# Dynamic Linear Models for analysing time series data in coastal environmental monitoring

Dominique Soudant ( ✉ dominique.soudant@ifremer.fr )

French Research Institute for Exploitation of the Sea

Tania Hernández-Fariñas

French Research Institute for Exploitation of the Sea

Additional Declarations: No competing interests reported.

# Dynamic Linear Models for analysing time series data in coastal environmental monitoring

Dominique Soudant[1*] and Tania Hernández-Fariñas[2]

[1*]Service Valorisation de l'Information et Gestion Intégrée des données de la surveillance, Ifremer, rue de l'île d'Yeu, BP 21105, Nantes Cedex 03, 44311, France .
[2]Laboratoire Environnement et Ressources de Normandie, Ifremer, Av. Du Général de Gaulle, BP 32, Port en Bessin, 14520, France .

*Corresponding author(s). E-mail(s): dominique.soudant@ifremer.fr;

**Abstract**

Global changes have led to a renewed interest in time series of environmental monitoring. In France, for example, the French Research Institute for the Exploitation of the Sea (Ifremer) has been managing for 40 years several networks with hundreds of active sites, with annual to fortnightly sampling frequencies, measuring dozens of variables. These long-term datasets are difficult to analyse due to their characteristics (e.g. missing data, outliers, changes in sampling frequency, shifts). For this large number of time series, this paper proposes a semi-automatic procedure based on Dynamic Linear Models, detailed from data pre-processing (e.g. time unit definition, aggregations, transformations), through model specification, automatic and manual intervention, outlier and shift handling, to model hypothesis testing. When applied to three time series combining the above features, the results showed that missing data and changes in sampling frequency were adequately handled. Outliers and structural breaks were identified automatically, but also added manually. Highlighted

shifts were identified as artefactual (e.g. probe drift), anthropogenic (e.g. ministerial decree) and ecological changes (e.g. storm impact). Finally, the presented treatment has been successfully applied routinely to more than 19,000 time series with a common and simple model structure. The broad theoretical framework offered by dynamic linear models opens up fruitful perspectives for improving and extending the results presented here, in particular for dealing with measurement quantification limits and time-varying observation variances.

# 1 Introduction

The 21st century has seen a growing interest in long-term time series (Koslow and Couture, 2013) for characterizing marine ecosystems in an era of global change and anthropogenic impacts. Observing programmes and networks collectively span a broad range of physical, biogeochemical and biological variables, supporting research on marine ecosystems functioning and changes (O'Brien et al, 2017; Benway et al, 2019). *In situ* sampling, remote sensors (e.g. satellite, airborne) and autonomous platforms (e.g. high frequency buoys) are being used to gather these data. Accordingly, the resulting long-term datasets encompass an increasing number of methods and parameters to survey the ocean. For example, in France, Ifremer (French Research Institute for the Exploitation of the Sea) has been operating several *in situ* monitoring networks at national and regional scales, over the past several decades. Sampling frequencies vary from annually to bi-monthly depending on the parameters measured. In this context, hundreds of sites have been sampled in continental France and overseas territories, and although some of them are no longer monitored, sustained observations have also been carried out across a network of sites all along the coastline. All this work, gathered over the years, has resulted in the existence of thousands of time series that are, owing to their unique

nature, challenging to analyse as long-term monitoring can result in data gaps or missing data, outliers, doubtful or false measurements. Additionally, sampling frequencies and methodologies may have changed through time, and both staff turnover and the evolution of their skills add another source of difficulty while exploring these datasets. These modifications may have induced different kinds of changes in the times series which should not be confused with those that are primarily being sought, i.e. ecological events and anthropogenic impacts. As a result, their analysis, starting with the basic distinction of signal and noise, the extraction of a trend and seasonality, is particularly difficult.

Many different methods are available to analyse time series. The simplest techniques are based on more descriptive approaches such as the moving average or a cumulative function (Ibanez et al, 1993; Legendre and Legendre, 2012). Methods based on smoothing techniques e.g. LOESS or LOWESS, (Cleveland and Devlin, 1988), General Additive Models (Hastie and Tibshirani, 1990), have been commonly implemented in marine sciences, as they are considered to be flexible tools that do not require any prior information on the shape of the trend (e.g. non-linear versus linear trends in Generalized Linear Models – GLM). However, one weakness of these techniques is the fact that they cannot be represented by a mathematical formula. Furthermore, a subjective smoothing parameter needs to be specified, given that it has the same inconvenience as the moving average, where a window length needs to be defined for which averages are calculated. These approaches are not time series analysis methods strictly speaking given that autocorrelations are not natively considered. The non-parametric Mann-Kendall trend test, specifically with a modification for autocorrelations (Hamed and Ramachandra Rao, 1998), can be used to test the existence of a monotonic trend and the Theil-Sen regression line (Theil, 1950; Sen, 1968) can be used to visualize a linear trend. These are very

useful approaches however they cannot capture the details of the time series evolution. The family of ARIMA (Autoregressive Integrated Moving Average) models has been specifically designed to analyse times series and is very versatile. However, its formalism is not easy to master and in practice, the computer packages that implement it require regular time series.

ARIMA models can be equivalently written in a state space representation. This latter could be considered close to the GLM formalism. Used in conjunction with the Kalman filter, this approach allows to manage missing data and changes in sampling frequency. Durbin and Koopman (2012) compared ARIMA and state space models: a key advantage of the latter is the structural analysis of the decomposition (i.e. trend, seasonality, noise), they are more general and encompass ARIMA models. Petris et al (2009) presented Dynamic Linear Models (DLM) as a special case of state space models. DLM have been highlighted as a promising approach for ecology (Levy et al, 2014; Auger-Méthé et al, 2021). Developed since the late 1950s (West, 2014), bayesian forecasting (aka DLM) is not a new method but a wide theoretical framework offering many fruitful perspectives. According to these arguments and our previous experiences with DLM (Soudant et al, 1997a,b), in 2012 we started to develop a process to analyse environmental time series collected by observing networks operated by Ifremer (cf. Hernández-Fariñas et al, 2013; Hernández-Fariñas et al, 2015). Our goal was to set up a process adapted to the data and not to adapt the data to a method. By applying this to all our available time series we have been able to make successive adjustments, until the process reaches its maturity and finally stability (cf. Ratmaya et al, 2019; Lheureux et al, 2022). The aim of this paper is to describe more precisely this process that can be used with all environmental monitoring time series. Three

|  | Teychan bis<br>44°40′25.0”N 1°09′30.9”W | | La Mouclière<br>45°58′15.5”N 1°06′08.8”W |
|---|---|---|---|
| Variable | Dissolved<br>oxygen | Micro-phytoplankton<br>abundances | Cadmium concentrations<br>in *Mytilus edulis* |
| Unit | mg/L | Cell./L | $\mu$g/kg |
| Period | 2007-2020 | 1987-2020 | 1979-2016 |
| Sampling<br>frequency | fortnightly | fortnightly | 1979-2002 quarterly<br>2003-2016 fall and winter |
| Measurement<br>method | Multiparameter<br>probe | Utermöhl | 1979-1986 FAA<br>1987-2004 FAA (Zeeman)<br>2005-2016 ICPMS |

**Table 1** Summary of sites and variables. FAA: Flameless Atomic Absorption, ICP-MS: Inductively Coupled Plasma Mass Spectrometry

examples with actual data illustrate its relevance and specifically the ability to identify artefactual, anthropogenic and ecological changes.

# 2 Data and Methods

## 2.1 Networks, sites, sampling and measurements

Data used in this paper have been collected as part of the REPHY (Observation and Monitoring program for Phytoplankton and Hydrology; Belin et al, 2021; REPHY, 2021) and ROCCH (Chemical Observation and Monitoring Network; Grouhel-Pellouin et al, 2022) networks monitored by Ifremer. Operating since 1987 and 1979 respectively, they count hundreds of sites, including over 450 active sites, and several tens of variables. Sites and variables have been chosen according to their illustrative value of the results of the presented methodology in terms of nature (i.e. outliers and shifts) and origin (i.e. artefactual, anthropogenic, ecological).

"Teychan bis" (cf. Fig. 1, Table 1), located in the Arcachon Bay, is a semi-closed lagoon. Water samples were taken all year round, at a fortnightly frequency and within two hours from high tide. Two variables have been selected for this site. The first one was dissolved oxygen as measured at the
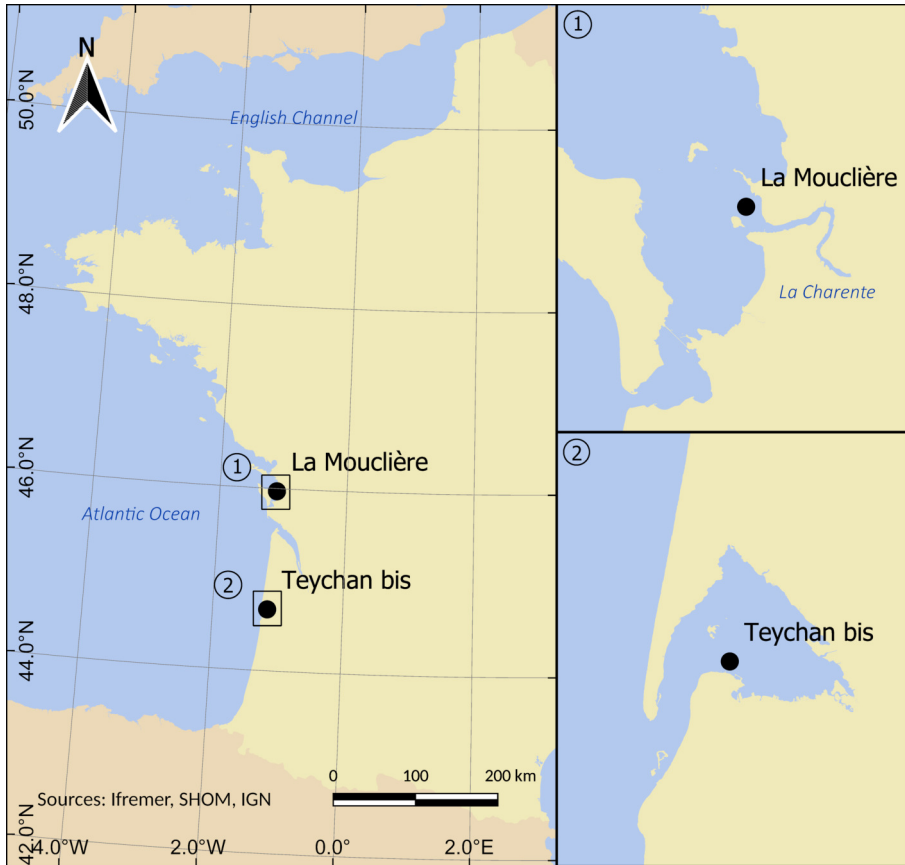
**Fig. 1** Study sites

water column bottom. Measurements have been taken since August 2007 using a multiparameter probe. The second variable is the micro-phytoplankton abundance estimated through the sum of diatom and dinoflagellate abundances from samples taken at the top of the water column: more precisely, up to and including 2007, samples were taken at a depth of 3 m, and then at the sub-surface (0-1 m). The absence of stratification in the Arcachon Bay (Neaud-Masson, pers. comm.) allows us to assume a homogenization over the entire water column and thus to treat the two depths as identical. Cell counts have been taken since 1987 as per the Utermöhl method (Utermöhl, 1958; Neaud-Masson, 2015).

"La Mouclière" (cf. Fig. 1, Table 1) is located at the mouth of the "La Charente" River and is thus under estuarine influence. Common mussels, *Mytilus edulis*, were sampled from 1979 to 2016. Until 2002, the mussels were sampled on a seasonal basis (winter, spring, summer and fall) and since then, only during fall and winter. For this study, the time series of cadmium concentrations found in freeze-dried *M. edulis* was used. Out of a total of 117 measurements, only four of them (i.e. fall measures from 2008 to 2011) were carried out by an external laboratory (i.e. not at Ifremer). Three analytical methods have been used: flameless atomic absorption until 1986, then the same method was applied with the Zeeman correction up to and including 2004 and lastly, Inductively Coupled Plasma Mass Spectrometry (ICP-MS) was used.

## 2.2 Data and pretreatment

Data from Ifremer's monitoring networks are entered into the Quadrige database[1]. Data have been extracted during April 2021 with a time window encompassing the first available results up to and including 2020. The time series approach implies the determination of an optimal temporal resolution for analysis. Thus, considering the theoretical frequency of observations and its actual application through the data, it is necessary to determine the duration of the time unit (e.g. day, week, fortnight, month) for which one observation is expected at most. The choice of this unit has two consequences: 1) time units will have to be created with a missing observation if they do not contain a sampling date, 2) time units with more than one observation must be reduced to one observation using an aggregation operator (e.g. median, maximum, minimum, mean). Operationally, an algorithm has been set up to determine the time unit for each series with the following ordered criteria: 1) minimizing the

---

[1]https://envlit.ifremer.fr/Quadrige-la-base-de-donnees

number of units for which a temporal aggregation was necessary 2) minimizing the time unit with missing data. The default aggregation operator was the median.

In order to stabilize the variances over time and thus to respect the homoscedasticity assumption of the model, data transformations may be needed. For cadmium, the logarithm of the concentrations is modelled. For cell counts, this is the decimal logarithm. Lastly, oxygen concentrations are treated within their original unit.

## 2.3  Methods

The model used here has two components: a local linear trend, in the form of a second order time series DLM (TSDLM), and a seasonal component. For DLMs, models are specified using two equations: an observation and an evolution equation. Data are described by the observation equation:

$$Y_t = \mu_t + \mathbf{FS}_t + \nu_t, \nu_t \sim N(0, V)$$

with

$Y_t$,      observation

$\mu_t$,      mean level or trend

$\mathbf{FS}_t$, seasonality

$\nu_t$,      error term or innovation

$V$,      observation variance

Here and hereafter, boldface text represents vectors and matrices. The observed signal is broken down into a mean level and seasonality. The vector $\mathbf{F}$ depends on the form, either factorial or trigonometric, and the time unit of the seasonality. Finally, an error term $\nu_t$, distributed according to a normal distribution of mean 0 and variance $V$ is added to represent the noise coming from all the variability in the data acquisition process (e.g. sampling strategy,

environmental variability, data entry) and factors that are not considered in the model. The sum of the mean level and seasonality represents an underlying unobservable process. These components are time-indexed, emphasizing that they can evolve over time. This evolution is described by the evolution equation, split into its components, below:

$$\mu_t = \mu_{t-1} + \beta_{t-1} + \omega_{\mu,t}, \ \omega_{\mu,t} \sim N(0,0)$$

$$\beta_t = \beta_{t-1} + \omega_{\beta,t}, \qquad\qquad \omega_{\beta,t} \sim N(0, W_{\beta,t})$$

$$\mathbf{S}_t = \mathbf{G}\mathbf{S}_{t-1} + \boldsymbol{\omega}_{\mathbf{S},t}, \qquad \boldsymbol{\omega}_{\mathbf{S},t} \sim N(0, \mathbf{W}_{\mathbf{S},t})$$

At time $t$, the mean level is equal to its value at time $t-1$ summed with $\beta_{t-1}$. This also means that $\mu_t - \mu_{t-1} = \beta_{t-1}$, i.e. the difference in the mean level between two units of time, and hence $\beta_{t-1}$ is the slope. In other words, the dynamic of $\mu_t$ includes a time-varying slope. In addition, an error term $\omega_{\mu,t}$ is added. Since the evolution equation describes the actual, unobservable underlying process, this error term is sometimes called *innovations*, as a source of change. But here, its variance is zero, which implies that the only source of variation in the mean level is the slope. This particular TSDLM, called Integrated Random Walk, is more parsimonious and efficient for extracting smooth trends (Durbin and Koopman, 2012). The last two equations describe, respectively, the changes in slope and seasonality equal to those of $t-1$ with *innovation* terms, $\omega_{\beta,t}$ and $\boldsymbol{\omega}_{\mathbf{S},t}$ with non-zero variance. The matrix $\mathbf{G}$ depends on the form, either factorial or trigonometric, and the time unit of seasonality.

Considering a theoretical bi-monthly sampling frequency and the usual observation of a spring and a fall bloom in temperate marine ecosystems (Cushing, 1959; Longhurst, 1995), a trigonometric form with two harmonics was chosen for micro-phytoplankton abundances and dissolved oxygen concentrations. For cadmium concentrations, the model implemented has

only one harmonic in accordance with the well-known annual variations in concentrations related to the physiology of bivalves (e.g. Amiard et al, 1986).

The parameters of the model, i.e. observation and evolution variances, are estimated using the maximum likelihood method. Initial values equal to the variance of the observed time series were chosen. Values at $t = 0$ for the mean level, slope and seasonality and their variances were chosen very non-informative: all means were set to 0 and all variances were set to $10^7$, which are the default values used by Petris et al (2009). Lastly, as the aim of the process was to carry out retrospective analyses, smoothed distributions were chosen as model estimations, rather than filtered distributions.

### 2.3.1 Interventions

An intervention is the name given to a change in a model parameter value, e.g. in order to consider exogenous information. Here, only changes in the mean level and outliers are considered. Previously, we pointed out that the slope is the only factor in the evolution of the trend, because evolution variance of the mean level is fixed at 0. When a change in the mean level is suspected at time $t$, a non-zero variance of the mean level evolution is specified for this time unit and estimated by the maximum likelihood method. For this latter, the initial value is the variance of the observed time series.

Outliers are defined as measurements with unusually high observational variances, regardless of their cause. If an outlier is suspected at time $t$, then an observation variance increase parameter is added to the model. This is a quantity greater than or equal to 1 and which acts as a multiplier of the routine observation variance. As before, the numerical relevance of these specifications is evaluated by the maximum likelihood method, with the initial value for the optimization process being 1.

### 2.3.2 Automatic identification of outliers and mean level changes

The statistical approach used to detect changes and outliers is based on the definition of outliers in a box–and–whisker plot. In this tool, outliers of a standard normal distribution are values higher than 2.7 or lower than -2.7, which correspond, respectively, to the 0.35% highest and the 0.35% lowest values for a total of 0.7% of the whole distribution. These threshold values, -2.7 and 2.7, are used in conjunction with the results of a DLM. For outliers, standardized errors are examined, the distribution of which is supposed to be a standard normal distribution. Therefore, values higher than 2.7 or lower than -2.7 potentially correspond to outliers and are thus candidates for appropriate treatment.

The same approach is used for mean level changes. The considered values are called auxiliary residuals (Harvey et al, 1999). For level changes, one must consider the auxiliary residuals of the mean level, i.e. the smoothed values of the mean level error term $\omega_{\mu,t}$. In the particular case of the Integrated Random Walk used here, since the variance of its error term is zero, its smoothed values are also zero and thus non-informative on level changes. However, auxiliary residuals are used to examine smoothed innovations, i.e. changes from one-time unit to another, i.e. first differences in the mean level $\mu_t - \mu_{t-1}$. Given that $\mu_t - \mu_{t-1} = \beta_{t-1}$, the smoothed values of $\beta_{t-1}$ standardized by their smoothed variances and centred on their mean carry similar information. By definition and for the sake of clarity, these values are referred to hereafter as "auxiliary level residuals". They can be compared to the threshold values 2.7 and -2.7, beyond which values of the slope are considered exceptional, i.e. data bring the model to the limits of its adaptative capacity, and suggest a change in the level. Unlike the case of outliers, for which the suggestion is unambiguously

associated with only one date, the exceptional values of the auxiliary residuals of the slope often constitute sequences: the measure used to suggest a change in level does not have a single value. Experience has shown that using the highest value is not always relevant because the event inducing the exceptional adjustment may have occurred before or after the highest value, and may even have occurred before or after the sequence of exceptional auxiliary slope residuals. Given this, the operational procedure for identifying the candidate measure for a level change suggestion is defined as follows. When the auxiliary residuals of the slope constitute an exceptional sequence (i.e. greater than 2.7 or less than -2.7), this latter is extended to the neighbouring values greater than 2 or less than -2 (i.e. the 2.5% highest and 2.5% lowest values) and, in this sequence of extended auxiliary residuals, a change point is identified using the changepoint package in R (Killick and Eckley, 2014), through the cpt.mean function with "at most one change (AMOC)" method and "none" as the penalty.

### 2.3.3 Semi-automated analysis strategy

In the approach defined above, parameter estimations and automatic identification of interventions are not performed jointly. Consequently, an analysis strategy must be defined. Since the statistical identification of outliers and level changes is based on the results of the model, there is a need to estimate the model parameters one time. Then, the suggested outliers and changes in mean level define a second model, etc. Based on experience in previous works, it appeared that the suggestions for outliers were always validated by thematic experts, whereas changes in levels could be more difficult to keep. In order to take this into account, instead of withholding suggestions for both outliers and level changes, for each model outliers are identified and treated first, and if there are no outliers, level changes are processed, if any. Thus, potentially,

the first model is the beginning of an iterative loop leading to other models, possibly in an infinite way. This pitfall is partially avoided by identifying outliers and level changes with thresholds defining them as very rare events. Furthermore, the maximum number of models has been arbitrarily limited to ten.

This analysis strategy implies that successive models are nested and hence a likelihood-ratio test can be performed to test the significance of the likelihood gain. In addition, comparisons between models are made using the Akaike Information Criterion (AIC). These measures can be used to retrospectively evaluate models produced by the automated process. Therefore, the most relevant model is not necessarily the last one. Furthermore, it is possible to exclude certain automatic interventions considering that, while they are numerically plausible, their contribution is limited from an explanatory point of view. It is also desirable, although not essential, to be able to justify automatic interventions retained in the model. Conversely, prior knowledge of a change that does not correspond to any of the automatic interventions can be used to suggest a new change. This possibility of removing automatic interventions and adding manual interventions justifies qualifying the strategy as semi-automated.

### 2.3.4 Diagnosis of the models

The two hypotheses checked are the normality and the independence of the standardized residuals. The first is assessed using a quantile-quantile diagram or Q-Q plot in conjunction with a formal Kolmogorov-Smirnov test. Independence of residuals is frequently assessed using the Ljung-Box test. Stoffer and Toloi (1992) proposed a modified version to handle missing data. This latter is used for its relevance according to the objective of an automated process applied to thousands of time series. The number of lags used was equal to $2m$
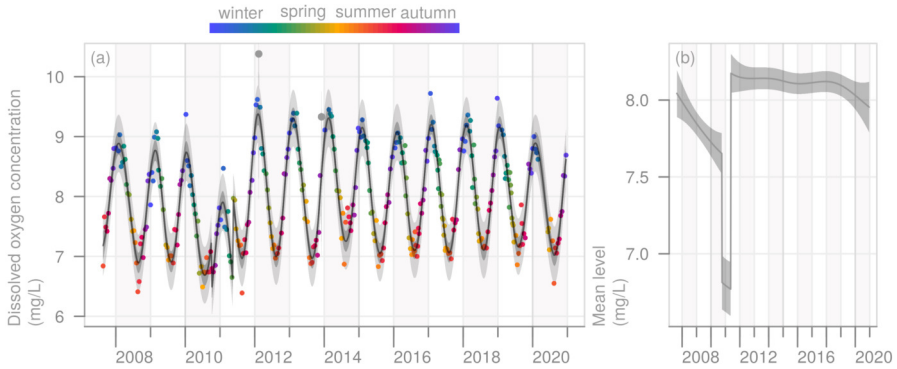
**Fig. 2** Dynamic linear model results for dissolved oxygen as measured at the bottom of the water column at "Teychan bis" site. a) Dots represent observations. Colours relate to the sampling date. Grey dots are observations treated as outliers. The solid line represents the model estimation, i.e. the mean level plus seasonality. The dark grey area corresponds to the 95% confidence interval. The light grey area is the 95% confidence interval of the observations. b) Mean level and its 95% confidence interval

as advocated by Hyndman and Athanasopoulos (2018), where $m$ is the period of the seasonality (e.g. $m = 52$ for the time unit "week") .

### 2.3.5 Software and hardware

All analyses and graphical representations were performed with the R software (R Core Team, 2022). Time series analyses with the DLM were performed with the dlm package for R (Petris, 2010). The use of a workstation, two processors, six cores each for a total of 24 processing threads with 64 GB of RAM made it possible to process the large number of time series at the same time.

## 3 Results

### 3.1 Dissolved oxygen

The initial time series included 359 observations. The time unit retained was the week, which generated aggregations of data for only 19 weeks and 354 weeks with missing data (i.e. 51% of weeks). Four successive models were produced through the automatic process. With the exception of the first one,
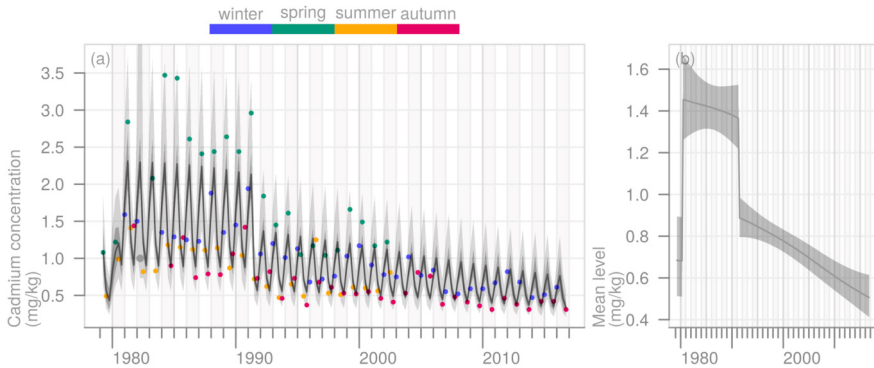
**Fig. 3** Dynamic linear model results for cadmium concentrations as measured in mussels sampled at the "La Mouclière" site. See Fig. 2 for details

they came with new automatic interventions that significantly increased the log-likelihood. The first outlier identified and treated as such is the highest value of the time series occurring in the third week of February 2012; the second one occurred in the first week of December 2013 with a value higher than expected in this season (Fig. 2a). Three automatic level changes were suggested: in the fourth week of January 2010, in the second week of October 2010 and in the last week of May 2011. In this last model, it appeared to us that the earliest level change suggestion was induced by a high, but not exceptional, value observed during the second week of the year. We subjectively decided to not retain this level change (Fig. 2b) as it did not provide deeper insight into understanding the time series, despite being a model with a significantly higher numerical log-likelihood. It should also be noted that despite specifying two harmonics, the model estimates only showed one. The standardized residuals of this final model satisfied the normality and independence hypotheses.

## 3.2 Cadmium

The observed time series was comprised of 118 observations; the selected time unit was the trimester with only one having aggregated data and 34 (i.e. 23%) with missing data. Twenty-eight of them (i.e. 82%) occurred starting
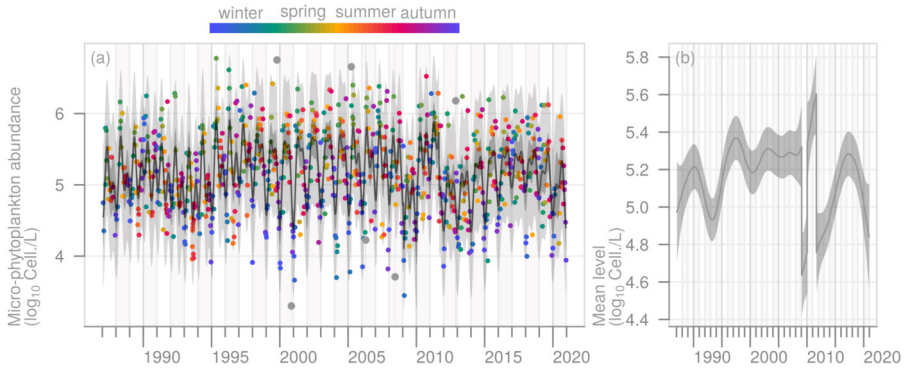
**Fig. 4** Dynamic linear model results for surface micro-phytoplankton abundances at the "Teychan bis" site. See Fig. 2 for details

from 2003, the year in which samples were taken twice a year, in the first and last trimester (Fig. 3a). The automatic process produced three successive models with significantly increasing log-likelihood. The last model included no outliers and three level changes: during the summer of 1980 and 1991 and the spring of 1982. It seemed to us that this latter was non-significant based on an overlap of approximately 20% of the 95% confidence intervals of the mean levels before and after the change. When this automatic intervention was discarded, the change in log-likelihood was not significantly different. Lastly, the standardized residuals of this latter model suggested that the lowest of the spring observations (i.e. 1982) was an outlier. The final model (Fig. 3) with this intervention significantly improved the log-likelihood. It should also be noted that the measurement method changes that occurred in 1986 and 2004 did not induce an automatic intervention nor did motivate a manual one. Its standardized residuals satisfied the normality but not the independence hypothesis.

## 3.3 Micro-phytoplankton abundances

Starting in 1987 and up to and including 2020, 825 samples have been analysed at the "Teychan bis" site. The selected time unit was the week. Only four

weeks were subject to data aggregation. In total, 47% of the weeks did not
have any measurements. Five successive models with significantly increasing
log-likelihood have been produced by the automatic process. There were two
suggested level changes during the fourth week of January 2010 and the first
week of September 2011. In addition, six outliers were identified and treated as
such: during 1999, 2000, 2005, 2006, 2008 and 2012. We noticed that the aux-
iliary residuals of the last model reached a value of -2.65 in mid 2008, which
was very close to the threshold of -2.7 at which an automatic search for a level
change is performed. We also visually perceived that the 2009 concentrations
appeared to be lower than those of the surrounding years. Furthermore, 16
consecutive residuals were negatives from the beginning of February to the end
May, suggesting out a lack of fit of the model. Lastly, from the third week of
January to the first week of February, the observed micro-phytoplankton con-
centration decreased from 350,000 Cell./L to 8,700 Cell./L. Thus, we decided
to add a manual intervention for a level change during the first week of
February 2009. The resulting final model (Fig. 4) had a significantly higher
log-likelihood. It should also be noted that there was no shift suggested at the
end of 2007, i.e. when the sampling depth changed. The standardized residuals
satisfied the normality assumption but rejected the independence hypothesis.

# 4 Discussion

The results showed the ability of our process to treat coastal environmental
time series, particularly with regards to irregular frequencies and missing data.
Automatic interventions for outliers and mean level changes induced numer-
ically enhanced models. Manual interventions resulting from statistical and
thematic review by experts allowed to obtain final models including an under-
taken subjectivity. The exogenous information supporting their relevance is

discussed below. Lastly, elements concerning model hypotheses checking are presented and some perspectives provided by the broad theoretical framework of Dynamic Linear Models are outlined.

## 4.1 Irregular frequencies and missing data

The algorithm determining the time unit of treatment was primary focused on minimizing data aggregation. The goal was to analyse the data as closely as possible to how they were collected. For example, for the 38 years of cadmium concentrations, the retained time unit was the trimester, which was also the sampling frequency during the first 24 years; after that, mussels were collected twice a year. The results proved the model's ability to handle time series with sampling frequency changes. For dissolved oxygen concentrations and micro-phytoplankton abundances, samples were taken every fortnight, however a weekly time unit was used for both of them. In fact, the two times series were missing roughly 50% of data. With long-term surveys, the planned frequency is a goal that is sometimes challenged by unforeseen events (e.g. a storm mean that a boat cannot be used, a micro-phytoplankton bloom may motivate an extra sample being taken). Hence, the sampling frequency was essentially every fortnight but with irregularities and shifts in sample weeks. This explained why the weekly time unit was chosen. DLM showed their capacity to manage such data. Lastly, analysing a time series with a time unit smaller than the sampling frequency means that information can be inferred even over time units with no data.

## 4.2 Automatic interventions

Successive models introducing new automatic interventions always significantly increased the log-likelihood.

### 4.2.1 Outliers

In our process, interventions for outliers were considered first. For a given outlier, the intervention consists of estimating a specific observation variance greater than the standard observation variance. This approach has several consequences. The value remained as an observation of the time series, but its weight is less than that of the other data. Thus, its influence is reduced in the estimation of the other parameters of the model but also in the filtering and smoothing processes of the observations. The observation variance is primarily impacted by the estimation of a specific variance for an outlier. By definition, the exceptional value often presents a large deviation from the other values and from the model. This deviation contributes significantly to the observation variance. Once the data is treated as exceptional, the deviation from the model has less weight and thus the estimated observation variance is reduced. The other parameters of the model are also mechanically affected by the treatment, but to a lesser extent. The reduction in the observation variance with the appropriate treatment of the exceptional data induces a noise that is lower than initially estimated. Since the information in the time series is, in the model, split between the observation-related noise and the underlying unobservable signal (i.e. structural part corresponding to the sum of the mean level and seasonality), a decrease in the noise part leads to an increase in the signal part. The signal-to-noise ratio is also modified. This ratio directly controls the adaptability of the model: a weak signal leads to a model with a long memory and low adaptability, whereas a strong signal leads to a model with a short memory and high adaptability. Thus, the reduction of the observation variance with the appropriate treatment of exceptional data allows, on the one hand, to tend towards a more realistic value estimation and, on the other hand, to obtain a signal-to-noise ratio and an adaptability in adequacy with

the standard process, i.e. excluding exceptional data. At the same time, the log-likelihood increases, reflecting the better adjustment of the model to the data, both because of the special treatment of exceptional values and because of the adaptability granted to the standard process. All these considerations clearly illustrate the importance of identifying and appropriately treating outliers, as well as the relevance of the "intervention" approach in the context of dynamic linear models.

By identifying exceptional values at time units for which the standardized errors belong to the lower 0.35% or the upper 0.35% of their distribution, the procedure unambiguously points to out-of-range results. However, confusion persists between the terms "exceptional" and "false/doubtful". The former should apply to a measurement for which the acquisition process is not questionable and/or for which there is exogenous information justifying the out-of-range nature of the measurement, for example, a one-time accidental pollution for chemical contaminants, a very large algal bloom for microphytoplankton abundances. Such a result is part of the history of the time series, must be shown in graphics and must be treated appropriately in the statistical analyses in order to limit its influence. On the other hand, a false or questionable value results from a deficient application of the data acquisition process or a deficient process, whether revealed through the existence of exogenous information or expertise pointing to it as incompatible with experience. In this case, the data should not appear in the representations of the time series and should not be included in the sets subjected to statistical analysis.

Lastly, it is worth mentioning that in recent years there has been a growing interest in detecting and analysing outliers as the main subject of the research, instead of being considered merely problematic and excluding such data from analysis (see Violle et al, 2017; Benhadi-Marín, 2018; Cook et al,

2021). Nevertheless, it is not always easy to fully understand the ecological origin of exceptional events, once the spurious nature is ruled out (e.g. six outliers identified on micro-phytoplankton time series). However, re-integrating rare or exceptional events into ecological science seems to be crucial for a better understanding of the extent of natural variations and our interpretation of biological processes (Cook et al., 2021).

### 4.2.2 Mean level shifts

Level changes are the second type of intervention implemented; their effects on the model are also important. By allowing the mean level to make a structural jump when justified, the deviations of the observations from the model are reduced, and with them the variance of the observations. This affects the signal-to-noise ratio and thus the adaptability and finally the fit of the model to the data, which induced an increase in the log-likelihood. Beyond this mechanical aspect of the intervention, the change in level constitutes a major event in the history of the time series, which needs to be argued. The identification of events synchronized with the suggested structural changes does not pretend to establish causal relationships. At best, the congruence highlighted allows to point out possible explanations to be explored. Indeed, an event (e.g. mild winter, storm) is the meeting of a set of particular conditions that contribute in unequal parts to a structural shock. In fact, apparently similar events do not necessarily have the same effects. Changes in level induced a different problem for thematic expertise. On the one hand, the increase in likelihood and the significant nature of this increase constitute arguments in favour of accepting interventions on the level. On the other hand, the difficulty in identifying exogenous elements potentially linked to these changes and the scientific implication, or even the economic, legal and societal consequences, to designate them as causal are obstacles to the validation of the shifts suggested in

the series. Lastly, the numerically optimal model may present an evolution to which the expert cannot relate his knowledge and experience. This pitfall, which often comes under the heading of "over parametrization", goes against the principle of parsimony and must be avoided. It is therefore a matter of the statistician and the expert having a dialogue so as to establish a model with undertaken subjectivity that is enlightening on the history of the time series, i.e. useful in Box's sense: "*All models are wrong, but some are useful*" (Box, 1976).

## 4.3 Statistical and thematic expertises

It follows that, for a given time series, the last model produced by the automatic process is not necessarily the finally accepted one and the three examples presented before illustrated how models can be modified. In the dissolved oxygen example, a level change has been not accepted. A level change should be motivated by at least two observations at the new level, preferably successive, otherwise it is an outlier. Here, the automatic level change intervention was only induced by one measurement and no particular exogenous information was found to maintain it. Hence, we decided to not keep this change, significantly losing in numerical likelihood but gaining in parsimony and assuming more easily the story told by the model. However, the two other level changes were induced by several data and appeared more relevant. A deeper investigation showed that the probe used to measure dissolved oxygen could not have been calibrated from 2010 October to 2011 May, which were the years and months of the two automatic interventions. Hence, this level change was most probably a probe drift highlighted by the model. This was identified as an artefactual induced change.

In the cadmium concentration example, the automatic intervention for the mean level in the spring of 1982 was abandoned and replaced by a manual intervention for an outlier. This model adjustment and the dissolved oxygen one pointed out that automatic process may lead to misdiagnosis and hence results should always be carefully examined and eventually models could be amended. There was also a level change in the summer of 1980. We learned (ROCCH, pers. comm.) that, despite a clearly defined quarterly sampling strategy during the first years of the network, the actual sampling date could be far from mid-trimesters, which was the actual target. The seasonality of the cadmium concentration is related to the mussel's physiological state and according to this point, early January is clearly different from end of March. Furthermore, it appeared that the median dry matter percentage rose from about 22% to 25% in 1985. This is related to the freeze-drying step of the measurement process and could have an impact on measurements. These two elements pointed out that the setup of a survey network may require several years to adjust the sampling strategy and measurement protocol. Hence, the hypothesis that the level change occurring in 1980 was linked to the early stages of ROCCH could not be rejected. But the shift in 1991 did not belong to this period. The investigation led to a two-step explanation. Firstly, in the Seine estuary, the cadmium concentration in mussels was related to factories used for the fertilizer industry (Chiffoleau et al, 2001). It has been found that there was and still is an active fertilizer factory on "La Charente" River (cf. coordinates 45°56′50.8″N 0°56′06.6″W), in the city of Tonnay-Charente, 15 km from the Charente estuary. Secondly, on January 23, 1991, a ministerial decree stipulated that the discharges from these factories must be at most equivalent to those induced by onshore storage (Lalonde, 1991). Lastly, it can be seen on aerial photos[2] that no such onshore storage existed in June 1989 but were built in July 1991.

---

[2]Available on https://remonterletemps.ign.fr.

Hence, this synchronicity between the level change suggested by the method and the ministerial decree added to the new onshore storage pleaded for a causal relationship. This was identified as an anthropogenic induced change.

In the micro-phytoplankton abundances example, a level change has been manually added the first week of February 2009 according to a close examination of data, residuals and auxiliary residuals. It turned out that on 24 January 2009, an exceptional storm named Klaus went through the south-west of France (Liberato et al, 2011). This event lasted for three days with hurricane-force winds and a path of highest winds precisely over Arcachon Bay. Here again, the concordance in time of a large decrease in micro-phytoplankton abundances (i.e. 340,000 Cell./L, 97%) and this major meteorological event suggested that there is a causal link between them, despite the fact that it had not been specifically explained yet. Nonetheless, storms and associated high winds are well-known for the disturbances generated over micro-phytoplankton communities, which depends on the characteristics of the storm (e.g. winds, precipitation) and the previous physical-chemical conditions of the water column (Wetz and Paerl, 2008). Decreases in abundance and biomass have been previously documented, especially through light limitation after sediment resuspension (Havens et al, 2011; Stockwell et al, 2020). This is a time-limited impact, the duration of which depends on the residence time of the water masses of the system; resilient communities usually return to pre-storms levels. The relevance of this manually added intervention illustrated that an automatic process based on a threshold had to be backed up by a careful examination of the data. The next level change occurred in the fourth week of January 2010 and, as previously mentioned it appeared to be associated with the process of ecosystem recovery, although the date of this change could open to discussion or even adjusted. The last automatic intervention was a

huge decrease in micro-phytoplankton abundance during September 2011. This change appeared to be relevant according to the data, which explained why it has been kept despite the absence of exogenous information to justify it. All these were identified as ecological induced changes.

## 4.4 What has not been seen, may be seen

In the oxygen example, two harmonics have been specified but only one has been shown in model estimates. Hence the model was obviously over-parametrized, but it should also be retained that specifying two harmonics does not force results with two harmonics, which is a recurring question for thematic experts. Similarly, the trend has been specified as a second order local linear, i.e. quadratic. However, this does not imply that the trend must be quadratic, but rather that the local trend can be up to quadratic, which includes a linear or constant trend.

In the example of cadmium (resp. micro-phytoplankton), changes in the measurement method (resp. sampling depth) did not induce an automatic intervention nor did they motivate a manual intervention. This lack of intervention is a result in itself. This does not mean that the changes did not impact the time series, but rather that according to the model in its structure, no changes seem necessary to explain the data. Thus, the process could also be used to plead in favour of an absence of change, e.g. when a site has been moved for operational reasons to a length considered insignificant according to the thematic.

## 4.5 Model hypotheses checking

The dynamic linear model approach is built upon three hypotheses: normality, independence and homoscedasticity of error terms. Trend, seasonality and confidence interval estimations are based on these assumptions. As DLM is an

inferential approach, we are founded in considering estimations to be correct to the extent that hypotheses are not violated. It follows three questions: 1) How to test? 2) Which impact of hypotheses violation on estimations 3) What needs to be done?

1. There is a very large amount of literature about normality testing. However, the use of the Q-Q plot is very common and its subjective interpretation may be completed by a formal test, such as the non-parametric Kolmogorov-Smirnov (KS) test. Independence is also very often evaluated through the Ljung-Box test, that should be replaced by the Stoffer-Toloi (ST) approach in order to appropriately consider missing data. Autocorrelation calculated with this latter test can be used to plot an estimation of the autocorrelated function (i.e. ACF). Among the three examples, all model residuals satisfied the normality assumption but only one satisfied the independence hypothesis. Our experience with thousands of time series is as follows. Firstly, it occurs that with time series with few data (e.g. 60 measures), the Q-Q plot (resp. ACF) may appear arced or "S" shaped (resp. with autocorrelations over the significant limit) but the KS (resp. TS) test concludes to normality (resp. independence). Conversely, for time series with a lot of data, the Q-Q plot and ACF may lead us to conclude to normality and independence, even though the tests were significant. This seems to us to be linked to the deviation from the null hypothesis that the test is able to highlight depending on the number of data. Hence the actual question should be about the amplitude of the deviation from the null hypothesis that we want to be able to highlight.

2. The question of the impact of hypotheses violation is often not addressed. Classically, biases in estimation are mentioned. In particular, this is an issue for variance estimations which are directly related to confidence intervals

and thus to any conclusion about the equality or difference of an estimated parameter with a given value.

3. The identification and appropriate treatment of outliers often improved normality and sometimes independence. Despite the fact that hypotheses are about residuals, when observations are patently non-Gaussian, a data transformation of observations may also lead to more favourable results. Lastly, with regard to independence, it may be suggested to add a low-order autoregressive (AR) process. However, our prior experiences with adding an AR term highlight the increase in model complexity without any important gain and without providing more explanatory elements that can be ecologically interpreted. This is particularly true when the order of autocorrelation added is greater than 1, which cannot be easily interpretable from an ecological point of view. On this particular point, we come close to the position of West and Harrison (1997, p. 349): "*It is sometimes tempting to explain more global movement in the series by such noise models when in fact they should be attributed to changes in trend or other components of the basic DLM*".

The homoscedasticity assumption has not been formally evaluated. Instead, data transformations have been performed. The log-normal distribution is ubiquitous in ecology: data are positive, distributions are often skewed and with a mean-variance relationship (Limpert et al, 2001). It is well known for species abundances that the log-transformation provide a stabilization of the variance. It is also a classical approach for chemical contaminants. Working with many dissolved oxygen concentration time series has led us to conclude that there is no need for any transformation.

As mentioned earlier, several tens of physical, chemical and biological variables, as well as derived ones (e.g. species richness and abundance, dominance,

nutrients ratios), were analysed using this modelling framework, resulting in different choices of data transformations. Square root transformation was used for counting data such as species richness and logit transformation for percentages (e.g. percentage of dinoflagellates related to total abundance). Ratios of log-normal variables (e.g. nutrients, dinoflagellates related to diatoms) are log-transformed: given that $log(A/B) = log(A) - log(B)$ and, by definition of a log normal distribution, $log(A)$ (resp. $log(B)$) is a normal distribution, and the difference of two normal distributions remains a normal distribution, it follows that the log-ratio is also normally distributed.

## 4.6 Future developments

Several methodological approaches are possible to improve and extend the results presented here. Seasonality has not been analysed in this paper but Ratmaya et al (2019) and Lheureux et al (2022) used it with micro-phytoplankton and nutrient concentrations. In a DLM, seasonality is time varying. Hence the seasonality factor is not a pattern identically repeated for each year, but a value adapted to the way seasons are deployed through time. This behaviour corresponds to the lived experience according to which, from year to year, seasons (e.g. winter) have the same main quality (e.g. cold) but are not identical in their intensities and stories. Thus, DLMs allow to extract phenological elements (e.g. duration, amplitude). As seasonality is modelled by a trigonometric function, annulations of its first derivative give dates at which a minimum and a maximum are attained. The use of two harmonics means that, eventually, dates for two minima and maxima are found. Annulations of the second derivative are inflection points, i.e. dates on which the increase or decrease in the observed variable is at tis maximum. Many other parameters may be defined

based upon this seasonality function (cf. Guallar et al, 2017; Karasiewicz and Lefebvre, 2022).

The transformations applied to the variables to stabilise the variances and thus satisfy the model hypotheses imply a quadratic relationship between mean and variance. This is a strong assumption that could be relaxed by using a time-varying variance model. Furthermore, this approach would allow for variations due to, for example, changes in methods or analytical laboratories. Another possibility to get around the transformation of the dependent variable would be the implementation of non-Gaussian models, e.g. setting the observation as log-normal. This approach remains compatible with the possibility of a time-varying variance. This last refinement requires the estimation to be performed through simulation methods (e.g. Gibbs sampler) or even using the Integrated Nested Laplacian Equation (INLA) approach (Rue et al, 2009). Then, these techniques open the door to more complex models. In particular, the current approach proceeds sequentially by fitting successive models; the results of one model are used, if necessary after a choice that may contain a subjective element, to specify the next one. Petris et al (2009) described a model for outliers and structural breaks that allows for the identification and treatment of outliers and structural changes for all model components including seasonality in a single estimation process. Fúquene et al (2015) published the results of a robust version of this type of model.

All measurement methods have quantification limits. In a time series, if a method changes over time, the quantification limit may also vary with time. These characteristics are very important (Helsel, 2009) and have an impact on estimated variances. Hence, the signal-to-noise ratio is affected and with it, the adaptive capacity of the DLM and thus the mean level. Allik et al (2016) have addressed this methodological issue in the context of the Kalman filter

and with a predictive purpose. This is a significant methodological investment for a gain in generality and relevance. Furthermore, the question arises as to what data should be entered: the limit of quantification, the measured value even if it is below the limit of quantification, or both.

It is possible to treat sites jointly to take advantage of the correlations that may exist in the evolution of the parameters studied. The methodological framework for this approach is provided by the SUTSE models, i.e. Seemingly Unrelated Time Series Equations (Petris et al, 2009). With these models, a correlation is not introduced for the noisy raw time series but for model parameters, primarily the mean level and seasonality. Moreover, by adding a hierarchical component to the specification of such models, in the form of a single random variable common to all locations, it is possible to distinguish in the observation variance a part common to all sites and a part specific to each of them. When a single laboratory is in charge of the measurements for the different sites considered, then the part common to all sites can be interpreted as the variance induced by the laboratory, i.e. only the sampling and laboratory operations, whereas the site-specific parts reflect the environmental variability specific to each monitoring site.

# 5 Conclusion

The process based upon Dynamic Linear Models described in this paper is currently used to provide preliminary analysis for about 20,000 time series (i.e. defined as at least two measurements performed on two different years on the same site) of variables and derived variables. The temporal granularity is as close as possible to the actual execution of the sampling strategy, and can even be smaller. DLM can consider missing data and changes in sampling frequency. The process appropriately suggests and treats outliers and level

changes. Resulting models could be amended according to statistical diagnostics and/or exogenous information. In our three examples, these latter showed that highlighted level changes could be artefactual, anthropogenic or ecological.

Although these results are relevant, they are only unrefined products. The opportunity offered to explore time series phenology remains widely open. Models may be enhanced through various ways including time-varying variance, non-gaussian models. More complex models may require the use of simulation method estimations but are in return a way to shift to a full Bayesian expression. Quantification limits, which are ubiquitous in ecology, need specific developments. Integrating the spatial dimension to take into account correlations across sites is already available through different ways. Lastly, this methodological approach is not new to ecology and will likely see further developments and fruitful results.

# Declarations

**Conflict of interest.**   The authors have no conflicts of interest to declare.

**Ethics approval.**   Not applicable.

**Consent to participate.**   Not applicable.

**Consent for publication.**   Not applicable.

**Availability of data.**   REPHY and ROCCH data are available at, respectively   https://doi.org/10.17882/47248   and   https://doi.org/10.17882/79255. The datasets generated during and/or analysed during the current study are available from the corresponding author on reasonable request.

**Code availability.**   The code can be obtained from the corresponding author.

**Authors' contributions.**   The authors have contributed equally to this work.

# References

Allik B, Miller C, Piovoso MJ, et al (2016) The Tobit Kalman filter: An estimator for censored measurements. IEEE Transactions on Control Systems Technology 24(1):365–371. https://doi.org/10.1109/TCST.2015.2432155

Amiard JC, Amiard-Triquet C, Berthet B, et al (1986) Contribution to the ecotoxicological study of cadmium, lead, copper and zinc in the mussel *Mytilus edulis*. Marine Biology 90:425–431. https://doi.org/10.1007/BF00428566

Auger-Méthé M, Newman K, Cole D, et al (2021) A guide to state–space modeling of ecological time series. Ecological Monographs 91(4):e01,470. https://doi.org/10.1002/ecm.1470

Belin C, Soudant D, Amzil Z (2021) Three decades of data on phytoplankton and phycotoxins on the French coast: Lessons from REPHY and REPHYTOX. Harmful Algae 102:101,733. https://doi.org/10.1016/j.hal.2019.101733

Benhadi-Marín J (2018) A conceptual framework to deal with outliers in ecology. Biodiversity and Conservation 27(12):3295–3300. https://doi.org/10.1007/s10531-018-1602-2

Benway HM, Lorenzoni L, White AE, et al (2019) Ocean time series observations of changing marine ecosystems: An era of integration, synthesis, and societal applications. Front Mar Sci 6:163–164. https://doi.org/10.3389/fmars.2019.00393

Box GEP (1976) Science and statistics. Journal of the American Statistical Association 71(356):791–799. https://doi.org/10.1080/01621459.1976.10480949

Chiffoleau JF, Auger D, Chartier E, et al (2001) Spatiotemporal changes in cadmium contamination in the Seine estuary (France). Estuaries 24(6):1029–1040

Cleveland WS, Devlin SJ (1988) Locally weighted regression: An approach to regression analysis by local fitting. Journal of the American Statistical Association 83(403):596–610. https://doi.org/10.1080/01621459.1988.10478639

Cook CN, Freeman AR, Liao JC, et al (2021) The philosophy of outliers: Reintegrating rare events into biological science. Integrative and Comparative Biology 61(6):2191–2198. https://doi.org/10.1093/icb/icab166

Cushing DH (1959) The seasonal variation in oceanic production as a problem in population dynamics. ICES Journal of Marine Science 24(3):455–464. https://doi.org/10.1093/icesjms/24.3.455

Durbin J, Koopman SJ (2012) Time Series Analysis by State Space Methods. Oxford University Press, https://doi.org/10.1093/acprof:oso/9780199641178.001.0001

Fúquene J, Álvarez M, Raúl Pericchi L (2015) A robust Bayesian dynamic linear model for Latin-American economic time series: "the Mexico and Puerto Rico cases". Latin American Economic Review 24(1):6. https://doi.org/10.1007/s40503-015-0020-z

Grouhel-Pellouin A, Verin F, Maheux F, et al (2022) Rocch dataset : chemical contaminants levels for shellfish area quality management. 2020-2021 data. https://doi.org/10.17882/79255

Guallar C, Bacher C, Chapelle A (2017) Global and local factors driving the phenology of *Alexandrium minutum* (Halim) blooms and its toxicity. Harmful Algae 67:44–60. https://doi.org/10.1016/j.hal.2017.05.005

Hamed KH, Ramachandra Rao A (1998) A modified Mann-Kendall trend test for autocorrelated data. Journal of Hydrology 204(1):182–196. https://doi.org/10.1016/S0022-1694(97)00125-X

Harvey A, Jan Koopman S, Penzer J (1999) Messy times series: a unified approach. In: Fomby TB, Carter Hill R (eds) Messy Data – Missing Observations, Outliers, and Mixed-Frequency Data, Advances in Econometrics, vol 13. Emerald Group Publishing Limited, p 103–143

Hastie T, Tibshirani R (1990) Generalized Additive Models. Chapman & Hall/CRC Monographs on Statistics & Applied Probability, Taylor & Francis

Havens KE, Beaver JR, Casamatta DA, et al (2011) Hurricane effects on the planktonic food web of a large subtropical lake. Journal of Plankton Research 33(7):1081–1094. https://doi.org/10.1093/plankt/fbr002

Helsel D (2009) Much Ado About Next to Nothing: Incorporating Nondetects in Science. The Annals of Occupational Hygiene 54(3):257–262. https://doi.org/10.1093/annhyg/mep092

Hernández-Fariñas T, Soudant D, Barillé L, et al (2013) Temporal changes in the phytoplankton community along the French coast of the eastern English Channel and the southern Bight of the North Sea. ICES Journal of Marine Science 71(4):821–833. https://doi.org/10.1093/icesjms/fst192

Hernández-Fariñas T, Bacher C, Soudant D, et al (2015) Assessing phytoplankton realized niches using a French national phytoplankton monitoring network. Estuarine, Coastal and Shelf Science 159:15–27. https://doi.org/10.1016/j.ecss.2015.03.010

Hyndman R, Athanasopoulos G (2018) Forecasting: principles and practice, 2nd edn., OTexts, Melbourne, Australia, chap Residual diagnostics. URL http://OTexts.com/fpp2, Accessed on 27 November 2022

Ibanez F, Fromentin J, Castel JC (1993) Application de la méthode des sommes cumulées à l'analyse des séries chronologiques en océanographie. Comptes rendus de l'Académie des sciences Sciences de la vie 316:745–748

Karasiewicz S, Lefebvre A (2022) Environmental impact on harmful species *Pseudo-nitzschia* spp. and *Phaeocystis globosa* phenology and niche. Journal of Marine Science and Engineering 10(2). https://doi.org/10.3390/jmse10020174

Killick R, Eckley IA (2014) changepoint: An r package for changepoint analysis. Journal of Statistical Software 58(3):1–19. https://doi.org/10.18637/jss.v058.i03

Koslow JA, Couture J (2013) Ocean science: Follow the fish. Nature 502(7470):163–164. https://doi.org/10.1038/502163a

Lalonde B (1991) Arrêté du 23 janvier 1991 relatif aux rejets de cadmium et d'autres substances dans les eaux en provenance d'installations classées pour la protection de l'environnement. Journal officiel "Lois et Décrets", URL https://www.legifrance.gouv.fr/jorf/id/JORFTEXT000000691416

Legendre P, Legendre L (2012) Numerical Ecology. Developments in Environmental Modelling, Elsevier Science

Levy O, Ball BA, Bond-Lamberty B, et al (2014) Approaches to advance scientific understanding of macrosystems ecology. Frontiers in Ecology and the Environment 12(1):15–23. https://doi.org/10.1890/130019

Lheureux A, David V, Del Amo Y, et al (2022) Bi-decadal changes in nutrient concentrations and ratios in marine coastal ecosystems: The case of the Arcachon bay, France. Progress in Oceanography 201:102,740. https://doi.org/10.1016/j.pocean.2022.102740

Liberato MLR, Pinto JG, Trigo IF, et al (2011) Klaus – an exceptional winter storm over northern Iberia and southern France. Weather 66(12):330–334.

https://doi.org/10.1002/wea.755

Limpert E, Stahel WA, Abbt M (2001) Log-normal Distributions across the Sciences: Keys and Clues: On the charms of statistics, and how mechanical models resembling gambling machines offer a link to a handy way to characterize log-normal distributions, which can provide deeper insight into variability and probability—normal or log-normal: That is the question. BioScience 51(5):341–352. https://doi.org/10.1641/0006-3568(2001) 051[0341:LNDATS]2.0.CO;2

Longhurst A (1995) Seasonal cycles of pelagic production and consumption. Progress in Oceanography 36(2):77–167. https://doi.org/10.1016/ 0079-6611(95)00015-1

Neaud-Masson N (2015) Observation et dénombrement du phytoplancton marin par microscopie optique photonique – spécifications techniques et méthodologiques appliquées au rephy. document de méthode. Report, Ifremer, Rue de l'île d'Yeu, BP 21105, 44311, Nantes Cedex 03, URL https: //archimer.ifremer.fr/doc/00292/40293/

O'Brien TD, Lorenzoni L, Isensee K, et al (2017) What are marine ecological time series telling us about the ocean? a status report. IOC Technical Series 129, IOC-UNESCO

Petris G (2010) An r package for dynamic linear models. Journal of Statistical Software 36(12):1–16. https://doi.org/10.18637/jss.v036.i12

Petris G, Petrone S, Campagnoli P (2009) Dynamic Linear Models with R. Use R!, Springer New York

R Core Team (2022) R: A Language and Environment for Statistical Computing. R Foundation for Statistical Computing, Vienna, Austria

Ratmaya W, Soudant D, Salmon-Monviola J, et al (2019) Reduced phosphorus loads from the Loire and Vilaine rivers were accompanied by increasing eutrophication in the Vilaine bay (south Brittany, France). Biogeosciences 16(6):1361–1380. https://doi.org/10.5194/bg-16-1361-2019

REPHY (2021) Rephy dataset – french observation and monitoring program for phytoplankton and hydrology in coastal waters. metropolitan data. https://doi.org/10.17882/47248

Rue H, Martino S, Chopin N (2009) Approximate bayesian inference for latent gaussian models by using integrated nested laplace approximations. Journal of the Royal Statistical Society: Series B (Statistical Methodology) 71(2):319–392. https://doi.org/10.1111/j.1467-9868.2008.00700.x

Sen PK (1968) Estimates of the regression coefficient based on Kendall's tau. Journal of the American Statistical Association 63(324):1379–1389. https://doi.org/10.1080/01621459.1968.10480934

Soudant D, Beliaeff B, Thomas G (1997a) Dynamic linear bayesian models in phytoplankton ecology. Ecological Modelling 99(2):161–169. https://doi.org/10.1016/S0304-3800(97)01949-2

Soudant D, Beliaeff B, Thomas G (1997b) Explaining dinophysis cf. acuminata abundance in Antifer (Normandy, France) using dynamic linear regression. Marine Ecology Progress Series 156:67 – 74. https://doi.org/10.3354/meps156067

Stockwell JD, Doubek JP, Adrian R, et al (2020) Storm impacts on phytoplankton community dynamics in lakes. Global Change Biology 26(5):2756–2784. https://doi.org/10.1111/gcb.15033

Stoffer DS, Toloi CM (1992) A note on the Ljung—Box—Pierce portmanteau statistic with missing data. Statistics & Probability Letters 13(5):391–396. https://doi.org/10.1016/0167-7152(92)90112-I

Theil H (1950) A rank-invariant method of linear and polynomial regression analysis. i, ii, iii. Nederl Akad Wetensch, Proc 53:386–392, 521–525, 1397–1412

Utermöhl H (1958) Zur Vervollkommnung der quantitativen Phytoplankton-Methodik. Mitt Internationale Ver Theoretische und Angewandte Limnologie 9:1–38

Violle C, Thuiller W, Mouquet N, et al (2017) Functional rarity: The ecology of outliers. Trends in Ecology & Evolution 32(5):356–367. https://doi.org/doi.org/10.1016/j.tree.2017.02.002

West M (2014) Bayesian forecasting. In: Balakrishnan N, Colton T, Everitt B, et al (eds) Wiley StatsRef: Statistics Reference Online. John Wiley & Sons, Ltd, https://doi.org/10.1002/9781118445112.stat00219

West M, Harrison J (1997) Bayesian Forecasting and Dynamic Models (Springer Series in Statistics). Springer-Verlag

Wetz MS, Paerl HW (2008) Estuarine phytoplankton responses to hurricanes and tropical storms with different characteristics (trajectory, rainfall, winds). Estuaries and Coasts 31(2):419–429. https://doi.org/10.1007/s12237-008-9034-y