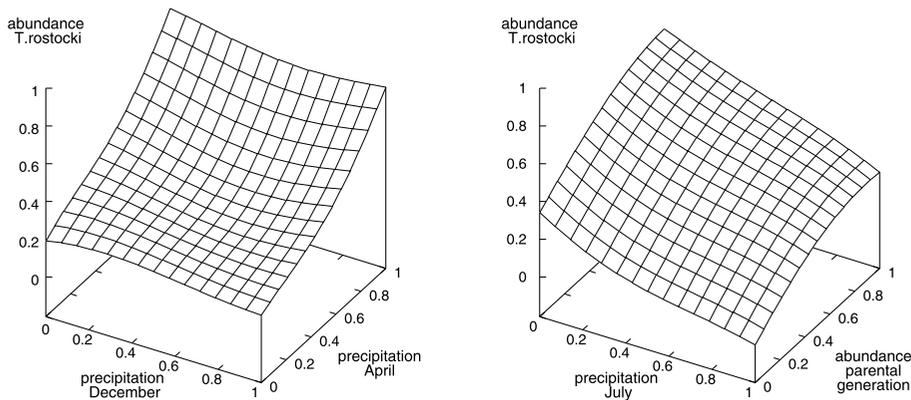


Fig. 4. Visualisation of the dependence of the annual abundance of *T. rostocki* (1969–1998) on two variables, using an RBF SOM model (data mapped onto [0,1]).

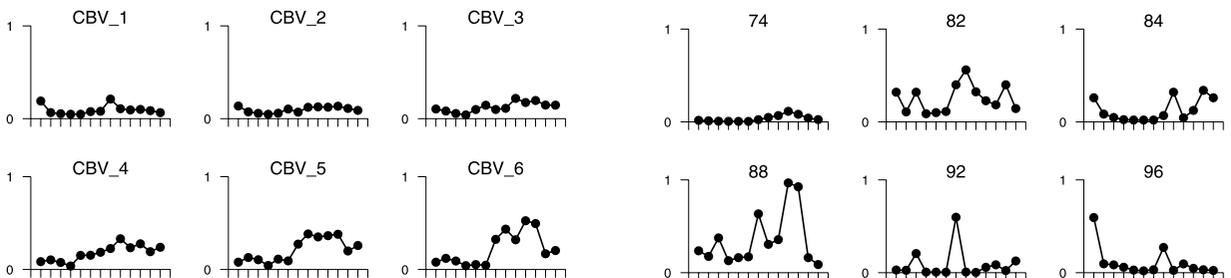


Fig. 5. Linear rescaled prototypes (codebook vectors, CBV) of discharge pattern classes during the lifecycle of *A. fimbriata* (left); input patterns of individual years (right); x-axis is a 13 month period, June to June; e.g. 74 = June 1973 to June 1974.

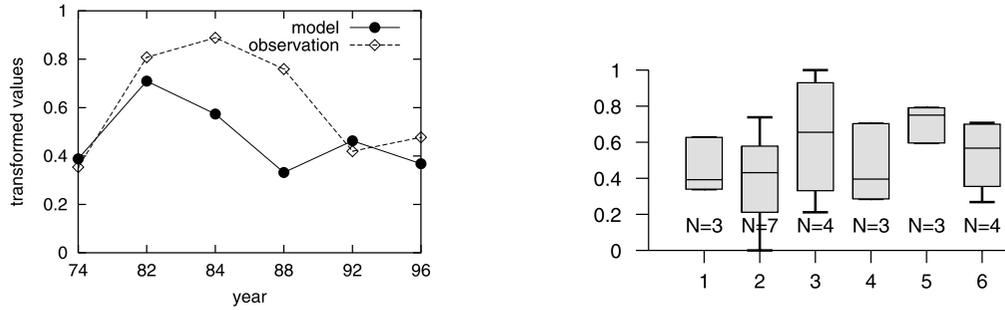


Fig. 6. Prediction of logarithmic rescaled annual abundance of *A. fimbriata*; comparison between observation and prediction (left); variability of data mapped mapped on six nodes (with CBV_1,..., CBV_6) of a SOM (right; compare Fig. 5).

Table 2

Output of RBFSOM; Testcase (TC), prediction of RBFSOM (Model), the difference between prediction and observation (Observ.), maximum activation of RBF neurones (MaxAct), class of input vector according to the SOM (Class), minimum distance (Mindist.), the number of training data assigned to that class serving as support of the network’s prediction (Supp.), the minimum (Min) and the maximum (Max) of the output values which respective inputs are assigned to that class, and their difference (Range) as a measure of variability of the data; mean of each class (MeanCls)

TC	Model	Observ.	Diff.	MaxAct	Class	Mindist.	Supp.	Range	Min	Max	MeanCls
1	0.38	0.35	0.03	0.88	2	0.17	7	0.73	0	0.73	0.40
2	0.70	0.80	-0.09	0.84	5	0.37	3	0.20	0.59	0.79	0.71
3	0.57	0.88	-0.31	0.83	3	0.26	4	0.78	0.21	10.63	
4	0.33	0.75	-0.42	0.74	6	0.57	4	0.43	0.27	0.71	0.52
5	0.46	0.42	0.04	0.69	2	0.43	7	0.73	0	0.73	0.40
6	0.36	0.47	-0.11	0.74	1	0.31	3	0.28	0.34	0.63	0.45

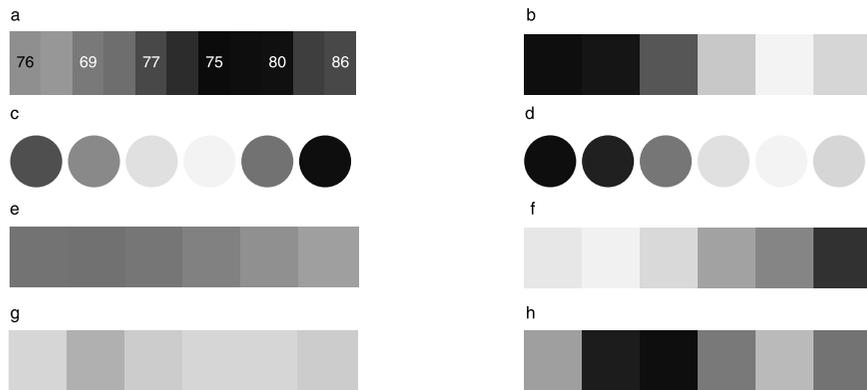


Fig. 7. Visualisation of the SOM nodes associated with CBV_1 to CBV_6 (left to right) in the first hidden layer of an RBFSOM network exemplified with data of *A. fimbriata*; U-matrix display (a); response surface of the RBFSOM (b); plane of November discharge (c); plane of February discharge (d); distance to the test year 1992 (e); activities of the RBF over the SOM for test year 1988 (f); support of the RBF centres (g); variability of the abundance of *A. fimbriata* for each class (h).

environmental data set that contained only maximum discharge. Mean or minimum values or other variables were neglected. Certainly, the variability of data in the classes can be further reduced by an optimised choice of net parameters.

Fig. 7 shows various possibilities for visualisation of the RBFSOM model. The U-matrix (Ultsch, 1993) displays the distances between the codebook vectors of neighbouring SOM nodes using grey shade gradients (a). (19)76 means the respective year fits to the codebook vector CBV_1 in Fig. 5. The darker the colour between the nodes, the greater the difference between the respective CBVs, e.g. patterns of (19)69 and (19)76 (CBV_1 and CBV_2) are more similar than of (19)80 and (19)86 (CBV_5 and CBV_6). The display of the response surface provides expected outputs in grey levels for prototype vectors (= input) of the respective nodes of the SOM (b); planes (c–d) use a grey scale to illustrate the values of the components of the codebook vectors, e.g. November discharge (component 6) is maximal with pattern four (1975), February discharge (component 9) with pattern five (1980); Fig. 7(e) displays the distances of a test case (discharge pattern 1992) to the SOM codebook vectors (CBV); activities of the RBF neurones for test year 1988 depend on the distances between prototypes (RBF centres) and the test pattern (f); support of the RBF centres is the number of training cases mapped on the respective SOM node (g), N in Fig. 6; range of the abundance of *A. fimbriata* for each class (h), whiskers in Fig. 6. The darker the node, the lower the value (c–h).

Table 2 shows the numerical output and quality measures provided by the RBFSOM model on test data.

Table 3

Error measures for the prediction of the abundance of *T. rostocki* in May, June and July with ANN types and the long term mean (P: phase) of the training data

ANN type	E_{ltm}	R_{ANN}	R_{P}	R_{May}	R_{Jun}	R_{Jul}	E_{ltmMay}	E_{ltmJun}	E_{ltmJul}
GRNN1	0.930	0.165	0.177	0.063	0.046	0.274	0.530	0.371	1.080
GRNN2	0.752	0.133	0.177	0.103	0.118	0.169	0.870	0.948	0.667
LNN	5.808	1.028	0.177	0.830	1.400	0.721	7.035	11.270	2.836
MLP	2.119	0.375	0.177	0.234	0.320	0.514	1.982	2.574	2.024

Root mean squared error R and phase error E_{ltm} are the three months mean; E_{ltm} was defined in Eq. (1).

3.4. Prediction of monthly abundances of *T. rostocki*

For the prediction of monthly abundances of *T. rostocki* during May, June, and July we compared different ANN types and error measures (Table 3).

LNN and MLP models have highest root mean squared errors R and E_{ltm} for the 3-month period and for individual months. Highest values for LNN indicate non-linear dependencies.

RMSE of the MLP is lower than of the LNN, but does not meet our requirements ($E_{\text{ltm}} > 1$). GRNNs fit our demands with $E_{\text{ltm}} < 1$. Both differed in the selection of relevant predictors with best results for May and June by GRNN1 compared with GRNN2. However, the low accuracy for the prediction of July significantly increased the total error E_{ltm} of GRNN1. The lowest R was found for the prediction of June abundance with GRNN1, i.e. an accuracy of less than 100 specimens.

4. Discussion

Prediction of ecosystems is an interesting field in research on ANNs and in ecology. Several difficulties arise with ecological systems. One is the restricted number of data, another is novelty, e.g. by environmental change or by unreliable data. However, a basic requirement is a meaningful ecological interpretation of the results.

The self-organised mapping of annual EPT data provides a comprehensive overview of similarities among abundance patterns (Fig. 2). The species are grouped in different areas of the dia-

gram. It is assumed that adjacent species have similar environmental requirements and life cycles. Functional feeding groups were partly visualised, although biological interactions may change or temporarily shift patterns. This is an indication that the variation of yearly abundance was determined by the same environmental predictors, and therefore leads to similar network models.

Prediction of monthly or yearly abundances of populations of aquatic insects is promising. Nevertheless, a basic difficulty is the overall variance that differs between taxonomic and environmental approaches. Comparing the methods used, the application of this alternative approach to the sliding window techniques reduce the amount of information by a factor of 12 (360 months to 30 years).

Environmental variability and species traits can be matched directly, but species traits sometimes are surprising trade-offs that violate assumptions of modelling efforts (e.g. Resh et al., 1994). Thus, a given habitat does not act as a single templet for all instars of a species. This may be one reason for the limited success of some models described.

Comparing the models in Fig. 3, the importance of using several measures of quality is evident. The variability in the test data may be high or low, but has to be compared meaningful (in 1986 a heavy input of insecticide into the stream (Zwick, 1993) negatively influenced species abundance, compare *A. fimbriata* in (Fig. 1)).

At first sight, results for *A. standfussi* are 'better' than for *T. rostocki*. Taking into account not only the absolute but also the relative errors, similar qualities of the models were estimated. Thus, it is proposed not to use only the 'best' error measure, but to assess model quality using relative measures in addition. Table 3 provides information on absolute and relative error measures for yearly and monthly abundance predictions of *T. rostocki*. It is evident that both types of errors provide different information. One is that the prediction for July is more difficult than for May and June. This has a negative effect on the error measure for the yearly data. It is necessary to compare monthly and yearly predictions with the long-term mean. Predictions based on the

mean are trivial; as shown above, ANN models should be distinctly better. Rescaling (e.g. logarithmic) of the data during pre-processing is convenient to influence errors.

In addition to error measures, several options exist to visualise data and the resulting models. Input data can be displayed using SOMs (Fig. 2). Improved SOMs were successfully applied to cluster input data, to visualise prototype patterns, and to compare them with test data. These processes are useful to combine self-organised learning nets and supervised trained RBF networks. This provides insight in the pattern recognition by RBF networks (Fig. 7).

The influence of two out of several input variables on the network output is visualised by three-dimensional surface plots (Fig. 4). The correlations among the three variables are obvious. However, it must be stated that all other variables are kept constant, i.e. other combinations of variables may alter the picture. The independence of input variables justifies the selected type of diagram.

If every RBF centre (SOM codebook vector) was supported by a sufficient number of training data, then local variability was successfully visualised by Box & Whisker plots (Fig. 6). The distinction of groups of patterns with good predictions and others with nonsensical predictions was meaningful for the assessment of the model's reliability and the interpretation of results. Box & Whisker plots provide more information than measures given in, e.g. Table 2, because in addition to the (outlier-sensitive) range of the output, its local distribution on input data clusters is displayed. However, due to low data density this quality measure is more adequate than that suggested by Leonard et al. (1992), and the Bayesian approach (Bishop, 1995). Class means (Table 2) as well as class medians (Fig. 6) may serve as benchmarks for comparison with the RBF predictions. They resemble predictions of Motoric Maps (Ritter et al., 1990). We visualised similarities among prototypes either as rectangular maps (3×2) or as 'chains' (6×1). The advantage of 'chains' over two-dimensional SOMs is the wider variability of prototype vectors, given a low number of nodes. The 24 training patterns limit the

application of larger SOMs. Small SOMs were used to guarantee a sufficient amount of data mapped on every single node.

The influence of the parent generation of *T. rostocki* in the linear model Eq. (2) is predominant for the success of the species in the given year. Other relevant factors describe environmental influence during larval development (July–May). Low temperatures in November are advantageous as well as high precipitation in September and April. There is a small but significant correlation among precipitation and discharge if monthly data are compared. This correlation measure probably will increase if shorter intervals and additional measures (e.g. mean, minimum, deviation of actual from long-term data) are used, or if a precipitation-discharge model is developed. Low discharge in September (month with overall lowest mean long-term discharge and standard deviation) limits available habitat, and moderate discharge in April inhibits desiccation of pupae.

Compared with *T. rostocki*, the low RMSE of the *A. standfussi* model depends on the low variance in the test data, and only one extreme observation in the training data. However, only the RMSE in relation to the long-term mean is adequate to compare the error measures of both models. It is easier to predict data with low variation close to the long-term mean than highly variable observations.

Remaining patterns in Fig. 6 (right) widely overlap. A future aim is to compute patterns among which the overlap is minimal. One strategy is to increase the numbers of SOM nodes with the risk of too few data per node. The selection of relevant predictors as proposed by the *T. rostocki* model may be advantageous.

5. Conclusion

The prediction of aquatic insect abundance with ANN models is promising. Quality measures were applied to objectify results and to weigh up the pros and cons of the different models. Reliability of the models was tested against the monthly or yearly long-term mean. The restricted number of input variables may have limited the

quality of the results. It is necessary to optimize the sampling rate and to increase the number of variables, e.g. mean and minimum discharge, and temperature. An increasing number of predictors may improve results. GRNNs with stepwise selection of important predictors are convenient for similar models, e.g. RBF networks.

Several methods of visualisation were presented, e.g. abundance patterns on a Sammon-map and comparison of prototypes of SOMs with test patterns. Box & Whisker plots are used to visualise the local variability of the output. The activation of neurones in the RBF layer, the distances of RBF centres as well as single components of vectors were mapped onto the SOM grid with scaled grey shades. Response surfaces visualised network output depending on two out of a number of variables.

Acknowledgements

The authors gratefully acknowledge the support by the German Research Foundation (Deutsche Forschungsgemeinschaft DFG), grant No. BO 1012/5-3.

References

- Bishop, C., 1995. Neural Networks for Pattern Recognition. Oxford University Press, Oxford.
- Borchardt, D., Schleiter, I.M., Werner, H., Dapper, T., Schmidt, K.-D., 1997. Modellierung ökologischer Zusammenhänge in Fließgewässern mit Neuronalen Netzwerken. Wasser und Boden, 49 (8), 38, 47–50.
- Borchardt, D., Dapper, T., Obach, M. et al., 2001. Abschlußbericht zum Forschungsvorhaben We 959/5-2, given to the German Research Foundation (Deutsche Forschungsgemeinschaft), submitted for publication.
- Chon, T.-S., Park, Y.S., Moon, K.H., Cha, E.Y., 1996. Patternizing communities by using an artificial neural network. Ecol. Modell. 90, 69–78.
- Cox, E.J., 1990. Studies on the algae of a small softwater stream. I: Occurrence and distribution with particular reference to the diatoms. Arch. Hydrobiol. Suppl. 83, 525–552.
- Dapper, T., 1998. Dimensionsreduzierende Vorverarbeitungen für Neuronale Netze mit Anwendungen in der Gewässerökologie. Dissertation im Fachbereich Mathematik/Informatik der Universität Gh Kassel. Berichte aus der Informatik D 34, Shaker Verlag, Aachen, p. 234.

- Foody, G.M., 1999. Applications of the self-organising feature map neural network in community data analysis. *Ecol. Modell.* 120, 97–107.
- Goldberg, D.E., 1989. *Genetic Algorithms in Search Optimization and Machine Learning*, first ed. Addison-Wesley, Reading MA, p. 412.
- Kohonen, T., 1995. *Self-Organizing Maps*. Springer, Heidelberg.
- Leonard, J.A., Kramer, M.A., Ungar, L.H., 1992. Using radial basis functions to approximate a function and its error bounds. *IEEE Trans. Neural Netw.* 3, 624–627.
- Obach, M., 1998. *Anwendung Statistischer Methoden und Künstlicher Neuronaler Netzwerke im Vergleich*. Diplomarbeit im Fachbereich Mathematik/Informatik der Universität Gesamthochschule Kassel.
- Palmer, M.A., Poff, N.L., 1997. The influence of environmental heterogeneity on patterns and processes in streams. *J. N. Am. Benthol. Soc.* 16, 169–173.
- Poff, N.L., 1992. Why disturbance can be predictable: a perspective on the definition of disturbance in streams. *J. N. Am. Benthol. Soc.* 11, 86–92.
- Poggio, T., Girosi, F., 1990. Networks for approximation and learning. *Pro. IEEE* 78 (9), 1481–1497.
- Reice, S.R., 1985. Experimental disturbance and the maintenance of diversity in a stream community. *Oecologia* 67, 90–97.
- Resh, V.H., Hildrew, A.G., Stutzner, B., Townsend, C.R., 1994. Theoretical habitat templets, species traits, and species richness: a synthesis of long-term research on the Upper Rhône River in the context of concurrently developed ecological theory. *Freshw. Biol.* 31, 539–554.
- Ringe, F., 1974. Die Chironomiden-Emergenz 1970 in Breitenbach und Rohrwiesenbach. *Arch. Hydrobiol.* 45 (Suppl.), 212–304.
- Ritter, H., Martinetz, T., Schulzen, K., 1990. *Neuronale Netze: Eine Einführung in die Neuroinformatik selbstorganisierter Netzwerke*. Addison-Wesley, Bonn.
- Sammon, J.W., 1969. A nonlinear mapping for data structure analysis. *IEEE Trans. Comput.* C-185, 401–409.
- Schleiter, I.M., Borchardt, D., Wagner, R., Dapper, T., Schmidt, K.-D., Schmidt, H.-H., Werner, H., 1999. Modelling water quality bioindication and population dynamics in lotic ecosystems using neural networks. *Ecol. Modell.* 120, 271–286.
- Schleiter, I.M., Obach, M., Borchardt, D., Werner, H., 2001a. Bioindication of chemical and hydromorphological habitat characteristics with benthic macro-invertebrates based on artificial neural networks. *Aquat. Ecol.*, in print.
- Schleiter, I.M., Obach, M., Borchardt, D., Werner, H., 2001b. Prediction of running water properties using radial basis function self-organising maps combined with input relevance detection, Second International Conference on Applications of Machine Learning to Ecological Modelling, Adelaide, Australia, 27th Nov.–1st Dec. 2000, in preparation.
- Southwood, T.R.E., 1988. Tactics, strategies and templets. *Oikos* 52, 3–18.
- Specht, D.F., 1991. A general regression network. *IEEE Trans. Neural Netw.* 6, 568–576.
- Townsend, C.R., 1989. The patch dynamics concept of stream community ecology. *J. N. Am. Benthol. Soc.* 8, 36–50.
- Townsend, R., Hildrew, A., 1994. Species traits in relation to a habitat templet for river systems. *Freshw. Biol.* 31, 265–275.
- Ultsch, A., 1993. Self organized feature maps for monitoring and knowledge acquisition of a chemical process. In: Gielen, S., Kappen, B. (Eds.), *Proceedings of the International Conference on Artificial Neural Networks ICANN*, vol. 93. Springer, London, pp. 864–867.
- Wagner, R., Schmidt, H.-H., Marxsen, J., 1994. The hyporheic habitat of the Breitenbach, spatial structure and physico-chemical conditions as a basis for benthic life. *Limnologica* 23, 285–294.
- Wagner, R., Dapper, T., Schmidt, H.-H., 2000. The influence of environmental variables on the abundance of aquatic insects—a comparison of ordination and artificial neural networks. *Hydrobiologia*, 422/423, pp. 143–152.
- Werner, H., Obach, M., 2001. New neural network types estimating the accuracy of response for ecological modelling. Second International Conference on Applications of Machine Learning to Ecological Modelling, Adelaide, Australia, 27th Nov.–1st Dec. 2000, *Ecol. Model.* 145.
- Zwick, P., 1993. Fließgewässergefährdung durch Insektizide. *Naturwissenschaften* 79, 437–442.