ELSEVIER

# A quality-control method for physical and chemical monitoring data. Application to dissolved oxygen levels in the river Loire (France)

F. Moatar[a],[*], J. Miquel[b], A. Poirel[b]

[a]*Laboratoire d'études des Transferts en Hydrologie et Environnement (UMR 5564: CNRS INPG IRD UJF), BP 53, 38041 Grenoble Cedex 9, France*
[b]*Electricité de France (Division Technique Générale) 21, Av. de l'Europe, BP 41, 38 040 Grenoble Cedex 9, France*

## Abstract

A quality-control method is proposed for examining continuous physical and chemical measurements, including temperature, dissolved oxygen, pH and electrical conductivity. Firstly, measurement consistency is evaluated by various modelling approaches: internal series structure, inter-variable relations or relations with external variables, spatial coherence and deterministic models. Secondly, outliers or systematic errors are detected using classical statistical tests. The method was evaluated for dissolved oxygen concentrations (DO) in the river Loire at Dampierre over a 5-year period (1990–1994), using data records containing fictitious errors, and raw data for the year 1995. The results demonstrate the effectiveness and advantages of a multi-model approach. In the case of dissolved oxygen for example, slow continuous drifts are always detected in under 4 days. © 2001 Published by Elsevier Science B.V.

*Keywords*: Water quality; Continuous monitoring; Quality control data; Dissolved oxygen; pH

## 1. Introduction

Instruments capable of making continuous in situ measurements of water quality are useful in detecting short-term changes in water composition. Water temperature, electrical conductivity, pH and dissolved oxygen (DO) are the most commonly monitored variables (Ranalli, 1998). Benefits of this continuous in situ monitoring include trend analysis. Considerable research has been devoted to water temperature, no doubt due to its long and reliable record (Webb, 1996). Studies of acidic streams consider variations in electrical conductivity and pH in terms of chemical and hydrological processes (Robson et al., 1993). Studies of eutrophic rivers use diurnal DO and pH measurements to estimate stream reaeration, primary production, and respiration rates (Chapra and Di Toro, 1991).

Data collected from in situ monitoring, by the electricity generating authority (EDF) at the nuclear power station, for example, are only useful if the measured values represent the in situ values. It is therefore essential to examine data critically and validate them by checking measurement consistency, in order to distinguish measurement anomalies from environmental changes. A strategy for critically examining and validating this type of data has been

---

\* Corresponding author. Fax: +33-4-76855286.
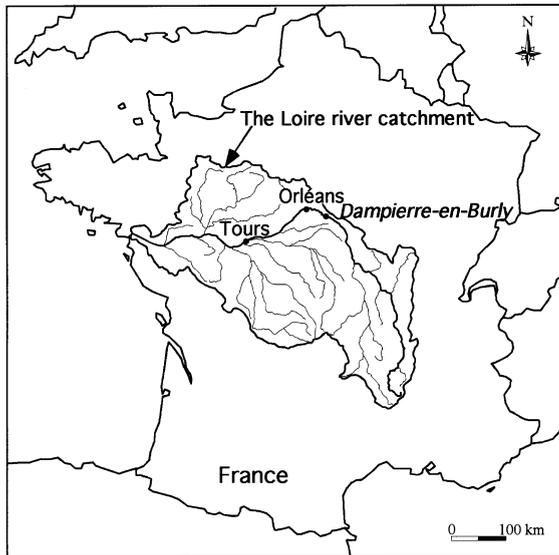  *E-mail address:* moatar@hmg.inpg.fr (F. Moatar).

Fig. 1. Location of the Dampierre en Burly study site.

developed (Moatar, 1997) and is discussed in this article.

## 2. Description of site and monitoring system

The river Loire is about 1012 km long, with a catchment area of 115,000 km$^2$ in central and western France (Fig. 1). The Dampierre site, considered in this study, is 550 km from the source and drains an area of 35,500 km$^2$.

The monitoring system consists of a floating platform with a sensor that measures water temperature at a depth of 20 cm, and a pumping device that sends a small flow of water (approx. 0.5 l/s) to three electrodes: pH (range 0–14 pH unit), DO (range 0–20 mg/l) and electrical conductivity (Cdv) (range 0–1000 μS/cm). Instantaneous measurements, taken every 5 s, are archived as hourly means. There are three monitoring stations: one upstream, one at the outflow of the power plant, and one 5 km downstream. Accuracy of the measurements (including electrode, transmission, and calibration), tested in situ by comparison with laboratory measurements, is estimated to be ±0.3°C, ±0.3 pH units, ±8% mg/l DO, ±5% μS/cm Cdv. Accuracy is defined as two times the standard deviation (S.D.) of the control sample readings. In spite of the control procedures, the data still contain anomalies, including outliers, gaps and systematic errors (discontinuities and drifts) (cf. Fig. 2a and b). Fig. 2a shows a measurement drift due to clogging of the sensor. Fig. 2b shows modification of the sensor and transmitter rating curves resulting in linear drifts in the DO measurements, with deviations ranging from −5 to +50%.

## 3. Quality control and validation method

The proposed method is based on three stages: (i) modelling, (ii) statistical analysis of model residuals, (iii) detection, diagnostic and correction.

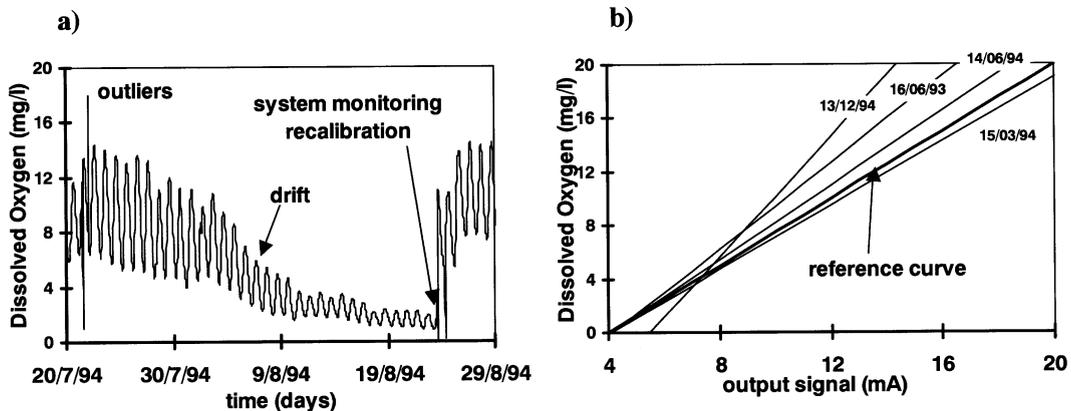(i) *Modelling* involves evaluating the probable



Fig. 2. Examples of errors: (a) outliers and drift due to clogging of the sensor; (b) rating curves.

value of the variables. It must be rapid, robust and reliable. A set of deterministic and stochastic models is chosen, based on analysis of the variables (internal structure), determination of inter-variable and inter-site relations, and assessment of relations with external variables that are reliable and available daily. *The internal structure of the series $X(t)$ is analysed using time series analysis methods* (Box and Jenkins, 1976). Data series with seasonal change only in the mean (temperature, conductivity, minimum oxygen value) are centred $U(t) = X(t) - Sm(t)$, while those with a seasonal change in the inter-annual standard deviation (daily range of DO and pH) are standardised. Harmonic analysis is performed to estimate these seasonal features:

$$Sm(t) = a_0 + \sum_{i=1}^{N} a_i \cos\left(\frac{2\pi t}{T_i}\right) + \sum_{i=1}^{N} b_i \sin\left(\frac{2\pi t}{T_i}\right) \quad (1)$$

where $a_0$ is the interannual mean, $T_i$ the harmonic period $i$ (days) and $a_i$ and $b_i$ are the harmonic terms. The stationary residual series $U(t)$ enables the structure of the series, and particularly its 'memory', to be studied. Autocorrelograms for the various series enable selection of autoregressive models AR($n$), with different orders of $n$ depending on the variable involved.

*Inter-variable relations* involve evaluating a variable (or component) as a function of another variable (or component) measured at the same station. Only relations between daily range of DO and daily mean of pH, and daily minimum DO value and daily mean of water temperature prove to be significant. *Inter-site relations* prove to be an effective means of detection, owing to the spatial consistency between the sites. Measurements at the upstream and downstream stations are highly correlated or even identical, owing to the minimal influence of the outfall on the physico-chemical characteristics of the water.

*Relations with external variables* mainly involve hydrometeorological variables. Air temperature is used to model water temperature. Solar radiation, discharge and water temperature are used to evaluate DO; solar radiation and discharge to estimate pH, and major anions and cations to estimate electrical conductivity. The most appropriate stochastic modelling methods are applied for each variable: Box and Jenkins transfer functions (Box and Jenkins, 1976) for

linear relations (water temperature, daily minimum of DO), neural networks for more complex and non-linear relations (Moatar et al., 1999a) and multiple correlation for variables with non-equidistant observations uncorrelated in time (electrical conductivity). For daily DO ranges, an empirical model was chosen, involving non-linear optimisation of conceptual equations parameters.

For the *deterministic models*, specific features of the middle Loire enable simplifying assumptions to be made in developing the models. This stretch of river normally shows a balance between thermal and ecological features and local meteorological conditions. There are no major tributaries affecting either hydraulic head or pollutant load. The water is relatively shallow and well mixed. Temporal variations in the variables are very pronounced. Consequently, equations defining changes in water temperature and DO can ignore upstream boundary conditions (Gilbert et al., 1986; Chapra and Di Toro, 1991).

(ii) *Statistical analysis of model residuals* involves comparing measurements of the variable, $C_0$, with the model forecasts, $M_i$ ($i = 1$ to $k$ models), using a series of 'classical' statistical tests (test of mean and test of gradient using sliding windows and Page–Hinkley cumulative test) described by Ragot et al., 1990. These tests are performed for normal variables with independent observations (Barnett and Lewis, 1995). It was therefore decided to work on the residuals between measurements and model forecasts, $\epsilon_i(t)$, and to decorrelate these series by autoregressive filters to extract the independent residual series $\mu_i(t)$.

*To detect a change in the mean signal* $\mu(t)$, the mean, $m_i$, is calculated using the current mobile window $F_i$, size $n$, and the test variable $u_i$: $u_i = (m_i - m_0)/(\sigma_0/\sqrt{n})$, where the reference mean, $m_0$, and the reference S.D., $\sigma_0$, are assumed to be known. These reference values are evaluated using a different and much larger sample, theoretically containing no measurement errors. In the absence of any change, the test variable $u_i$ has a normal centred and reduced distribution. If there are no reference values for testing the calculated values, these are evaluated using two windows (size $n_1$ and $n_2$) to enable a comparison to be made. In this case, the test variable $t_i$ follows a Student law with $n_1 + n_2 - 2$ degrees of freedom.

*To detect a discontinuity*, it is also possible to

calculate the slope of the regression curve $y_{i+j}$ on the $i$th mobile window of size $n$, fitting it as closely as possible to the signal.

$$y_{i+j} = b_{1i}j + b_{0i}, \qquad j = 1, ..., n \tag{2}$$

A variable for testing the significance of coefficient $b_{1i}$ is then estimated by applying the classical results concerning simple linear regression. The reduced deviation: $t_i = b_{1i}/s(b_{1i})$ is distributed according to Student's law with $(n-2)$ degrees of freedom ($s(b_{1i})$ being the variance of the parameter $b_{1i}$).

*The Page–Hinkley test* (Baseville, 1986), performed as a cumulative test sum, consists in fixing a minimum jump, $\delta$ to be detected, and running two tests in parallel, the 'direction' of the jump not being known a priori. The variable to be tested is redefined as

$$U_n = \sum_{i=1}^{n} \left( \mu(t) - m_0 - \frac{\delta}{2} \right), \qquad n > 0 \text{ and } U_0 = 0 \tag{3}$$

$$m_n = \min_{0 \le k \le n} (U_k)$$

$$T_n = \sum_{i=1}^{n} \left( \mu(t) - m_0 + \frac{\delta}{2} \right), \qquad n > 0 \text{ and } T_0 = 0 \tag{4}$$

$$M_n = \max_{0 \langle k \langle n} (T_k)$$

The detector signals the first time, $n$, at which $U_n - m_n > \lambda$ in the case of an increasing mean, and the first time, $n$, at which $M_n - T_n > \lambda$ for a decreasing mean. For the limit, $\lambda$, Ragot et al. (1990), suggest the expression: $\lambda = 2h\sigma_0/\delta$ where $h$ is equal to 2 for normal distributions and $\sigma_0$ is the S.D. of the signal.

(iii) *Detection, diagnostic and correction.* Wherever an anomaly is found, measurement consistency is checked. Measurements, model results, and residuals are analysed to identify the origin of the anomaly, which could be a drift, error, incident during measurement, accidental external phenomenon (pollution), limit of model validity, etc. When the measurement is suspect, a correction is proposed.

The method was validated by performing efficiency tests based on simulated error detection. Two types of perturbation were systematically introduced into the minimum values series and daily range of DO: linear drifts, characteristic of a modification in the rating curve (cf. Fig. 2b), and damped type drifts (exponential), characteristic of the sensor becoming clogged (cf. Fig. 2a). Both the efficiency of the proposed tests and speed of detection could thus be evaluated (cf. Section 4.3).

## 4. Application of the method to dissolved oxygen

Data recorded between 1990 and 1994 at the Dampierre station were used for finalising the method, as follows:

- Development and calibration of models covering 4 years (for which data quality could be checked) and cross-validation over one year: each of the 5 years was chosen in turn as the validation period, the other four being used as calibration. This allowed the robustness of the models and the quality of the data to be checked. Iterations were necessary between the first calibration, problem detection, elimination of doubtful periods and a second calibration.
- Introduction of artificial perturbations into series considered to be 'correct', i.e. homogeneous periods with stationary measurement/model residuals and small standard deviations.
- Raw data for the year 1995 were then used for overall validation of the quality-control technique.

### 4.1. Modelling of dissolved oxygen

Changes in DO can be represented by two components which have specific behaviour patterns: *Daily minimum values* (denoted $DO_{min}$) close to saturation are relatively predictable from the water temperature, and *daily ranges* (difference between $DO_{min}$ and $DO_{max}$, denoted $\Delta DO$) reflect the hydrological regime and seasonal photosynthesis activity of phytoplankton (Moatar et al., 1999b). Stochastic modelling was used, based on these two components.

### 4.1.1. Internal structure

*The minimum value of DO* for a given day is the sum of the annual component $DO_{min}[S(t)]$ and short-term component $DO_{min}[U(t)]$, which explains the daily
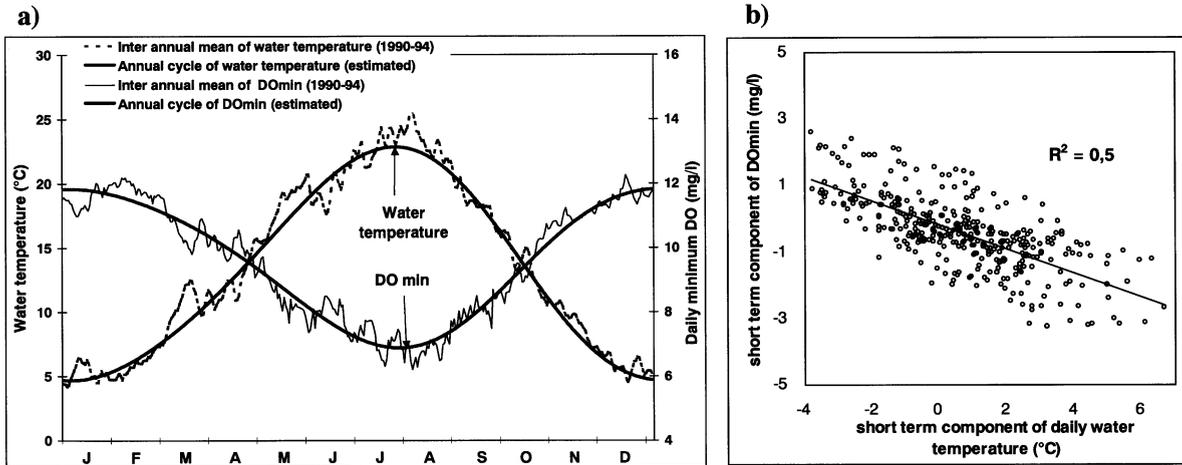
**a)**



**b)**



Fig. 3. (a) Annual cycle of daily minimum DO and of water temperature. (b) Short term evolution of $DO_{min}$ versus short term evolution of water temperature (cf. Moatar et al. 1999b).

fluctuations.

$$DO_{min}(t) = DO_{min}[Sm(t)] + DO_{min}[U(t)] \quad (5)$$

Spectral density analysis shows that the first two harmonics are sufficient for a proper representation of the seasonal component $DO_{min}[Sm(t)]$ ($R^2 = 0.72$). The coefficients of the harmonics (Eq. (1)) are as follows: $a_0 = 9.55$; $T_1 = 365$, $a_1 = 2.32$, $b_1 = 0.73$; $T_2 = 182$, $a_2 = -0.05$, $b_2 = -0.28$. The short-term component is an autoregressive process of order 1:

$$DO_{min}[U(t)] = 0.91DO_{min}[U(t-1)] + \epsilon(t) \quad (6)$$

As *the daily ranges* of DO displayed a S.D. that varied over the year, the series was standardised. The model for the short-term component $\Delta DO[U(t)]$ is an AR(3):

$$\Delta DO[U(t)] = 0.47\Delta DO[U(t-1)]$$

$$+0.20\Delta DO[U(t-2)] + 0.19\Delta DO[U(t-3)] + \epsilon(t) \quad (7)$$

### 4.1.2. Inter-variable relations

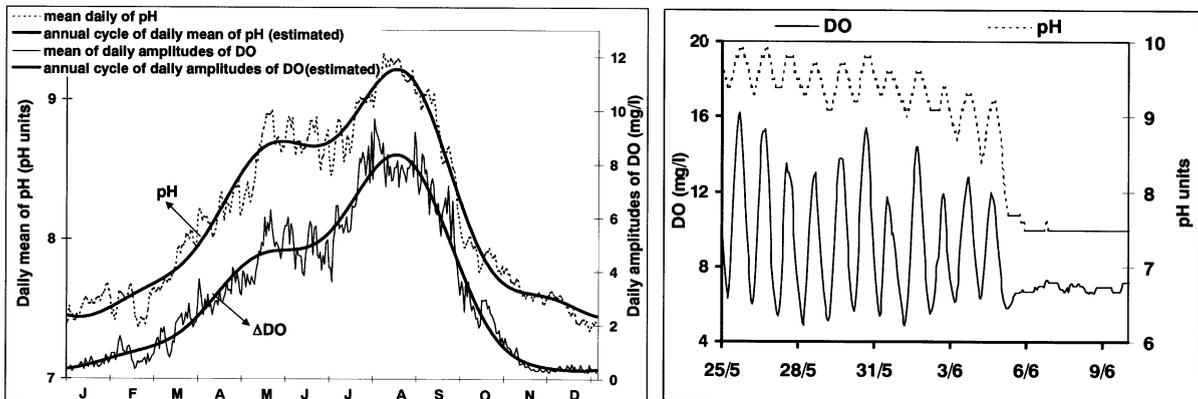The annual cycles of $DO_{min}$ and water temperature are remarkably in phase (Fig. 3a). The lowest values





Fig. 4. (a) Annual cycle of daily minimum DO and of water temperature. (b) Short term evolution of $DO_{min}$ versus short term evolution of water temperature (cf. Moatar et al., 1999b).

of $DO_{min}$ are observed in July and August, concomitant with the highest temperature values. Similarly, the highest values of $DO_{min}$ are concomitant with the lowest temperature values in December and January. Comparison of the short term components, obtained by removing the seasonal trend of $DO_{min}[S(t)]$ and water temperature $T[S(t)]$, shows a reasonably strong linear relationship (Fig. 3b). According to the pattern in the cross-correlation function (exponential decrease from lag 1) (Vandaele, 1983) and after testing several models, the model with the minimum Akaike's information criterion was chosen for the short-term component of the minimum value of $DO(DO_{min}[U^*(t)])$ :

$$DO_{min}[U^*(t)] = -0.26T[U(t)] + \epsilon(t) \qquad (8)$$

where $T[U(t)]$ is the short-term component of the daily mean water temperature $T(t)$.

The ranges of DO and pH are influenced by photosynthesis and phytoplankton respiration via $CO_2$. The annual cycle of the DO range shows a similar pattern to that of the daily mean of pH (Fig. 4a). During periods of relatively constant and low flow in summer, the daily maximum of pH occurs systematically half to 1 h later than the maximum of DO. Conversely, during wintertime and summer flooding, DO and pH are controlled by physical and chemical rather than biological processes, and daily cycles level off significantly (Fig. 4b) (Moatar et al., 1999b). After eliminating the seasonal effect in both variables ($\Delta$DO and daily pH), the following relationship was determined by linear regression:

$$\Delta DO[U^*(t)] = 3.89 \, pH[U(t)] \qquad (9)$$

### 4.1.3. External variables

The mass balance (O'Connor and Di Toro, 1970) (photosynthesis, respiration and reaeration) has been adapted to local and daily conditions (Moatar, 1997). The daily DO range could therefore be expressed as the sum of two factors: the first representing the effect of temperature ($T_{h\max}$ and $T_{h\min}$) on DO saturation concentration, and the second (denoted $\Delta$DO_B) the ratio between photosynthesis (P) and reaeration

rate ($k_a$).

$$DO_{max} - DO_{min} = [DOsat(T_{h\max}) - DOsat(T_{h\min})] + \frac{P}{k_a} \qquad (10)$$

$T_{h\max}$ and $T_{h\min}$ represent the temperature corresponding to maximum and minimum DO timing. In the Loire, during photosynthesis periods (favourable meteorological conditions) and where flow is constant and rather low, the DO peak is observed around 3 or 4pm, and the minimum before sunrise (5 or 6am). However, during medium or high flow, these cycles shift progressively until they eventually invert, $DO_{min}$ being observed during the day (approx. 2pm) and the maximum around midnight. Photosynthesis is weak during these periods, as DO production is controlled by the physical equilibrium between water and atmosphere (Moatar et al., 1999b).

Photosynthesis is represented as $P = f \times CP \times PHY$ where $f$ is a calibration parameter representing the conversion of algal growth into oxygen potential (mgDO/$\mu$gchla$_a$), CP is algal growth rate and PHY is living phytoplankton (mgchla$_a$/m$^3$). Algal growth rate is expressed as a classical multiplicative relation (Bowie, 1985) between light effect $h(I_0, H)$ and temperature effect $g(T)$ : $CP = c_{max}h(I_0, H)g(T)$, where $c_{max}$ is a calibration parameter (day$^{-1}$) representing maximum growth rate (when $(I_0, H) = g(T) = 1$). Steele's equation (1962), integrated over flow depth $H$, represents the effect of light attenuation on growth $h(I_0, H)$. Lassiter and Kearns' formula (1973) was used to represent the effect of temperature on metabolic growth processes $g(T)$, where $T$ is the daily mean temperature.

The reoxygenation coefficient $k_a$ is defined by O'Connor and Dobbins' equation (1958). It is expressed as a function of discharge via the empirical formulae established for the middle Loire ($H = 0.134Q^{0.4124}$; $V = 0.165Q^{0.2175}$). The light extinction coefficient $k_e$ in Steele's formula (after integration and averaged along the vertical) and the biomass PHY are considered to be a power as a function of the discharge ($Q$): $k_e = \alpha Q^{\beta}$; $PHY = m/Q^n$. $m$ and $n$ are calibration parameters. These parameters were initially estimated from Secchi disk reading measurements and from chlorophyll concentration at the study site. These estimates served as initial conditions for the optimisation algorithm. The

Table 1
Root mean square error, RMSE and correlation coefficients $R^2$ for dissolved oxygen models

|  | Autoregressive relation | | Inter variable relation | | Inter-site relation | | External variables relation | | Deterministic model | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
|  | RMSE (mg/l) | $R^2$ | RMSE (mg/l) | $R^2$ | RMSE (mg/l) | $R^2$ | RMSE (mg/l) | $R^2$ | RMSE (mg/l) | $R^2$ |
| $DO_{min}$ | 0.50 | 0.92 | 0.74 | 0.83 | 0.79 | 0.81 | – | – | 0.91 | 0.75 |
| $\Delta DO$ | 0.96 | 0.90 | 1.50 | 0.79 | 0.86 | 0.93 | 1.15 | 0.87 | 1.82 | 0.69 |

above considerations led to the following equation for the factor $\Delta DO\_B$:

$$\Delta DO\_B = f c_{max} \frac{m}{Q^n} \frac{Q^{0.48}}{32.2} * \frac{1}{\alpha Q^\beta 0.134 Q^{0.412}} * \left[ e^{1 - \frac{I_0}{I_s} e^{- \alpha Q^\beta 0.134 Q^{0.412}}} - e^{1 - \frac{I_0}{I_s}} \right] * e^{a(T - T_{opt})} \left[ \frac{(T_{max} - T)}{(T_{max} - T_{opt})} \right]^{a(T_{max} - T_{opt})}$$

(11)

$$\underbrace{\qquad}_{\text{PHY } 1/k_a} \quad \underbrace{\qquad\qquad}_{h(I_0, H)} \quad \underbrace{\qquad\qquad\qquad}_{g(T)}$$

A series of factors was fixed, based on values evaluated on site (Champ, 1980): $T_{opt} = 25.6°C$; $T_{max} = 36°C$; $I_s = 250$ W/m². Others were optimised using the Levenberg Marquardt algorithm from Matlab (The Math Works Inc, 1996): $fc_{max}m = 447.2$; $n = 0.57$; $\alpha = 0.01$; $\beta = 0.56$; $a = 0.14$.

### 4.1.4. Deterministic model

The BIOMOX model (Gosse, 1989) was used to determine DO and phytoplankton (PHY) on an hourly basis. The DO equation was adapted and simplified according to specific features of the Loire. The benthic compartment was ignored (the river bed consisting of soft sand). As BOD5 in the section studied depends mainly on the phytoplankton biomass (Khalanski, 1989), allochthonic contributions and organic matter not connected with algae were also ignored. The oxygen consumption terms for degradation of organic plant matter and corresponding nitrification were therefore replaced by a term that depends directly on the quantity of phytoplankton ($c_{dbo}$PHY) ($c_{dbo}$ being the coefficient representing organic matter degradation (mgDO/μgchla):

$$\frac{dDO}{dt} = f(CP - RP)PHY$$
$$+ k_a 1.025^{(T-20)}(DO_{sat} - DO) - c_{dbo}PHY$$

(12)

where CP is the growth rate (day$^{-1}$), expressed by the same multiplication relation $c_{max}h(I_0, H)g(T)$, and RP the respiration rate (day$^{-1}$) considered as a calibration parameter.

O'Connor and Dobbins' formula (1958) was used to determine the coefficient of reoxygenation at 20°C (day$^{-1}$). $I_s$, $T_{opt}$, $T_{max}$ were fixed on the basis of values evaluated on site (Champ, 1980) and the parameters $f$, $c_{max}$, RP, $c_{dbo}$ were manually calibrated: $f = 0.3$ mgDO/μgchla$_a$, $c_{max} = 1.5$ day$^{-1}$; RP = 0.05 day$^{-1}$; $c_{dbo} = 0.1$ mgDO/μgchla.

### 4.2. Results for each model

Table 1 presents the root mean square errors (RMSE) and correlation coefficients ($R^2$) between the measured and calculated DO values for minimum values and daily range. Fig. 5a and b shows the time series of $DO_{min}$ and daily ranges, $\Delta DO$, for 1990. The effectiveness of the various models is quite good (with $R^2$ varying between 0.75 and 0.92), except possibly for the representation of daily ranges by the deterministic model ($R^2 = 0.69$, RMSE = 1.82 mg/l). This model simulates minimum values better than range. This is due to under-estimation of maximum values of DO, due in turn to under-estimation of the phytoplankton biomass at certain periods. The mathematical formulation adopted includes a single phytoplankton compartment, inadequate for simulating the
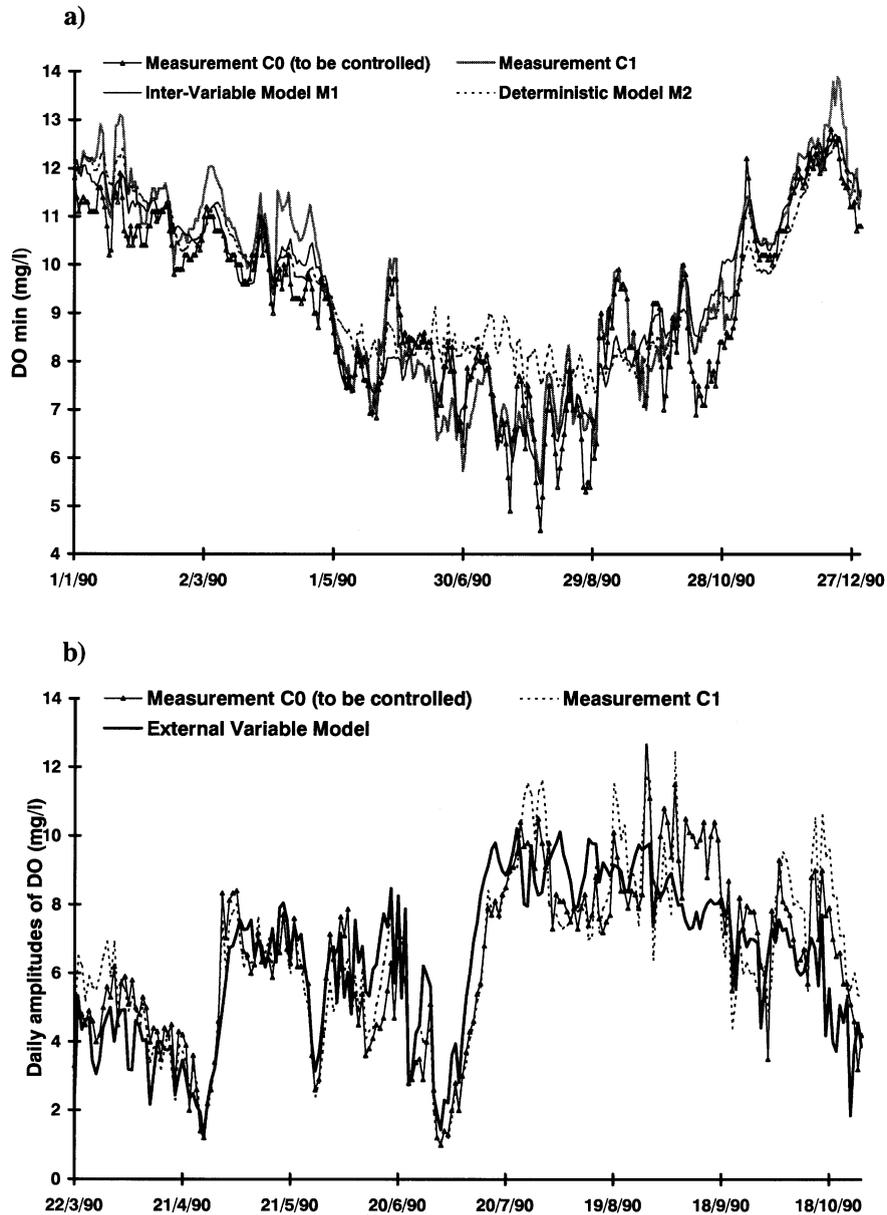
a)



b)



Fig. 5. (a) Minimum daily values of DO (year 1990) (b) Daily ranges of DO (year 1990).

spring-time growth of algae observed in this stretch of river.

The autoregressive model is shown to be the most efficient, due to the strong autocorrelation existing for the daily time step. This type of model could therefore be highly effective in detecting outliers, but not slow continuous drifts. For daily minimum values, the 'inter-variable', 'upstream/downstream' and deterministic models are identical. The efficiency of these models differs according to season (on average, RMSE $= 0.65$ mg/l and $R^2 = 0.70$ in winter, as against RMSE $= 1.2$ mg/l and $R^2 = 0.5$ in summer). The poor summer representation can be explained by the difficulty in measuring DO content in summer
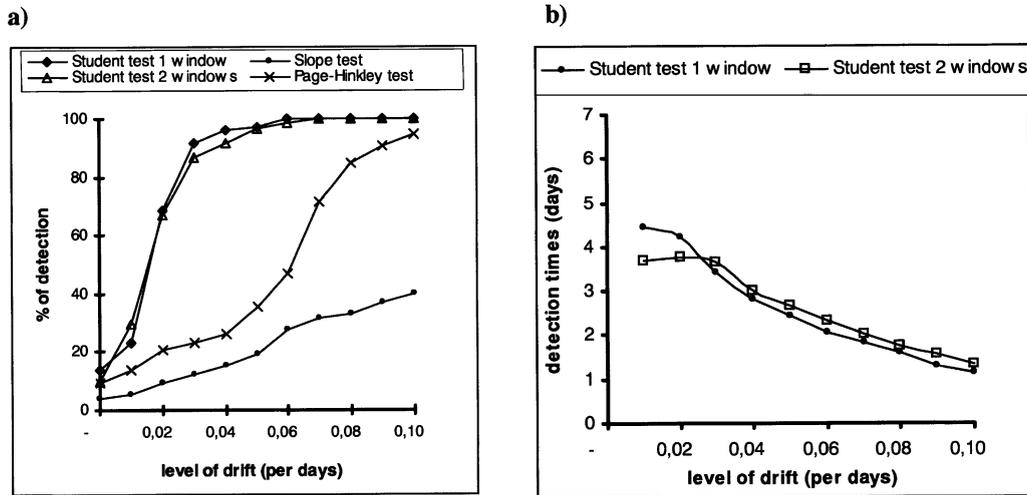
a)

b)



Fig. 6. (a) Percentage of detection of introduced drifts: Detection in relation to the inter-variable model $M_1$. (b) Detection times.

(because of risks of drifts and saturation due to the development of algae). Also, the DO model does not take into account biodegradation unconnected with phytoplankton (allochthonic pollution).

With regard to the daily range of DO, the performance of the inter-variable model depends to a great extent on the quality of the pH data. For example, from mid-January to mid June 1994, the external variable model reproduced $\Delta$DO quite well, while the inter-variable model systematically overestimated them because of erroneous pH values. In fact, the method of critical data analysis applied to pH shows that this time period is suspect (Moatar et al., 1999a). Otherwise, the model provides much better results with RMSE = 1.07 mg/l and $R^2 = 0.85$.

### 4.3. Results of tests with artificial perturbations

*Linear drifts* of variable levels (nd) and lasting one week were introduced systematically into the measurement series of daily minimum DO values at the upstream site $C_0$ (Fig. 5a). Drift detection capacity was then examined by using the statistical tests and comparing them with the downstream site $C_1$, and from the inter-variable model, $M_1$, and deterministic model, $M_2$. The deviations $\Delta C_1(t)$, $\Delta M_1(t)$ and $\Delta M_2(t)$ were decorrelated by AR(1) (cf. Eqs. (13)–(15) and the tests were applied to their residuals

$\mu C_1(t)$, $\mu M_1(t)$ and $\mu M_2(t)$ :

$$\mu_{C_1}(t) = \Delta_{C_1}(t) - 0.88\Delta_{C_1}(t - 1) \tag{13}$$

$$\mu_{M_1}(t) = \Delta_{M_1}(t) - 0.84\Delta_{M_1}(t - 1) \tag{14}$$

$$\mu_{M_2}(t) = \Delta_{M_2}(t) - 0.91\Delta_{M_2}(t - 1) \tag{15}$$

For each drift level, the percentage detected by each test was calculated for a confidence level $(1 - \alpha)$ of 95% for the three different reference cases (Fig. 6a). The percentage detected was comparable in the three cases. A drift of 4% per day on average was detected in more than 90% of cases, with mean comparison tests on 1 and 2 mobile windows. The slope test and Page–Hinkley test were very poor at detecting drifts. Regarding detection time, the mobile mean tests detected the drift in under 3 days on average (Fig. 6b), which has interesting operational implications.

*Exponential drifts* were simulated for daily DO ranges by the equation: $D(t) = C_0(t) - C_0(t)nd(k)^{(t-t_0)}$, where $t_0$ is the beginning of the 7 days and nd(k) is the level of drift (nd(k) = 0.02k, k = 0, ..., 20). Mean comparison tests using mobile windows were shown to be the most efficient in detecting linear drifts introduced into the minimum values. Detection by comparison with $C_1$ was slightly better, with a drift level of 0.1% per day being detected in 80% of cases on average, against 70% of cases with $M_1$ and $M_2$. Detection took 3.6 days on
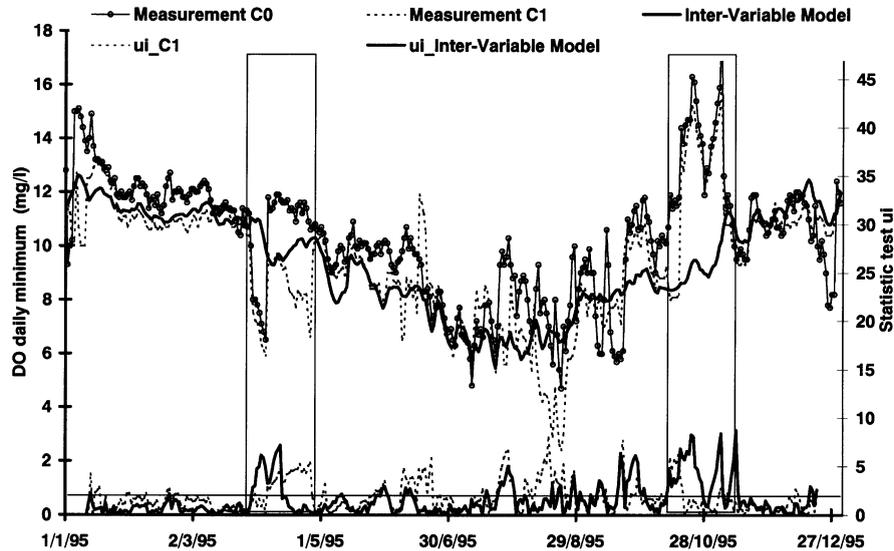
Fig. 7. Minimum daily values of DO (measurement to be checked $C_0$ and two references: $C_1$ and 'inter-variable' model), as well as test variable $u_i$ (year 1995).

average. Consequently, for real DO ranges of 6 mg/l, sensor clogging was detected as soon as the measured range dropped below 4 mg/l. Fig. 2a shows how using this method could have avoided 20 days of erroneous measurements.

Similar observations can be made for the other variables. However, the sophistication of the models, results and applicability depend on the variable involved. Water temperature is a robust, exact measurement, provided an appropriate measurement site is chosen. When artificial anomalies were introduced to test efficiency of detection, the estimate was shown to be representative, both according to other measurements, and to the properties of the models. Detection by inter-site comparison was naturally very good. Where a double measurement was not available, the reference temperature, estimated from models, showed greater disparity with the data (S.D. of the residuals for the External Variables model = 0.9°C and S.D. of the residuals of the deterministic model = 1.2°C). A sudden discontinuity of 0.4°C was detected in all cases in inter-site comparison and in 40% of cases with a 25% chance of false alarms, when compared with the models. A change of 1°C was detected by models in 60% of the cases.

The other variables — DO, pH and electrical conductivity — involve more sensitive measuring systems, requiring calibration, compensation and regular maintenance, and errors could be due to drifts caused by the equipment. The models for DO and pH enabled this kind of drift to be detected. However, the results for electrical conductivity were inconclusive, as the major anions and cations used by the models were only available 6–24 days per year.

## 4.4. Results of application of the method to the year 1995

The raw results of DO measurements recorded by the automatic station in 1995 presented obvious 'errors', i.e. outliers or gaps due to poor operation or interruptions in the measuring sequence. Before any work could be carried out on these series (e.g. simple transformation for calculating rates and daily means), outliers had to be eliminated and the series reconstituted as accurately as possible. The AR model applied to short-term daily minimum values $DO_{min}[U(t)]$ was used to detect outliers (cf. Section 4.1.2., Eq. (8)). The alarm threshold was set when the value of the residuals was greater than four S.D. of the series $\epsilon(t)$. Outliers were replaced by linear interpolations, or by Fourier series in the case of daily cycles.

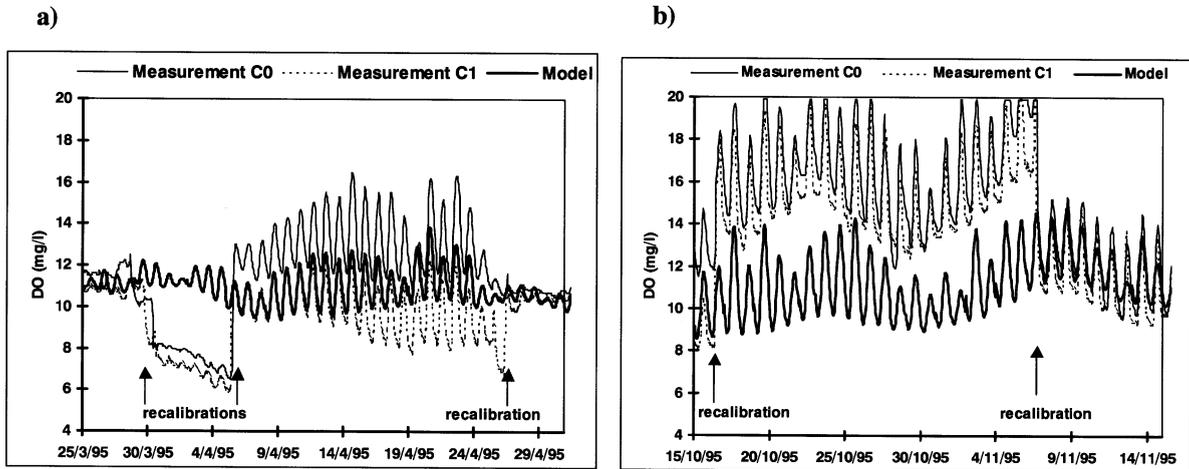a)                                                           b)

Fig. 8. Periods detected as being 'doubtful'; (a) 25 March–May 1995; (b) 15 October–15 November.

The method for detecting continuous or sudden changes in the DO measurement series at the upstream station ($C_0$) was applied first to the minimum values and then to the daily ranges. Fig. 7 shows the $DO_{min}$ measured at the stations upstream ($C_0$) and downstream ($C_1$) and the values 'predicted' by the inter-variable model. The graph also shows the test variable $u_i$ for comparing the mean of the residuals with a 7-day sliding window and the control limit corresponding to a confidence level of 95%. For minimum values, detection tests revealed 33% of anomalies in measured values compared with 'probable values' on the 'inter-variable' model, and 26% when compared with the measured downstream value $C_1$. For daily ranges, there were 9% of errors when compared with the 'external variables' model, and 14% compared with the downstream measurement $C_1$. Fig. 8a and b describes two periods detected as 'doubtful'.

Fig. 8a shows the shifts observed in the measured and calculated values, and the dates of three inspections by the automatic station maintenance team. The 'probable' hourly values are reconstituted from several models (inter-variable model for minimum values, external variables model for daily range and hourly reconstitution using a Fourier series). The model provides very well-calibrated data before, during and after the suspect period. The ranges and variations are similar to the two signals.

## 5. Conclusions

Various deterministic and stochastic models have been proposed in the literature for water temperature, electrical conductivity, pH and dissolved oxygen. All these models are geared essentially to environmental forecasting and impact studies, but rarely to data quality control. Quality-control methods using control charts exist for industrial processes (Steiner, 1984). The control limits are generally determined by product and process specifications. With regard to natural water quality, these values vary over time due to uncontrollable external factors (hydrometeorological, physico-chemical and biological). It was this aspect that was factored in when developing the methods discussed here.

The principles of this method could be generalised to other types of data. For example, Box–Jenkins transfer/noise models for spatial interpolation of groundwater head series have been proposed (Van Geer and Zuur, 1997). Coupling these with statistical tests would allow errors to be automatically detected and critically examined, and missing data to be filled in.

## References

Barnett, V., Lewis, T., 1995. Outliers in Statistical Data. . 3rd ed.Wiley, England.

Basseville, B., 1986. On line detection of jumps in mean. Lect. Notes Control Inf. Sci. 77, 12–26.

Bowie, G.L., Mills, W.B., Porcella, D.B., Campbell, C.L., Pagenkopf, J.R., Rupp, G.L., Johnson, K.M., Chan, P.W.H., Gherini, S.A., 1985. Rates, Constants, and Kinetic Formulations in Surface Water Quality Modelling. 2nd ed.. EPA/600/3-85/040US Environmental Protection Agency, Office of Research and Development, Athens, GA.

Box, G., Jenkins, G., 1976. Time Series Analysis; Forecasting and Control. Holden-Day, San Francisco.

Champ, P., 1980. Biomasse et production du phytoplancton de la Loire en amont et en aval de la centrale nucleaire de Saint-Laurent-des-Eaux. rapport EDF, HE. 31/80/01.

Chapra, S.C., Di Toro, D.M., 1991. Delta method for estimating primary production, respiration, and reaeration in streams. J. Environ. Engng. 117 (5), 640–655.

Gilbert, A., Gras, R., Roult, D., 1986. Numerical computation of natural river temperature. International Conference on Water Quality Modelling in the Inland Natural Environment. Bournemouth, England: 10–13 June, Paper M1, 1986, pp. 457–472.

Gosse, P., 1989. Monitoring the quality of water in the Doubs river (France). Hydroécol. Appl. (France) 1/2, 85–116.

Khalanski, M., 1989. Impact hydrobiologique du C.P.N. de Saint-Laurent-des-Eaux. Bilan sur la periode. 1977–1988 HE/32/87/17.

Lassiter, R., Kearns, D., 1973. Phytoplankton population changes and nutrient fluctuations in a simple aquatic ecosystem model. Rates, constants and kinetic formulations in surface water quality modeling. EPA, 1978US Environmental Protection Agency, Office of Research and Development, Athens, GA, pp. 256–257.

Moatar, F., 1997. Modélisations statistiques et déterministes des paramètres physico-chimiques utilisées en surveillance des eaux de rivières. Application à la validation des séries de mesures en continu (Cas de la Loire Moyenne). PhD thesis Institut National Polytechnique de Grenoble, France.

Moatar, F., Fessant, F., Poirel, A., 1999a. pH modelling by neural networks. Application of control and validation data series in the Middle Loire river. Ecol. Modell. 120, 141–156.

Moatar, F., Obled, Ch., Poirel, A., 1999b. Analyse de séries temporelles de mesures de l'oxygène dissous et du pH sur la Loire au niveau du site nucléaire de Dampierre (Loiret). 1. Compréhension des variations temporelles des teneurs en oxygène dissous et du pH en relation avec des données hydrométéorologiques. Hydroécol. Appl. Tome 11 (1/2), 127–151.

O'Connor, D.J., Dobbins, W.E., 1958. Mechanism of reaeration in natural stream. Trans. Am. Soc. Civ. Engrs. 123, 641–684.

O'Connor, D.J., Di Toro, D.M., 1970. Photosynthesis and oxygen balance in streams. J. Sanit. Engng. Div., ASCE 96 (2), 547–571.

Ragot, J., Darouach, M., Maquin, D., Bloch, G., 1990. Validation de données et diagnostic. Hermès, Paris.

Ranalli, A., 1998. An evaluation of in-situ measurements of water temperature, specific conductance, and pH in low ionic strength streams. Water, Air, Soil Pollution 104, 423–441.

Robson, A.J., Neal, C., Hill, S., Smith, C.J., 1993. Linking variations in short-and medium-term stream chemistry to rainfall inputs – some observations at Plynlimon Mid-Wales. J. Hydrol. 144, 291–310.

Steele, J.H., 1962. Environmental control of photosynthesis in the sea. Limnol. Océanogr. 7, 137–150.

Steiner, E.H., 1984. Statistical methods in quality control. In: Herschdoerfer, S.M. (Ed.). 2nd ed.. Quality Control in the Food Industry, vol. 1. Academic Press, London, pp. 169–298.

The Math Works Inc., 1996. Matlab Optimisation Toolbox User's Guide.

Van Geer, F.C., Zuur, A.F., 1997. An extension of Box-Jenkins transfert/noise models for spatial interpolation of groundwater head series. J. Hydrol. 192, 65–80.

Vandaele, W., 1983. Applied Time Series and Box-Jenkins Models. Academic Press, San Diego, CA.

Webb, B.W., 1996. Trends in stream and river temperature. Hydrol. Processes 10, 205–226.