

Statistical Validation of Water Quality Data

Lieven Clement, Olivier Thas, Peter Vanrolleghem and
Jean-Pierre Ottoy

Department of Applied Mathematics, Biometrics and
Process Control, Ghent University, Belgium

Outline

- Background
- Problem definition
- Methodology
- Validation in a simulation study
- Evaluation of deviating measurements
- Case study
- Implementation
- Conclusions

Background:

European Water Framework Directive

- Goal: maintain and improve aquatic environment
- Essential step: registration of the water quality
 - In Flanders, task performed by Flemish Environmental Protection Agency (VMM)
 - VMM established several monitoring networks

Problem Definition:

Data validation of the physico-chemical monitoring network

- 1317 measuring points located in different catchment area's of Flanders
- Every location is evaluated 12-26 times a year on a basic spectrum of physico-chemical variables
- Hugh amount of data to validate and to interpret
- It is important to detect possible anomalies, which can indicate water quality losses
- Data validation is carried out by human experts

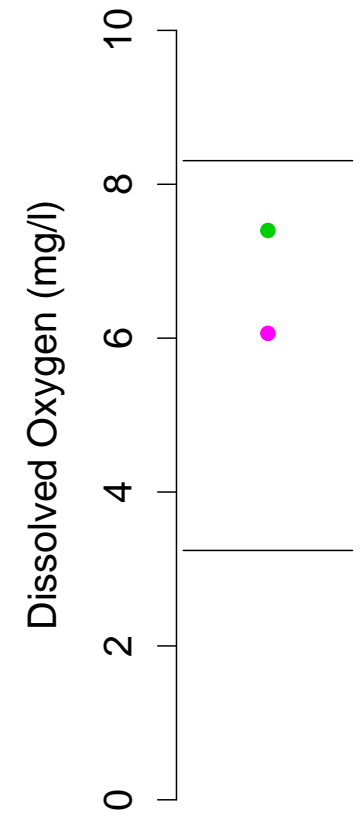
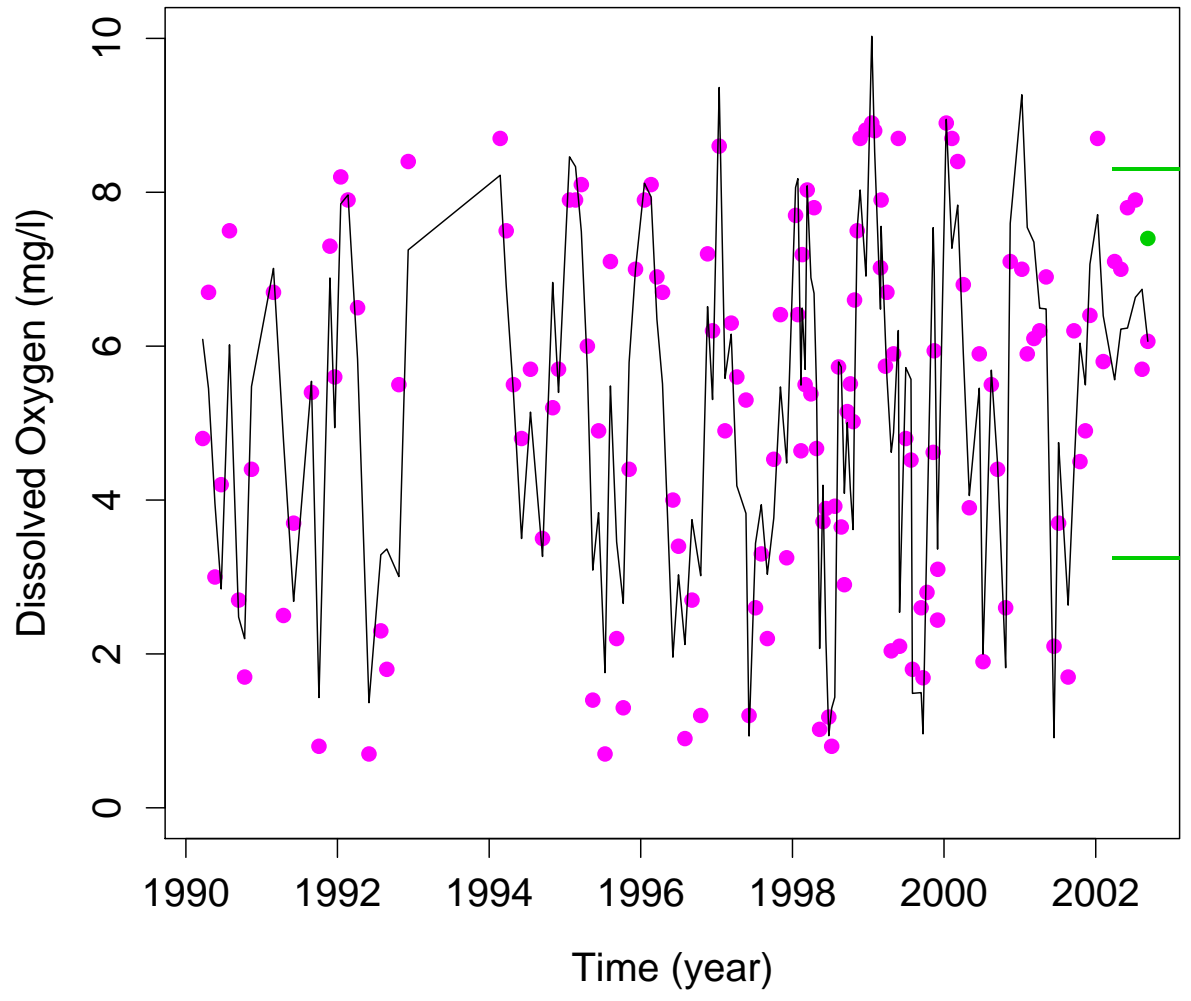
Need for semi-automatic data validation procedure

- Automatic validation: data accepted \Rightarrow addition to Dbase
- Deviation data is passed on to expert for further evaluation

Data validation: the statistical translation

- Fit model with historical data $t=1..n$
- Construct prediction and prediction interval at $t=n+1$
- Evaluate if new measurement is covered by prediction interval

Prediction Interval Dissolved Oxygen (mg/l)



Methodology:

- The GAM-model Family
- Model selection
- Bootstrap prediction intervals

Methodology:

GAM-Family

- Flexible model needed:
 - Adaptation to changing environment
 - Valid for different rivers and different locations
 - Able to model non-linear relations
 - No intervention of the user needed
- A flexible model family was chosen: GAM-models

GAM-Family: Model Structure

- Model structure

$$Y = \beta_0 + \sum_{j=1}^p f_j(X_j) + \delta$$

- f_j : Local linear regression smoother
- No analytical solution: back-fitting algorithm is used

GAM-Family: Model Structure Characterization

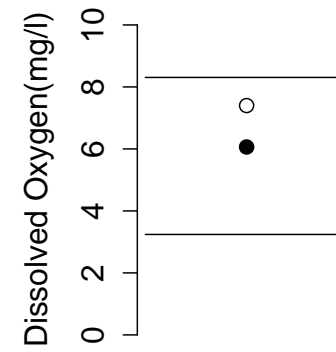
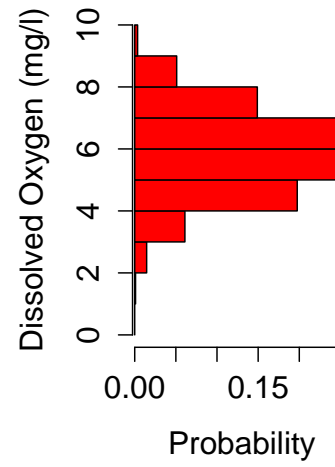
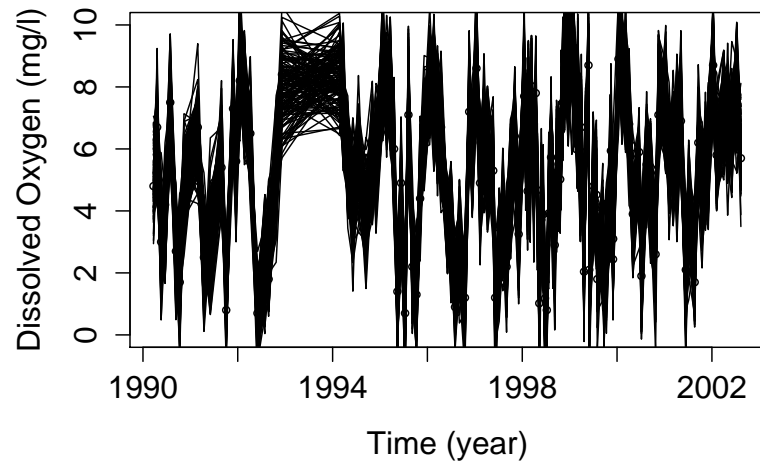
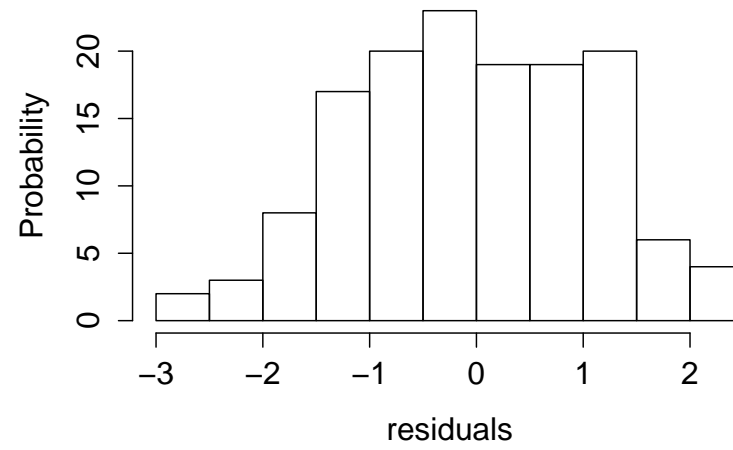
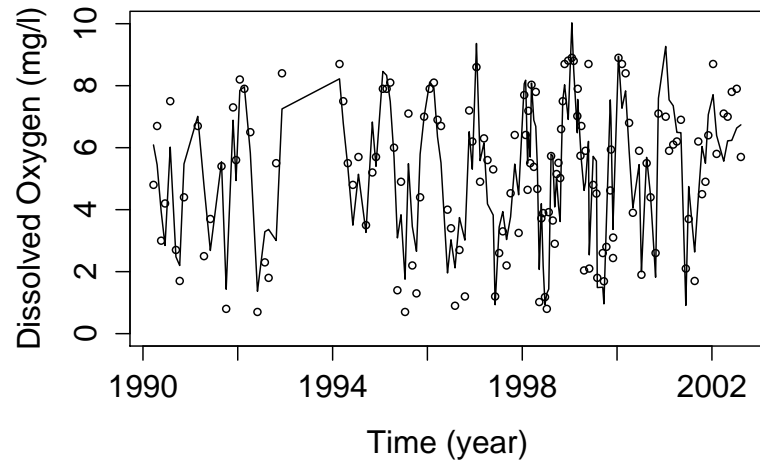
- Model structure: Long term trend, Seasonal trend, several physico-chemical variables are selected
- Model selection is performed at every location and at each time-step
- Selection Statistic:

$$AIC = \frac{1}{n} \sum_{i=1}^n D(y_i, \hat{y}_i) + 2 \frac{df \phi}{n}$$

Methodology:

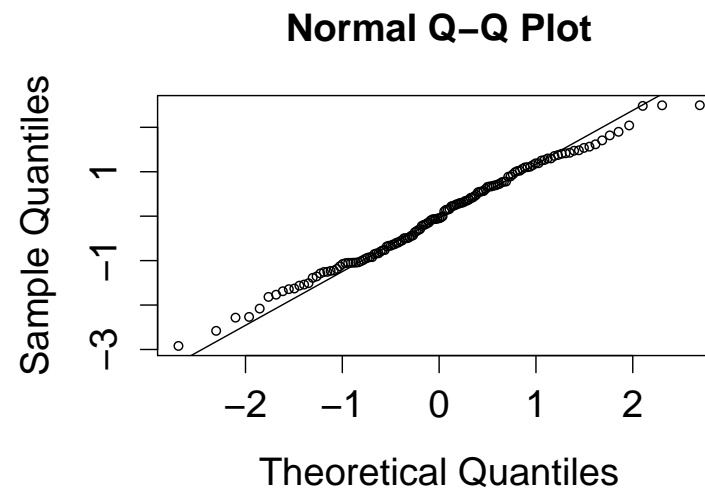
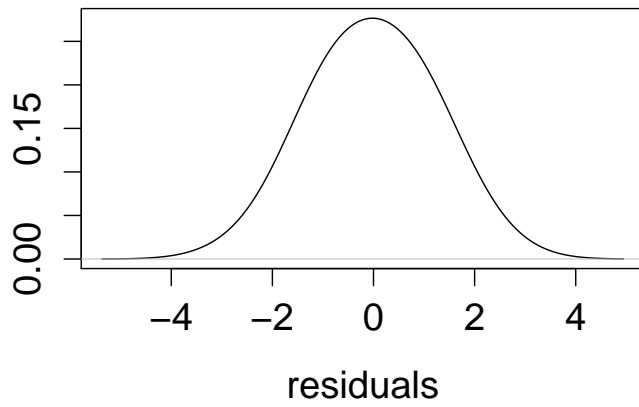
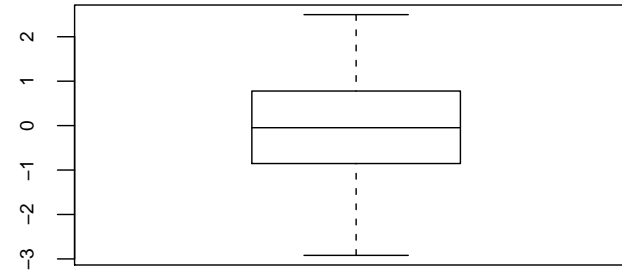
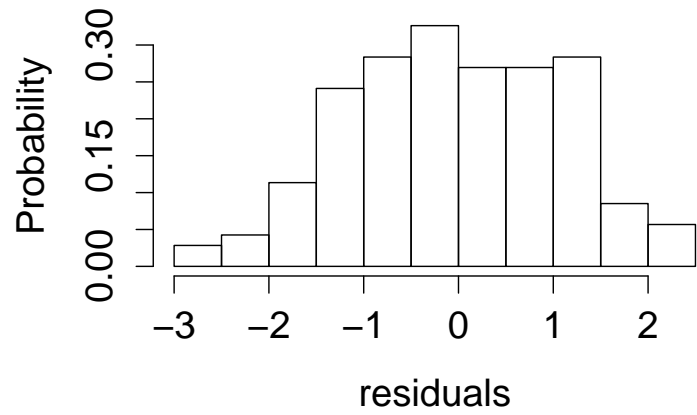
Bootstrap Prediction Intervals

- GAM-model: no analytical solution of confidence intervals
- Bootstrap methodology is used
 - New dataset: sample multivariate data with replacement
 - Fit GAM-model
 - Perform prediction
 - add random residual to the prediction
 - repeat this procedure

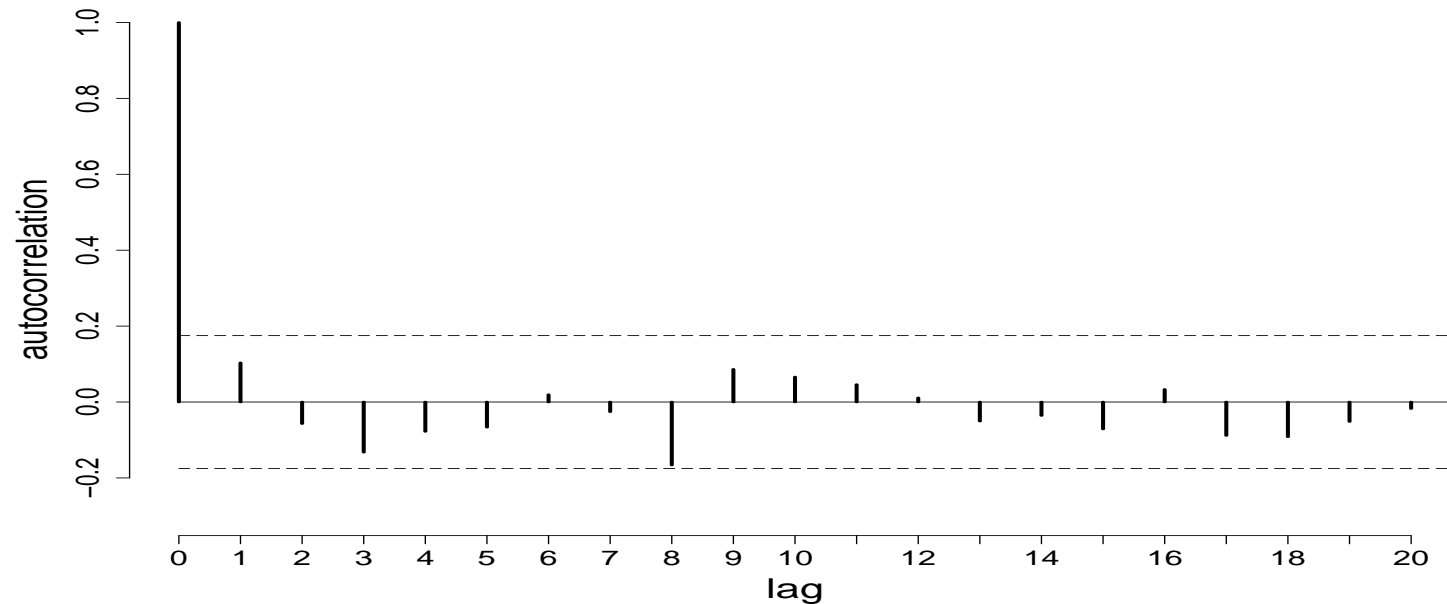


Validation: Simulation study

- Data generation
 - Model Fit: obtain \hat{Y}
 - Estimate variance of errors: s^2
 - 10000 new datasets: predictor variables and new response
 $Y = \hat{Y} + N(0, s^2)$
- Observed coverage of 95% bootstrap prediction intervals = 95.7%.



- No problems with serial correlation occurred
no moving block bootstrap was needed



- Possible explanations: Serial correlation could be partially covered in predictor variables
 - Seasonal trend and long term trend in the model
 - Other physico-chemical variables

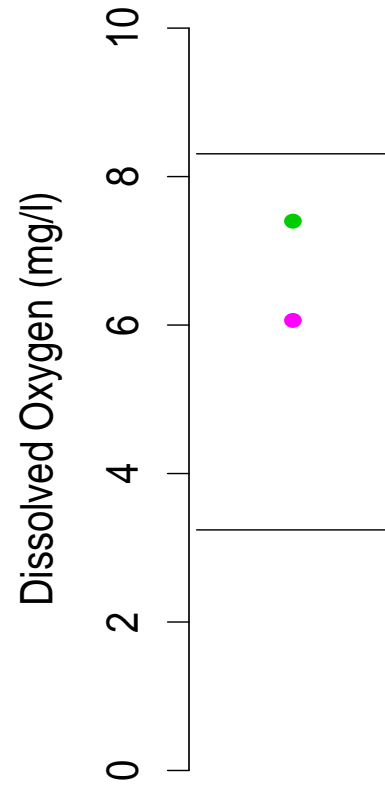
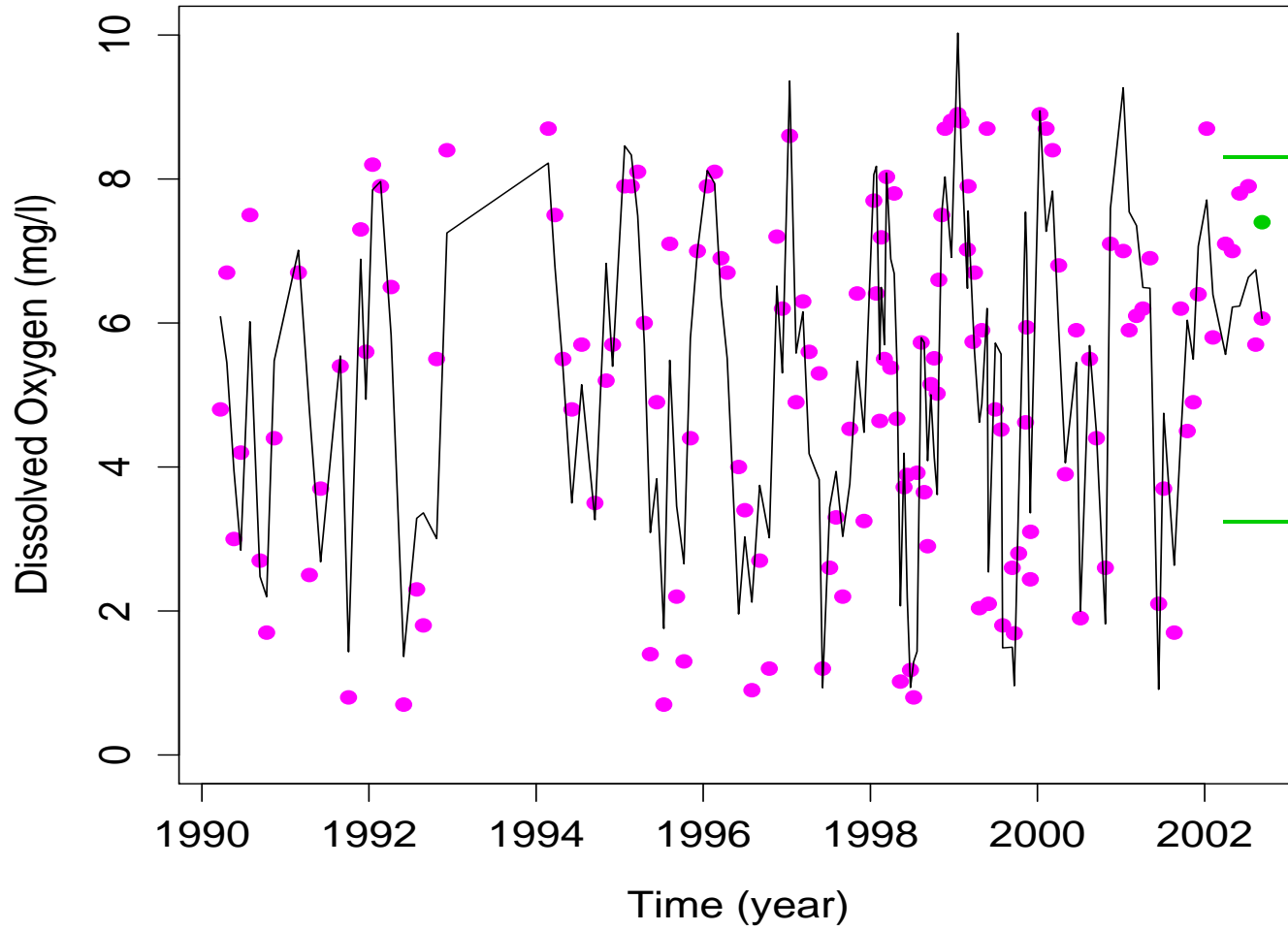
Evaluation of deviating measurements

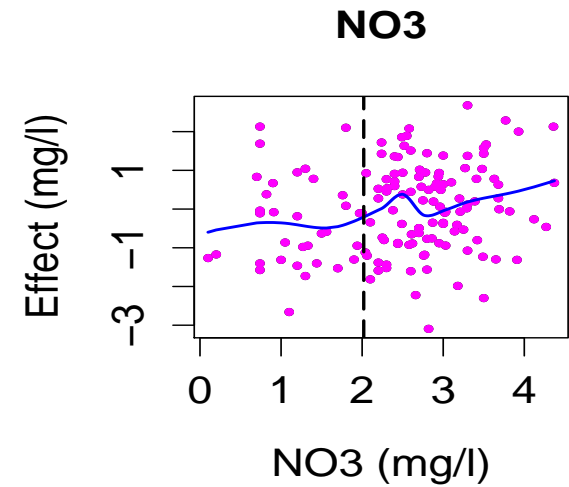
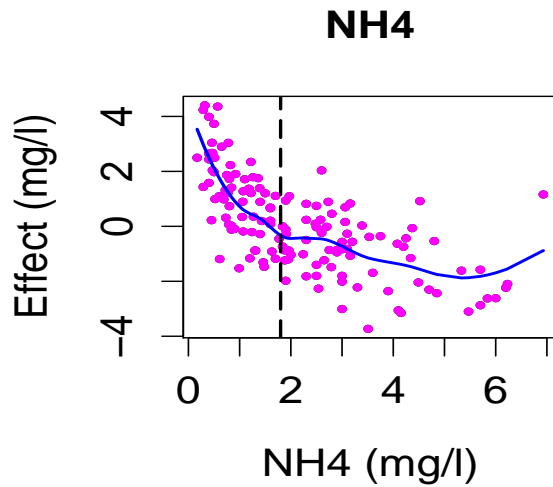
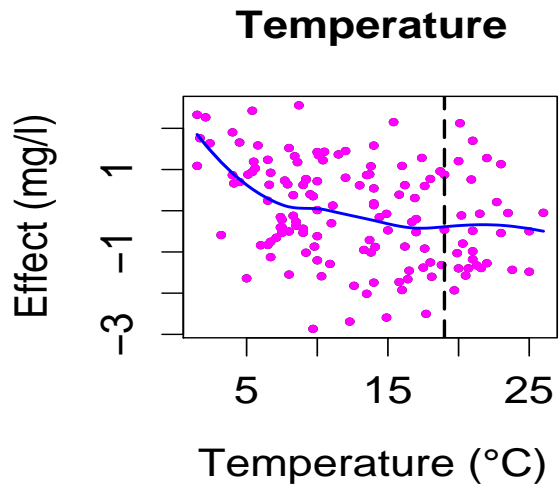
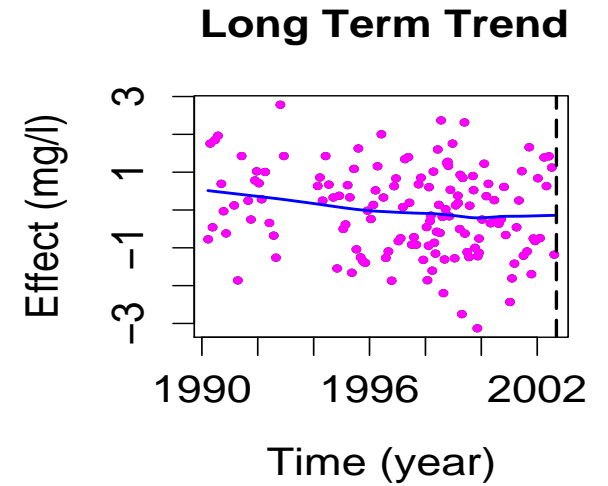
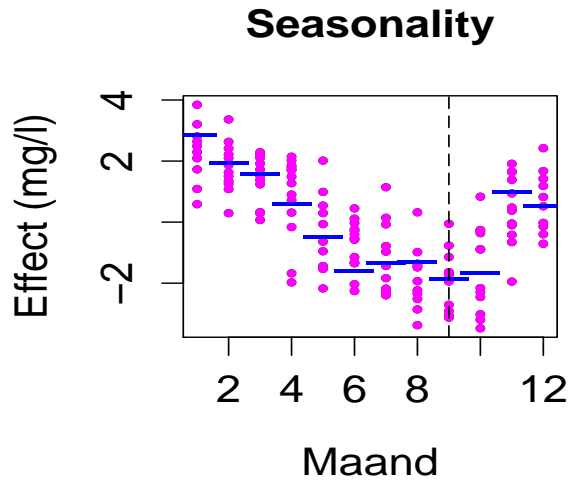
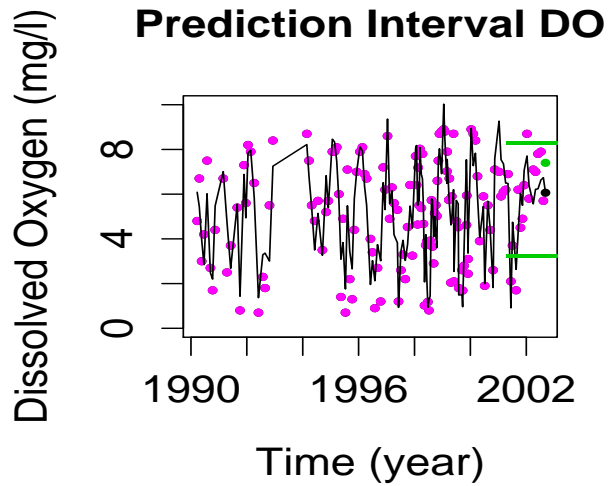
- Possible causes of deviating measurements
 - Due to the measurement
 - Due to a change in the system
 - Due to deviations in the predictor variables
- Detection of deviations in the predictor variables
 - Leave predictor variables one by one out the model
 - Construct prediction interval
 - Evaluate if measurement is covered by the new interval

Case study

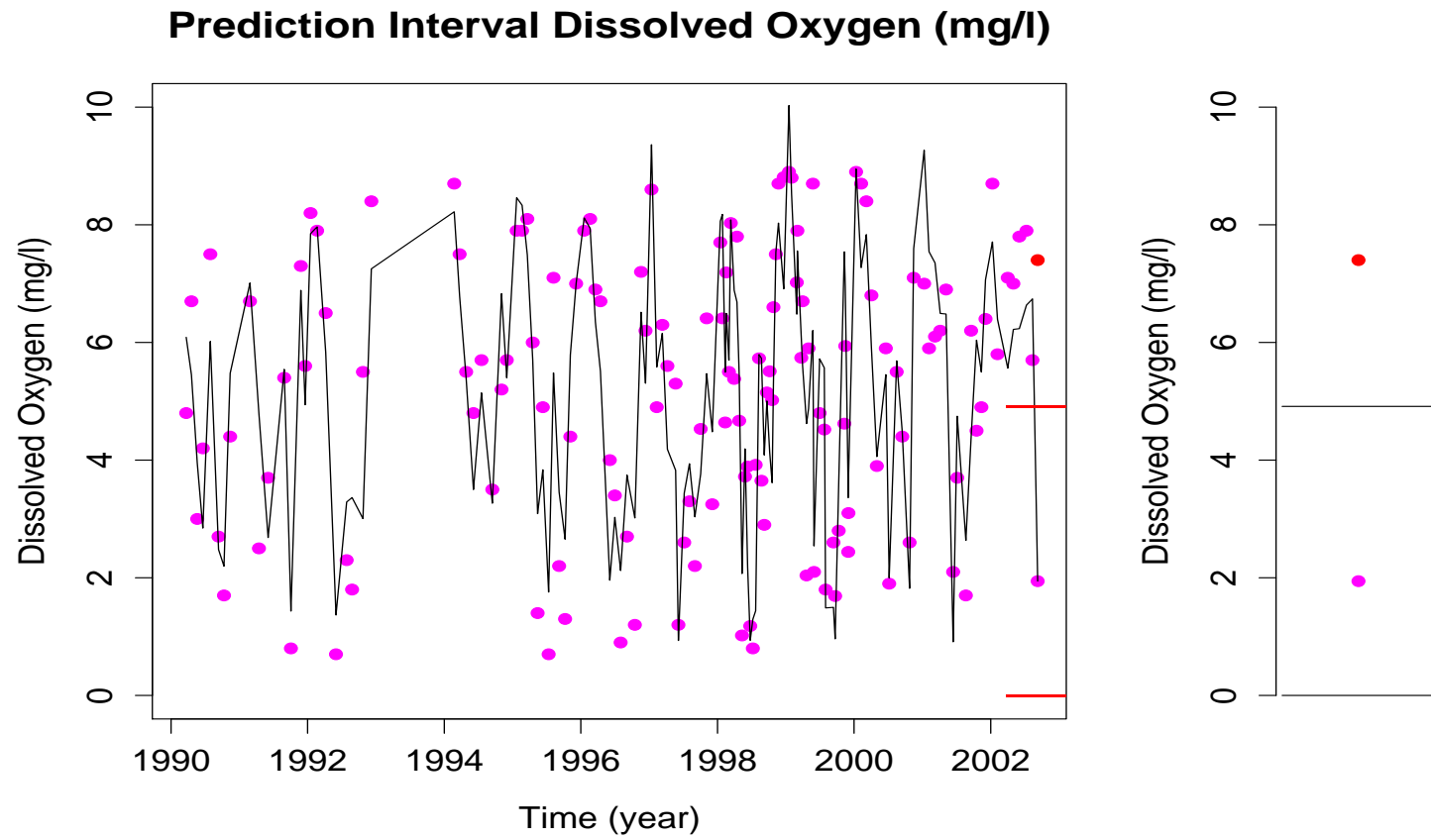
- Data validation of Dissolved Oxygen
- Introduction of error in prediction variable

Prediction Interval Dissolved Oxygen (mg/l)

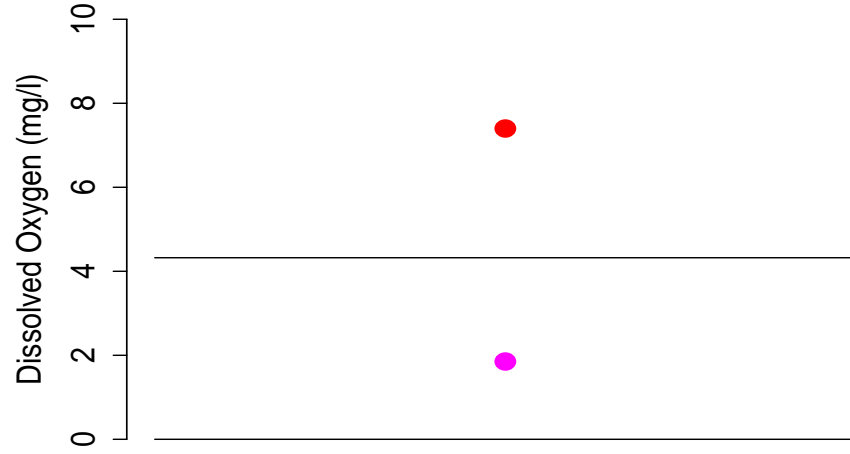




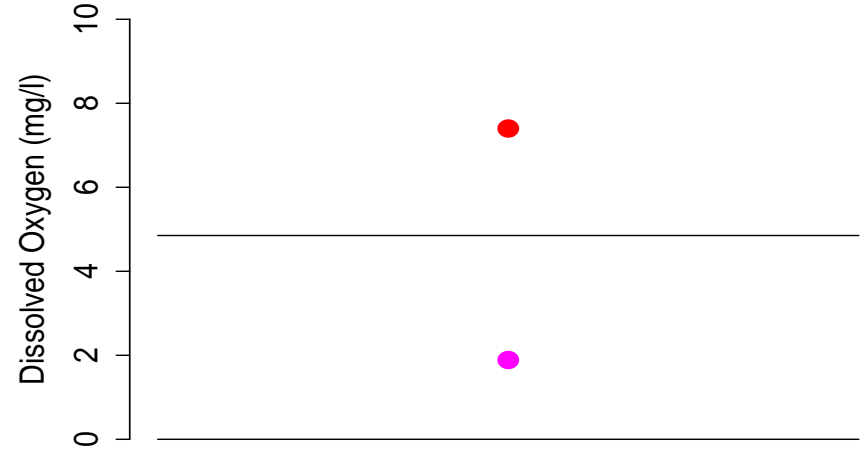
error introduced in NH4 concentration



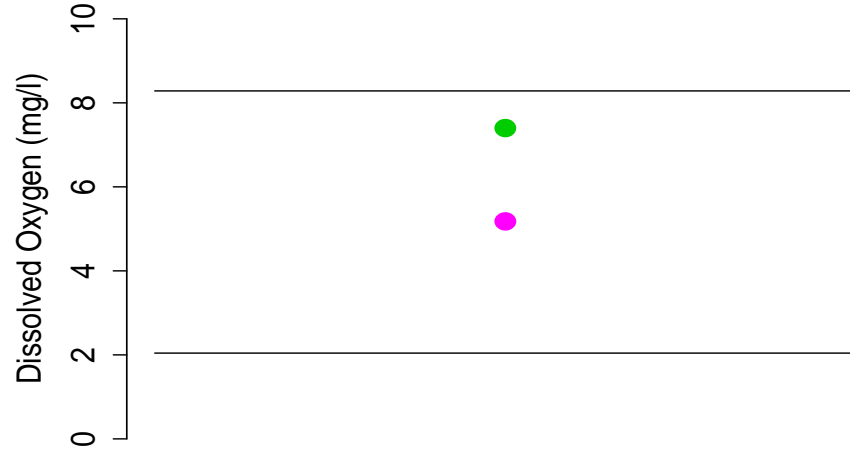
Prediction Interval without Long Term Trend



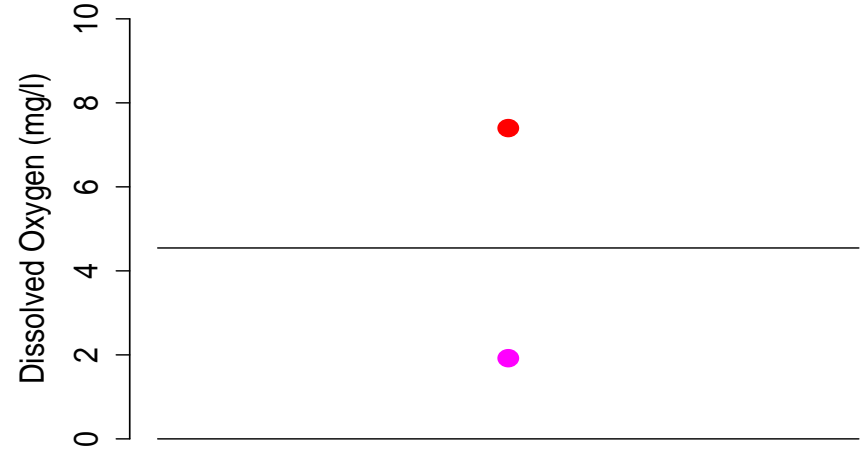
Prediction Interval without Temperature



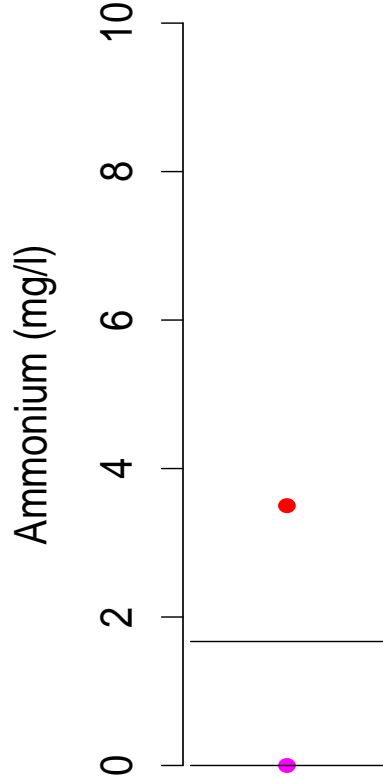
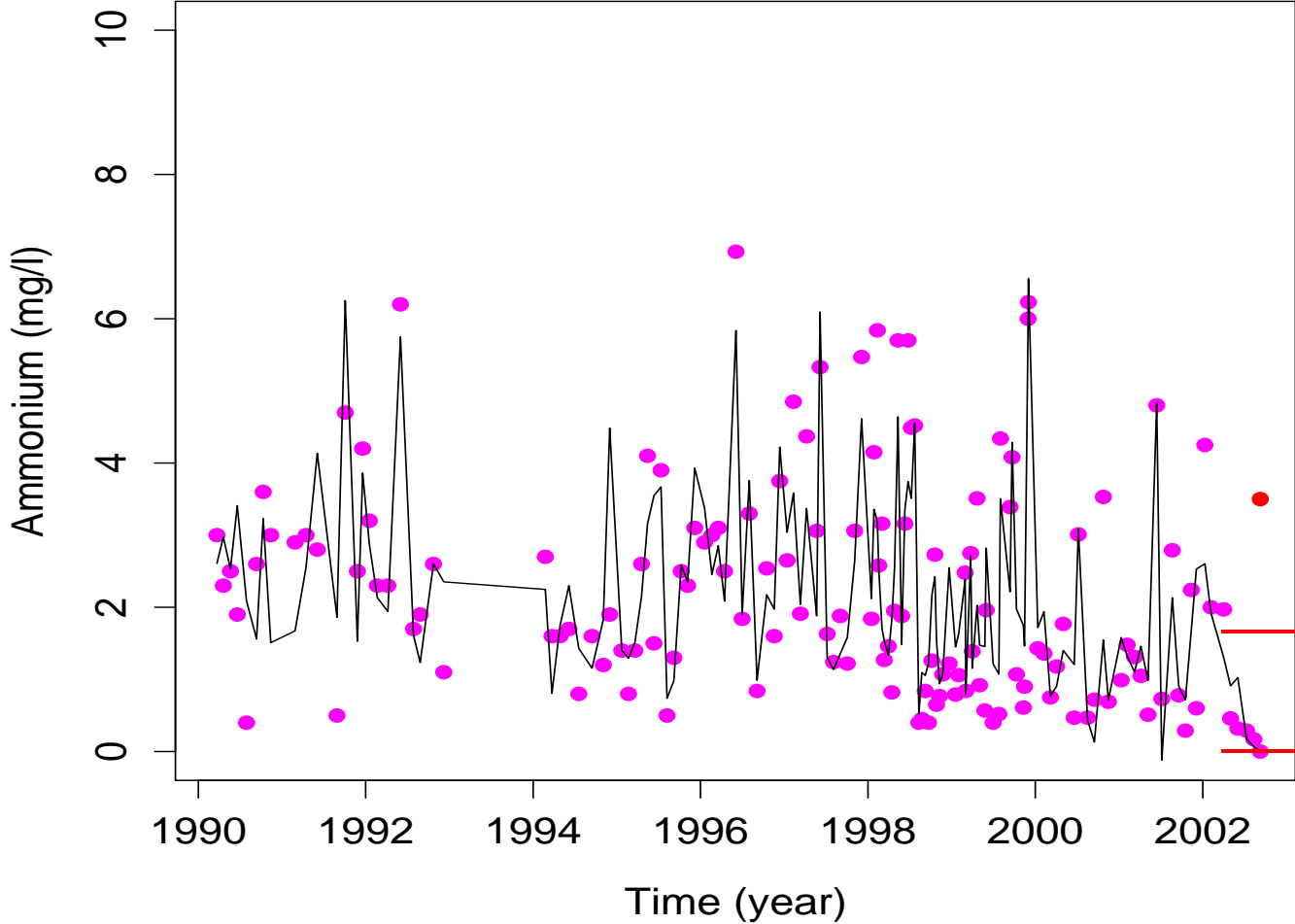
Prediction Interval without NH4



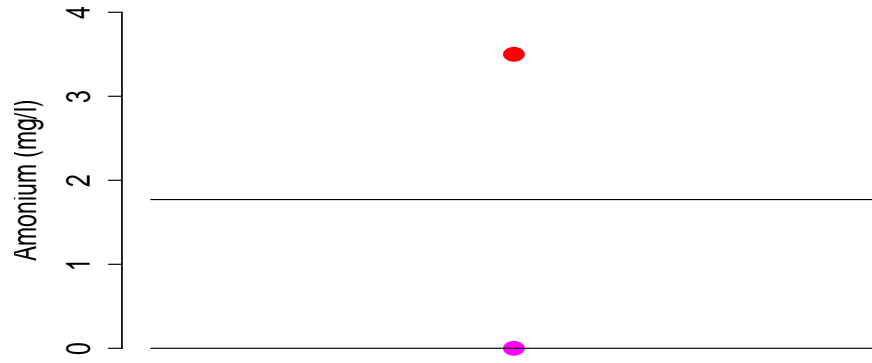
Prediction Interval without NO3



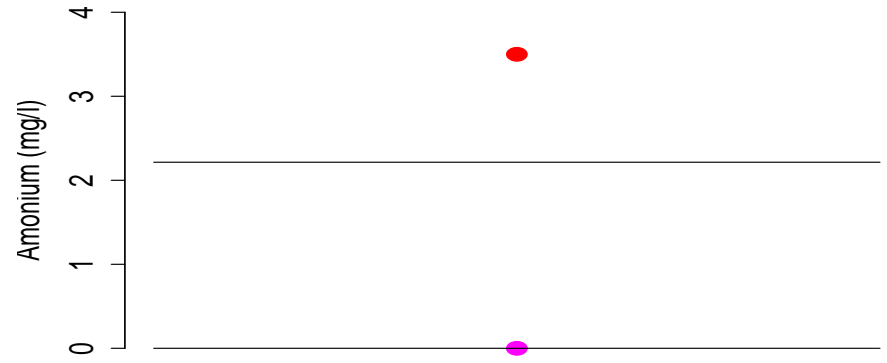
Prediction Interval Ammonium (mg/l)



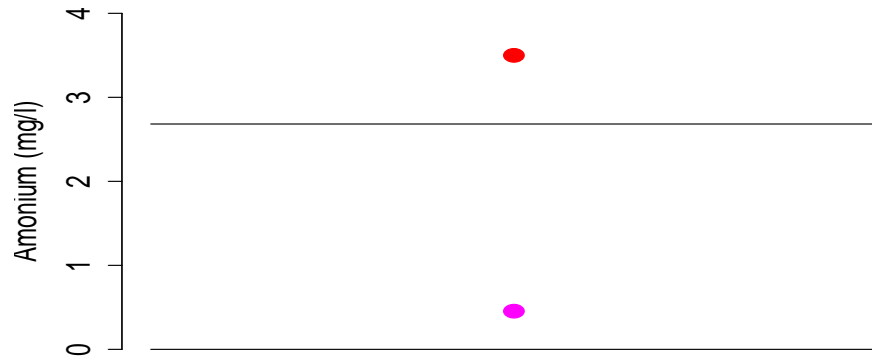
Prediction Interval without Long Term Trend



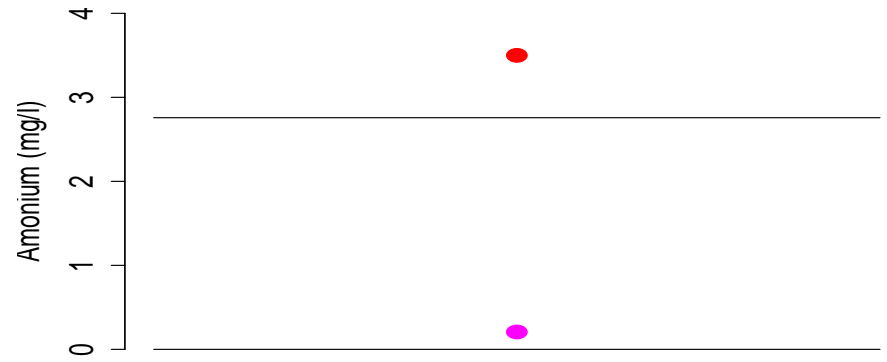
Prediction Interval without Temperature



Prediction Interval without NO2



Prediction Interval without NO3



Practical Implementation

- The coming months the procedure will be implemented with the VMM
- Following loop will be implemented at each location:
 - Model structure characterization
 - Construction Prediction Intervals
 - Evaluation of new measurement

Conclusions:

- GAM Model Family:
 - Provide the flexibility needed for the semi-automatic procedure
 - Trends and relations between physico-chemical variables can be deduced
- Good coverage of the Prediction Intervals
- Methodology can give an indication why deviating measurements occur