

FLOWNORM 2.0

- **a Visual Basic program for computing riverine loads of substances and extracting anthropogenic signals from time series of load data**

User's Manual

2004-04-13

**Anders Grimvall
Department of Mathematics
Linköping University**

Introduction

Monitoring of river water quality and discharge plays a key role in many water management programmes. In particular, such monitoring can provide important information about the loads of nutrients and toxic substances that are carried by rivers to the sea. However, increases or decreases in riverine loads must be interpreted with care, because the natural variation between years can be very large. Both the leaching of dissolved substances from soil to water and the transport of particulate matter can increase significantly during periods of high precipitation. Moreover, runoff, temperature and light conditions can strongly influence processes in the river itself. This calls for efficient procedures to separate human impact from natural variation in the collected data.

The Visual Basic programme FLOWNORM has been developed to facilitate:

- **calculation of riverine loads of substances;**
- **extraction of anthropogenic signals from time series of riverine loads.**

FLOWNORM, Version 2.0, consists of five Visual Basic macros having the names and functions listed below. Two Excel worksheets containing concentration and flow data for an arbitrary number of sampling sites form the starting point for an analysis. The final result consists of time series of **normalised loads**, i.e. riverine loads that have been adjusted to remove natural fluctuations and clarify anthropogenic impacts.

Macro	Function
auditdailydata	Identify illegal entries in raw data
matchflowandcons	Define pairs of flow and concentration data for load calculations
computeloads	Compute monthly and annual loads from flow and concentration data ordered by date
definenormalisationmodels	Define response variables and normalisation models
flownormalise	Select normalisation models by cross-validation and compute monthly and annual normalised loads

Each macro operates on predefined worksheets for inputs and outputs. The table below shows which worksheets that are used for the different macros. Further details are given in the documentation of each macro.

Macro	Input worksheets	Output worksheets
auditdailydata	"Concentration by date" or "Flow by date"	The same as input, and "Concentration data summary" or "Flow data summary"
matchflowandcons	"Concentration by date" and "Flow by date"	"Matched pairs"
computeloads	"Matched pairs", "Flow by date", and "Concentration by date"	"Annual totals" and "Seasonal totals"
definormalisationmodels	"Seasonal totals"	"Normalisation models"
flownormalise	"Normalisation models"	"Normalisation models", "Normalised seasonal totals", and "Normalised annual totals"

Theoretical background

Estimation of monthly and annual riverine loads

Monthly and annual riverine loads are calculated by first expanding the time series of observed concentration and flow data to complete series of daily data and then summing daily values of the product of concentration and water discharge. The expanded values are computed by connecting observed values with straight lines.

Normalisation of riverine loads

Monthly values of riverine loads are normalised by employing parametric and semiparametric regression models to remove or suppress the temporal variation that can be attributed to fluctuations in water discharge, water temperature, salinity, or other indicators of natural variation.

The semiparametric normalisation model has the general form

$$y_{ij} = \alpha_{ij} + \beta_{1,j}x_{1,ij} + \dots + \beta_{p,j}x_{p,ij} + \varepsilon_{ij}, \quad i = 1, \dots, n \quad j = 1, \dots, m$$

where y_{ij} is the observed response for the j th month of the i th year, $x_{k,ij}$, $k=1, \dots, p$ represent contemporaneous values of p explanatory variables, and ε_{ij} is a random error term with mean zero. The slope parameters ($\beta_{k,j}$, $k=1, \dots, p$) are permitted to vary with the season (j) under consideration, and the intercept (α_{ij}) is permitted to vary with both season (j) and year (i). However, rapid changes in the intercept are controlled by so-called roughness penalty factors (λ_1 and λ_2), and the intercept and slope parameters are estimated by minimising the expression

$$S(\alpha, \beta) = \sum_{i,j} (y_{ij} - \alpha_{ij} - \beta_{1,j} x_{1,ij} - \dots - \beta_{p,j} x_{p,ij})^2 + \lambda_1 \sum_{i,j} (\alpha_{ij} - \frac{\alpha_{i+1,j} + \alpha_{i-1,j}}{2})^2 + \lambda_2 \sum_{i,j} (\alpha_{ij} - \frac{\alpha_{i,j-1} + \alpha_{i,j+1}}{2})^2,$$

where the first sum ranges over all values of i and j for which both the response variable and the explanatory variables have been observed. A univariate form of this model was first used by Stålnacke and co-workers (1999), and detailed information about algorithms for parameter estimation has been published by Stålnacke and Grimvall (2001).

The parametric normalisation model is an ordinary multiple regression model of the general form

$$y_{ij} = \alpha + \beta_{1,j} x_{1,ij} + \dots + \beta_{p,j} x_{p,ij} + \varepsilon_{ij}, \quad i = 1, \dots, n \quad j = 1, \dots, m$$

where y_{ij} , $x_{k,ij}$, and ε_{ij} have the same meaning as above. The slope parameters ($\beta_{k,j}$, $k=1, \dots, p$) are permitted to vary with the season (j) under consideration, whereas the intercept is assumed to be constant.

Selection of roughness penalty factors and assessment of the predictive ability of the tested normalisation models

The penalty factors λ_1 and λ_2 are determined by cross-validation. With this technique, the entire data set is separated into an estimation set (or training set) and a test set. The model is first fitted to the estimation set and is subsequently used to predict the observations in the test set, that is, the values that have been left out of the estimation step. If the observation period covers m years, we define m estimation sets M_i , $i=1, \dots, m$ by leaving out one-year-long blocks of observations, and then we compute a so-called PRESS-value (i.e., a sum of squared prediction errors):

$$S(\lambda_1, \lambda_2) = \sum_i \sum_{(i,j) \notin M_i} (y_{ij} - \hat{\alpha}_{ij} - \hat{\beta}_{1,j} x_{1,ij} - \dots - \hat{\beta}_{p,j} x_{p,ij})^2.$$

Finally, the factors λ_1 and λ_2 are selected in such a way that $S(\lambda_1, \lambda_2)$ is minimised, and the corresponding Root Mean PRESS value

$$\min \left\{ \sqrt{\frac{1}{N} S(\lambda_1, \lambda_2)}; \lambda_1 > 0, \lambda_2 > 0 \right\}$$

is used as a measure of the predictive ability of the normalisation model under consideration.

Literature references

Stålnacke P, Grimvall, A., Sundblad, K., and Wilander, A.: 1999, 'Trends in nitrogen transport in Swedish rivers', *Environ. Monit. Assess.* **59**, 47-72.

Stålnacke, P., and Grimvall, A.: 2001, 'Semiparametric approaches to flow-normalisation and source apportionment of substance transport in rivers', *Environmetrics* **12**, 233-250.

Input data

The concentration and flow data that shall be analysed are pasted on the worksheets 'Concentration by date' and 'Flow by date', respectively. Figure 1 shows how the input data shall be organised for the worksheet 'Concentration by date'. For each sampling site, a date column shall be followed by an arbitrary number of columns containing the values of the monitored variables. Strings found in the row immediately above the first observations (row 3 in Figure 1) are interpreted as variable names. Sampling site names are extracted from the date columns two lines above the first observation. Missing values shall be entered as empty cells. 'Less-than-values' of the form <0.05 are permitted.

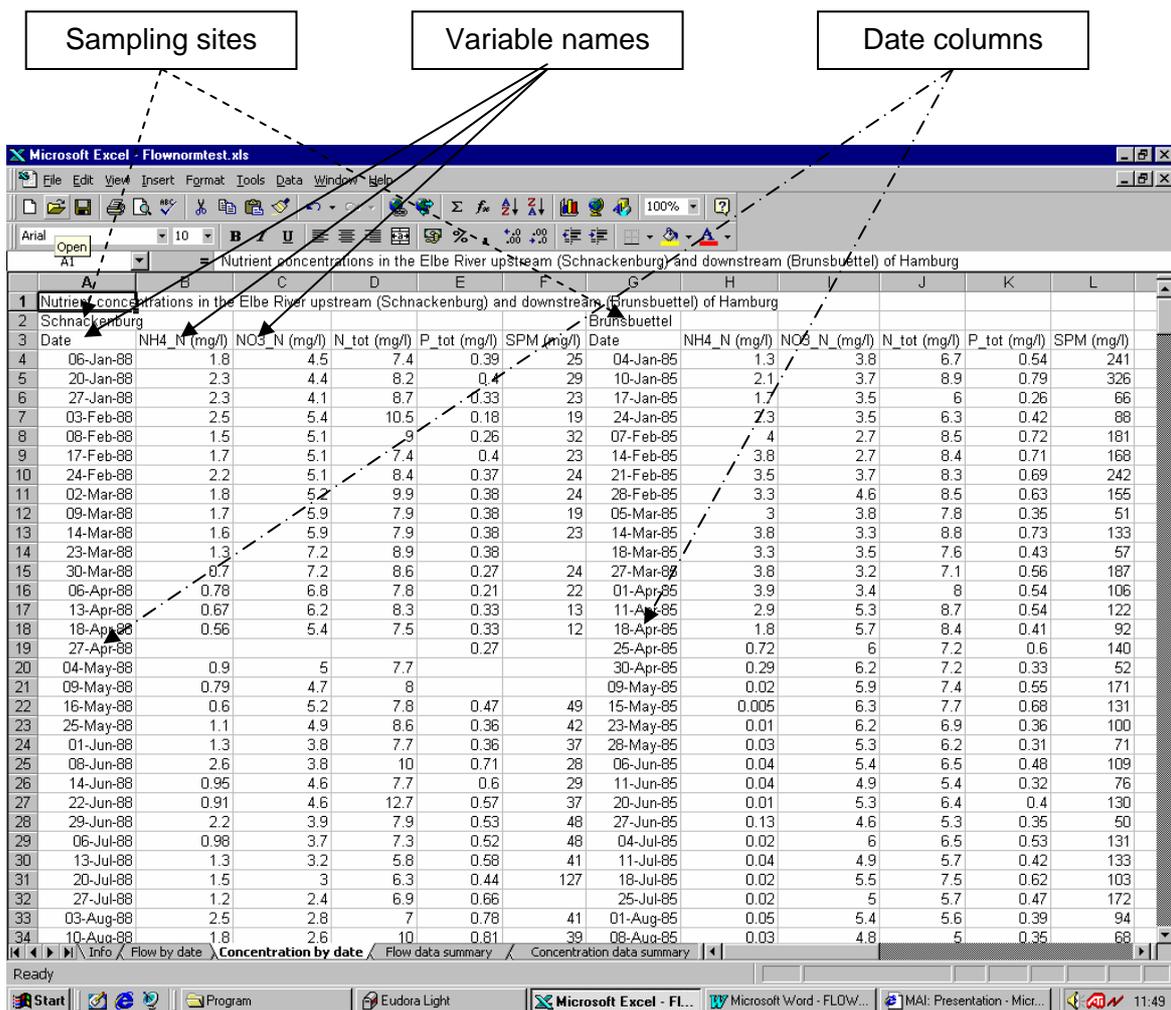


Figure 1. Worksheet 'Concentration by date' for input of concentration data.

Flow data are organised analogously. However, only one column of flow-data is permitted for each sampling site name.

Auditdailydata

This macro checks that the data pasted on the worksheets 'Concentration by date' and 'Flow by date' are of correct type and properly organised. Missing values shall be entered as empty cells. 'Less-than-values' are permitted.

When 'Auditdailydata' is run, the macro first identifies the name of active worksheet ('Concentration by date' or 'Flow by date'), and then it searches for cells containing date values. To be more precise, the macro searches for the first row containing date values and then uses the cell values in this row to determine which of the columns that shall be regarded as date columns. In figure 2, column A and E will be identified as date columns.

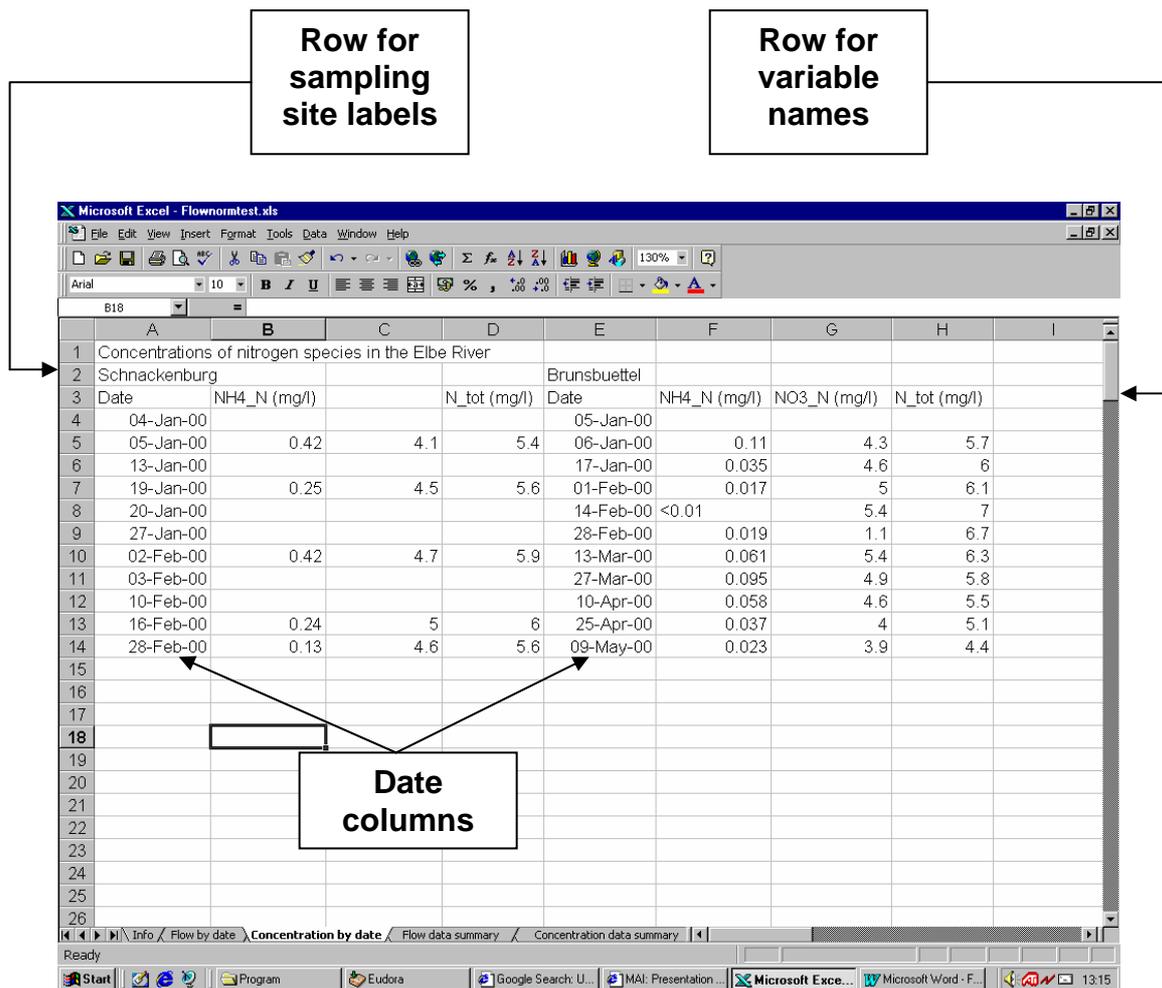


Figure 2. Worksheet 'Concentration by date' before the macro 'Auditdailydata' is run.

Variable names are extracted from the row immediately above the first row containing date values. If no variable name has been assigned to a variable (see cell C3 in Figure 2), the macro assigns a default name. Names of sampling sites are extracted from cells located in date columns and two rows above the first row with date values (cells A2 and E2 in figure 2). A maximum of two rows above the sampling site names can be used for comments.

Figure 3 shows the results obtained by letting 'Auditdailydata' operate on the worksheet in Figure 2. A default name for the variable in column C has been added, and an extra line for comments has been inserted. The number of observations for each of the identified variables is shown on the worksheet "Concentration data summary".

1	A	B	C	D	E	F	G	H	I
2	Concentrations of nitrogen species in the Elbe River								
3	Schnackenburg				Brunsbuettel				
4	Date	NH4_N (mg/l)	Var2	N_tot (mg/l)	Date	NH4_N (mg/l)	NO3_N (mg/l)	N_tot (mg/l)	
5	04-Jan-00				05-Jan-00				
6	05-Jan-00	0.42	4.1	5.4	06-Jan-00	0.11	4.3	5.7	
7	13-Jan-00				17-Jan-00	0.035	4.6	6	
8	19-Jan-00	0.25	4.5	5.6	01-Feb-00	0.017	5	6.1	
9	20-Jan-00				14-Feb-00	<0.01	5.4	7	
10	27-Jan-00				28-Feb-00	0.019	1.1	6.7	
11	02-Feb-00	0.42	4.7	5.9	13-Mar-00	0.061	5.4	6.3	
12	03-Feb-00				27-Mar-00	0.095	4.9	5.8	
13	10-Feb-00				10-Apr-00	0.058	4.6	5.5	
14	16-Feb-00	0.24	5	6	25-Apr-00	0.037	4	5.1	
15	28-Feb-00	0.13	4.6	5.6	09-May-00	0.023	3.9	4.4	
16									
17									
18									
19									
20									
21									
22									
23									
24									
25									
26									

Figure 3. Worksheet 'Concentration by date' after the macro 'Auditdailydata' has been run.

Matchflowandconc

This macro aims to facilitate the matching of flow and concentration data for load calculations. Based on the identified names of the sampling sites for flow and concentration the macro prints a preliminary list of matched pairs on the worksheet 'Matched pairs'. This list can then be edited prior to the load calculations.

Figures 4-6 show an example of input worksheets to the macro and the output worksheet before editing. As can be seen, the macro has identified two sets of flowdata (daily data and midweek data) from NeuDarchau on the Elbe. Furthermore, concentration data has been found for two sites: Schnackenburg and Brunsbuettel. Suppose that we are only interested in calculating loads using midweek discharge data. We can then edit worksheet 'Matched pairs' as shown in Figure 7, where the pairs (NeuDarchau_midweek, Schnackenburg) and (NeuDarchau_midweek, Brunsbuettel) have been listed.

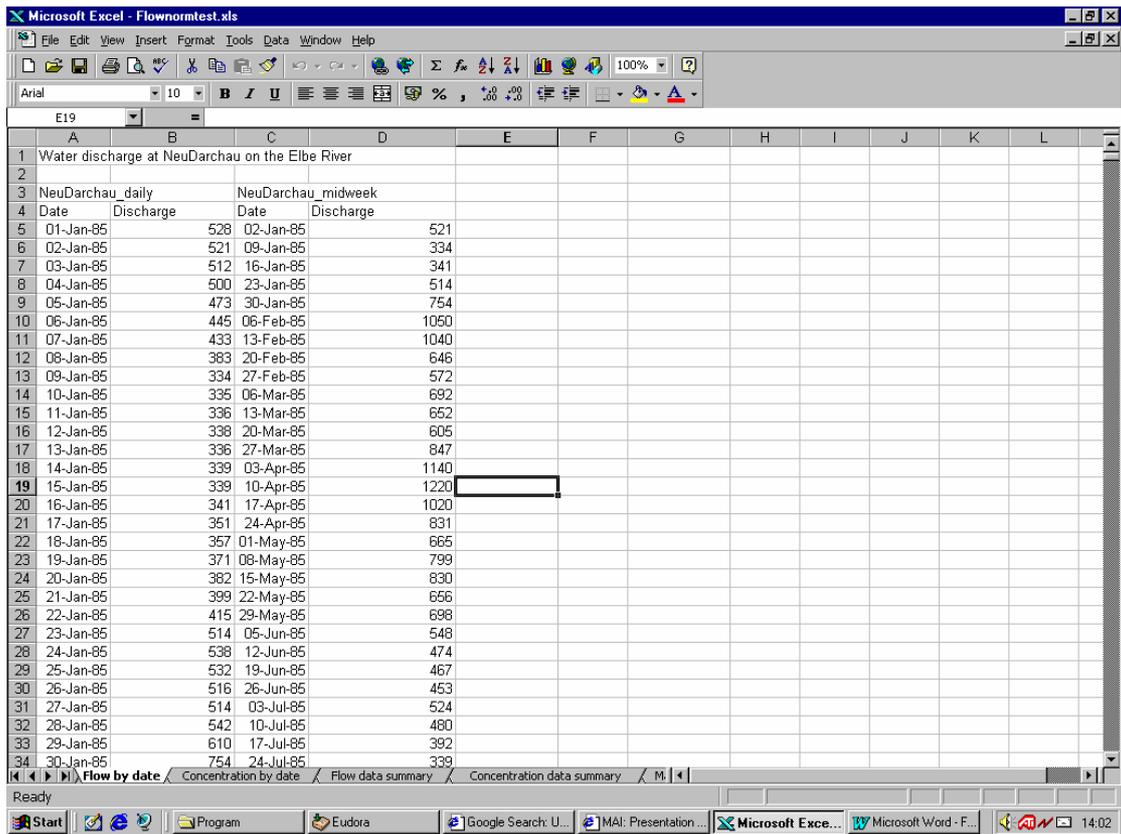


Figure 4. Flow data on worksheet 'Flow by date'

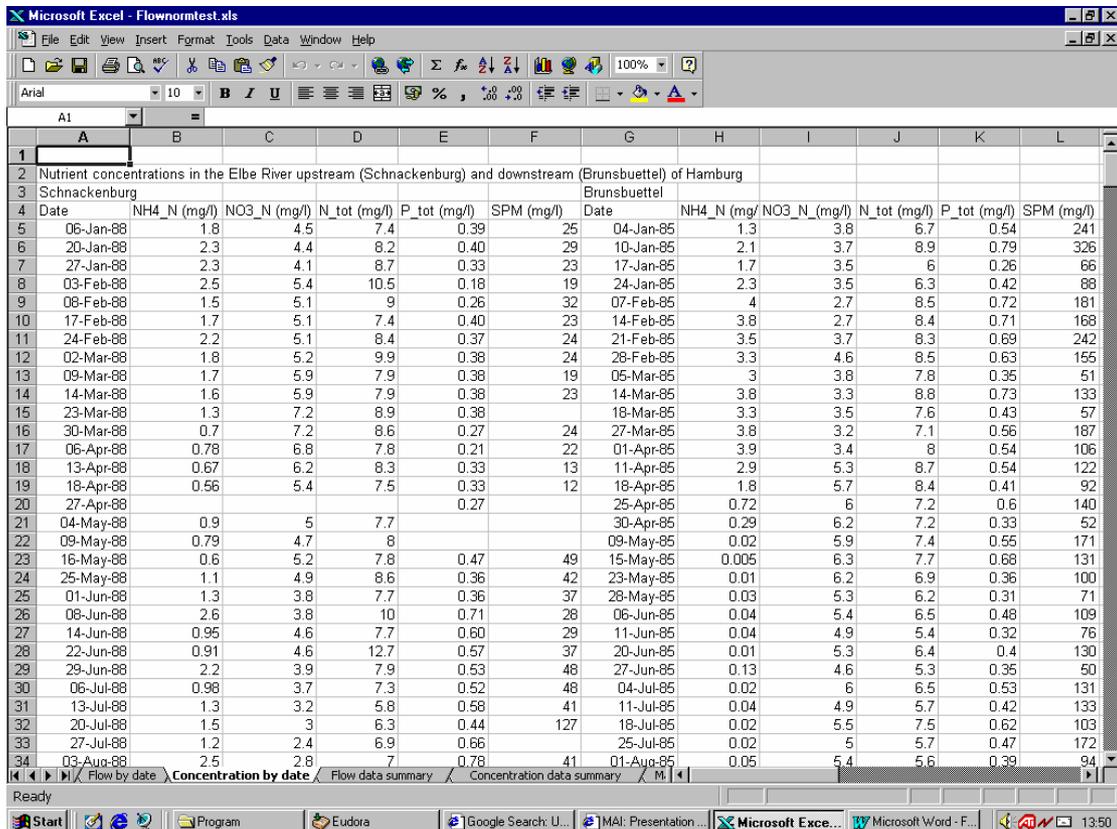


Figure 5. Concentration data on worksheet 'Concentration by date'

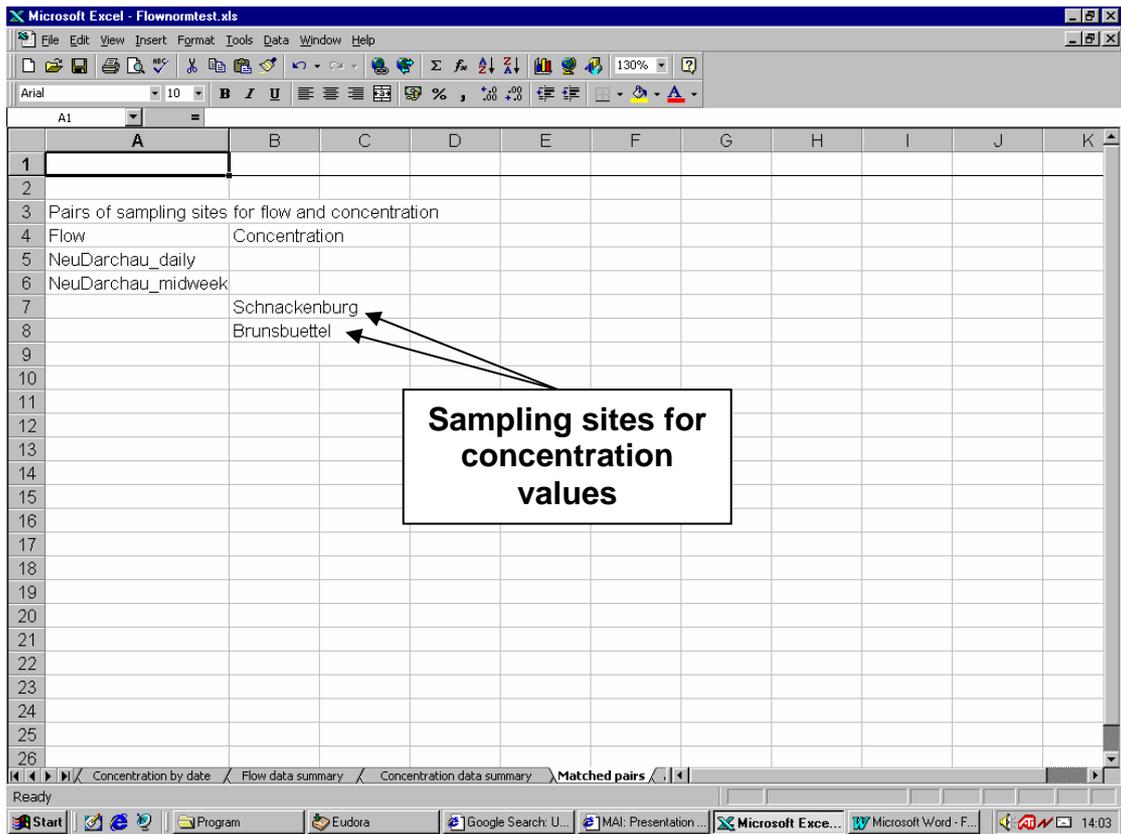


Figure 6. Identified sampling sites for flow and concentration.

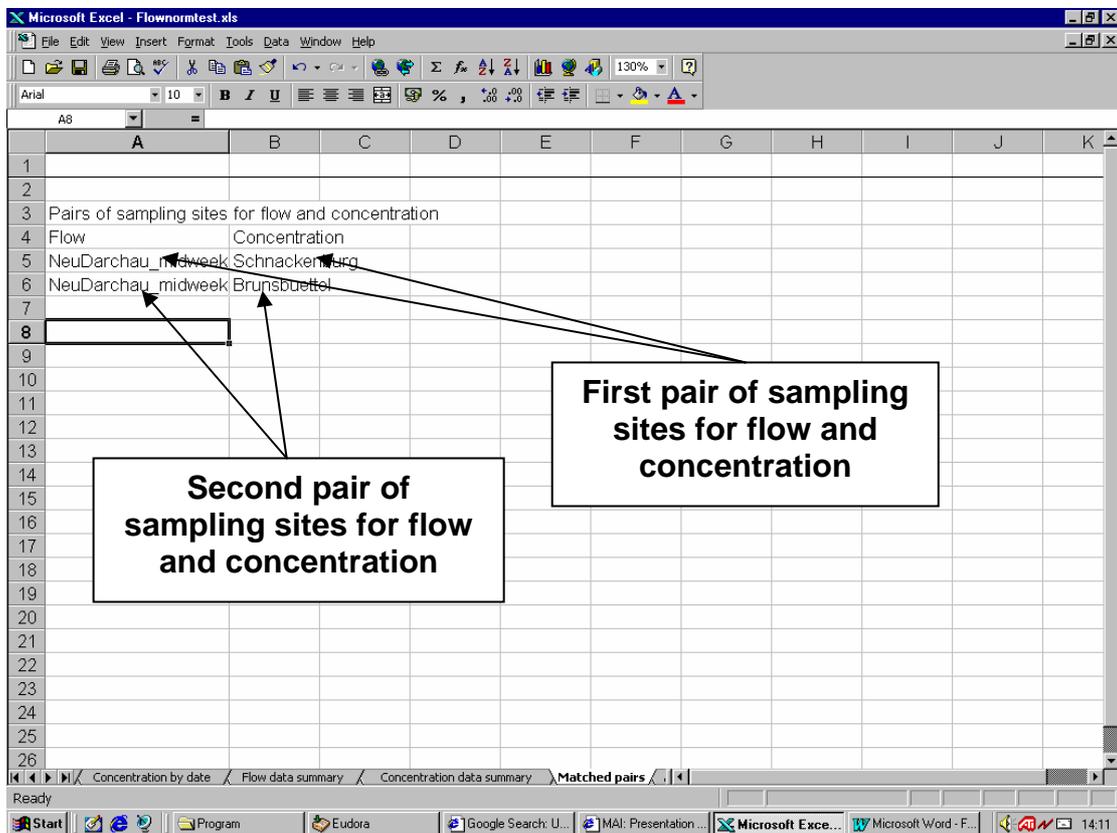


Figure 7. Worksheet 'Matched pairs' after editing

Computeloads

This macro operates on the worksheets 'Flow by date', 'Concentration by date' and 'Matched pairs'. The output worksheets are 'Annual totals' and 'Monthly totals'. Figure 8 and 9 illustrate the output obtained for the input shown in Figures 4, 5 and 7.

Because no concentration data were available for Schnackenburg in 1985-87, no substance loads were calculated for that period (see Figure 8). To avoid misleading extrapolations and interpolations of observed data, the following general rules are applied:

- the calculation of substance loads starts the first month for which both flow and concentration data are available and ends the last month for which such data exist;
- gaps in flow data are filled by interpolated values provided that the gap does not contain a full month without any observations;
- gaps in concentration data are filled by interpolated values provided that the gap contains a maximum of two full months without any observations;
- annual loads are only given for years with a complete set of monthly loads.

'Less-than-values' are replaced by a fixed percentage of the detection limit. The user is asked to enter the desired percentage when the macro is run.

If the water discharge is expressed in m^3/s and the concentrations in mg/l , the riverine loads will be expressed in ton/month or ton/year . The monthly and annual totals for the water discharge are expressed in 10^9 m^3 per month and year, respectively.

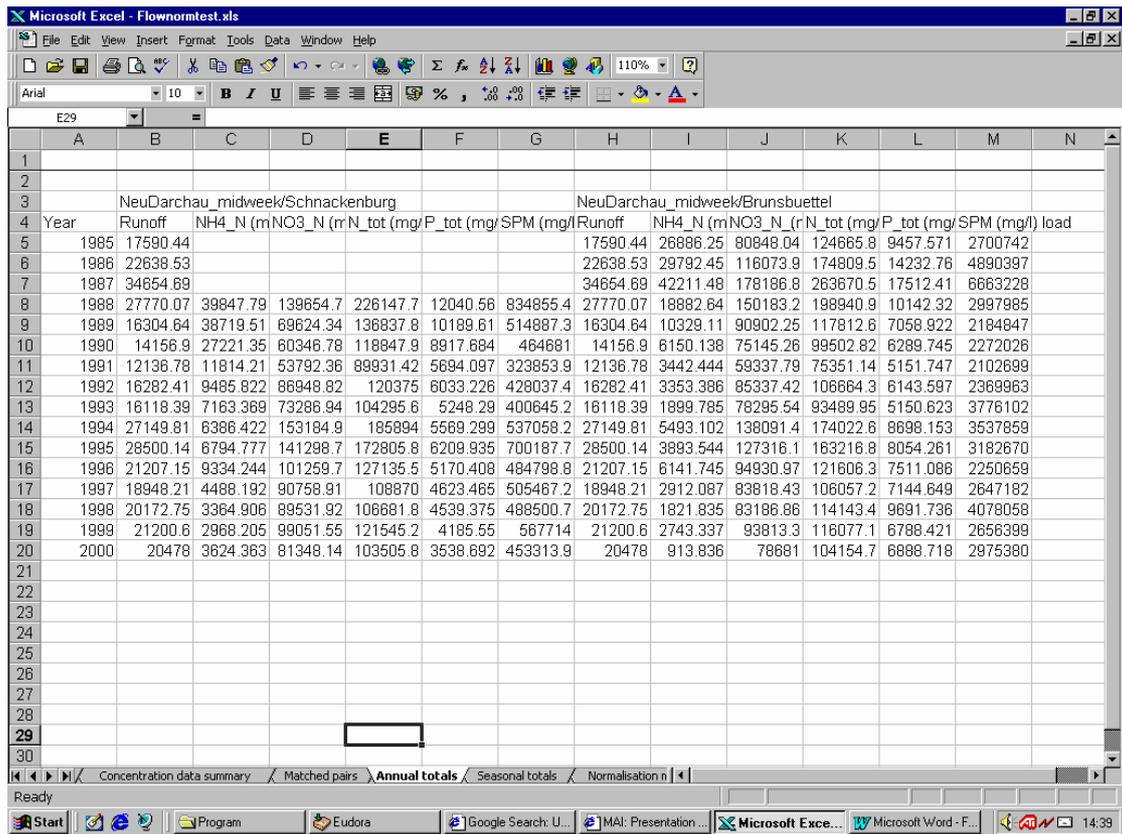


Figure 8. Annual loads of the substances monitored at Schnackenburg and Brunsbuettel.

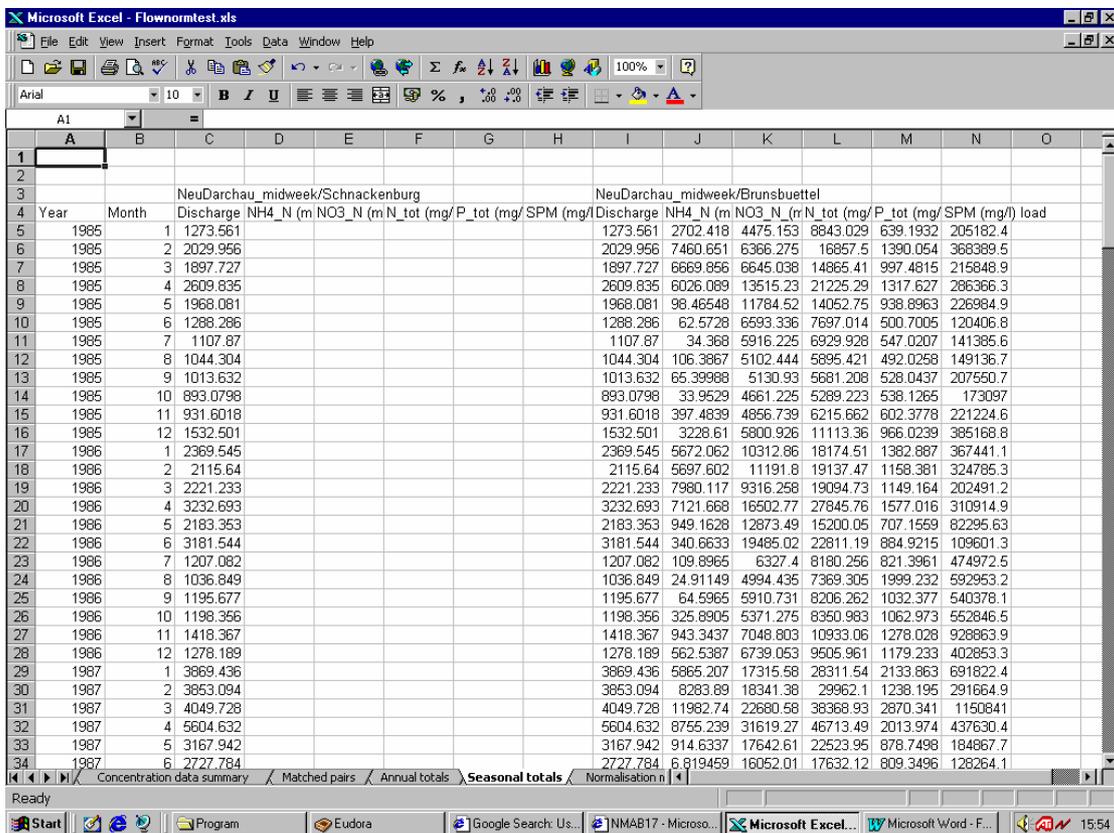


Figure 9. Monthly loads of the substances monitored at Schnackenburg and Brunsbuettel.

Definenormalisationmodels

This macro, which aims to facilitate the formulation of normalisation models, operates on the worksheet 'Seasonal totals'. Before the macro is run, response variables and explanatory variables shall be defined by typing 'y' or 'x' two rows above the variable names in this worksheet (see Figure 10).

Figure 11 shows a list of possible normalisation models. This list can be edited prior to the normalisation. For example we can decide to use the list in Figure 12 instead.

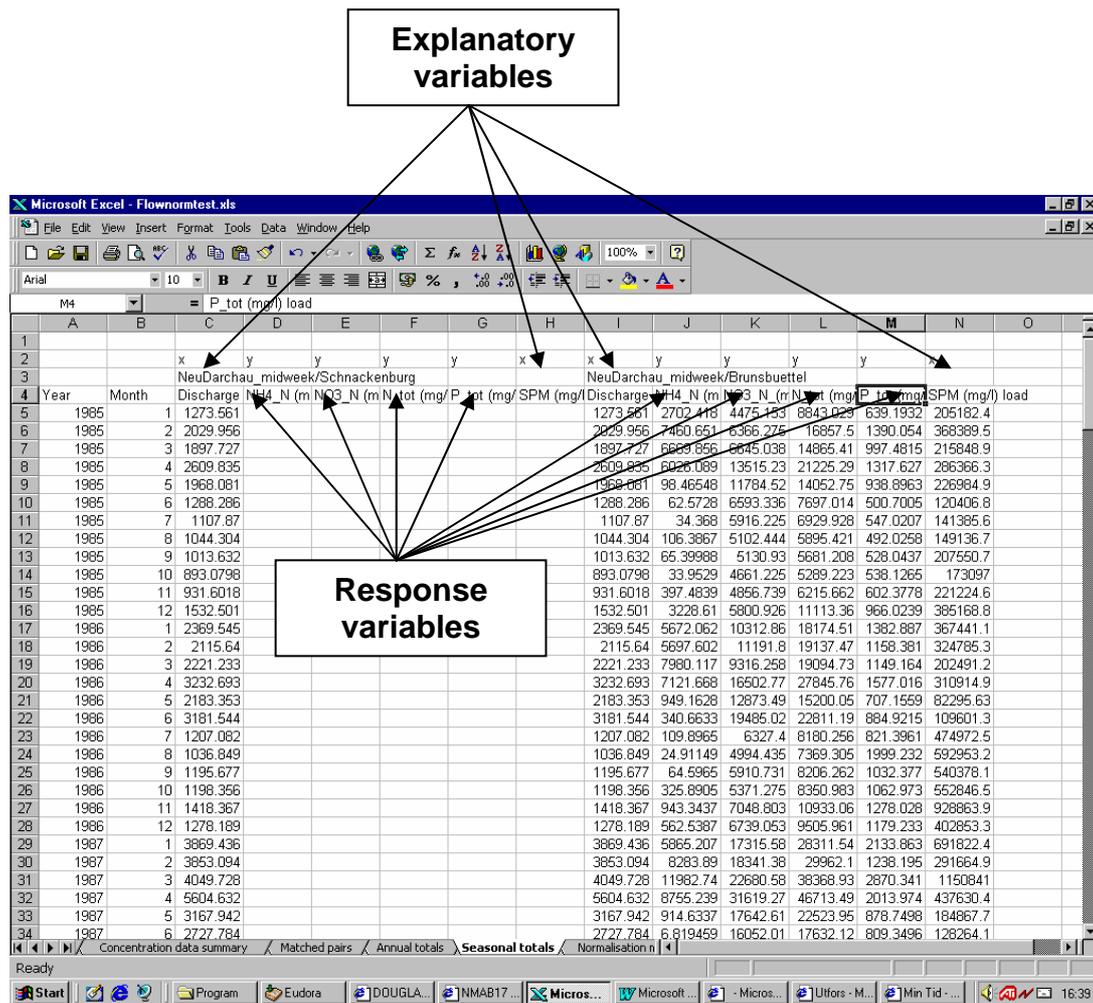


Figure 10. Response variables and explanatory variables indicated by 'y' and 'x', respectively, two rows above the variable names.

	A	B	C	D	E	F	G
1							
2							
3							
4	Sampling site	Response variable	Explanatory variable 1	Explanatory variable 2			
5	NeuDarchau_midweek/Schnackenburg	NH4_N (mg/l) load	Discharge				
6	NeuDarchau_midweek/Schnackenburg	NH4_N (mg/l) load	SPM (mg/l) load				
7	NeuDarchau_midweek/Schnackenburg	NH4_N (mg/l) load	Discharge	SPM (mg/l) load			
8	NeuDarchau_midweek/Schnackenburg	NO3_N (mg/l) load	Discharge				
9	NeuDarchau_midweek/Schnackenburg	NO3_N (mg/l) load	SPM (mg/l) load				
10	NeuDarchau_midweek/Schnackenburg	NO3_N (mg/l) load	Discharge	SPM (mg/l) load			
11	NeuDarchau_midweek/Schnackenburg	N_tot (mg/l) load	Discharge				
12	NeuDarchau_midweek/Schnackenburg	N_tot (mg/l) load	SPM (mg/l) load				
13	NeuDarchau_midweek/Schnackenburg	N_tot (mg/l) load	Discharge	SPM (mg/l) load			
14	NeuDarchau_midweek/Schnackenburg	P_tot (mg/l) load	Discharge				
15	NeuDarchau_midweek/Schnackenburg	P_tot (mg/l) load	SPM (mg/l) load				
16	NeuDarchau_midweek/Schnackenburg	P_tot (mg/l) load	Discharge	SPM (mg/l) load			
17	NeuDarchau_midweek/Brunsbuettel	NH4_N (mg/l) load	Discharge				
18	NeuDarchau_midweek/Brunsbuettel	NH4_N (mg/l) load	SPM (mg/l) load				
19	NeuDarchau_midweek/Brunsbuettel	NH4_N (mg/l) load	Discharge	SPM (mg/l) load			
20	NeuDarchau_midweek/Brunsbuettel	NO3_N (mg/l) load	Discharge				
21	NeuDarchau_midweek/Brunsbuettel	NO3_N (mg/l) load	SPM (mg/l) load				
22	NeuDarchau_midweek/Brunsbuettel	NO3_N (mg/l) load	Discharge	SPM (mg/l) load			
23	NeuDarchau_midweek/Brunsbuettel	N_tot (mg/l) load	Discharge				
24	NeuDarchau_midweek/Brunsbuettel	N_tot (mg/l) load	SPM (mg/l) load				
25	NeuDarchau_midweek/Brunsbuettel	N_tot (mg/l) load	Discharge	SPM (mg/l) load			
26	NeuDarchau_midweek/Brunsbuettel	P_tot (mg/l) load	Discharge				
27	NeuDarchau_midweek/Brunsbuettel	P_tot (mg/l) load	SPM (mg/l) load				
28	NeuDarchau_midweek/Brunsbuettel	P_tot (mg/l) load	Discharge	SPM (mg/l) load			
29							

Figure 11. List of normalisation models generated by the macro 'Definenormalisationmodels'.

	A	B	C	D	E	F	G
1							
2							
3							
4	Sampling site	Response variable	Explanatory variable 1	Explanatory variable 2			
5	NeuDarchau_midweek/Schnackenburg	NH4_N (mg/l) load	Discharge				
6	NeuDarchau_midweek/Schnackenburg	NO3_N (mg/l) load	Discharge				
7	NeuDarchau_midweek/Schnackenburg	N_tot (mg/l) load	Discharge				
8	NeuDarchau_midweek/Schnackenburg	P_tot (mg/l) load	Discharge				
9	NeuDarchau_midweek/Schnackenburg	P_tot (mg/l) load	SPM (mg/l) load				
10	NeuDarchau_midweek/Schnackenburg	P_tot (mg/l) load	Discharge	SPM (mg/l) load			
11	NeuDarchau_midweek/Brunsbuettel	NH4_N (mg/l) load	Discharge				
12	NeuDarchau_midweek/Brunsbuettel	NO3_N (mg/l) load	Discharge				
13	NeuDarchau_midweek/Brunsbuettel	N_tot (mg/l) load	Discharge				
14	NeuDarchau_midweek/Brunsbuettel	P_tot (mg/l) load	Discharge				
15	NeuDarchau_midweek/Brunsbuettel	P_tot (mg/l) load	SPM (mg/l) load				
16	NeuDarchau_midweek/Brunsbuettel	P_tot (mg/l) load	Discharge	SPM (mg/l) load			
17							
18							
19							
20							
21							
22							
23							
24							
25							
26							
27							
28							
29							

Figure 12. Edited list of normalisation models.

Flownormalise

This macro aims to remove or suppress the natural variation in monthly riverine loads. The response variable and the explanatory variables in each of the tested normalisation models are listed on the worksheet 'Normalisation models'. Input data are read from the worksheet 'Seasonal totals', and the outputs of the macro are printed on the worksheets 'Normalised annual totals' (see Figure 13) and 'Normalised seasonal totals' (not shown).

For each combination of response variable and explanatory variables, normalised values are computed by employing:

- a semiparametric regression model in which the intercept is permitted to vary slowly with season and year;
- an ordinary regression model with constant intercept.

The smoothing parameters in the semiparametric model can be selected by cross-validation or assigned user-defined values. When cross-validation is employed, the number of test sets is equal to the number of years in the observation period, and each test set consists of all observations made the same year.

The worksheet 'Seasonal totals' can be edited before the riverine loads are normalised. For example, it is possible to change the definition of seasons. However, the label 'Year' and the years below that label must be left unchanged, and the season numbers must be typed in the column to the right of the years.

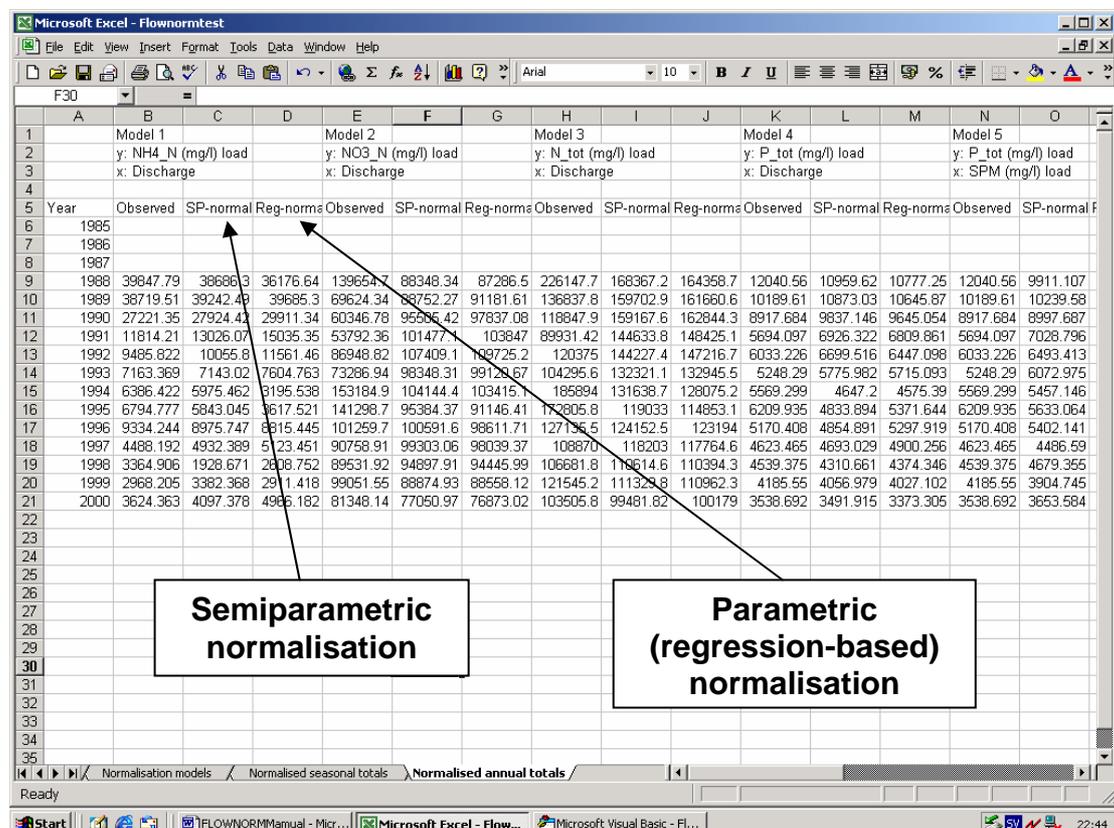


Figure 13. Normalised annual loads

The predictive ability of the different normalisation models can be assessed by computing the mean square prediction error, that shall be as low as possible. The MSPE-values in Figure 14 show that, for the tested data sets, semiparametric models perform better than the simpler parametric models. In addition, the results indicate that the load of total-P at Schnackenburg is most efficiently normalised by using a semiparametric model with load of suspended particulate matter (SPM) as explanatory variable.

The goodness-of-fit is determined by computing the mean square residual. The results obtained for total-P at Schnackenburg (Figure 14) show that the predictive ability of a normalisation model can decrease, even if the goodness-of-fit is improved when a new explanatory variable is included in the model.

The lambda-values in Figure 14 show the values of the smoothing parameters that were selected by cross-validation. To reduce the computational burden, the search for optimal smoothing parameters is restricted to a grid-set of values.

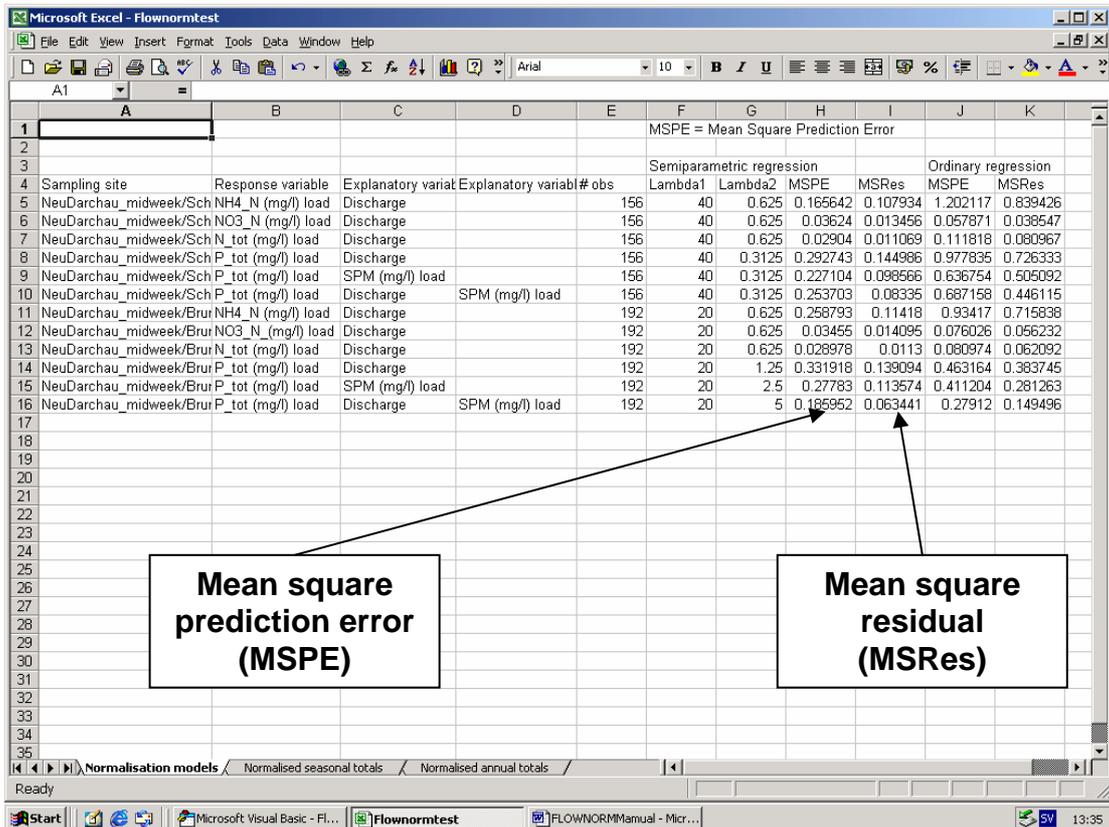


Figure 14. Predictive ability (mean square prediction error) and goodness-of-fit (mean square residual) for the tested normalisation models.