



# Estimation of phytoplankton pigments from ocean-color satellite observations in the Senegalo–Mauritanian region by using an advanced neural classifier

Khalil Yala<sup>1</sup>, N'Dèye Niang<sup>2</sup>, Julien Brajard<sup>1,4</sup>, Carlos Mejia<sup>1</sup>, Mory Ouattara<sup>2</sup>, Roy El Hourany<sup>1</sup>, Michel Crépon<sup>1</sup>, and Sylvie Thiria<sup>1,3</sup>

<sup>1</sup>IPSL/LOCEAN, Sorbonne Université (Université Paris 6, CNRS, IRD, MNHN), 4 Place Jussieu, 75005 Paris, France

<sup>2</sup>CEDRIC, CNAM, 292 rue Saint Martin, 75003 Paris, France

<sup>3</sup>UVSQ-LATMOS, 78035 Versailles, France

<sup>4</sup>Nansen Center, Thormøhlensgate 47, 5006 Bergen, Norway

**Correspondence:** Michel Crepon (crepon@locean-ipsl.upmc.fr)

Received: 15 February 2019 – Discussion started: 20 March 2019

Revised: 23 December 2019 – Accepted: 18 February 2020 – Published: 24 April 2020

**Abstract.** We processed daily ocean-color satellite observations to construct a monthly climatology of phytoplankton pigment concentrations in the Senegalo–Mauritanian region. Our proposed new method primarily consists of associating, in well-identified clusters, similar pixels in terms of ocean-color parameters and in situ pigment concentrations taken from a global ocean database. The association is carried out using a new self-organizing map (2S-SOM). Its major advantage is allowing the specificity of the optical properties of the water to be taken into account by adding specific weights to the different ocean-color parameters and the in situ measurements. In the retrieval phase, the pigment concentration of a pixel is estimated by taking the pigment concentration values associated with the 2S-SOM cluster presenting the ocean-color satellite spectral measurements that are the closest to those of the pixel under study according to some distance. The method was validated by using a cross-validation procedure. We focused our study on the fucoxanthin concentration, which is related to the abundance of diatoms. We showed that the fucoxanthin starts to develop in December, presents its maximum intensity in March when the upwelling intensity is maximum, extends up to the coast of Guinea in April and begins to decrease in May. The results are in agreement with previous observations and recent in situ measurements. The method is very general and can be applied in every oceanic region.

## 1 Introduction

Phytoplankton are the basis of the ocean food web and consequently drive ocean productivity. They also play a fundamental role in climate regulation by trapping atmospheric carbon dioxide (CO<sub>2</sub>) through gas exchanges at the sea surface and consequently lowering the rate of anthropogenic increase in the atmosphere of CO<sub>2</sub> concentration by about 25 % (Le Quéré et al., 2018). With the growing interest in climate change, one may ask how the different phytoplankton populations will respond to changes in ocean characteristics (temperature, salinity, acidity) and nutrient supply, which presents an important societal impact with respect to both climate and fisheries, with a possible effect on fish that graze phytoplankton via the marine food chain.

Methods for identifying phytoplankton have greatly progressed during the last 2 decades. Phytoplankton were first described by microscopy. Microscopy is time-consuming and unable to identify picoplankton. Imaging flow cytometry (IFC) has renewed microscopic methods, thanks to the speed at which they are able to characterize phytoplankton in a water sample (IOCCG, 2014). An alternative method is the analysis of seawater samples by high-performance liquid chromatography (HPLC), which is widely used to categorize broad phytoplankton groups such as phytoplankton functional type (PFT) or phytoplankton size class (PSC) (Jeffreys et al., 1997; Brewin et al., 2010; Hirata et al., 2011). HPLC enables the identification of 25 to 50 pigments within a single

analysis, which is much easier and faster to conduct than microscopic observations (Sosik et al., 2014). Each phytoplankton group is associated with specific diagnostic pigments, and a conversion formula, the so-called diagnostic pigment analysis, can be derived to estimate the percentage of each group from the pigment measurements (Vidussi et al., 2001; Uitz et al., 2010). HPLC measurements are now recognized as the standard for calibrating and validating satellite-derived chlorophyll *a* (chl *a* in the following) concentration and for mapping groups of phytoplankton (IOCCG, 2014).

The use of satellite ocean-color sensor measurements has permitted researchers to map the ocean surface at a daily frequency. Satellite sensors measure the sunlight, at several wavelengths, backscattered by the ocean. The downwelling sunlight interacts with the seawater through backscattering and absorption in such a manner that the upwelling radiation transmitted to the satellite (“water-leaving” reflectance) contains information related to the composition of the seawater. The light transmitted to the satellite depends on the phytoplankton cell shape (backscattering), its pigments (absorption) and the dissolved matter (e.g., CDOM).

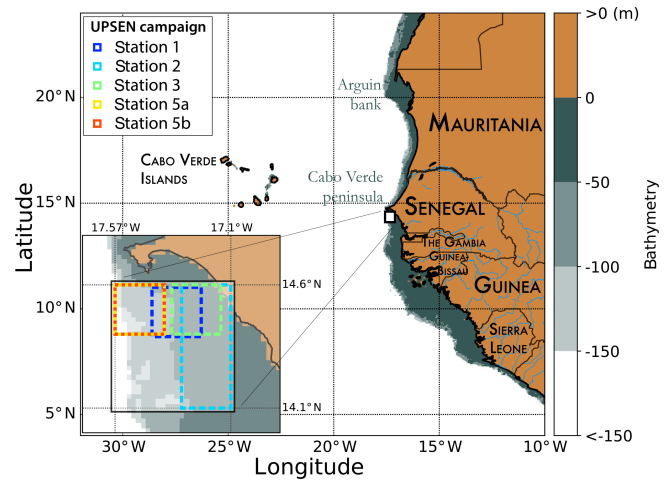
This upwelling radiation, the so-called remotely sensed reflectance  $\rho_w(\lambda)$ , is determined by the spectral absorption  $a$  and backscattering ( $b_b$ ;  $\text{m}^{-1}$ ) coefficients of the ocean (pure water and various particulate and dissolved matter) using the simplified formulation (Morel and Gentili, 1996)

$$\rho_w(\lambda) = G b_b(\lambda) / (a(\lambda) + b_b(\lambda)), \quad (1)$$

where ( $a$ ;  $\text{m}^{-1}$ ) is the sum of the individual absorption coefficients of water, phytoplankton pigments, colored dissolved organic matter and detrital particles; ( $b_b$ ;  $\text{m}^{-1}$ ) depends on the shape of the phytoplankton species.  $G$  is a parameter mainly related to the geometry of the situation (sensor and solar angles) but also to environmental parameters (wind, aerosols).

In the open ocean far from the coast (in case 1 waters), the light seen by the satellite sensor mainly contains information on phytoplankton abundance and diversity. Ocean-color measurements have been used intensively to estimate chlorophyll *a* concentration in the surface waters of the ocean and marginal seas and lakes (Longhurst et al., 1995; Antoine et al., 1996; Behrenfeld and Falkowski, 1997; Behrenfeld et al., 2005; Westberry et al., 2008).

It has been shown that it is also possible to extract additional information such as phytoplankton size classes (PSCs) by using some relationship between chlorophyll concentration and PSC (Uitz et al., 2006; Ciotti and Bricaud, 2006; Hirata et al., 2008; Mow and Yoder, 2010). These algorithms try to establish a relationship between the chl *a* concentration and the chl *a* concentration fractions associated with each of the three PSCs. Some of them (Uitz et al., 2006; Aiken et al., 2009) break down the chl *a* abundance into several ranges for each of which a specific relationship is computed. Others (Brewin et al., 2010; Hirata et al., 2011) are based on a continuum of chl *a* abundance. Studies have also been done to



**Figure 1.** Mauritania and Senegal coastal topography. The land is in brown, and the ocean depth is represented in meters by the color scale on the right side of the figure. The UPSÉN stations are shown at the bottom left of the figure.

estimate the phytoplankton groups (PFTs) by taking into account spectral information (Sathyendranath et al., 2004; Alvain et al., 2005, 2012; Hirata et al., 2011; Ben Mustapha et al., 2014; Farikou et al., 2015). This is of fundamental interest to the understanding of phytoplankton behavior and to modeling its evolution.

Due to highly nonlinear relationship linking the multispectral ocean-color measurements with the pigment concentrations, we proposed a neural network clustering algorithm (2S-SOM) able to deal with multi-variables linked by complex relationships. The 2S-SOM algorithm is well adapted to this complex task by weighting the different inputs. The clustering algorithm was calibrated on a restricted database composed of remotely sensed observations collocated with measurements taken in the global ocean.

In the present paper, we propose the retrieval of the major pigment concentrations from satellite ocean-color multispectral sensors in the Senegalo–Mauritanian upwelling, which is an oceanic region off the coast of West Africa where a strong seasonal upwelling occurs (Fig. 1).

The Senegalo–Mauritanian upwelling is one of the most productive eastern boundary upwelling systems (EBUSs) with strong economic impacts on fisheries in Senegal and Mauritania. Since the region has been poorly surveyed in situ, we have chosen to extract pertinent biological information from ocean-color satellite measurements. The region has been intensively studied through analysis of SeaWiFS (Sea-Viewing Wide Field-of-View Sensor) ocean-color data and AVHRR sea surface temperature as reported in Demarcq and Faure (2000), Sawadogo et al. (2009), Farikou et al. (2013, 2015), Ndoye et al. (2014), and more recently by Capet et al. (2017) with in situ observations.

The paper is organized as follows: in Sect. 2, we present the data we used (in situ and remote sensing observations). The mathematical aspect of the clustering method (2S-SOM) is detailed in Sect. 3. In Sect. 4 we present the methodological results. The spatiotemporal variability of the fucoxanthin and chl *a* concentration in the Senegalo–Mauritanian upwelling region are presented in Sect. 5, as are the results of the oceanic UPSEN campaigns. In Sect. 6 we discuss the results and the method. A conclusion is presented in Sect. 7.

## 2 Materials

In this study we used three distinct datasets: the first was used to calibrate the method, the second to conduct a climatological analysis of the Senegalo–Mauritanian upwelling region and the third was obtained during the oceanographic UPSEN campaign. These datasets are composed of satellite remote sensing observations and in situ measurements.

### 2.1 The calibration database (DPIG)

The calibration database (DPIG) comprises in situ pigment measurements collocated with satellite ocean-color observations by the SeaWiFS (Sea-Viewing Wide Field-of-View Sensor).

This DPIG is composed of 515 matched satellite observations and in situ measurements made in the global ocean (mainly in the North Atlantic and the equatorial ocean; Ben Mustapha et al., 2014). The matchup criteria were quite severe: we used satellite pixels situated at a distance of less than 20 km from the in situ measurement in a time window of  $\pm 12$  h. The geographic distribution of the 515 coincident in situ and satellite measurements is shown in Fig. 2. The matchup procedure between in situ and satellite observations is a crucial question to estimate remote sensing algorithms. If the parameters of the procedure are too severe, the number of collocated data points dramatically decreases. If the parameters are too large, it is the accuracy of the matching that decreases. We accordingly chose some compromise. Usually people use a matchup window of  $3 \times 3$  pixels (Alvain et al., 2005), which corresponds to a distance of less than 20 km between the satellite pixel and in situ measurement, since we deal with level 3 satellite observations whose pixel size is of the order of  $9 \times 9$  km. This criterion refers to the typical length of ocean variability (Lévy et al., 2012; Lévy, 2003).

In Fig. 3 we present the  $R^2$  coefficient between the in situ chl *a* and the SeaWiFS chl *a* computed by using the OC4V4 algorithm (O'Reilly et al., 2001) for the DPIG collocated observations. We remark that the two measurements are in good agreement at global scale. Each data point of DPIG is a vector having 17 components (five ocean reflectance  $\rho_w(\lambda)$  and  $Ra(\lambda)$  at five wavelengths (412, 443, 490, 510 and 555 nm), SeaWiFS chl *a*, five in situ pigment ratios, and in situ chl *a*

concentration). The in situ chl *a* concentration ranges between 0.007 and 3 mg m<sup>-3</sup> (see Table 1).

The five  $Ra(\lambda)$  are defined following Alvain et al. (2012):

$$Ra(\lambda) = \rho_w(\lambda) / \rho_{wref}(\lambda, chl\ a), \quad (2)$$

where the parameter  $\rho_{wref}(\lambda, chl\ a)$  is an average reflectance depending on the chl *a* concentration only that was computed according to the procedure reported in Farikou et al. (2015).  $Ra(\lambda)$  is a nondimensional parameter that depends on the chl *a* abundance at second order and is mainly sensitive to the secondary pigments (Alvain et al., 2012).

The DPIG database thus provides information on the existing links between the pigment composition and the SeaWiFS measurements. The pigment composition is defined by the pigment ratios, which are nondimensional variables of the form in the present study:

$$\text{pigment ratio} = DP / Tchl\ a, \quad (3)$$

which is defined as the ratio of the diagnostic pigment (DP) versus the total chl *a* ( $Tchl\ a = chl\ a + \text{divinyl chl } a$ ), according to Alvain et al. (2005).

The pigments of the DPIG and their statistical characteristics are given in Table 1. The statistical tests presented in Fig. 3 ( $R^2$  and RMSE) and in Table 1 (MEAN, SD, MIN, MAX) were computed in milligrams per cubic meter (mg m<sup>-3</sup>).

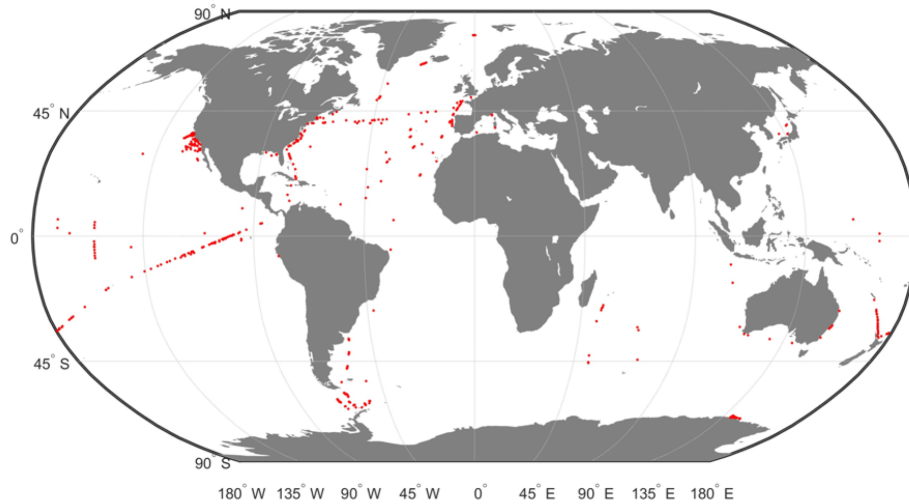
### 2.2 The Senegalo–Mauritanian upwelling satellite data (DSAT)

The satellite dataset we processed to retrieve the pigment concentration consists of five  $\rho_w(\lambda)$  and five  $Ra(\lambda)$  at five wavelengths (412, 443, 490, 510 and 555 nm), as well as the SeaWiFS chl *a* concentration observed in the Senegalo–Mauritanian upwelling region (8–24° N, 14–20° W; Fig. 3) during 11 years (1998–2009) by SeaWiFS. This dataset is denoted here as DSAT.

The satellite observations ( $\rho_w(\lambda)$  and chl *a* concentration) were provided by NASA with a resolution of 9 km. Due to the presence of Saharan dust in this region, very few estimations of satellite  $\rho_w(\lambda)$  and in situ chl *a* were available, and some satellite estimations of chl *a* could present strong overestimations (Gregg et al., 2004). For this reason, we reprocessed the  $\rho_w(\lambda)$  and chl *a* data with an atmospheric correction algorithm developed specifically for Saharan dust (Diouf et al., 2013; <http://poacc.locean-ipsl.upmc.fr/>, last access: 4 March 2020) in order to improve the satellite observations.

### 2.3 The UPSEN database

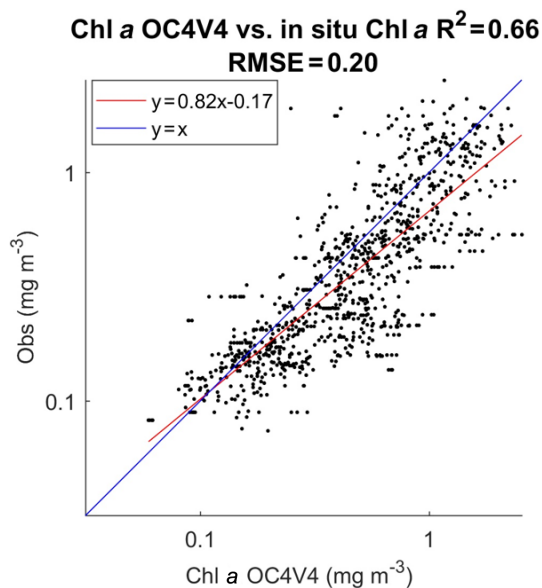
Recently, some HPLC measurements were made in the Senegalo–Mauritanian region during two oceanographic cruises (UPSEN campaigns) of the oceanographic ship



**Figure 2.** Geographic positions of the 515 in situ and satellite collocated measurements of the DFIG database.

**Table 1.** Pigments of the DFIG and their statistical characteristics: SD (standard deviation), MIN (minimum value), MAX (maximum value).

	Rdiviny A	Rperid	Rfuco	R19HF	Rzeax	Chl <i>a</i> in situ
Mean	0.1414	0.0272	0.1248	0.1859	0.1696	0.5292
SD	0.1584	0.0196	0.0971	0.0996	0.2063	0.5720
MIN	0.0037	0.0035	0.0053	0.0066	0.0027	0.007
MAX	0.8889	0.2027	0.8514	0.7654	1.5574	2.9980



**Figure 3.** Dispersion diagram of DFIG chl *a* computed from the SeaWiFS observations using the OC4V4 algorithm versus in situ chl *a*. The coefficient of determination  $R^2$  and the RMSE (root mean square error) were computed in milligrams per cubic meter ( $\text{mg m}^{-3}$ ).

*Le Suroit* from 7 to 17 March 2012 and from 5 to 26 February 2013 as reported in Ndoye et al. (2014) and Capet et al. (2017). The goal was to study the dynamics and the biological variability of the Senegalo–Mauritanian upwelling. During these campaigns, in situ HPLC measurements were carried out. We expected to be able to collocate them with the ocean-color VIIRS (Visible Infrared Imaging Radiometer Suite) sensor observations, whose wavelengths are close to those of the SeaWiFS. Unfortunately, we were only able to process satellite observations made on 21 February 2013 due to the presence of clouds and Saharan aerosols the other days. We processed the satellite observations provided by the VIIRS sensor at four wavelengths (443, 490, 510, 555 nm) for pixels in the vicinity of the ship stations (within a distance of 20 km) observed in a time window of  $\pm 12$  h and for which the satellite chl *a* was less than  $3 \text{ mg m}^{-3}$ , which is the limit of validity of our method imposed by the range of chl *a* observed in DFIG (mean of  $0.52 \text{ mg m}^{-3}$ ). Only five stations off the Cabo Verde peninsula fit these requirements (see Fig. 1 for their positions).

### 3 The proposed method (2S-SOM)

Classification methods were applied to retrieve geophysical parameters from large databases in several studies including weather forecasting (Lorenz, 1969; Kruizinga and Murphy,

1983), short-term climate prediction (Van den Dool, 1994), downscaling (Zorita and von Storch, 1999), reconstruction of oceanic pCO<sub>2</sub> (Friedrichs and Oschlies, 2009) and chl *a* concentration under clouds (Jouini et al., 2013). In the present study, we used a new neural network classifier, which is an extension of the SOM algorithms.

### 3.1 The SOM clustering

The SOM algorithms (Kohonen, 2001) constitute powerful nonlinear unsupervised classification methods. They are unsupervised neural classifiers that have been commonly used to solve environmental problems (Cavazos, 2000; Hewitson and Crane, 2002; Richardson et al., 2003; Liu and Weisberg, 2005; Liu et al., 2006; Niang et al., 2003, 2006; Reusch et al., 2007). The SOM aims at clustering vectors  $z_i \in \mathbb{R}^N$  of a multidimensional database  $D$ . Clusters are represented by a fixed network of neurons (the SOM), each neuron  $c$  being associated with the so-called referent vector  $w_c$  representing a cluster. The self-organizing maps are defined as an undirected graph, usually a rectangular grid of size  $p \times q$ . This graph structure is used to define a discrete distance (denoted by  $\delta$ ) between two neurons of the  $p \times q$  rectangular grid that presents the shortest path between two neurons. Each vector  $z_i$  of  $D$  is assigned to the neuron whose referent  $w_c$  is the closest in the sense of the Euclidean distance:  $w_c$  is called the projection of the vector  $z_i$  on the map. A fundamental property of an SOM is the topological ordering provided at the end of the clustering phase: close neurons on the map represent data that are close in the data space. The estimation of the referent vectors  $w_c$  of an SOM and the topological order is achieved through a minimization process in which the referent vectors  $w$  are estimated from a learning dataset (the DFIG database in the present case). The cost function is shown in Appendix A.

The SOMs have frequently been used in the context of completing missing data (Jouini et al., 2013), so the projected vectors  $z_i$  may have missing components. Under these conditions, the distance between a vector  $z_i \in D$  and the referent vectors  $w_c$  of the map is the Euclidean distance that considers only the existing components (the truncated distance or TD hereinafter).

### 3.2 The 2S-SOM classifier

In the present case, we used the 2S-SOM algorithm, a modified version of the SOM, which is very powerful in the case of a large number of variables. It automatically structures the variables having some common characteristics into conceptually meaningful and homogeneous blocks. The 2S-SOM takes advantage of this structuration of  $D$  and the variables into different blocks, which permits an automatic weighting of the influence of each block and consequently of each variable. The block weighting facilitates the clustering procedure by considering the most pertinent variables. The vectors of

DFIG defined in Sect. 2 can be decomposed into four blocks. The essence of this decomposition into blocks is that each of the 17 components of the DFIG vectors gathers information with a different physical influence in the classification phase. The composition of each block is done as follows.

The 2S-SOM is able to deal with a large quantity of variables, choosing those that are the most significant for the classification and neutralizing those that are the least significant. This is done by estimating weights on the blocks and the variables. We fully describe the 2S-SOM algorithm in Appendix A. In the following we use a simplified version of 2S-SOM in which only the blocks are weighted.

### 3.3 The calibration phase

Similarly to the standard SOM, the 2S-SOM is determined through a learning phase by using a more complex cost function (see Appendix A) that estimates for each neuron, in addition to the referent vector, a weight ( $\alpha$ ) for each block. For a neuron  $c$ , we define the weights of each block  $b$  ( $b = 1 \dots 4$ ).

At the end of the calibration phase, each element  $z_i$  of the dataset DFIG is associated with a referent  $w_c$  whose components are partitioned into four blocks. In the present study, the 2S-SOM is represented by a two-dimensional ( $9 \times 18 = 162$ ) grid that represents the partition of the DFIG dataset into different classes. Each class provided by the 2S-SOM is associated with a so-called referent vector  $w_c$  with  $c \in \{1 \dots 162\}$ . The size of the map has been determined by using the procedure provided by the SOM software available at <http://www.cis.hut.fi/projects/somtoolbox/download/> (last access: 4 March 2020).

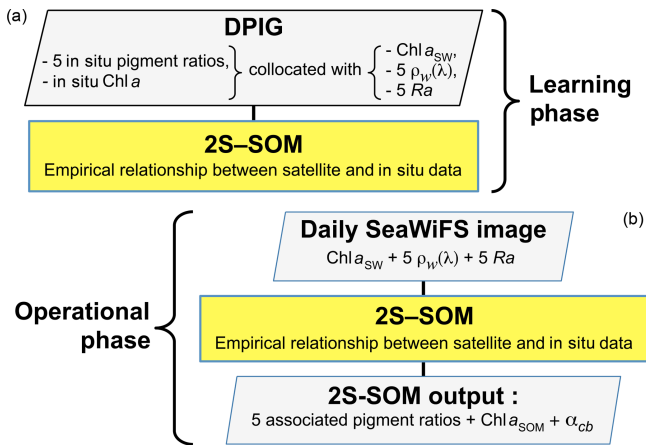
### 3.4 The pigment retrieval

In the second phase, which is an operating phase, we estimated the pigment concentration ratios of a pixel from its satellite ocean-color sensor observations only. The 11 ocean-color satellite observations ( $5\rho_w(\lambda)$ ,  $5Ra(\lambda)$  and chl *a*) of pixel  $PX_m$  were projected onto the 2S-SOM using the truncated Euclidean distance (Sect. 3.1). We select the neuron  $c$  associated with a referent vector whose 11 ocean-color parameters are the closest to those observed by the satellite sensor. The pigment ratios  $PX_m$  are those associated with the neuron  $c$ . At the end of the assignment phase, each pixel  $PX_m$  of a satellite image is associated with a referent vector  $w_c$ , which has six pigment concentration ratios among its 17 components. The flowcharts of the method (2S-SOM learning and pigment retrieval) are presented in Fig. 4.

## 4 Methodological results

### 4.1 Statistical validation of the method

The validation of the method was focused on the retrieval of the fucoxanthin ratio, which is a characteristic of diatoms,



**Figure 4.** Flowchart of the method: (a) learning phase; (b) operational phase that consists of pigment retrieval and the determination of the block parameters.

but the same procedure could be applied to any pigment. The hyper-parameter  $\mu$  (see Appendix A) was optimized in order to retrieve that ratio, while  $\eta$  was set as constant since only the blocks were weighted in the present study. Due to the small amount of data in the DFIG, we estimated the accuracy of the fucoxanthin retrieval by a cross-validation procedure, which is a powerful procedure in statistics. The principle is the following: we learned 30 2S-SOMs using 30 different learning datasets  $L_i$  constituting 90 % of DFIG taken at random, and then we computed a statistical estimator on the retrieved quantities using 30 test datasets (10 % of DFIG). The algorithm was as follows.

Starting with  $i = 1 \dots 30$ :

1. determination at random of a learning dataset  $L_i$  (90 % of DFIG) and a test dataset  $TL_i$  (10 % of DFIG);
2. training of a 2S-SOM  $M_i$  using  $L_i$  (see Sect. 3.2 and 3.3);
3. validation using  $TL_i$  according to the procedure described in Sect. 3.4; and
4. estimation of the  $RMSE_i$  and  $R_i^2$  on  $TL_i$  between the estimated and observed fucoxanthin ratios.

The flowchart of the cross-validation procedure is presented in Fig. 5 for the computation of the mean RMSE and  $R^2$  ( $R^2$ ,  $RMSE = \frac{1}{30} \sum_{i=1}^{30} R_i^2$ ,  $RMSE_i$ ).

Statistical parameters ( $R^2$  coefficients, RMSE and  $P$  values) of the cross-validation between the DFIG in situ pigments and the pigments given by the 2S-SOM averaged for the 30 2S-SOM realizations, which are presented in Table 2, show the good performance of the method.

**Table 2.** Statistical parameters ( $R^2$  coefficients, RMSE and  $P$  values) of the cross-validation between the DFIG in situ pigments and the pigments given by the 2S-SOM averaged for the 30 2S-SOM realizations.

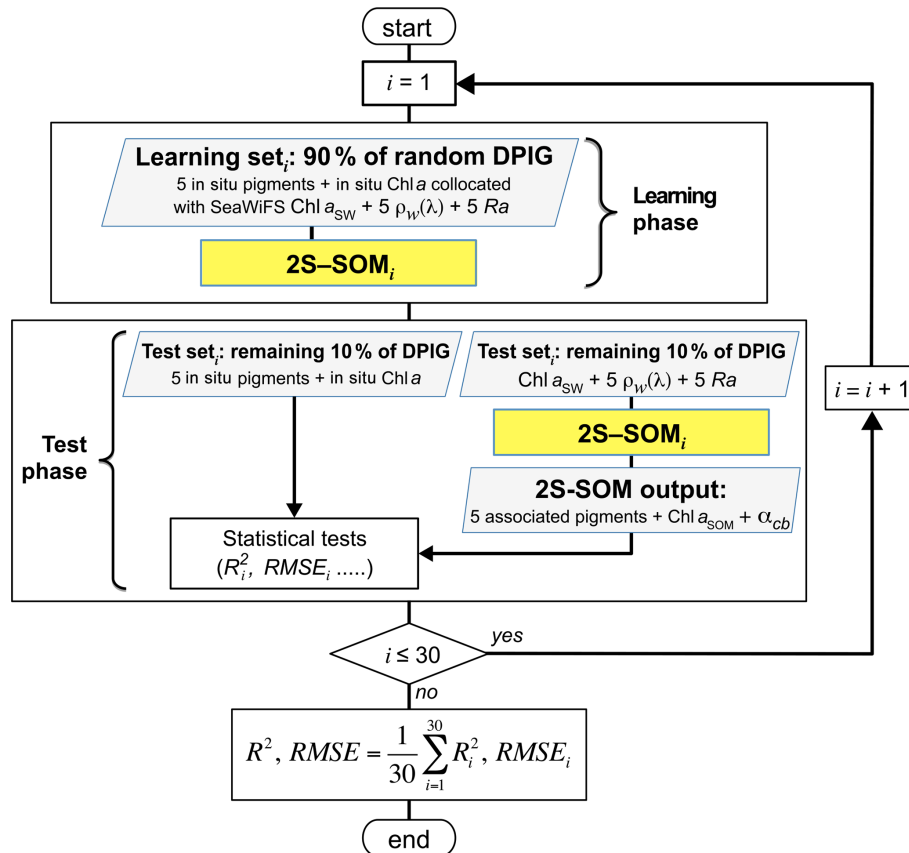
	$R^2$	RMSE ( $\text{mg m}^{-3}$ )	PVAL
Chl <i>a</i> SOM	0.84	0.22	0.001
DV Chl <i>a</i>	0.60	0.02	0.001
Fuco	0.87	0.02	0.001
Perid	0.81	0.01	0.001

#### 4.2 Analysis of the topology of the 2S-SOM

As explained in Sect. 3.2 and 3.3, the referent vector components ( $\mathbf{w}_c \in R^{17}$ ), which are estimated during the learning phase, are partitioned into four blocks B1, B2, B3 and B4. The hyper-parameter  $\mu$  was tuned in order to favor the accuracy of the retrieval of the fucoxanthin ratio. We recall that all the pigment ratios are estimated during the calibration phase, but in the present paper attention was focused on the fucoxanthin ratio when selecting the parameter  $\mu$ . In Fig. 6, we present six of the referent vector components of the 2S-SOM. These components are  $\rho_w(490)$ ,  $Ra(490)$ , SeaWiFS chl *a*, and the ratios of fucoxanthin, which is a specific diatom pigment, and of peridinin and divinyl. They exhibit a coherent topological order, with the components having values that are close together on the topological map. The remaining 11 components (not shown) exhibit the same coherent topological order. One can observe a very good topological order for the fucoxanthin ratio that was favored by the determination of the hyper-parameter  $\mu$ . Moreover, the bottom right region in the 2S-SOM (Fig. 6) may correspond to the diatoms with a good confidence since high fucoxanthin is associated with a high chlorophyll concentration and low peridinin. This is confirmed in Sect. 5 by looking at the geographical location of the different pigment concentrations (Figs. 8, 10, 11). Another important remark is that the value of each component presents a large range of variation of the same order as the range of variation found in the DFIG variables. This means that the 2S-SOM has captured most of the variability of the dataset.

Figure 6 shows a strong link between the values of the referent vectors for fucoxanthin and chl *a* (high fucoxanthin and chl *a* values at the bottom right of the 2S-SOM), while fucoxanthin is high and chl *a* low for the referent vectors at the bottom left of the 2S-SOM. Additional information will be provided by the  $Ra(490)$  values when the fucoxanthin is less closely linked to the chlorophyll.

In addition, for each neuron, the 2S-SOM provides a weight for each block ( $\alpha_{cb}$ ) and each variable ( $\beta_{cbj}$ ). For a given neuron  $c$  the weights ( $\alpha_{cb}$ ) of the blocks are normalized, their sum being 1. A value of 1 for one block (and there-



**Figure 5.** Flowchart of the cross-validation procedure for 30 partitions of the DPIG database.

for a value of 0 for the other blocks) indicates that the data in the neuron are gathered with respect to that block only because there is too much noise in the variables in the other blocks. By examining the weights on the map, one can see which block most influences the link between the satellite measurements and the pigment ratios.

In Fig. 7, we present the  $\alpha_{cb}$  values estimated during the learning phase of the four blocks (B1, B2, B3, B4). For some neurons, only the blocks related to the reflectance and the reflectance ratio are used for the definition of the neuron, while the weights for the two other blocks (pigments and chl  $a$ ) are null, indicating that for these neurons, in situ observations and SeaWiFS chl  $a$  are more noisy than the reflectance. These neurons correspond to very small chl  $a$  concentrations, which are estimated with large error. We remark that high  $\alpha$  values for chl  $a$  correspond to high chl  $a$  concentration values (bottom right of the chl  $a$  panel in Figs. 7 and 6). For these cases, the clustering assembled data that mainly depend on chl  $a$  concentration.

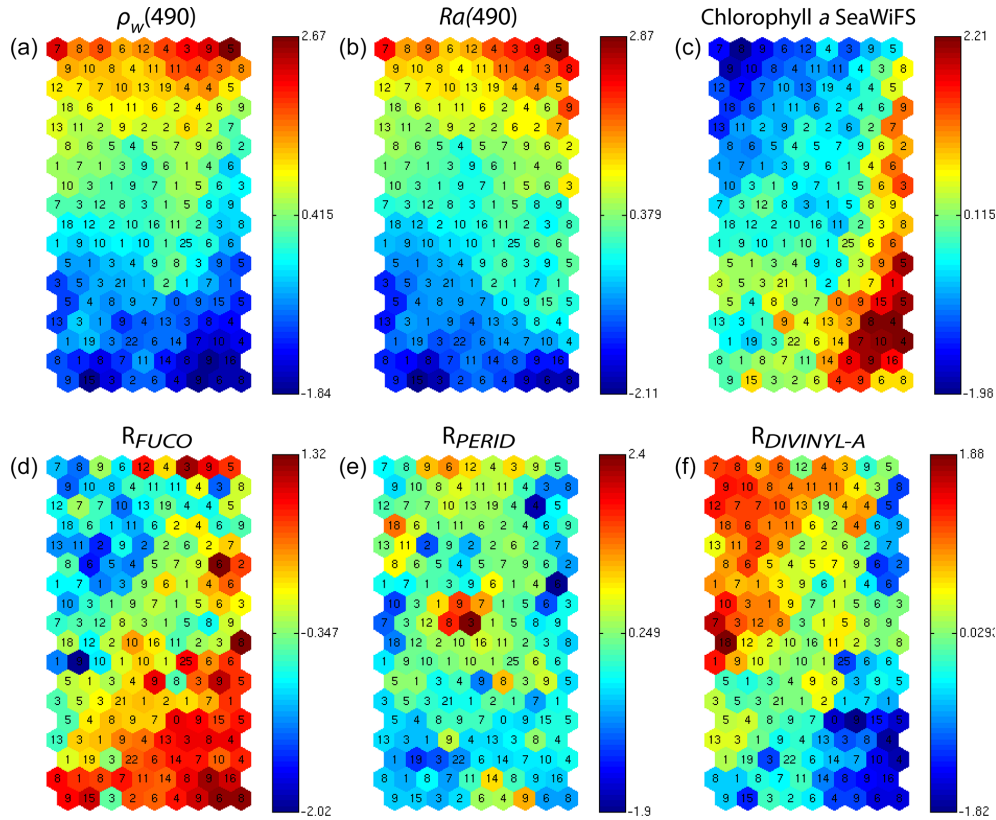
## 5 Geophysical results

In the present study, we apply the 2S-SOM (Sect. 3), which explicitly makes weighted use of the data according to their

specificity (ocean-color signals or in situ observations) to retrieve the fucoxanthin concentration from remotely sensed data in the Senegalo–Mauritanian upwelling region where in situ measurements are lacking. According to the good results of the cross-validation method as shown in Sect. 4.1, we expect that the 2S-SOM will provide pertinent results in a region that has been poorly surveyed.

### 5.1 The pigment estimation from SeaWiFS observations in the Senegalo–Mauritanian upwelling region

We decoded the DSAT database (Sect. 2.3) using the 2S-SOM for 11 years (1998–2009) of SeaWiFS data observed in the Senegalo–Mauritanian upwelling region (8–24° N, 14–20° W). This study was done according to the retrieval phase described in Sect. 3.4. For each day, we projected the 11 SeaWiFS observations ( $5\rho_w(\lambda)$ ,  $5Ra(\lambda)$  and chl  $a$ ) of each pixel on the 2S-SOM. At the end of the assignment phase, each pixel of a satellite image was associated with six pigment concentration ratios. The underlying assumption is that the link between the remote sensing information and the pigment ratios of a pixel is provided by the selected referent  $w_c$ . Thanks to the topological order provided by the 2S-SOM,



**Figure 6.** 2S-SOM. From left to right and top to bottom, values of the referent vectors for (a)  $\rho_w(490)$ , (b)  $R_a(490)$ , (c) SeaWiFS chl  $a$ , and the (d) fucoxanthin, (e) peridinin and (f) divinyl ratios. The number in each neuron indicates the amount of DPIG data captured at the end of the learning phase; the values indicated by the color bars are centered–reduced nondimensional values.

we expected that the best neurons chosen during the retrieval would give accurate concentration ratios. In Figs. 8, 10 and 11 we present the fucoxanthin concentration ratio estimation for 3 different days and the associated SeaWiFS chlorophyll images (1 and 6 January and 28 February 2003). Due to the limited size of the DPIG, the range of the ratio learned for fucoxanthin is between 0.3 % and 20 % with a mean of 10 %, and the chl  $a$  content is between 0.5 and 3  $\text{mg m}^{-3}$ . The statistical estimator we used cannot extrapolate what has not been learned, and for that reason we flagged the pixels in the SeaWiFS images that have a chl  $a$  concentration greater than 3  $\text{mg m}^{-3}$ .

Regarding the images obtained for 1 January 2003 in the Senegalo–Mauritanian region (Fig. 8a–d), we observe that the chl  $a$  (Fig. 8a) is very high at the coast and decreases offshore in accordance with the upwelling intensity as shown in the sea surface temperature (SST) image (Fig. 9). Moreover, we observed a persistent well-marked chl  $a$  pattern south of the Cabo Verde peninsula in the form of a W, which is the signature of a baroclinic Rossby wave (Sirven et al., 2019).

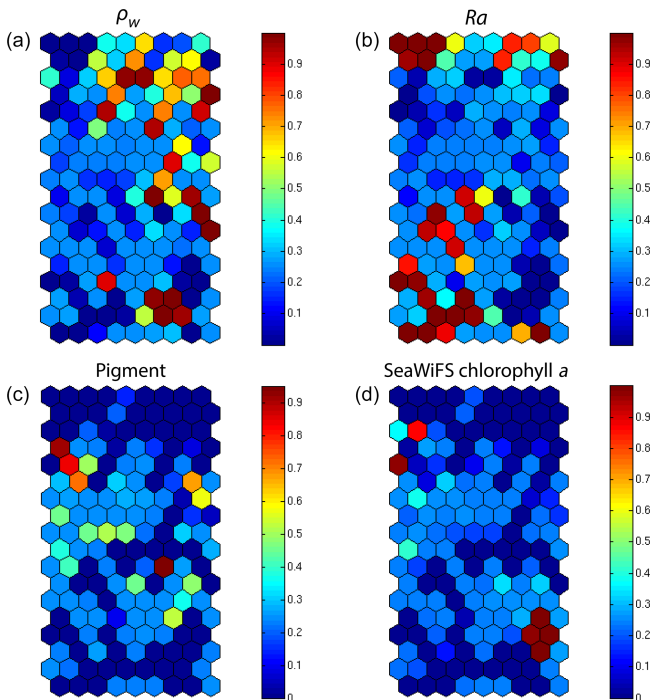
Except in the southern part of the region, the AOT (aerosol optical thickness) is low; this means that the atmospheric correction of the reflectance is quite small, which gives confidence in the ocean-color data products. The fucoxanthin con-

centration is maximum at the coast and decreases offshore as does the chl  $a$  concentration, in agreement with the works of Uitz et al. (2006, 2010). Fucoxanthin presents coherent spatial patterns. The peridinin concentration is somewhat complementary to that of fucoxanthin, with the low fucoxanthin concentration area corresponding to the high peridinin concentration area (northern part of Fig. 8b, d). This behavior is also observed in Fig. 10 (6 January 2003) and in Fig. 11 (28 February 2003), supporting the analysis shown in Fig. 8.

For 28 February, we selected two square box regions (Fig. 11), one near the coast (NSB, long.  $[-20^\circ, -18^\circ]$ , lat.  $[12^\circ, 14^\circ]$ ) and the other about 800 km offshore (OFB, long.  $[-28^\circ, -26^\circ]$ , lat.  $[12^\circ, 14^\circ]$ ). NSB waters correspond to upwelling waters, while OFB waters correspond to oligotrophic waters. We projected the 11 ocean-color parameters of the NSB and OFB pixels on the 2S-SOM.

Figure 12 presents the reflectance spectra (in blue) captured by three neurons of the 2S-SOM corresponding to pixels located in the NSB region (panels a–c) and those captured by three neurons corresponding to pixels located in the OFB region (panels d–f). The reflectance spectra of the associated referent vectors  $w$  are in yellow. The satellite reflectance spectra match the referent vector spectra; moreover, the fucoxanthin ratio varies inversely with the mean value of the





**Figure 7.** 2S-SOM. Weights ( $\alpha_{cb}$ ) of the four block parameters determined at the end of the learning phase; (a)  $\rho_w$ , (b)  $Ra$ , (c) pigment, (d) SeaWiFS chl  $a$ . The color bars show the percent of the weight estimated by 2S-SOM, with a value of 1 or 0 indicating that the data in the neuron are assembled with respect to that block only.

spectrum: the higher the fucoxanthin ratio, the smaller the mean value of the spectrum. The pigment concentration is greater near the coast.

We note a strong difference between the shape and the intensity of the nearshore (NSB) and offshore (OFB) spectra. The OFB spectra present mean values higher than those of the NSB spectra. This is due to the fact that NSB spectra were observed in a region where diatoms are abundant, as shown by the high value of the fucoxanthin concentration in this region (Figs. 8, 10 and 11), which is a proxy for diatoms along with a higher chl  $a$  concentration. In Fig. 12, we note the lower values of the coastal spectra at 443 nm, which can be interpreted as a predominant effect of spectral absorption by phytoplankton pigments and CDOM. The different spectra are close together in the OFB region and more disperse in the NSB region. This can be explained by the fact that the OFB region corresponds to case 1 waters, while the NSB region waters are close to case 2 waters and are influenced by the variability of nearshore process like turbidity or the presence of dissolved matter and dynamical instabilities.

We analyzed the weights of the blocks for the neurons selected in the analysis of the coastal (NSB) and offshore (OFB) boxes. Figure 13 presents the box plot of the weight  $\alpha_{cb}$  corresponding to the neurons belonging to the four blocks (B1, B2, B3, B4), with the constraint that the sum

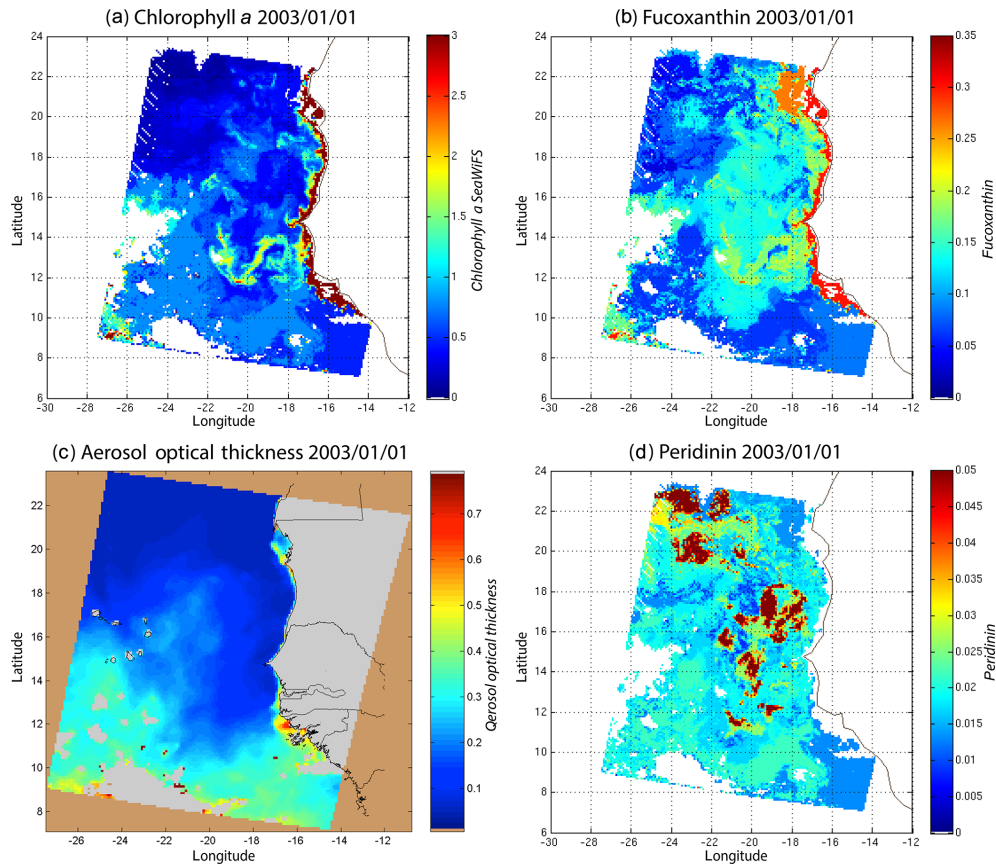
of the weights of a neuron is 1; a weight  $\alpha$  larger than 0.25 indicates the predominance of a block in the learning for the classification (see Sect. 3.5). It is clear that the weights for pixels near the coast (Fig. 13a) are different from those for offshore pixels (Fig. 13b). As already mentioned in Sect. 4.3 and also shown in Fig. 7, the weights of the 2S-SOM play a significant role in the 2S-SOM topology and consequently in the pigment retrieval. The weights of blocks B1 and B4 that take into account the influence of the pigment ratios and the chlorophyll content in the retrieval are very low for the offshore (OFB) oligotrophic region and more important for the coastal (NSB) region. The weights of the blocks B2 and B3, which take into account the influence of the reflectance ( $\rho_w(\lambda)$ ,  $Ra(\lambda)$ ), dominate for the offshore regions. In coastal waters, the weights of all the blocks are used, with a smaller influence of B3, which is associated with  $Ra$ . This gives information on the role played by the different variables in the classification in waters having different phytoplankton concentrations and compositions. It also shows the automatic adaptation of the 2S-SOM to the environment in order to optimize the clustering efficiency with respect to a classical SOM.

In order to study the seasonal variability of the fucoxanthin concentration with some statistical confidence in the Senegalo–Mauritanian upwelling region, we constructed a monthly climatology for an 11-year period (1998–2009) of the SeaWiFS observations by summing the daily pixels of the month under study. The resulting climatology is presented in Fig. 14 for December (Fig. 14a), March (Fig. 14b) and May (Fig. 14c), which correspond to the most productive period (Fig. 14c). The fucoxanthin concentration, and consequently the associated diatoms, presents a well-marked seasonality. Fucoxanthin starts to develop in December north of  $19^\circ$  N, presents its maximum intensity in March when the upwelling intensity is maximum, extends up to the coast of Guinea ( $12^\circ$  N) in April and begins to decrease in May when it is observed north of the Cabo Verde peninsula ( $15^\circ$  N) in agreement with the observations reported by Farikou et al. (2015) and Demarcq and Faure (2000).

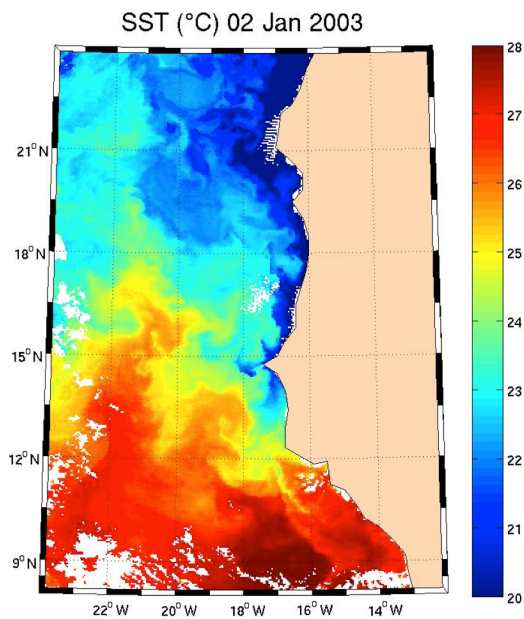
Figure 15 shows the fucoxanthin (in green) and the chl  $a$  (in blue) concentrations computed from satellite observations for an 11-year period of SeaWiFS observations in the NSB region. There is a good correlation in phase between these two variables but not in amplitude (a good coincidence of peak occurrence but weak correlation in peak amplitude), showing that the relationship between fucoxanthin and chl  $a$  is complex as mentioned by Uitz et al. (2006). In particular, there is a weak peak in fucoxanthin in October 2001, which is not correlated with a chl  $a$  peak.

## 5.2 Analysis of the UPSEN campaigns

Figure 16 shows, for each UPSEN station 1, 2, 3, 5a and 5b (see Fig. 1 for their geographical position), the averaged in situ UPSEN spectrum (in blue) and the referent spectrum (in



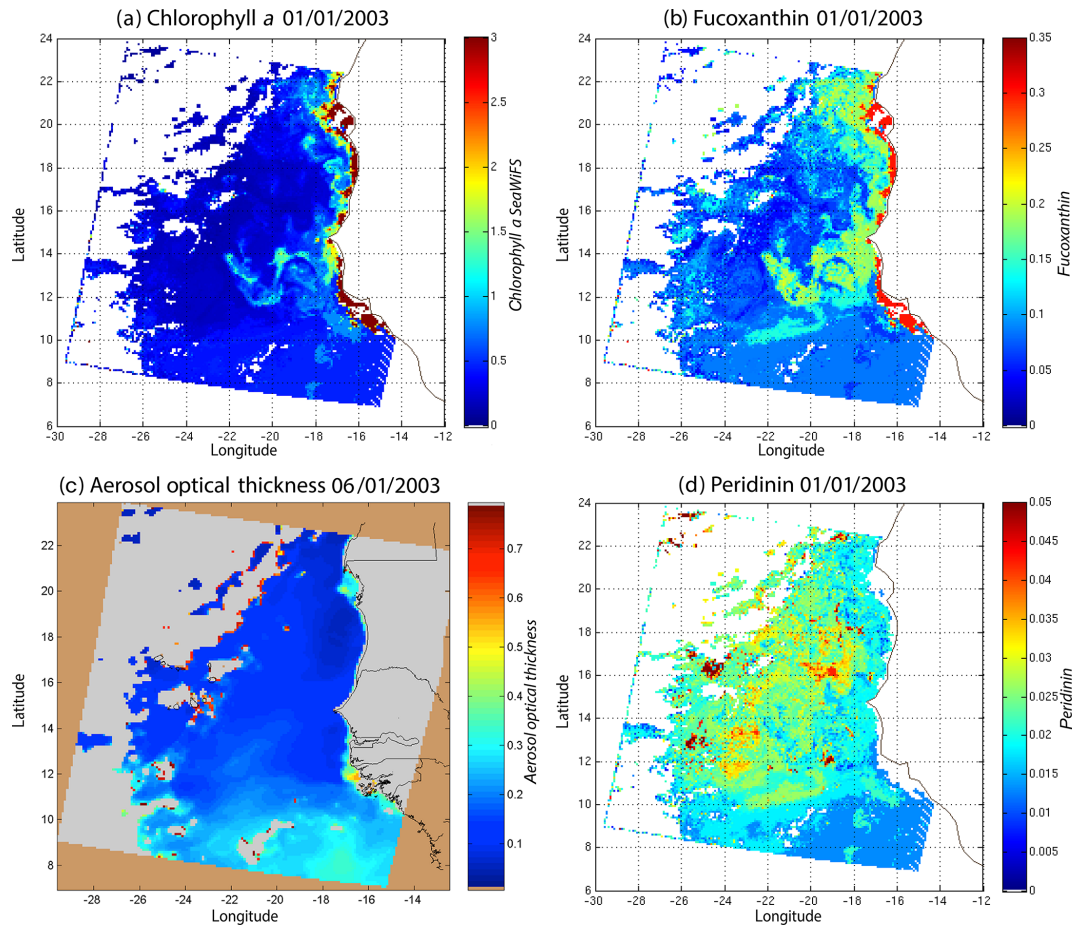
**Figure 8.** (a) Chl *a* concentration, (b) fucoxanthin ratio and (c) aerosol optical thickness, (d) peridinin for 1 January 2003. Panels (b, d) show that second-order information was retrieved, which is correlated with the chl *a* concentration (a) but not equivalent. The aerosol optical thickness (c) does not seem to contaminate the estimated parameters (fucoxanthin and peridinin ratios).



**Figure 9.** SST for 2 January 2003. Note the well-marked upwelling (cold temperature) north of 13° N.

red) of the 2S-SOM neuron captured by the collocated satellite VIIRS sensor observations. The referent spectrum is the mean of the different spectra captured by that neuron during the learning phase. Among these different spectra, there is one (black curve in Fig. 16) that is the closest to the UPSEN spectrum. Obviously, the black curve is closer to the blue curve than the red one that is flattened due to the averaging process. These three spectra are close together, showing the good functioning of the 2S-SOM.

Their shapes are close to those observed in the NSB region (Fig. 12) but their intensity is lower, meaning that their waters are more absorbing than the NSB waters due to a higher pigment concentration. In fact, the UPSEN stations were located close to the coast (Fig. 1) in the Hann bight south of the Cabo Verde peninsula, which is very rich in phytoplankton pigments. In Table 3, we present the fucoxanthin ratios associated with the referent vectors ( $R_{\text{fuco}_{2\text{S-SOM}}}$ ), the closest DPIG fucoxanthin ratios captured by the neuron of the referents and the fucoxanthin ratios measured during the UPSEN campaign. We note that the fucoxanthin ratios of the in situ measurements are in the range of the DPIG (see Table 1), which allows for the good functioning of the 2S-SOM esti-



**Figure 10.** (a) Chl *a* concentration, (b) fucoxanthin ratio, (c) aerosol optical thickness and (d) peridinin for 6 January 2003. Panels (b, d) show that second-order information was retrieved, which is correlated with the chl *a* concentration (a) but is not equivalent. It is found that the aerosol optical thickness (c) does not contaminate the estimated parameters (fucoxanthin and peridinin ratios).

matom. The pigment ratios obtained from ocean-color observations through the 2S-SOM are close to pigment concentrations measured at the ship stations, which confirms the validity of the method we have developed. We remark that the best 2S-SOM estimate of the fucoxanthin ratio with respect to the UPSEN in situ measurement is given at station 5b, which is the farthest off the coast. These results support the climatological study of the Senegalo–Mauritanian upwelling region we have done with the 2S-SOM (Sect. 5.1).

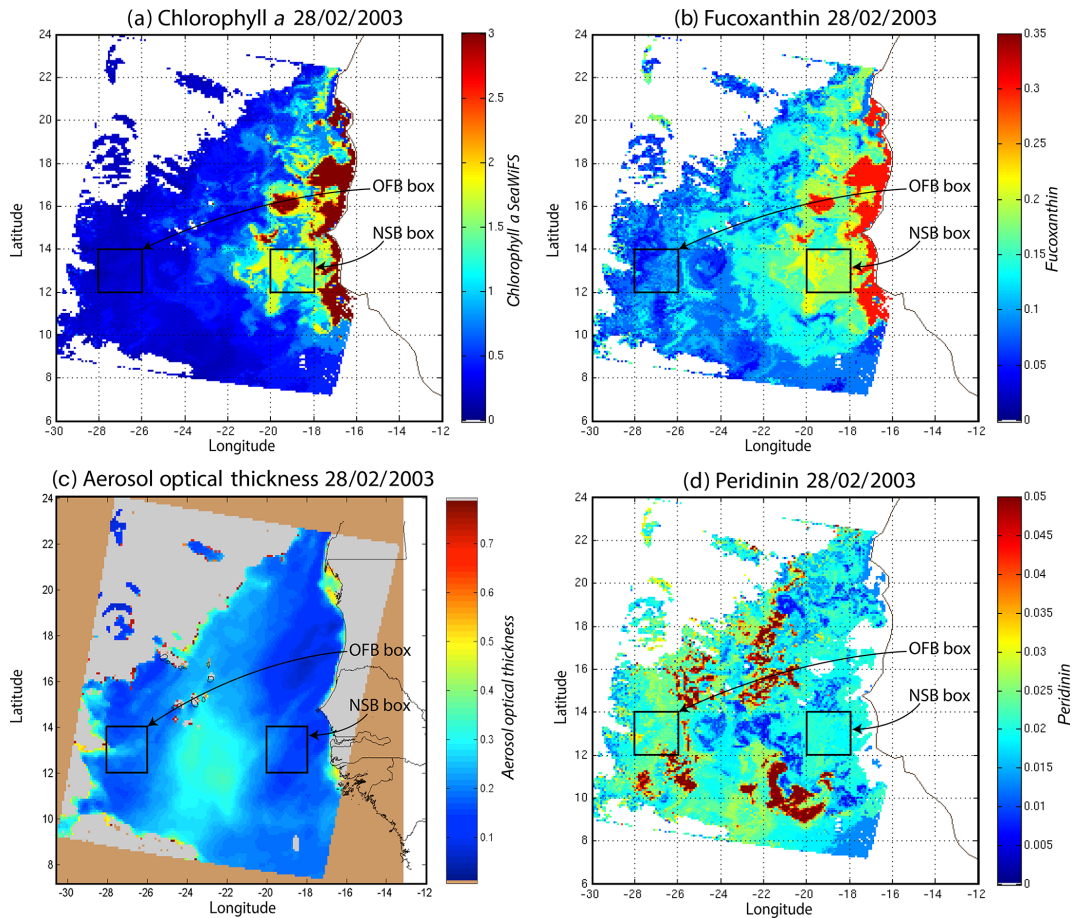
The 2S-SOM method gives pigment concentrations that are close to those obtained by in situ observations. The method could be applied to a large variety of other parameters in the context of studying and managing the planet Earth. The major constraint to obtaining accurate results is to deal with a learning dataset that statistically reflects all the situations encountered in the observations processed. Due to its construction, the method cannot be used to find values beyond the range of the learning dataset.

**Table 3.** For ship stations 1, 2, 3, 5a and 5b of the UPSEN campaign, we show the referent captured by the VIIRS observations, the fucoxanthin ratio associated with this referent (Rfuco-2S-SOM), the fucoxanthin ratio of the closest DFIG fucoxanthin ratio captured by the neuron of the referent and the fucoxanthin ratio measured in situ during the UPSEN campaign.

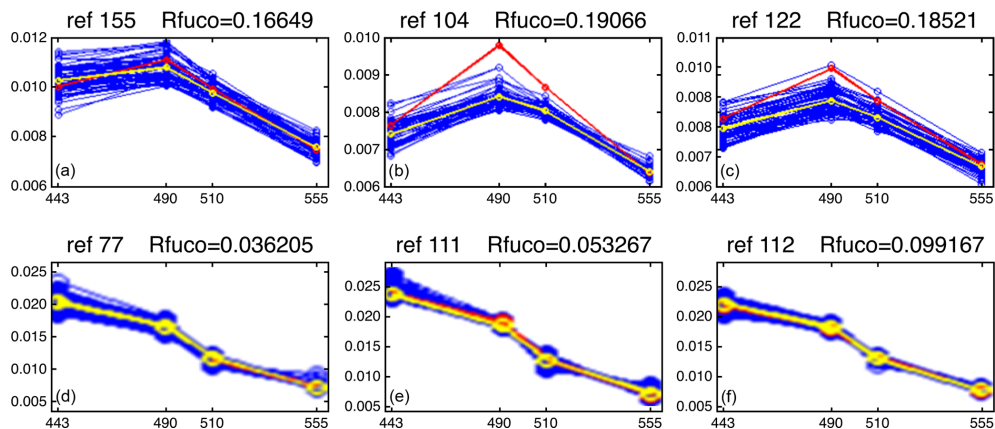
UPSEN station	Referent No.	Rfuco		
		2S-SOM	DFIG	UPSEN
STAT 1 17.3° E 14.5° N	126	0.213	0.236	0.378
STAT 2 17.2° E 14.4° N	126	0.213	0.236	0.391
STAT 2 17.2° E 14.5° N	126	0.213	0.236	0.436
STAT 5A 17.5° E 14.5° N	126	0.213	0.171	0.299
STAT 5B 17.5° E 14.5° N	143	0.242	0.258	0.295

## 6 Discussion

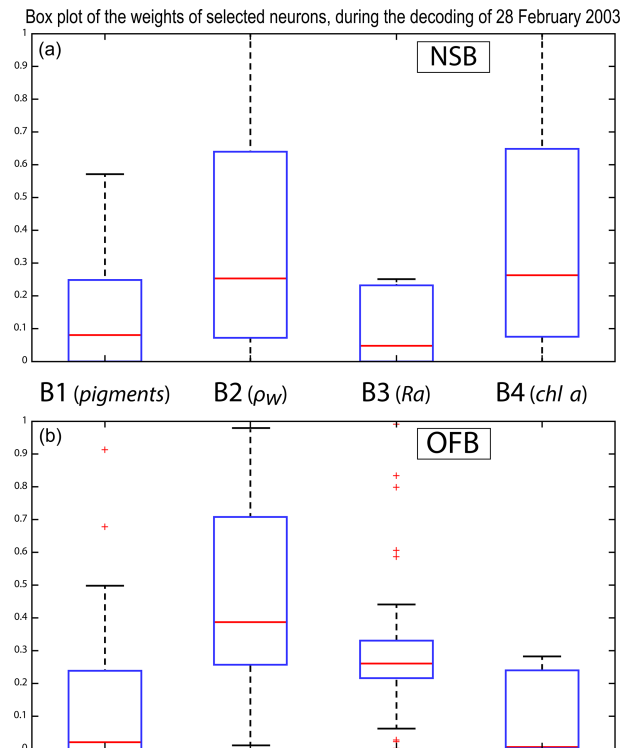
Machine-learning methods are powerful methods to invert satellite signals as soon as we have an adequate database to support the calibration. Several techniques have been used



**Figure 11.** (a) Chl *a* concentration, (b) fucoxanthin ratio, (c) aerosol optical thickness and (d) peridinin for 28 February 2003. Panels (b, d) show that second-order information was retrieved, which is correlated with the chl *a* concentration (a) but is not equivalent. It is found that the aerosol optical thickness (c) does not contaminate the estimated parameters (fucoxanthin and peridinin ratios). The positions of the NSB and OFB are outlined by black square boxes.



**Figure 12.** Reflectance spectra (in blue) captured on 28 February by six neurons whose referent vector spectra are in yellow: (a–c) pixels in the NSB region (long.  $[-20^\circ, -18^\circ]$ , lat.  $[12^\circ, 14^\circ]$ ); (d–f) pixels in the OFB region (long.  $[-28^\circ, -26^\circ]$ , lat.  $[12^\circ, 14^\circ]$ ).



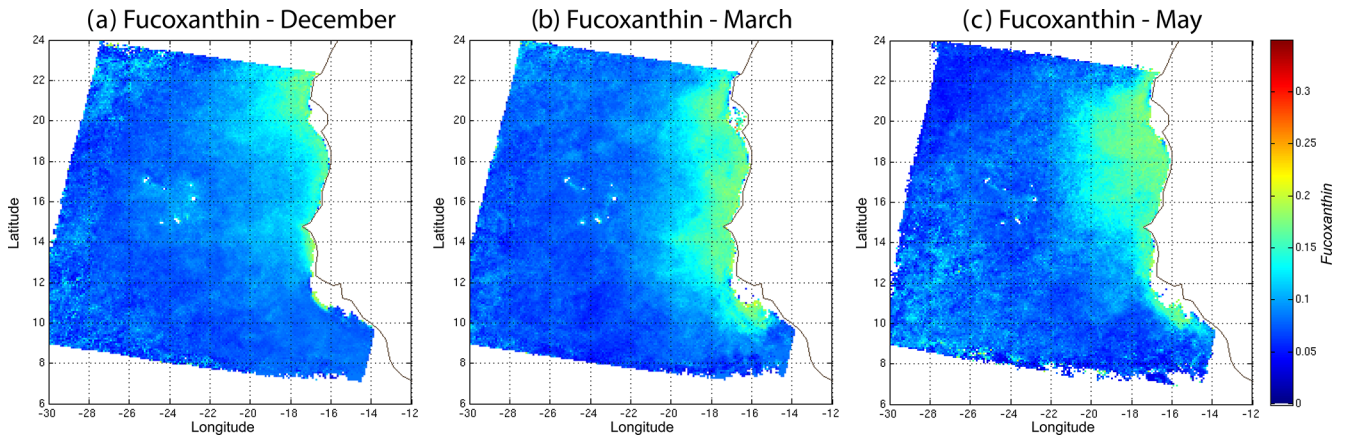
**Figure 13.** Box plot of the weights of the selected neurons during the decoding of the 28 February data. From left to right are the weights of blocks B1, B2, B3 and B4 (a) in the NSB region (long.  $[-20^\circ, -18^\circ]$ , lat.  $[12^\circ, 14^\circ]$ ) and (b) in the OFB region (long.  $[-28^\circ, -26^\circ]$ , lat.  $[12^\circ, 14^\circ]$ ).

for retrieving biological information from ocean-color satellite observations. First, studies have employed multilayer perceptrons (MLPs), which are a class of neural networks suitable to model transfer function (Thiria et al., 1993). Gross et al. (2000, 2004) retrieved the  $chl\ a$  concentration from SeaWiFS, Bricaud et al. (2006) modeled the absorption spectrum with MLP, and Raitos et al. (2008) and Palacz et al. (2013) introduced additional environmental variables in their MLPs such as SST in the retrieval of PSC and PFT from SeaWiFS, which improved the skill of the inversion. Another suitable procedure was to embed NN in a variational inversion, which is a very efficient way when a direct model exists (Jamet et al., 2005; Brajard et al., 2006a, b; Badran et al., 2008). Statistical analysis of the absorption spectra of phytoplankton and pigment concentrations was conducted by Chazottes et al. (2006, 2007) using an SOM.

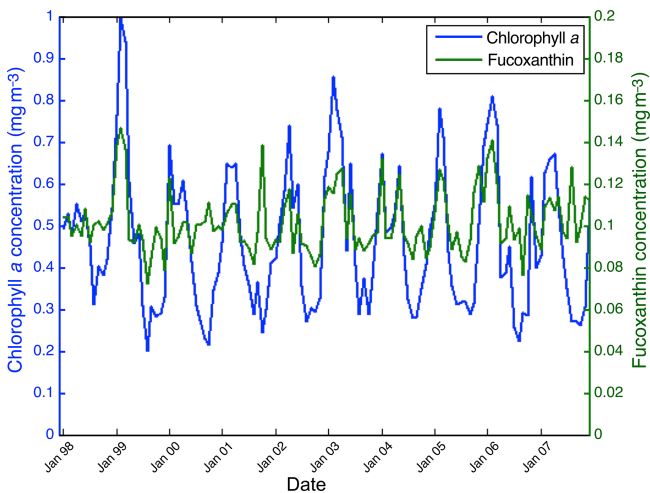
In the present study, due to the fact that the learning dataset was quite small (515 elements), we used an unsupervised neural network classification method, which is an extension of the SOM method well adapted to dealing with a small database whose elements are very inhomogeneous. We clustered available satellite ocean-color reflectance at five wavelengths and their derived products, such as chlorophyll concentration and the associated in situ pigment ratios.

The major points of this study are as follows.

1. The clustering was carried out by developing a new neural classifier, the so-called 2S-SOM, which presents several advantages with respect to the classical SOM. As in the SOM, we defined clusters that assemble vectors that are close together in terms of a specified distance. This classifier was learned from a worldwide database (DPIG) whose vectors are ocean-color parameters observed by satellite multispectral sensors and associated pigment concentrations measured in situ. In the operational phase, SeaWiFS images are decoded, allowing for the estimation of the pigment concentration ratios. The major advantage of 2S-SOM with respect to the classical SOM is to cluster variables having similar physical significance into blocks having specific weights. The weights attributed to the four blocks are computed during the learning phase and vary with the quality of the variables and with respect to their location in the ocean (near the coast or offshore). This permits us to modulate the variable influence in the cost function, which makes the clustering more informative than that provided by the SOM. The block decomposition provides useful scientific information. For offshore, the weight analysis allowed us to show that more influence is given to the reflectance ratios  $R_a(\lambda)$  and less to the  $chl\ a$  and pigment concentrations; in contrast, near the coast the weights indicate a more active use of the pigment composition and the  $chl\ a$  concentration. Therefore, the resulting 2S-SOM clustering at best takes into account the information that belongs to the specific water content.
2. The 2S-SOM decomposes the DPIG into a large number of significant ocean-color classes, allowing for the reproduction of the different possible situations encountered in the dataset we analyze. We assume that the relationship between the pigment concentration and the remotely sensed ocean-color observations is independent of the location, which is justifiable since the relationship depends on the optical properties of ocean waters through well-defined physical laws that are region-independent. This also supports the fact that we used a global database to retrieve pigments in a definite region. In contrast, the different phytoplankton species vary from one region to another, making the relationship between the pigment ratio and phytoplankton species strongly dependent on the region. This justifies the fact we focused our study on the pigment retrieval rather than on the PSC or PFT, as mentioned above. Moreover, most of the recent phytoplankton in situ identifications have been made using pigment measurements with the HPLC method (Hirata et al., 2011). It is therefore more natural to retrieve the pigment concentration, which is the quantity we measured, than the associated PSC or PFT, which are estimated from the pigment observations through complex nonlinear and region-dependent



**Figure 14.** Monthly fucoxanthin concentration averaged over 11 years (1998–2009) for December (a), March (b) and May (c).



**Figure 15.** Chl *a* (in blue) and fucoxanthin (in green) concentrations for nearshore pixels (in the NSB region).

algorithms (Uitz et al., 2006). Due to the characteristics of the DPIG, the method can retrieve pigment concentration patterns over a large range ( $0.02\text{--}2\text{ mg m}^{-3}$ ).

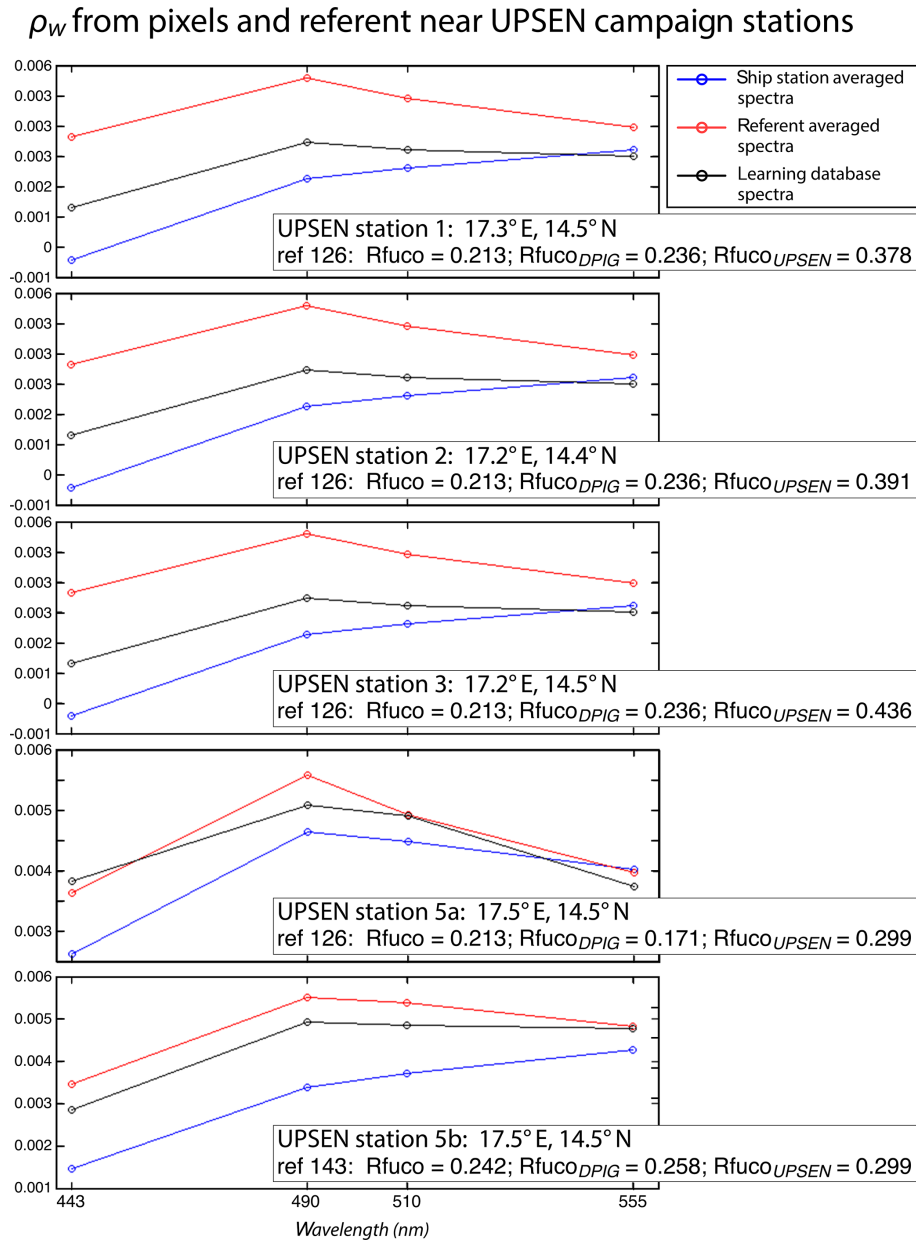
3. We were able to analyze the pigment concentration in the Senegalo–Mauritanian region by processing satellite ocean-color observations with the 2S-SOM. We found an important seasonal signal of fucoxanthin concentration with a maximum occurring in March. We found evidence of a large offshore gradient of fucoxanthin concentrations, the nearshore waters being richer than the offshore ones. We showed that the offshore region waters correspond to case 1 waters, while the nearshore waters are close to case 2 waters and are influenced by the variability of nearshore process like turbidity or the presence of dissolved matter. The UP-SEN measurements show that the pigment ratios of the Senegalo–Mauritanian region are in the range of the

DPIG database used to calibrate the method, which justifies the use of the 2S-SOM algorithm to investigate this region.

4. We used daily satellite observations to construct a monthly climatology of pigment concentrations of the Senegalo–Mauritanian upwelling region, which has been poorly surveyed by oceanic cruises. Due to the highly nonlinear character of the algorithms for determining the pigment concentrations from satellite measurements, it is mathematically more rigorous to apply these algorithms to daily satellite data and average this daily estimate for the climatology period under study than to estimate them from the satellite data climatology, as many authors have done (Uitz et al., 2010; Hirata et al., 2011). We found that fucoxanthin starts developing in December north of  $19^\circ\text{ N}$ , presents its maximum intensity in March when the upwelling intensity is maximum, extends up to the coast of Guinea ( $12^\circ\text{ N}$ ) in April and begins to decrease in May.

Another important aspect of our study concerns the validity of our results. The 2S-SOM method has been validated by focusing the retrieval accuracy on the fucoxanthin ratio by using a cross-validation procedure. These results were qualitatively confirmed by two other independent studies.

- We first applied a cross-validation procedure (see Sect. 4.1), which is a powerful technique for validating models (Kohavi, 1995; Varma and Simon, 2006). We learned 30 different 2S-SOMs using 30 different learning dataset determined at random from the DPIG dataset (each learning dataset representing 90 % of DPIG) and 30 test datasets (10 % of DPIG). By averaging the results, we found that the 2S-SOM method retrieves the fucoxanthin concentration with a good score (see the statistical parameters in Table 2), which confirms the pertinence of the method.



**Figure 16.** For ship stations 1, 2, 3, 5a and 5b, we show the averaged spectrum of the in situ spectra of the UPSEN stations in blue and the spectrum of the referent vector (in red) of the 2S-SOM neuron that has captured the closest satellite observations to the UPSEN station. Among the different spectra constituting the referent spectrum, the spectrum of the learning database (DPIG) that is the closest to the averaged satellite spectra is shown in black. In the rectangular boxes, we show the position of the UPSEN station, the number of the neuron of the 2S-SOM that has captured the satellite observation, the Rfuc<sub>0</sub> of the referent vector, the Rfuc<sub>DPIG</sub> of the closest DPIG and the in situ Rfuc<sub>UPSEN</sub>.

- We then found that our fucoxanthin climatology is in agreement with in situ observations of phytoplankton reported in Blasco et al. (1980) in March to May 1974 off the coast of Senegal during the JOINT I experiment. These authors analyzed 740 water samples collected with Niskin bottles at 136 stations extending along a line at 21°40' N (in the northern part of the

studied region) from 0 to 100 km offshore. The samples were taken at several depths (mostly at 100, 50, 30, 15, 5 m). Phytoplankton cells were counted and identified by the Utermöhl inverted microscope technique (Blasco, 1977). These authors found that diatoms reach their maximum concentration in April–May and are the most abundant group in that period, whereas the other

cells predominate in March. Similar microscope observations were reported in the ocean area south of Dakar by Dia (1985) during several ship surveys in February–March 1982–1983.

- Our method is also in agreement with the monthly 11-year climatology presented in Farikou et al. (2015), who used a modified PHYSAT method to retrieve the PFT in the Senegalo–Mauritanian region.
- The pigment concentrations provided by the 2S-SOM from the VIIRS sensor observations are in qualitative agreement with the in situ measurements done at five stations during the two UPSEN campaigns in 2012 and 2013, showing that the method is able to function in waters where the pigment concentrations are quite high (fucoxanthin ratios of the order of 0.4).

## 7 Conclusions

We developed a new neural network clustering method, the so-called 2S-SOM algorithm, to retrieve phytoplankton pigment concentration from satellite ocean-color multispectral sensors. The 2S-SOM algorithm is an SOM specifically designed to deal with a large number of heterogeneous components such as optical and chemical measurements. The major advantage of 2S-SOM with respect to the classical SOM is to cluster variables having similar significance into blocks having specific weights. The weights attributed to the blocks during the learning phase vary with the quality of the variables in the classification. This permits us to modulate the variable influence in the cost function, which makes the clustering more informative than that provided by the SOM. The block weighting provides useful information on the functioning of the classification by permitting us to identify the variables that control it. It also allows us to better understand the dynamics of the phytoplankton communities.

The 2S-SOM method is efficient and rapid as soon as the calibration is done, since it uses elementary algebraic operations only. The 2S-SOM method is like a piecewise regression that takes advantage of the unsupervised classification of the SOM. We decomposed the DFIG database into quite a large number of partitions ( $9 \times 8 = 162$ ) when comparing our study to other studies (Uitz et al., 2006). The validity of the method has been controlled through a cross-validation procedure and confirmed by three qualitative studies. Statistical parameters ( $R^2$  coefficients, RMSE and  $P$  values) of the cross-validation between the DFIG in situ pigments and the pigments given by the 2S-SOM averaged for the 30 2S-SOM realizations presented in Table 2 show the good performance of the method. It must be noted that the performance mainly depends on the size of the learning set used to calibrate the 2S-SOM. This set must include all the situations encountered in the pigment retrieval. The larger the learning set, the better the method performs. Due to its generic character and its

flexibility, the method could be used to determine a large variety of measures with satellite remote sensing observations.

In this work, the method was applied to study the seasonal variability of the fucoxanthin concentration in the Senegalo–Mauritanian upwelling region. We showed a large offshore gradient of fucoxanthin, the higher concentration being situated near the shore. We were able to construct a monthly climatology for an 11-year period (1998–2009) of the SeaWiFS observations by summing the daily pixels of the month under study in a region that was poorly surveyed by oceanic cruises. The fucoxanthin concentration, and consequently the associated diatoms, presents a well-marked seasonality (Fig. 10). Fucoxanthin starts developing in December north of  $19^\circ$  N, presents its maximum intensity in March when the upwelling intensity is maximum, extends up to the coast of Guinea ( $12^\circ$  N) in April and begins to decrease in May when it is observed north of the Cabo Verde peninsula ( $15^\circ$  N), in agreement with the observations reported by Farikou et al. (2015) and Demarcq and Faure (2000). The UPSEN campaign results confirm the validity of the study of the Senegalo–Mauritanian upwelling region done with the 2S-SOM.



Appendix A

A1 Cost function of the SOM

Let us recall the following notation:

$\mathbf{D} = \{z_1, \dots, z_i, \dots, z_K\}$  is the dataset composed of  $K$  vectors  $z_i \in \mathbb{R}^N$ , and

$\mathbf{W} = \{w_1, \dots, w_c, \dots, w_C\}$  is the set of weights  $w_c \in \mathbb{R}^N$ , where  $C = p \times q$  is the size of the SOM.

The  $w_c$  of the SOM is estimated by minimizing a cost function of the form

$$J_{\text{SOM}}^T(\chi, \mathbf{W}) = \sum_{i=1}^K \sum_{c=1}^{p \times q} K^T(\delta(c, \xi(z_i))) \|z_i - w_c\|^2, \quad (\text{A1})$$

where  $c$  indices are the neurons of the SOM,  $\xi$  is the allocation function that assigns each element  $z_i$  of  $\mathbf{D}$  to its referent vector  $w_c$ , which is of the form  $\xi(z_i) = \arg \min_c \|z_i - w_c\|$ ,  $\delta(c, \xi(z_i))$  is the discrete distance on the SOM between a neuron if index  $c$  and the neuron are allocated to observation  $z_i$ , and  $K^T$  is a kernel function parameterized by  $T$  that weights the discrete distance on the map and decreases during the minimization process.  $T$  acts as a regularization term (Kohonen, 2001; Niang et al., 2003). In the present case  $K^T$  is of the form

$$K^T(\delta) = (1/T)K(\delta/T),$$

where  $K$  is the Gaussian function of mean 0 and standard deviation 1.

The cost function (A1) takes into account the proper inertia of the partition of the dataset  $\mathbf{D}$  and ensures that its topology is preserved.

A2 Definition of the algorithm 2S-SOM

The 2S-SOM algorithm is an extension of the self-organizing maps (SOMs; Kohonen, 2001) based on the K-mean method (Ouattara, 2014). It automatically structures the variables having some common characters into conceptually meaningful and homogeneous blocks during the learning phase. The 2S-SOM takes advantage of this structuration of  $\mathbf{D}$  and the variables into  $B$  different blocks, which permits an automatic weighting of the influence of each block and consequently of each variable in the classification phase. The 2S-SOM is based on a modification of the cost function of the SOM algorithm. For a neuron of index  $c$ , we define the weights  $\alpha_{cb}$  of each block  $b(b = 1, \dots, B)$  and the weights  $\beta_{cbj}$  of the variables  $j(j = 1, \dots, P_b)$  in this block, where  $P_b$  is the number of variables in the block indexed by  $b$ . The vectors of weights are denoted

$$\alpha = \{\alpha_{cb}\}_{1 \leq c \leq C, 1 \leq b \leq B} \text{ and } \beta = \{\beta_{cbj}\}_{1 \leq c \leq C, 1 \leq b \leq B, 1 \leq j \leq P_b}.$$

The new cost function is

$$J_{\text{2S-SOM}}^T(\chi, \mathbf{W}, \alpha, \beta) = \sum_c \left( \sum_{b=1}^B \left( \sum_{z_i} \alpha_{cb} K^T(\delta(c, \xi(z_i))) d_{\beta_{cb}}(i) + J_{cb} \right) + I_c \right), \quad (\text{A2})$$

with

$$d_{\beta_{cb}}(i) = \sum_{j=1}^{P_b} \beta_{cbj} (z_{ib}^j - w_{ib}^j), \quad (\text{A3})$$

where  $c$  indices are the neurons of the 2S-SOM under the two constraints

$$\sum_{b=1}^B \alpha_{cb} = 1, \alpha_{cb} \in [0, 1] \forall c, 1 \leq c \leq C \quad (\text{A4})$$

and

$$\sum_{j=1}^{P_b} \beta_{cbj} = 1; \beta_{cbj} \in [0, 1], \forall c, 1 \leq c \leq C; \forall b, 1 \leq b \leq B. \quad (\text{A5})$$

$I_c$  and  $J_{cb}$  are used to regularize the weights  $\alpha$  and  $\beta$ . They are defined as negative entropies weighted for the blocks and for the variables of each block:

$$I_c = \mu \sum_{b=1}^B \alpha_{cb} \log(\alpha_{cb}), \quad (\text{A6})$$

and

$$J_{cb} = \eta \sum_{j=1}^{P_b} \beta_{cbj} \log(\beta_{cbj}). \quad (\text{A7})$$

The topological conservation properties of 2S-SOM are influenced by the weights  $\alpha_{cb}$  and  $\beta_{cbj}$  in the classification through the hyper-parameters  $\mu$  and  $\eta$  as well as the neighborhood parameter  $T$ .

The weights  $\alpha_{cb}$  and  $\beta_{cbj}$  respectively indicate the relative importance of blocks and variables in the neurons. Thus, the greater the weight of a block  $b$  or a variable  $j$ , the more the block or the variable contributes to the definition of the class (or neuron) in the sense that it makes it possible to reduce the variability of the observations in the cell and in its close neighborhood. For a high value of  $\eta$  and a fixed one for  $\mu$ , the  $\beta_{cbj}$  values in a block are equal to  $1/P_b$ . In this case, only the blocks are modified according to their capacity to define the neurons. In this context, the 2S-SOM then makes it possible to weight the different blocks for each neuron.

- For high values of  $\mu$ ,  $I_c$  is large. The minimization of  $J_{cb}$  forces all its coefficients to become equal. For a fixed value of  $\eta$ , the  $\alpha_{cb}$  values associated with the blocks are all equal to  $1/B$ . In this case, only the  $\beta_{cbj}$  values of the variables inside the blocks weight the neurons.
- When  $\mu$  and  $\eta$  tend to very large values, the blocks are equiprobable as are the variables. Thus, the 2S-SOM algorithm is comparable to the SOM.

### A3 How the 2S-SOM algorithm works

For fixed  $\mu$  and  $\eta$ , the learning of the 2S-SOM algorithm is as follows.

- *Step 0.* Initialization with the iteration of the algorithm SOM by setting  $\alpha$  and  $\beta$  to homogeneous values.

The optimization is carried out through an iterative process composed of three steps (1, 2 and 3) presented below.

- *Step 1.* The  $w_c$  referents and the weights  $\alpha$  and  $\beta$  are known and fixed, and the observations are assigned to the neurons by respecting the assignment function

$$c(z_i) = \chi(z_i) = \arg \min_{r \in C} \left( \sum_{r \in C} K^T(\delta(r, c)) \left( \sum_{b=1}^B \alpha_{cb} d_{\beta_{cb}}(i) \right) \right). \quad (\text{A8})$$

- *Step 2.* Updating the neuron centers (the  $w_c$  referents) according to the formula of the SOM algorithm.
- *Step 3.* The assignment function and the referents  $w_c$  being fixed,  $\alpha$  and  $\beta$  are determined according to Eqs. (A9)–(A12) by minimizing the cost function with respect to  $\alpha$  and  $\beta$  under the following constraints (Eqs. A4 and A5):

$$\alpha_{cb} = \frac{\exp\left(\frac{-\psi_{cb}}{\mu}\right)}{\sum_{b=1}^B \exp\left(\frac{-\psi_{cb}}{\mu}\right)} \quad (\text{A9})$$

with

$$\psi_{cb} = \sum_{z_i \in \mathcal{D}} K^T(\delta(\chi(z_i), c)) d_{\beta_{cb}}(i) \quad (\text{A10})$$

and

$$\beta_{cbj} = \frac{\exp\left(\frac{-\phi_{cbj}}{\eta}\right)}{\sum_{b=1}^{p_b} \exp\left(\frac{-\phi_{cbj}}{\eta}\right)} \quad (\text{A11})$$

with

$$\psi_{cbj} = \sum_{z_i \in \mathcal{D}} \alpha_{cb} K^T(\chi(z_i), c) \left( z_{ib}^i w_{cb}^j \right)^2. \quad (\text{A12})$$

This algorithm is repeated by sampling the hyper-parameters  $\mu$  and  $\eta$  until convergence.

Finally, at the convergence, the 2S-SOM provides a topological map allowing us to visualize the data and a weight system for the neurons of the map allowing us to interpret the role of the different variables, choose those that are the most significant for the classification and neutralize those that are the least significant.

*Code and data availability.* The satellite data (ocean color and SST) are available at the following website: <http://poacc.locean-ipsl.upmc.fr/> (last access: 4 March 2020, Diouf et al., 2013).

The DFIG database was kindly provided by Séverine Alvain (severine.alvain@univ-littoral.fr).

The UPSEN data are available at [alban.lazar@locean-ipsl.upmc.fr](mailto:alban.lazar@locean-ipsl.upmc.fr).

The 2S-SOM code is available on request at [carlos.mejia@locean-ipsl.upmc.fr](mailto:carlos.mejia@locean-ipsl.upmc.fr).

*Author contributions.* N'DN and MO provided the 2S-SOM code, KY processed the data and did the computations with the 2S-SOM, ST, MC and JB analyzed the results, and CM and REH did the statistical tests presented in tables and Fig. 13. ST conceived and supervised the study.

*Competing interests.* The authors declare that they have no conflict of interest.

*Acknowledgements.* The study was supported by the CNES (Centre National d'Etudes Spatiales) (project nos. CNES-TOSCA 2013-2014 and 2014-2015). The water-leaving reflectances were obtained from the SeaWiFS daily reflectances,  $\rho_{\text{obsTOAw}}(\lambda)$ , provided by NASA/GSFC/DAAC observed at the top of the atmosphere (TOA) and processed with the SOM-NV algorithm (Diouf et al., 2013) from 1998 to 2010. They are available at the following website: <http://poacc.locean-ipsl.upmc.fr/> (last access: 4 March 2020). The DFIG database was kindly provided by Séverine Alvain. We thank Alban Lazar and Eric Machu for providing in situ data measured during the UPSEN experiments as well as stimulating discussions for their interpretation. We also thank Ray Griffiths for editing an earlier version of the paper.

*Financial support.* This research has been supported by the CNES (Centre National d'Etudes Spatiales) (project nos. CNES-TOSCA 2013-2014 and 2014-2015).

*Review statement.* This paper was edited by Oliver Zielinski and reviewed by two anonymous referees.

## References

Aiken, J., Pradhan, Y., Barlow, R., Lavender, S., Poulton, A., and Hardman-Mountford, N.: Phytoplankton pigments and functional types in the Atlantic Ocean: A decadal assessment, 1995–2005, *Deep-Sea Res. Pt II*, 56, 899–917, <https://doi.org/10.1016/J.DSR2.2008.09.017>, 2009.

Alvain, S., Moulin, C., Dandonneau, Y., and Breon, F. M.: Remote sensing of phytoplankton groups in case-I waters from global SeaWiFS imagery, *Deep-Sea Res. Pt. I*, 52, 1989–2004, 2005.

Alvain, S., Loisel, H., and Dessailly, D.: Theoretical analysis of ocean color radiances anomalies and implications for phytoplankton group detection, *Opt. Express*, 20, 1070–1083, 2012.

Antoine, D., André, J. M., and Morel, A.: Oceanic primary production: Estimation at global scale from satellite (Coastal Zone Color Scanner) chlorophyll, *Global Biogeochem. Cy.*, 10, 57–69, 1996.

Badran, F., Berrada, M., Brajard, J., Crepon, M., Sorrow, C., Thiria, S., Hermant, J. P., Meyer, M., Perichon, L., and Asch, M.: Inversion of satellite ocean colour imagery and geoacoustic characterization of seabed properties: Variational data inversion using a semi-automatic adjoint approach, *J. Marine Syst.*, 69, 126–136, 2008.

Behrenfeld, M. J. and Falkowski, P. G.: Photosynthetic rates derived from satellite base chlorophyll concentration, *Limnol. Oceanogr.*, 42, 1–20, 1997.

Behrenfeld, M. J., Boss, E., Siegel, D. A., and Shea, D. M.: Carbon-based ocean productivity and phytoplankton physiology from space, *Global Biogeochem. Cy.*, 19, GB1006, <https://doi.org/10.1029/2004GB002299>, 2005.

Ben Mustapha, Z. S., Alvain, S., Jamet, C., Loisel, H., and Desailly, D.: Automatic water leaving radiance anomalies from global SeaWiFS imagery: application to the detection of phytoplankton groups in open waters, *Remote Sens. Environ.*, 146, 97–112, 2014.

Blasco, D.: Red tide in the upwelling region of Baja California, *Limnol. Oceanogr.*, 22, 255–263, 1977.

Blasco, D., Estrada, M., and Jones, B.: Relationship between the phytoplankton distribution and composition and the hydrography in the northwest African upwelling region, near Cabo Corbeiro, *Deep-Sea Res.*, 27A, 799–821, 1980.

Brajard, J., Jamet, C., Moulin, C., and Thiria, S.: Atmospheric correction and oceanic constituents retrieval with a neuro-variational method, *Neural Networks*, 19, 178–185, 2006a.

Brajard, J., Jamet, C., Moulin, C., and Thiria, S.: Neurovariational inversion of ocean color images, *Journal of Atmospheric Space Research*, 38, 2169–2175, 2006b.

Brewin, R. J. W., Sathyendranath, S., Hirata, T., Lavender, S. J., Barciela, R., and Hardman-Montford, N. J.: A three-component model of phytoplankton size class for the Atlantic Ocean, *Ecol. Model.*, 22, 1472–1483, 2010.

Bricaud, A., Mejia, C., Blondeau Patissier, D., Claustre, H., Crepon, M., and Thiria, S.: Retrieval of pigment concentrations and size structure of algal populations from absorption spectra using multilayered perceptrons, *Appl. Optics*, 46, 1251–1260, 2006.

Capet, X., Estrade, P., Machu, E., Ndoye, S., Grelet, J., Lazar, A., Marié, L., Dausse, D., and Brehmer, P.: On the Dynamics of the Southern Senegal Upwelling Center: Observed Variability from Synoptic to Superinertial Scales, *J. Phys. Oceanogr.*, 47, 155–180, 2017.

Cavazos, T.: Using Self-Organizing Maps to Investigate Extreme Climate Events: An Application to Wintertime Precipitation in the Balkans, *J. Climate*, 13, 1718–1732, 2000.

Chazotte, A., Crepon, M., Bricaud, A., Ras, J., and Thiria, S.: Statistical analysis of absorption spectra of phytoplankton and of pigment concentrations observed during three POMME cruises using a neural network clustering method, *Appl. Optics*, 46, 3790–3799, 2007.

- Chazottes, A., Bricaud, A., Crepon, M., and Thiria, S.: Statistical analysis of a data base of absorption spectra of phytoplankton and pigment concentrations using self-organizing maps, *Appl. Optics*, 45, 8102–8115, 2006.
- Ciotti, A. and Bricaud, A.: Retrievals of a size parameter for phytoplankton and spectral light absorption by colored detrital matter from water-leaving radiances at SeaWiFS channels in a continental shelf region off Brazil, *Limnol. Oceanogr.-Meth.*, 4, 237–253, 2006.
- Demarcq, H. and Faure, V.: Coastal upwelling and associated retention indices from satellite SST. Application to Octopus vulgaris recruitment, *Oceanol. Acta*, 23, 391–407, 2000.
- Dia, A.: Biomasse et biologie du phytoplancton le long de la petite côte sénégalaise et relations avec l'hydrologie, Rapport interne No. 44 du CRODT, Réf: OC000798, 1981–1982, available at: <http://www.sist.sn/gsd/collect/publi/index/assoc/HASH2127.dir/doc.pdf> (last access: 4 March 2020), 1985.
- Diouf, D., Niang, A., Brajard, J., Crepon, M., and Thiria, S.: Retrieving aerosol characteristics and sea-surface chlorophyll from satellite ocean color multi-spectral sensors using a neural-variational method, *Remote Sens. Environ.*, 130, 74–86, <https://doi.org/10.1016/j.rse.2012.11.002>, 2013.
- Farikou, O., Sawadogo, S., Niang, A., Brajard, J., Mejia, C., Crépon, M., and Thiria, S.: Multivariate analysis of the Sénégal-Mauritanian area by merging satellite remote sensing ocean color and SST observations, *Research Journal of Environmental and Earth Sciences*, 12, 756–768, 2013.
- Farikou, O., Sawadogo, S., Niang, A., Diouf, D., Brajard, J., Mejia, C., Dandonneau, Y., Gasc, G., Crepon, M., and Thiria, S.: Inferring the seasonal evolution of phytoplankton groups in the Senegalo-Mauritanian upwelling region from satellite ocean-color spectral measurements, *J. Geophys. Res.-Oceans*, 120, 6581–6601, 2015.
- Friedrich, T. and Oschlies, A.: Basin-scale pCO<sub>2</sub> maps estimated from ARGO float data: A model study, *J. Geophys. Res.*, 114, C10012, <https://doi.org/10.1029/2009JC005322>, 2009.
- Gregg, W. W., Casey, N., and McClain, C.: Recent trends in global ocean chlorophyll, *Geophys. Res. Lett.*, 32, L03606, <https://doi.org/10.1029/2004GL021808>, 2005.
- Gross, L., Thiria, S., Frouin, R., and Mitchell, B. G.: Artificial neural networks for modeling transfer function between marine reflectance and phytoplankton pigment concentration, *J. Geophys. Res.*, 105, 3483–3949, 2000.
- Gross, L., Frouin, R., Dupouy, C., Andre, J. M., and Thiria, S.: Reducing biological variability in the retrieval of chlorophyll *a* concentration from spectral marine reflectance, *Appl. Optics*, 43, 4041–4054, 2004.
- Hewitson, B. C. and Crane, R. G.: Self organizing maps: application to synoptic climatology, *Clim. Res.*, 22, 13–26, 2002.
- Hirata, T., Aiken, J., Hardman-Mountford, N., Smyth, T. J., and Barlow, R. G.: An absorption model to determine phytoplankton size classes from satellite ocean color, *Remote Sens. Environ.*, 112, 3153–3159, 2008.
- Hirata, T., Hardman-Mountford, N. J., Brewin, R. J. W., Aiken, J., Barlow, R., Suzuki, K., Isada, T., Howell, E., Hashioka, T., Noguchi-Aita, M., and Yamanaka, Y.: Synoptic relationships between surface Chlorophyll-*a* and diagnostic pigments specific to phytoplankton functional types, *Biogeosciences*, 8, 311–327, <https://doi.org/10.5194/bg-8-311-2011>, 2011.
- IOCCG: Phytoplankton Functional Types from Space, in: Reports of the International Ocean-Colour Coordinating Group, edited by: Sathyendranath, S., IOCCG, Dartmouth, Canada, IOCCG Report No. 15, 156 pp., 2014.
- Jamet, C., Thiria, S., Moullin, C., and Crepon, M.: Use of a neural inversion for retrieving Oceanic and Atmospheric constituents for Ocean Color imagery: a feasibility study, *J. Atmos. Ocean. Tech.*, 22, 460–475, <https://doi.org/10.1175/JTECH1688.1>, 2005.
- Jeffreys, S. W. and Vesk, M.: Introduction to marine phytoplankton and their pigment signatures, in: *Phytoplankton pigments in oceanography: guidelines to modern methods*, edited by: Jeffery, S. W., Mantoura, R. F. C., and Wright, S. W., UNESCO, Paris, 33–84, 1997.
- Jouini, M., Lévy, M., Crépon, M., and Thiria, S.: Reconstruction of ocean color images under clouds using a neuronal classification method, *Remote Sens. Environ.*, 131, 232–246, 2013.
- Kohavi, R.: A study of cross-validation and bootstrap for accuracy estimation and model selection, in: *Proceedings of the Fourteenth International Joint Conference on Artificial Intelligence*, San Mateo, CA, Morgan Kaufmann Publishers Inc., 2, 1137–1143, 1995.
- Kohonen, T.: *Self-organizing maps*, 3rd edn., Springer, Berlin Heidelberg New York, 2001.
- Kruizinga, S. and Murphy, A.: Use of an analogue procedure to formulate objective probabilistic temperature forecasts in the Netherlands, *Mon. Weather Rev.*, 111, 2244–2254, 1983.
- Le Quéré, C., Andrew, R. M., Friedlingstein, P., Sitch, S., Hauck, J., Pongratz, J., Pickers, P. A., Korsbakken, J. I., Peters, G. P., Canadell, J. G., Arneeth, A., Arora, V. K., Barbero, L., Bastos, A., Bopp, L., Chevallier, F., Chini, L. P., Ciais, P., Doney, S. C., Gkritzalis, T., Goll, D. S., Harris, I., Haverd, V., Hoffman, F. M., Hoppema, M., Houghton, R. A., Hurtt, G., Ilyina, T., Jain, A. K., Johannessen, T., Jones, C. D., Kato, E., Keeling, R. F., Goldewijk, K. K., Landschützer, P., Lefèvre, N., Lienert, S., Liu, Z., Lombardozzi, D., Metzl, N., Munro, D. R., Nabel, J. E. M. S., Nakaoka, S., Neill, C., Olsen, A., Ono, T., Patra, P., Peregon, A., Peters, W., Peylin, P., Pfeil, B., Pierrot, D., Poulter, B., Rehder, G., Resplandy, L., Robertson, E., Rocher, M., Rödenbeck, C., Schuster, U., Schwinger, J., Séférian, R., Skjelvan, I., Steinhoff, T., Sutton, A., Tans, P. P., Tian, H., Tilbrook, B., Tubiello, F. N., van der Laan-Luijkx, I. T., van der Werf, G. R., Viovy, N., Walker, A. P., Wiltshire, A. J., Wright, R., Zaehle, S., and Zheng, B.: *Global Carbon Budget 2018*, *Earth Syst. Sci. Data*, 10, 2141–2194, <https://doi.org/10.5194/essd-10-2141-2018>, 2018.
- Lévy, M.: Mesoscale variability of phytoplankton and of new production: Impact of the large-scale nutrient distribution, *J. Geophys. Res.*, 108, 3358, <https://doi.org/10.1029/2002JC001577>, 2003.
- Lévy, M., Iovino, D., Resplandy, L., Klein, P., Madec, G., Tréguier, A.-M., Masson, S., and Takahashi, K.: Large-scale impacts of submesoscale dynamics on phytoplankton: Local and remote effects, *Ocean Model.*, 43–44, 77–93, 2012.
- Liu, Y. and Weisberg, R. H.: Patterns of ocean current variability on the West Florida Shelf using the self-organizing map, *J. Geophys. Res.*, 110, C06003, <https://doi.org/10.1029/2004JC002786>, 2005.
- Liu, Y., Weisberg, R. H., and He, R.: Sea surface temperature patterns on the West Florida Shelf using growing hierarchical self-organizing maps, *J. Atmos. Ocean. Tech.*, 23, 325–338, 2006.

- Longhurst, A. R., Sathyendranath, S., Platt, T., and Caverhill, C.: An estimation of global primary production in the ocean from satellite radiometer data, *J. Plankton Res.*, 17, 1245–1271, 1995.
- Lorenz, E. N.: Atmospheric predictability as revealed by naturally occurring analogs, *J. Atmos. Sci.*, 26, 639–646, 1969.
- Morel, A. and Gentili, G.: Diffuse reflectance of oceanic waters. III. Implication of bidirectionality for the remote-sensing problem, *Appl. Optics*, 35, 4850–4862, 1996.
- Mouw, C. B. and Yoder, J. A.: Optical determination of phytoplankton size composition from global SeaWiFS imagery, *J. Geophys. Res.*, 115, C12018, <https://doi.org/10.1029/2010JC006337>, 2010.
- Ndoye, S., Capet, X., Estrade, P., Sow, B., Dagorne, D., Lazar, A., Gaye, A., and Brehmer, P.: SST patterns and dynamics of the southern Senegal-Gambia upwelling center, *J. Geophys. Res.-Oceans*, 119, 8315–8335, 2014.
- Niang, A., Gross, L., Thiria, S., Badran, F., and Moulin, C.: Automatic neural classification of ocean colour reflectance spectra at the top of atmosphere with introduction of expert knowledge, *Remote Sens. Environ.*, 86, 257–271, 2003.
- Niang, A., Badran, F., Moulin, C., Crépon, M., and Thiria, S.: Retrieval of aerosol type and optical thickness over the Mediterranean from SeaWiFS images using an automatic neural classification method, *Remote Sens. Environ.*, 100, 82–94, 2006.
- O'Reilly, J. E., Maritorena, S., Siegel, D. A., O'Brien, M. C., Toole, D., Mitchell, B. G., Kahru, M., Chavez, F. P., Strutton, P., Cota, G. F., Hooker, S. B., McClain, C. R., Carder, K. L., Muller-Karger, F., Harding, L., Magnuson, A., Phinney, D., Moore, G. F., Aiken, J., Arrigo, K. R., Letelier, R., and Culver, M.: Ocean color chlorophyll algorithms for SeaWiFS, OC2 and OC4: Version 4, in: *SeaWiFS postlaunch calibration and validation analyses: Part 3*, edited by: Hooker, S. B. and Firestone, E. R., NASA Goddard Space Flight Center, Greenbelt, MD, NASA Tech. Memo. 2000-206892, 11, 9–23, 2001.
- Ouattara, M.: Développement et mise en place d'une méthode de classification multi-blocs: application aux données de l'OQAI, PhD thesis, available at: <https://www.theses.fr/179489704>, last access: 4 March 2020.
- Palacz, A. P., John, M. A. St., Brewin, R. J. W., Hirata, T., and Gregg, W. W.: Distribution of phytoplankton functional types in high-nitrate, low-chlorophyll waters in a new diagnostic ecological indicator model, *Biogeosciences*, 10, 7553–7574, <https://doi.org/10.5194/bg-10-7553-2013>, 2013.
- Raitsos, D. E., Lavender, S. J., Maravelias, C. D., Haralambous, J., Richardson, A. J., and Reid, P. C.: Identifying phytoplankton functional groups from space: an ecological approach, *Limnol. Oceanogr.*, 53, 605–613, <https://doi.org/10.4319/lo.2008.53.2.0605>, 2008.
- Reusch, D. B., Alley, R. B., and Hewitson, B. C.: North Atlantic climate variability from a self-organizing map perspective, *J. Geophys. Res.*, 112, D02104, <https://doi.org/10.1029/2006JD007460>, 2007.
- Richardson, A., Risien, C., and Shillington, F.: Using self-organizing maps to identify patterns in satellite imagery, *Prog. Oceanogr.*, 59, 223–239, <https://doi.org/10.1016/J.POCEAN.2003.07.006>, 2003.
- Sathyendranath, S., Watts, L., Devred, E., Platt, T., Caverhill, C. M., and Maass, H.: Discrimination of diatom from other phytoplankton using ocean-colour data, *Mar. Ecol. Prog. Ser.*, 272, 59–68, 2004.
- Sawadogo, S., Brajard, J., Niang, A., Lathuilière, C., Crepon, M., and Thiria, S.: Analysis of the Senegalo-Mauritanian upwelling by processing satellite remote sensing observations with topological maps, in: *2009 International Joint Conference on Neural Networks (IJCNN)*, Atlanta, GA, USA, 14–19 June 2009, IEEE, 313–319, 2009.
- Sirven, J., Mignot, J., and Crépon, M.: Generation of Rossby waves off the Cape Verde Peninsula: the role of the coastline, *Ocean Sci.*, 15, 1667–1690, <https://doi.org/10.5194/os-15-1667-2019>, 2019.
- Sosik, H. M., Sathyendranath, S., Uitz, J., Bouman, H., and Nair, A.: In situ methods of measuring phytoplankton functional types, in: *Phytoplankton Functional Types from Space*, edited by: Sathyendranath, S., IOCCG, Dartmouth, NS, Canada, IOCCG report, No. 15, 21–38, 2014.
- Thiria, S., Mejia, C., Badran, F., and Crépon, M.: A neural network approach for modeling nonlinear transfer functions: application for wind retrieval from spaceborne scatterometer data, *J. Geophys. Res.*, 98, 22827–22841, 2003.
- Uitz, J., Claustre, H., Morel, A., and Hooker, S. B.: Vertical distribution of phytoplankton communities in open ocean: an assessment based on surface chlorophyll, *J. Geophys. Res.*, 111, C08005, <https://doi.org/10.1029/2005JC003207>, 2006.
- Uitz, J., Claustre, H., Gentili, B., and Stramski, D.: Phytoplankton class-specific primary production in the world's ocean: seasonal and interannual variability from satellite observations, *Global Biogeochem. Cy.*, 24, GB3016, <https://doi.org/10.1029/2009GB003680>, 2010.
- Van den Dool, H.: Searching for analogs, how long must we wait?, *Tellus A*, 46, 314–324, 1994.
- Varma, S. and Simon, R.: Bias in error estimation when using cross-validation for model selection, *BMC Bioinformatics*, 7, 91, <https://doi.org/10.1186/1471-2105-7-91>, 2006.
- Vidussi, F., Claustre, H., Manca, B. B., Luchetta, A., and Marty, J. C.: Phytoplankton pigment distribution in relation to upper thermocline circulation in the eastern Mediterranean sea during winter, *J. Geophys. Res.*, 106, 19939–19956, 2001.
- Westberry, T., Behrenfeld, M. J., Siegel, D. A., and Boss, E.: Carbon-based productivity modeling with vertically resolved photoacclimation, *Global Biogeochem. Cy.*, 22, GB2024, <https://doi.org/10.1029/2007GB003078>, 2008.
- Zorita, E. and von Storch, H.: The Analog Method as a Simple Statistical Downscaling Technique: Comparison with More Complicated Methods, *J. Climate*, 12, 2474–2489, 1999.