





35 May. The results are in agreement with previous observations and recent in situ measurements. The  
36 method is very general and can be applied in every oceanic region.

37

38

## 39 1 - INTRODUCTION

40

41 Phytoplankton is the basis of the ocean food web and consequently drives the ocean productivity. It  
42 also plays a fundamental role in climate regulation by trapping atmospheric carbon dioxide (CO<sub>2</sub>)  
43 through gas exchanges at the sea surface and consequently lowering the rate of anthropogenic  
44 increase in the atmosphere of CO<sub>2</sub> concentration by about 30% (*Behrenfield et al, 2005*). With the  
45 growing interest in climate change, one may ask how the different phytoplankton populations will  
46 respond to changes in ocean characteristics (temperature, salinity, acidity) and nutrient supply, which  
47 presents an important societal impact with respect to both climate and fisheries with a possible effect  
48 on fish grazing on phytoplankton.

49 Methods for identifying phytoplankton have greatly progressed during the last two decades.  
50 Phytoplankton was first described by microscopy. The number of samples that can be analyzed is  
51 very limited due to the necessary presence of an expert to discriminate the different taxa observed  
52 with the microscope.

53 Recently pigment analysis of seawater samples by high-performance liquid chromatography (HPLC)  
54 has been widely used to categorize broad phytoplankton size classes (PSC) (*Jeffreys et al, 1997,*  
55 *Brewin et al, 2010*) and even phytoplankton functional types (PFT; *Hirata et al, 2011*). HPLC  
56 enables identification of 25 to 50 pigments within a single analysis, which is much more easy and  
57 fast to conduct than microscopic observations. Each phytoplankton group (PSC and PFT) is  
58 associated with specific diagnostic pigments and a conversion formula can be derived to estimate the  
59 percentage of each group from the pigment measurements (*Vidussi et al, 2001; Uitz et al, 2010*).  
60 HPLC measurements are now widely used to determine phytoplankton species in situ.

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

68 This upwelling radiation, the so-called remotely sensed reflectance  $\rho_w(\lambda)$ , is determined by the



69 spectral absorption  $a$  and backscattering ( $b_b$  ( $\text{m}^{-1}$ )) coefficients of the ocean (pure water and various  
70 particulate and dissolved matters) using the simplified formulation (Morel and Gentili, 1996):

71

$$72 \quad \rho_w(\lambda) = G b_b(\lambda) / (a(\lambda) + b_b(\lambda)) \quad (1)$$

73

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

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

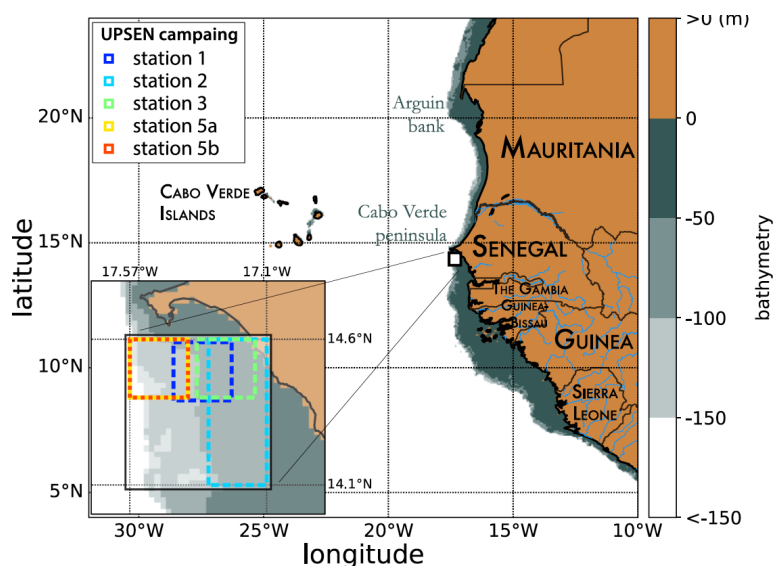
83 It has been shown that it is also possible to extract additional information such as phytoplankton size-  
84 classes (PSC) by using some relationship between chlorophyll concentration and PSC (Uitz *et al.*,  
85 2006; Ciotti and Bricaud, 2006; Hirata *et al.*, 2008; Mow and Yoder, 2010). These algorithms try to  
86 establish an algebraic relationship between the  $chl-a$  concentration and the PSC percentage. Some of  
87 them (Uitz *et al.*, 2006; Aiken *et al.*, 2009) break-down the  $chl-a$  abundance into several ranges for  
88 each of which a specific relationship is computed. Others (Brewin *et al.*, 2010; Hirata *et al.*, 2011) are  
89 based on a continuum of  $chl-a$  abundance. Studies have also been done to estimate the  
90 phytoplankton groups (PFT) by taking into account spectral information (Sathyendranath *et al.*,  
91 2004, Alvain *et al.*, 2005, 2012; Hirata *et al.*, 2011; Ben Mustapha *et al.*, 2013; Farikou *et al.*, 2015),  
92 which is of fundamental interest to the understanding of the phytoplankton behavior and to modeling  
93 its evolution.

94 Due to highly non-linear relationship linking the multispectral ocean color measurements with the  
95 pigment concentrations, we proposed a neural network clustering algorithm (2S-SOM) able to deal  
96 with multi variables related by complex relationships. The 2S-SOM algorithm is well adapted to this  
97 complex task by weighting the different inputs. The clustering algorithm was calibrated on a  
98 restricted database composed of remote sensed observations co-located with measurements taken in  
99 the global ocean.

100 In the present paper, we propose the retrieval of the major pigment concentrations from satellite  
101 ocean color multi-spectral sensors in the Senegalo-Mauritanian upwelling, which is an oceanic



102 region off the coast of West Africa where a strong seasonal upwelling occurs (Figure 1).  
103



104  
105 Figure 1 : Mauritania and Senegal coastal topography. The land is in red and the ocean depth is  
106 represented in meters by the color scale on the right side of the figure. The UPSEN stations are  
107 shown in the bottom left cartoon of the figure.  
108

109  
110  
111 The Senegalo-Mauritanian upwelling is a one of the most productive eastern boundary upwelling  
112 system (EBUS) with strong economic impacts on fisheries in Senegal and Mauritania. Since this  
113 region has been poorly surveyed in situ, we have chosen to extract pertinent biological information  
114 from ocean-color satellite measurements. This region has been intensively studied by analysis of  
115 SeaWiFS ocean-color data and AVHRR sea-surface temperature as reported in *Demarcq and Faure*  
116 (2002), and more recently by *Sawadogo et al.* (2009); *Farikou et al.* 2013, 2015; *Ndoye et al.* 2014;  
117 *Capet et al.* 2017.

118 The paper is articulated as follows: in section 2, we present the data we used (in situ and remote  
119 sensing observations). The mathematical aspect of the clustering method (2S-SOM) is detailed in  
120 section 3. In section 4 we present the methodological results. The spatio-temporal variability of the  
121 fucoxanthin and chl-a concentration in the Senegalo-Mauritanian upwelling region are presented in  
122 section 5 as well as the results of the oceanic UPSEN campaigns. In section 6 we discuss the results  
123 and the method. A conclusion is presented in section 7.

124



125

126

## 127 2- MATERIALS

128

129 In this study we used three distinct datasets: the first was used to calibrate the method, the second to  
130 conduct a climatological analysis of the Senegalo-Mauritanian upwelling region and the third was  
131 obtained during the oceanographic UPSEN campaign. These datasets are composed of satellite  
132 remote sensing observations (datasets 1, 2, 3) and in-situ measurements (datasets 1, 3).

133

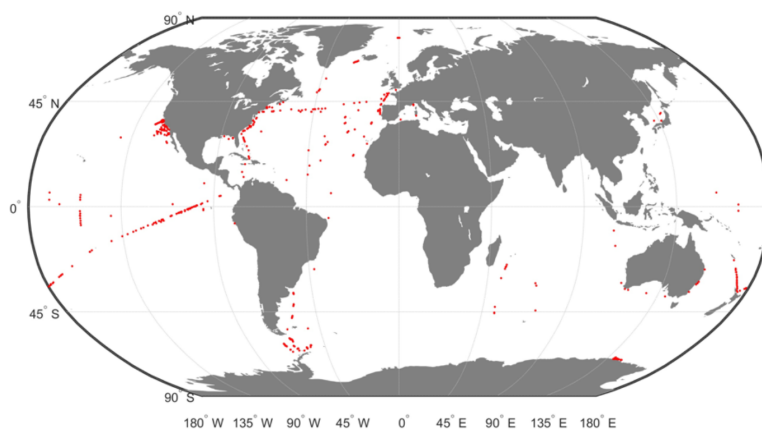
### 134 2.1 The calibration data base (DPIG)

135 The calibration database (DPIG) comprises in situ pigment measurements co-located with satellite  
136 ocean-color observations done by the SeaWiFS (sea-viewing, wide-field-of-view sensor).

137 This learning dataset (DPIG) is composed of 515 matched satellite observations and in situ  
138 measurements made in the global ocean (mainly in the North Atlantic and the equatorial ocean; *Ben*  
139 *Mustapha et al.*, 2014). The match-up criteria were quite severe: we used satellite pixel situated at a  
140 distance less than 20km of the in situ measurement in a time window of +/- 12h. The geographic  
141 distribution of the 515 coincident in situ and satellite measurements is shown in Fig. 2.

142

143



144

145

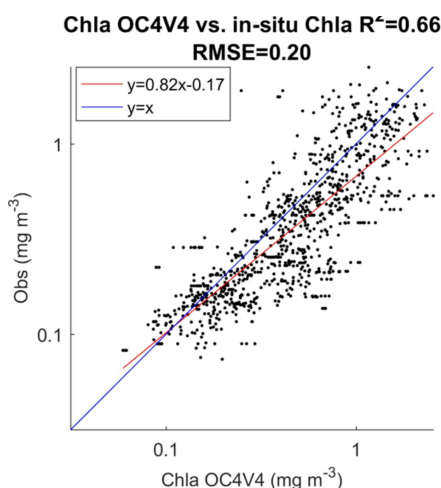
146 Figure 2 : *Geographic positions of the 515 in situ and satellite collocated measurements of*  
147 *the DPIG database.*

148

149



150 In Figure 3 we present the  $R^2$  coefficient between the in situ *chl-a* and the SeaWiFS *chl-a*  
 151 computed by using the OC4V4 algorithm for the DPIG collocated observations. We remark that the  
 152  
 153



154  
 155  
 156  
 157  
 158  
 159

Figure 3: Dispersion diagram of DPIG *chl-a* computed from the SeaWiFS observations using the OC4V4 algorithm versus in situ *chl-a*.

160 two measurements are in good agreement at global scale. Each data of DPIG is a vector having 17  
 161 components (five ocean reflectance ( $\rho_w(\lambda)$ ) at five wave lengths, five  $Ra(\lambda)$ , SeaWiFS *chl-a*, five in  
 162 situ pigment ratios and an in situ *chl-a* concentration). The in situ *chl-a* concentration ranges  
 163 between 0.007 and 3.  $\text{mg m}^{-3}$  (see Table 1).  
 164

	RDIVINY A	RPERID	RFUCO	R19HF	RZEAX	CHLORO IN SITU
MEAN	0.1414	0.0272	0.1248	0.1859	0.1696	0.5292
STD	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

165  
 166  
 167  
 168  
 169

Table 1 : Pigments of the DPIG and their statistical characteristics



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

$$171 \quad Ra(\lambda) = \rho_w(\lambda) / \rho_{wref}(\lambda, chl-a) \quad (2)$$

172 where the parameter  $\rho_{wref}(\lambda, chl-a)$  is an average reflectance depending on the *chl-a* concentration  
173 only which was computed according to the procedure reported in *Farikou et al*, 2015.  $Ra(\lambda)$  is a non-  
174 dimensional parameter which is independent of the *chl-a* abundance and sensitive the secondary  
175 pigments only (*Alvain et al* , 2012).

176

177 The DPIG database thus provides information on the existing links between the pigment composition  
178 and the SeaWiFS measurements. The pigment composition are defined by the pigment ratios which  
179 are non-dimensional variables of the form in the present study:

$$180 \quad \text{Pigment Ratio} = DP / Tchl-a \quad (3)$$

181 Defined as the ratio of the diagnostic pigment (DP) related to the total *chl-a* ( $Tchl-a = chl-a$   
182  $+Dyvinilchl-a$ ) (*Alvain et al.*, 2005).

183 The pigments of the DPIG and their statistical characteristics are given in Table 1.

184

## 185 **2.2 The Senegalo-Mauritanian upwelling satellite data (DSAT)**

186 Since the Senegalo-Mauritanian region has been poorly surveyed in situ, we have chosen to extract  
187 pertinent biological information from ocean-color satellite measurements. This region has been  
188 intensively studied by the analysis of SeaWiFS ocean-color data and AVHRR sea-surface  
189 temperature as reported in *Demarcq and Faure* (2002), *Sawadogo et al.* (2009), *Farikou et al.* (2013,  
190 2015) and more recently by *Ndoye et al*, 2014; *Capet et al*, 2017 with in situ observations. The  
191 satellite dataset we processed to retrieve the pigment concentration consist of five  $\rho_w(\lambda)$  and five  $Ra(\lambda)$   
192 at five wavelengths (412 nm, 443nm, 490nm, 510nm and 555nm), and the SeaWiFS *chl-a*  
193 concentration observed in the Senegalo-Mauritanian upwelling region (8°N-24°N, 14°W-20°W;  
194 Figure 3) during 11 years (1998-2009) by SeaWiFS. This data set is here below denoted *DSAT*.

195 The satellite observations ( $\rho_w(\lambda)$  and *chl-a* concentration) were provided by NASA with a resolution  
196 of nine kilometers. Due to the presence of Saharan dusts in this region, very few estimations of  
197 satellite  $\rho_w(\lambda)$  and in situ *chl-a* were available, and some satellite estimations of *chl-a* could present  
198 strong over-estimations (*Gregg et al*, 2004). For this reason, we reprocessed the  $\rho_w(\lambda)$  and *chl-a* data  
199 with an atmospheric correction algorithm developed specifically for Saharan dust (*Diouf et al*, 2013,  
200 <http://poac.locean-ipsl.upmc.fr>) in order to improve the satellite observations.

201

202

203



204

205

### 206 **2.3 The UPSEN database**

207 Recently, some HPLC measurements were made in the Senegalo-Mauritanian region during two  
208 oceanographic cruises (UPSEN campaigns) of the oceanographic ship “Le Suroit” from 7 to 17  
209 March 2012 and from 5 to 26 February 2013 as reported in *Ndoye et al*, 2014; *Capet et al*, 2017. The  
210 goal was to study the dynamics and the biological variability of the Senegalo-Mauritanian upwelling.  
211 During these campaigns, in-situ HPLC measurements were carried out. We expected to be able to  
212 co-locate them with the ocean-color VIIRS (visible infra-red imaging radiometer suite) sensor  
213 observations whose wavelengths are close to those of the SeaWiFS. Unfortunately, we were only  
214 able to process satellite observations made on 21 February 2013 due to the presence of clouds and  
215 aerosols the other days. We processed the satellite observations provided by the VIIRS sensor at four  
216 wavelengths (443, 490, 510, 555 nm) for pixels in the vicinity of the ship stations (within a distance  
217 of 20km) and observed in a time window of +/- 12h, and for which the satellite *chl-a* was less than  
218 3mg/m<sup>3</sup>, which is the limit of validity of our method imposed by the range of *chl-a* observed in  
219 DGIP (mean of 0,52 mg/m<sup>3</sup>). Only five stations off Cabo Verde peninsula fitted these requirements  
220 (see Figure 1 for their positions).

221

222

## 223 **3 - THE PROPOSED METHOD (2S-SOM)**

224

225 Classification methods were applied for retrieving geophysical parameters from large databases in  
226 several studies including weather forecasting (*Lorenz*, 1969; *Kruizinga and Murphy*, 1983), short-  
227 term climate prediction (*Van den Dool*, 1994), downscaling (*Zorita and von Storch*, 1999),  
228 reconstruction of oceanic pCO<sub>2</sub> (*Friedrichs and Oschlies.*, 2009), and of *chl-a* concentration under  
229 clouds (*Jouini et al*, 2013). In the present study we used a new neural network algorithm, which is an  
230 extension of the SOM algorithms (*Kohonen*, 2001)

### 231 **3-1 The SOM clustering**

232 SOM algorithms constitute powerful nonlinear unsupervised classification methods. They are  
233 unsupervised neural classifiers, which have been commonly used to solve environmental problems  
234 (*Cavazos*, 1999; *Hewitson et al*, 2002; *Richardson et al*, 2003; *Liu et al*, 2005, 2006; *Niang et al*,  
235 2006; *Reusch et al*, 2007). SOM aims at clustering vectors of a multidimensional database (**D**) into  
236 classes represented by a fixed network of neurons (the SOM map). The self-organizing maps are  
237 defined as an undirected graph, usually a rectangular grid of dimension  $p \times q$ . This graph structure is





238 used to define a discrete distance (denoted by  $\delta$ ) between the neurons of the map which presents the  
239 shortest path between two neurons. Moreover, SOM enables the partition of  $\mathbf{D}$  in which each cluster  
240 is associated with a neuron of the map and is represented by a prototype that is a synthetic  
241 multidimensional vector (the referent vector  $\mathbf{w}$ ). Each vector  $\mathbf{z}_i$  of  $\mathbf{D}$  will be assigned to the neuron  
242 whose referent  $\mathbf{w}$  is the closest, in the sense of the Euclidean Norm (EN), and will be called the  
243 projection of the vector  $\mathbf{z}_i$  on the map. A fundamental property of a SOM is the topological ordering  
244 provided at the end of the clustering phase: two close neurons on the map represent data that are  
245 close in the data space. The estimation of the referent vectors  $\mathbf{w}$  of a SOM and the topological order  
246 is achieved through a minimization process in which the referent vectors  $\mathbf{w}$  are estimated from a  
247 learning data set (The DFIG data base in the present case). The cost function is of the form:

$$248 \quad J_{SOM}^T(\chi, W) = \sum_{\mathbf{z}_i \in \mathbf{D}} \sum_{c \in SOM} K^T(\delta(c, \chi(\mathbf{z}_i))) \|\mathbf{z}_i - \mathbf{w}_c\|^2 \quad (4)$$

249 where  $c \in SOM$  indices the neurons of the SOM map,  $\chi$  is the allocation function that assigns each  
250 element  $\mathbf{z}_i$  of DFIG to its referent vector  $\mathbf{w}_{\chi(\mathbf{z}_i)}$  and  $\delta(c, \chi(\mathbf{z}_i))$  is the discrete distance on the SOM  
251 between a neuron  $c$  and the neuron allocated to observation  $\mathbf{z}_i$  and  $K^T$  a kernel function parameterized  
252 by  $T$  (where  $T$  stands for “temperature” in the scientific literature dedicated to SOM) that weights the  
253 discrete distance on the map and decreases during the minimization process.  $T$  acts as a  
254 regularization term.

255 This cost function takes into account the proper inertia of the partition of the data set  $\mathbf{D}$  and ensures  
256 that its topology is preserved.

257 SOMs have frequently been used in the context of completing missing data (Jouini *et al*, 2013), so  
258 the projected vectors  $\mathbf{z}_i$  may have missing components. Under these conditions, the distance between  
259 a vector  $\mathbf{z}_i \in \mathbf{D}$  and the referent vectors  $\mathbf{w}$  of the map is the Euclidean distance that considers only the  
260 existing components (the Truncated Distance or *TD* hereinafter).

261

### 262 **3-2 The 2S-SOM Classifier**

263 In the present case, we used 2S-SOM, which is a modified version of SOM, very useful in the case of  
264 a large number of variables. It automatically structures the variables having some common characters  
265 into conceptually meaningful and homogeneous blocks. The 2S-SOM takes advantage of this  
266 structuration of  $\mathbf{D}$  and variables into different blocks, which permits an automatic weighting of the  
267 influence of each block and consequently of each variable. Due to its capacity to weight the different  
268 parameters, 2S-SOM is able to deal with a large quantity of parameters, choosing those that are the  
269 most significant for the classification and neutralizing those which are less significant. This is done  
270 by using a more complex cost function that introduces a set of new parameters estimated during the



271 learning phase along with the referent vectors. For a neuron  $c$ , we define the weights  $\alpha_{cb}$  of each  
 272 block  $b$  and the weights  $\beta_{cbj}$  of the variables  $j$  in this block  $b$ . The new cost function is

$$273 \quad J_{2S-SOM}^T(\chi, W, \alpha, \beta) = \sum_c \left( \sum_{b=1}^B \left( \sum_{zi \in D} \alpha_{cb} K^T(\delta(c, \chi(z_i))) \right) d_{\beta_{cb}}(i) + J_{cb} \right) + I_c$$

274 (5)

$$275 \quad \text{with } d_{\beta_{cb}}(i) = \sum_{j=1}^{P_b} \beta_{cbj} \|z_{ib}^j - w_{ib}^j\|^2 \quad (6)$$

276 where  $c \in 2S-SOM$  indices the neurons of the 2S-SOM map,  $P_b$  is the number of variables in the  
 277 block  $b$ , and  $B$  is the number of blocks, under the constraints:

$$278 \quad \sum_{b=1}^B \alpha_{cb} = 1; \alpha_{cb} \in [0,1] \forall c \in [2S - SOM] \quad (7)$$

279 and

$$280 \quad \sum_{j=1}^{P_b} \beta_{cbj} = 1; \beta_{cbj} \in [0,1] \forall c \in [2S - SOM]; \forall b \quad (8)$$

281

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

$$284 \quad I_c = \mu \sum_{b=1}^B \alpha_{cb} \log(\alpha_{cb})$$

$$285 \quad \text{and } J_{cb} = \eta \sum_{j=1}^{P_b} \beta_{cbj} \log(\beta_{cbj})$$

286

287 The penalty coefficients  $\mu$  and  $\eta$  are two hyper-parameters that are determined at the end of the  
 288 learning phase depending on the problem under study. The 2S-SOM allows, during the learning  
 289 phase, an automatic weighting of the influence of each block and of each variable. After the learning  
 290 phase, the influence of each group and of each variable is different for each neuron, which permits  
 291 determination of a better classification that identifies the relevant variables for the different PSCs.

292 The learning procedure comprises three distinct phases, the third being iterated in order to choose the  
 293 two hyper-parameters:

- 294 • *First* : a standard SOM map ( $SOM_{init}$ ) is learned using the  $J_{SOM}^T(\chi, W)$  cost function
- 295 • *Second* : the different values of  $(\mu, \eta)$  are sampled
- 296 • *Third* : for each pair of  $(\mu, \eta)$
- 297 •  $SOM_{init}$  is used, as the initial condition of a new learning phase.
- 298 • All the parameters are estimated at the same time using  $J_{2S-SOM}^T(\chi, W, \alpha, \beta)$ .
- 299 • A  $2S-SOM(\mu, \eta)$  is determined

300 The best  $SOM(\mu, \eta)$  is chosen according to the problem under consideration.

301

302 The 2S-SOM algorithm is available on : [https://github.com/carmman/2S-SOM\\_versionCM](https://github.com/carmman/2S-SOM_versionCM)



303

### 304 **3.3 The calibration phase**

305 The vectors of DPIG defined in section 2 can be decomposed in four blocks. The essence of the  
306 decomposition of the components of the DPIG database vectors in blocks is that each of the 17  
307 components of the DPIG vectors gathered information with a different physical influence in the  
308 classification phase. We therefore consider different blocks (four in the present case) composed of  
309 variables having similar physical significance. The composition of each block is done as follows:

310 **First Block** (B1) comprises the five pigment in-situ concentration ratios (divinyl chlorophyll-a,  
311 peridinin, fucoxanthin, 19'hexanoyloxyfucoxanthin, zeaxanthin concentration ratios). The pigment  
312 ratios are defined in Eq. 3.

313 **Second Block** (B2) comprises the water-leaving reflectance  $\rho_w(\lambda)$  at the five SeaWiFS  
314 wavelengths

315 **Third Block** (B3) comprises the five  $Ra(\lambda)$ ,

316 **Fourth Block** (B4) comprises two variables: The in situ and the SeaWiFS *chl-a* concentrations.

317

318 At the end of the calibration phase, each element  $z_i$  of the dataset DPIG is associated with a referent  
319  $w_k$ , whose components are partitioned into four blocks, which is the closest  $w_k \in W$  in term of the  
320 weighted assignment function used in 2S-SOM. In the present study, the 2S-SOM map is represented  
321 by a two-dimensional (9x18=162) grid that represents the partition of the DPIG dataset into different  
322 classes. Each class provided by the 2S-SOM is associated with a so-called referent vector  $w_k$  with  $k$   
323  $\in \{1, \dots, 162\}$ . The size of the map has been heuristically determined.

324

### 325 **3.4 The Pigment retrieval**

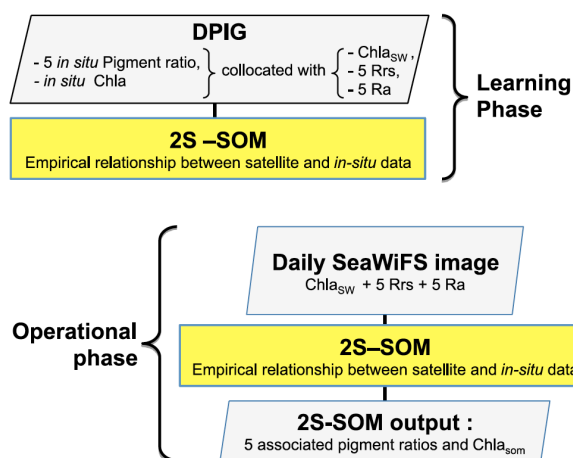
326 In the second, phase which is an operating phase, we estimate the pigment concentration ratios of a  
327 pixel  $P_j$  from its satellite ocean-color sensor observations only. The 11 ocean color satellite  
328 observations (5  $\rho_w(\lambda)$ , 5  $Ra(\lambda)$ , and *chl-a*) of pixel  $P_j$  are projected onto the 2S-SOM using the  
329 Truncated Euclidian Distance (section 3.1). We select the neuron  $k$  associated with a referent vector  
330 whose the 11 ocean-color parameters are the closest to those observed by the satellite sensor. The  
331 pigment ratios of  $P_j$  are those associated with the neuron  $k$ . At the end of the assignment phase, each  
332 pixel  $k$  of a satellite image is associated with a referent vector  $w_k$ , which has 6 pigment concentration  
333 ratios among its 17 components. The flowcharts of the method (2S-SOM learning and pigment  
334 retrieval) are presented in Figure 4.

335



336

337



338

339

340 Figure 4 : Flowchart of the method: top panel - Learning phase; bottom panel – operational phase  
 341 and pigment retrieval.

342

#### 343 4 - METHODOLOGICAL RESULTS

344

##### 345 4-1 Statistical validation of the method

346 The validation of the method was focused on the retrieval of the fucoxanthin ratio, which is a  
 347 characteristic of diatoms, but the same procedure could be applied to any pigment. The hyper-  
 348 parameters  $\mu$  and  $\eta$  were optimized in order to retrieve that ratio. Due to the small amount of data in  
 349 the DPIG, we estimated the accuracy of the fucoxanthin retrieval by a cross-validation procedure,  
 350 which is a powerful procedure in statistics. The principle is the following: we iterated the procedure  
 351 for one realization of the learning phase 30 times. The procedure was the following:

352  $i=1 \dots 30$

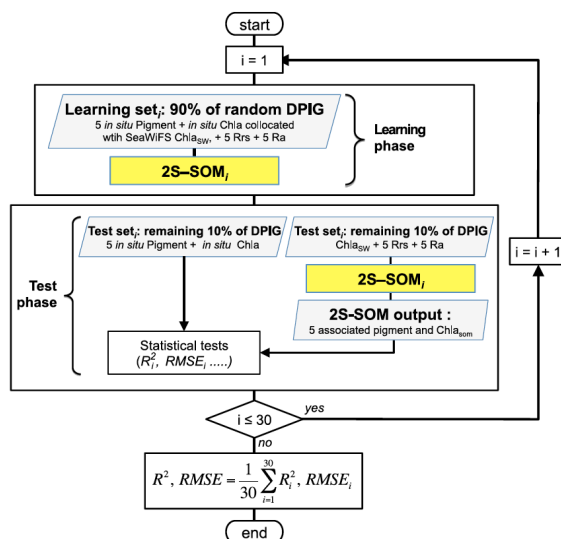
- 353 1. determination at random of a learning dataset  $L_i$  (90% of DPIG) and a test dataset  $T_i$  (10% of  
 354 DPIG)
- 355 2. training of a 2S-SOM map  $M_i$  using  $L_i$  (see section 3.2 and 3.3).
- 356 3. Validation using  $T_i$  according to the procedure described in section 3.4
- 357 4 Estimation of the  $RMSE_i$  and  $R^2_i$  on  $T_i$  between the estimated and observed fucoxanthin ratios
- 358 end

359 Computation of the mean RMSE and  $R^2$  ( $R^2, RMSE = \frac{1}{30} \sum_{i=1}^{30} R^2_i, RMSE_i$ )



360  
 361  
 362

The flowchart of the cross-validation procedure is presented in Figure 5.



363  
 364  
 365  
 366  
 367  
 368  
 369  
 370  
 371  
 372

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

Statistical parameters ( $R^2$  coefficients, RMSE and P-values) of the cross validation between the DPIG 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.

	$R^2$	RMSE [MG M <sup>-3</sup> ]	PVAL
CHLA SOM	0.84	0.22	0.001
DVCHLA	0.60	0.02	0.001
FUCO	0.87	0.02	0.001
PERID	0.81	0.01	0.001

373  
 374  
 375  
 376  
 377  
 378  
 379

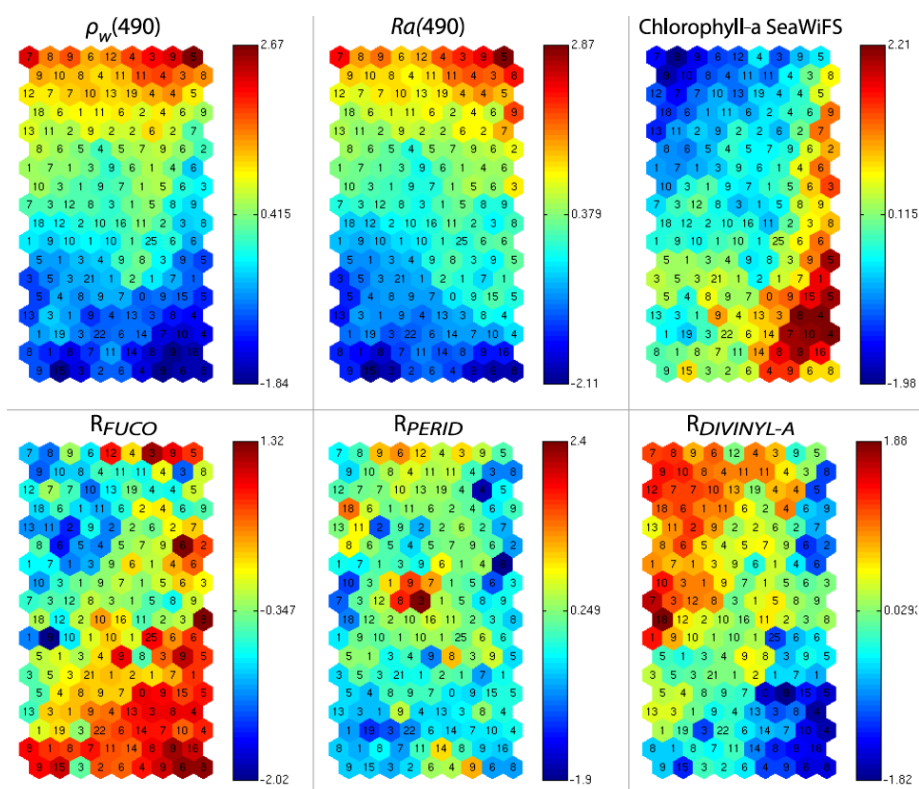
Table 2 : Statistical parameters ( $R^2$  coefficients, RMSE and P-values) of the cross validation between the DPIG in situ pigments and the pigments given by the 2S-SOM averaged for the 30 2S-SOM realizations



380

381 **4-2 Analysis of the topology of the 2S-SOM**

382 As explained in sections 3-2 and 3-3, the referent vector components ( $w_k \in R^{17}$ ), which are estimated  
 383 during the learning phase, are partitioned in four blocks B1, B2, B3 and B4. The hyper parameters  $\mu$   
 384 and  $\eta$  are tuned in order to favor the accuracy of the retrieval of the fucoxanthin ratio. We recall that  
 385 all the pigment ratios are estimated during the calibration phase, but in the present paper attention  
 386 was focused on the fucoxanthin ratio when selecting the two hyper parameters. In Figure 6, we  
 387  
 388



389

390

391 Figure 6 : 2S-SOM Map. From left to right and top to bottom, values of the referent vectors for  $\rho_w(490)$ ,  
 392  $Ra(490)$ , SeaWiFS chl-a, fucoxanthin, divinyl, peridinin. The number in each neuron indicates the amount  
 393 of DPIG data captured at the end of the learning phase, the values indicated by the color bars are centered-  
 394 reduced and non-dimensional values.

395



396 present six of the referent vector components of the 2S-SOM map. These components are  $\rho_w(490)$ ,  
397  $Ra(490)$ , SeaWiFS *chl-a*, and the ratios of fucoxanthin, which is a specific diatom pigment. They  
398 exhibit a coherent topological order, the components having close values being close together on the  
399 topological map. The remaining eleven components (not shown) exhibit the same coherent  
400 topological order. One can observe a very good topological order for the fucoxanthin ratio that was  
401 favored by the determination of the hyperparameters  $\mu$  and  $\eta$ . Moreover, the region in the 2S-SOM  
402 that characterizes the diatoms (high fucoxanthin ratio and high *chl-a*) seems to be well evidenced.  
403 Another important remark is that the value of each component presents a large range of variation of  
404 the same order as the range of variation found in the DPIG variables. It means that the SOM map has  
405 captured most of the variability of the dataset.

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

410 Besides, for each neuron, the 2S-SOM provides a weight for each block ( $\alpha_{cb}$ ) and each variable  
411 ( $\beta_{cbj}$ ). For a given neuron the weights ( $\alpha_{cb}$ ) of the blocks are normalized, their sum being 1. A value  
412 of 1 for one block (and therefore a value of 0 for the other blocks) indicates that the data in the  
413 neuron are gathered with respect to that block only because there is too much noise in the variables  
414 in the other blocks. By examining the weights on the map one can see which block most influences  
415 the link between the satellite measurements and the pigment ratios.

416 In Figure 7, we present the  $\alpha_{cb}$  values estimated during the learning phase of the 4 blocks (B1, B2,  
417 B3, B4). For some area on the map, only the blocks related to the reflectance and the reflectance  
418 ratio are used for the definition of the neuron, while the weights for the two other blocks (pigments  
419 and *chl-a*) are null, indicating that for these neurons, in situ observations and SeaWiFS *chl-a* are  
420 more noisy. For these cases, the clustering assembled the data that are similar with respect to the  
421 reflectance only.

422

423

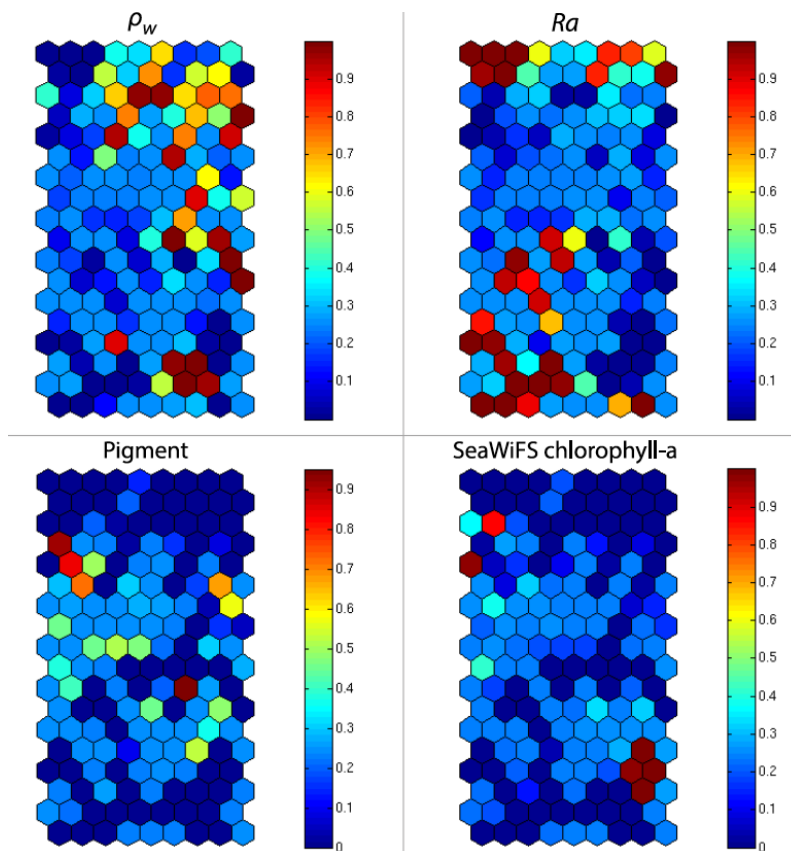
## 424 5 - GEOPHYSICAL RESULTS

425

426 In the present study, we apply the 2S-SOM (section 3), which explicitly makes a weighted use of the  
427 data according to their specificity (ocean-color signals or in situ observations) to retrieve the



428 fucoxanthin concentration from remote sensed data in the Senegalo-Mauritanian upwelling region  
429  
430



431  
432

433 Figure 7: 2S-SOM map. Weights of the four block parameters ( $\alpha_{cb}$ ) determined at the end of the learning  
434 phase; from left to right and top to bottom:  $\rho_w$ ,  $Ra$ , Pigment, SeaWifs chl-a. The color bars show the % of  
435 the weight estimated by 2S-SOM, a value of 1 or 0 indicating that the data in the neuron are assembled  
436 with respect to that block only.

437

438 where in situ measurements are lacking. According to the good results of the cross validation method  
439 as shown in section 4.1, we expect that the 2S-SOM will provide pertinent results in a region which  
440 has been poorly surveyed.

441

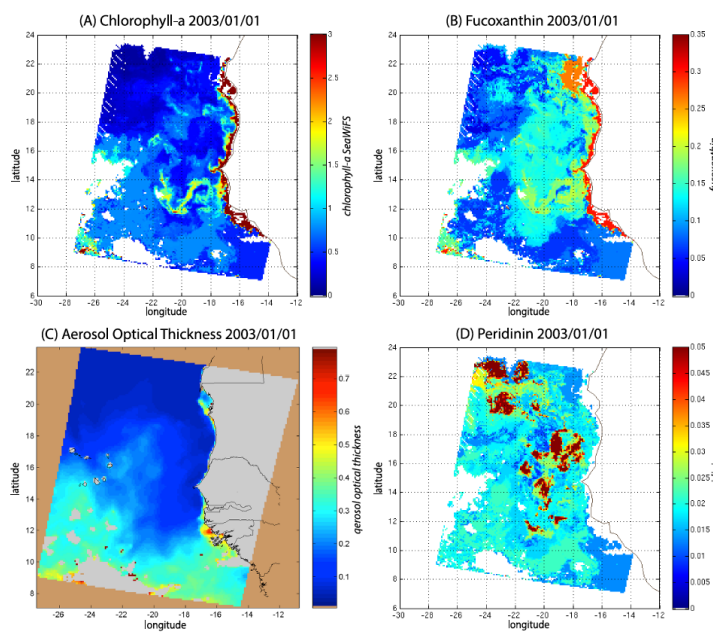
442





443 **5-1 The pigment estimation from SeaWiFS observations in the S enegal-Mauritanian upwelling**  
444 **region**

445 We decoded the DSAT database (section 2-3) using the 2S-SOM for 11 years (1998-2009) of  
446 SeaWiFS data observed in the Senegalo-Mauritanian upwelling region (8 N-24 N, 14 W-20 W).  
447 This study was done according to the retrieval phase described in section 3.4. For each day, we  
448 projected the 11 SeaWiFS observations (5  $\rho_w(\lambda)$ , 5  $Ra(\lambda)$  and  $chl-a$ ) of each pixel  $P_j$  on the 2S-SOM.  
449 At the end of the assignment phase, each pixel of a satellite image was associated with 6 pigment  
450 concentration ratios. The underlying assumption is that the link between the remote sensing  
451 information and the pigment ratios of a pixel is this provided by the selected referent  $w_k$ . Thanks to  
452 the topological order provided by the 2S-SOM, we expect that the best neurons chosen during the  
453 retrieval would give accurate concentration ratios. In Figures 8, 10 and 11 we present the  
454 fucoxanthin concentration ratio restitution for three different days and the associated SeaWiFS  
455 Chlorophyll images (1 and 6 January and 28 February 2003). Due to the limited size of the DPIG, the  
456



457

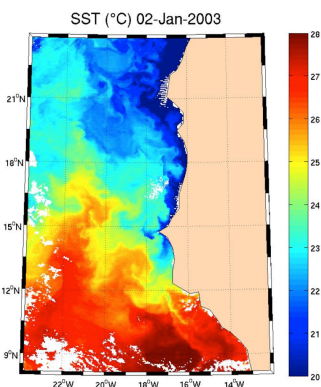
458

459 Figure 8 : *A) chl-a concentration, (B) fucoxanthin ratio, (C) aerosol optical thickness, (D) peridinin*  
460 *for 1 January 2003. Panels (B) and (D) show that a second-order information was retrieved, which*  
461 *is correlated with the chl-a concentration (A) but not equivalent. The aerosol optical thickness (C)*  
462 *does not seem to contaminate the estimated parameters (fucoxanthin and peridinin ratios).*

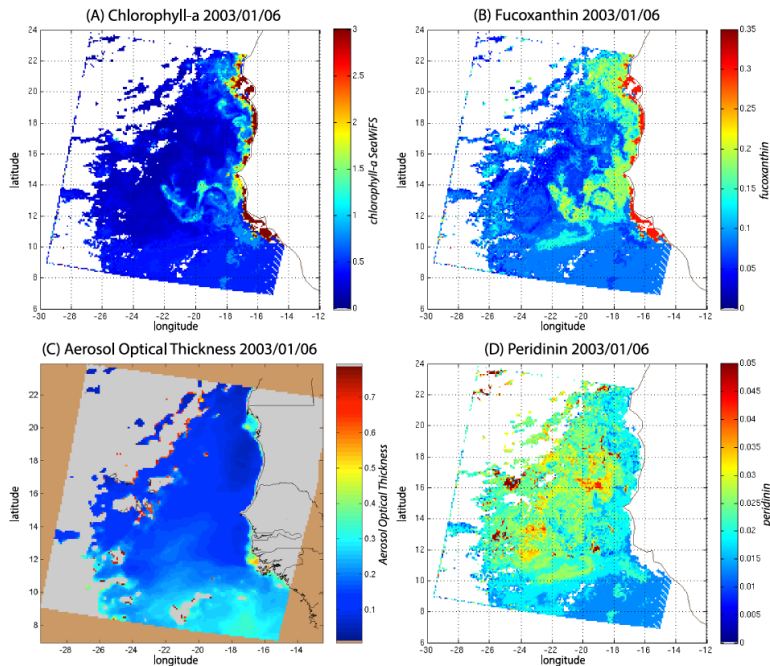
463



464 range of the ratio learned for the fucoxanthin is between 0.3% and 20% with a mean of 10% and the  
465 *chl-a* content is between  $0.5 \text{ mg m}^{-3}$  and  $3 \text{ mg m}^{-3}$ . The statistical estimator we used cannot  
466



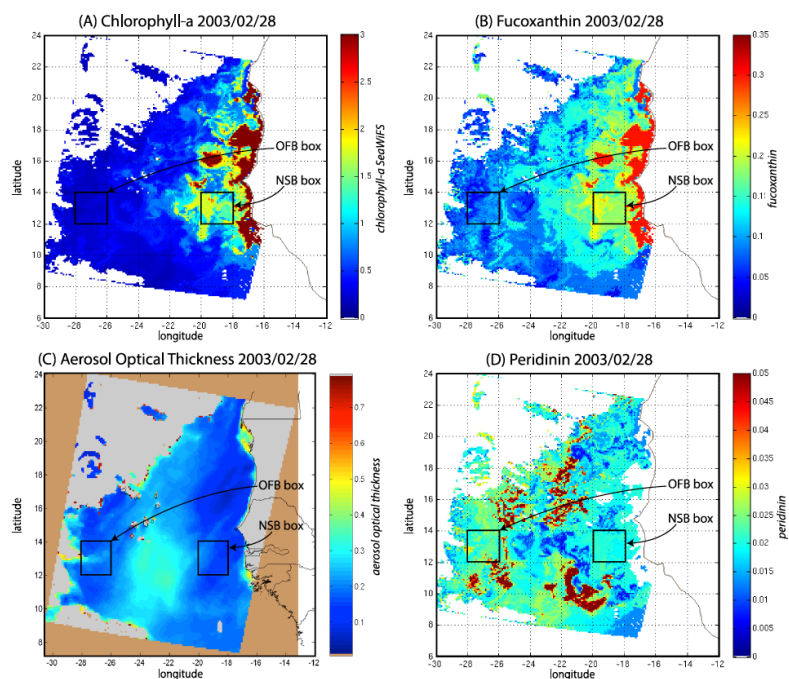
467  
468 Figure 9 : SST for 2 January 2003. Note the well-marked upwelling (cold temperature) north of 13°N.  
469



470  
471 Figure 10 : (A) *chl-a* concentration, (B) fucoxanthin ratio, (C) aerosol optical thickness, (D) peridinin for 6  
472 January 2003. Panels (B) and (D) show that a second-order information was retrieved, which is correlated  
473 with the *chl-a* concentration (A) but is not equivalent. It is found that the aerosol optical thickness (C) does  
474 not contaminate the estimated parameters (fucoxanthin and peridinin ratios).  
475



476 extrapolate what has not been learned and for that reason we flagged the pixels in the SeaWiFS  
477 images that have a *chl-a* concentration greater than  $3\text{ mg m}^{-3}$ .  
478



479  
480

481 Figure 11 : (A) *chl-a* concentration, (B) fucoxanthin ratio, (C) aerosol optical thickness, (D) Peridinin  
482 for 28 February 2003. Panels (B) and (D) show that a second order information was retrieved, which is  
483 correlated with the *chl-a* concentration (A) but is not equivalent. It is found that the aerosol optical  
484 thickness (C) does not contaminate the estimated parameters (fucoxanthin and peridinin ratios).  
485

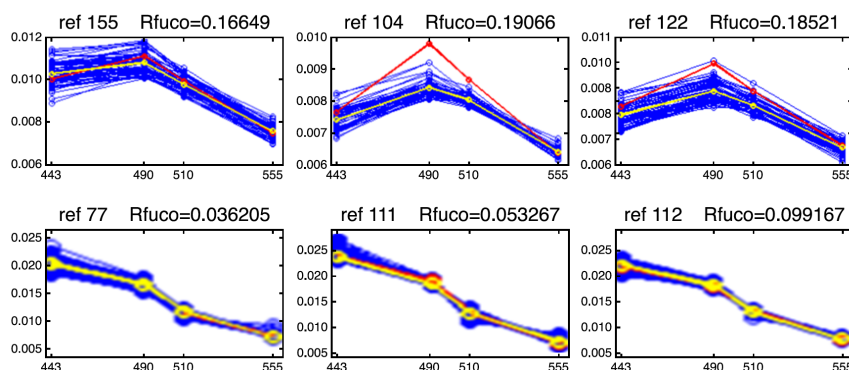
486

487 Regarding the images obtained for 1 January 2003 in the Senegalo-Mauritanian region  
488 (Fig 8A, B, C, D), we observe that the *chl-a* (Fig 8A) is very high at the coast and decreases offshore  
489 in accordance with the upwelling intensity as shown in the SST image (Fig 9). Moreover, we  
490 observed a persistent well-marked *chl-a* pattern south of the Cap Vert peninsula in form of a “W”.  
491 Except in the southern part of the region, the AOT (Aerosol Optical Thickness) is low, which means  
492 that the atmospheric correction of the reflectance is quite small, which gives confidence in the ocean-  
493 color data products. The fucoxanthin concentration is maximum at the coast and decreases offshore  
494 as does the *chl-a* concentration, in agreement with the works of Uitz *et al.*, (2006, 2010).  
495 Fucoxanthin presents coherent spatial patterns. Peridinin concentration is somewhat complementary



496 to that of fucoxanthin, with the low fucoxanthin concentration area corresponding to high peridinin  
497 concentration area (northern part of Figs 8B, D). This behavior is also observed in Figure 10 (6  
498 January 2003) and in Figure 11 (28 February, 2003) endorsing the analysis shown in Figure 8.  
499 For 28 February, we selected two square box regions (Fig. 11), one near the coast (NSB,  
500 long [-20°, -18°], lat [12°,14°]) and the other about 800 km offshore (OFB, long [-28°, -26°], lat  
501 [12°,14°]). NSB waters correspond to upwelling waters while OFB waters correspond to oligotrophic  
502 waters. We projected the eleven ocean color parameters of the NSB and OFB pixels on the 2S-SOM  
503 map.

504



505

506

507

508 Figure 12 : Reflectance spectra (in blue) of six referent vectors (red) selected during the decoding of  
509 28 February: top line, in the NSB region (long. [-20°, -18°], lat. [12°, 14°]); bottom line, in the OFB  
510 region (long. [-28°, -26°], lat. [12°, 14°]). The reflectance spectra of the captured pixels are in blue  
511 and those of the  $w$  in yellow.

512

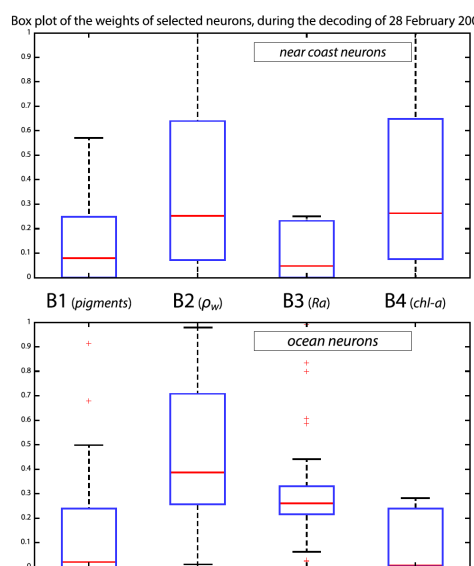
513

514 Figure 12 presents the reflectance spectra (in blue) captured by three neurons of the 2S-SOM  
515 corresponding to pixels located in the NSB region (top line) and those captured by three neurons  
516 corresponding to pixels located in the OFB region (bottom line). The reflectance spectra of the  
517 associated referent vectors  $w$  are in yellow. The satellite reflectance spectra match the referent vector  
518 spectra; moreover the fucoxanthin ratio varies inversely with the mean value of the spectrum: the  
519 higher the fucoxanthin ratio, the smaller the mean value of the spectrum. The pigment concentration  
520 is greater near the coast.

521 We note a strong difference between the shape and the intensity of the near-shore (NSB) and  
522 offshore (OFB) spectra. The OFB spectra present mean values higher than those of the NSB spectra.



523 This is due to the fact that NSB spectra were observed in a region where diatoms are abundant, as  
524 shown by the high value of fucoxanthin concentration in this region (Figs 8, 10, and 11), which is a  
525 proxy for diatoms along with higher *chl-a* concentration. In Figure 12, we note the lower values of  
526 the coastal spectra at 443 nm, which can be interpreted as a predominant effect of spectral absorption  
527 by phytoplankton pigments and CDOM.  
528 The different spectra are close together in the OFB region and more disperse in the NSB region. This  
529 can be explained by the fact that the OFB region corresponds to Case-1 waters while the NSB region  
530 waters are close to Case-2 waters and are influenced by the variability of near shore process like  
531 turbidity, or presence of dissolved matters.  
532



533  
534

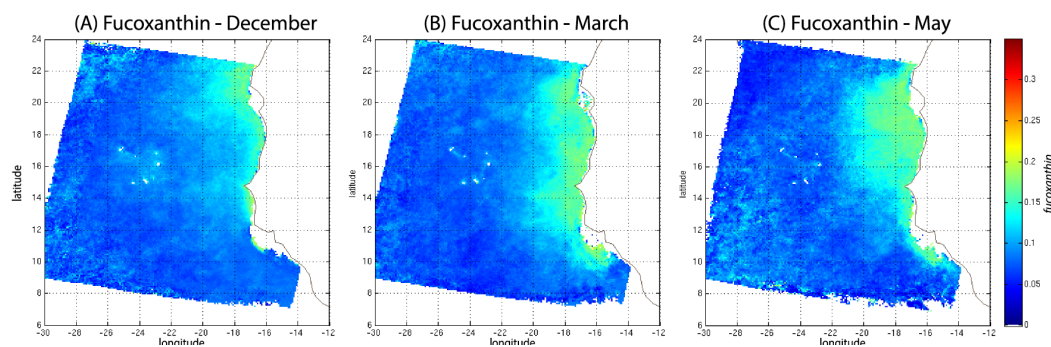
535 Figure 13 Box plot of the weights of the selected neurons during the decoding of the 28 February  
536 data. From left to right, weights of blocks B1, B2, B3, B4. Top panel, in the OFB region (long. [-20°,  
537 -18°], lat. [12°, 14°]); bottom panel, in the NSB region (long. [-28°, -26°], lat. [12°, 14°]).

538  
539  
540

541 We analyzed the weights of the blocks for all the neurons selected in the analysis of the coastal (NSB)  
542 and offshore (OFB) boxes. Figure 13 presents the box plot of the weight  $a_{cb}$  corresponding to the  
543 neurons belonging to the four blocks (B1, B2, B3, B4), with the constrain that the sum of the weights  
544 of a neuron is 1; a weight larger than 0.25 indicates the predominance of a block in the learning for  
545 the classification (see section 3.5). It is clear that the weights for pixels near the coast (Fig 13, top



546 panel) are different from those for offshore pixels (Fig. 13, bottom panel). As already mentioned in  
547 section 4.3 and also shown in Figure 7, the weights of the 2S-SOM play a significant role in the 2S-  
548 SOM topology and consequently in the pigment retrieval. The weights of blocks B1 and B4 that take  
549 into account the influence of the pigment ratios and the chlorophyll content in the retrieval are very  
550 low for the offshore (OFB) oligotrophic region and more important for the coastal (NSB) region. The  
551 weights of the blocks B2 and B3, which take into account the influence of the reflectances ( $\rho_w(\lambda)$ ,  
552  $R_a(\lambda)$ ), dominate for the offshore regions. In coastal waters the weights of all the blocks are used,  
553 with a smaller influence of B3, which is associated with  $R_a$ . This shows the automatic adaptation of  
554 the 2S-SOM to the environment in order to optimize the clustering efficiency with respect to a  
555 classical SOM.  
556



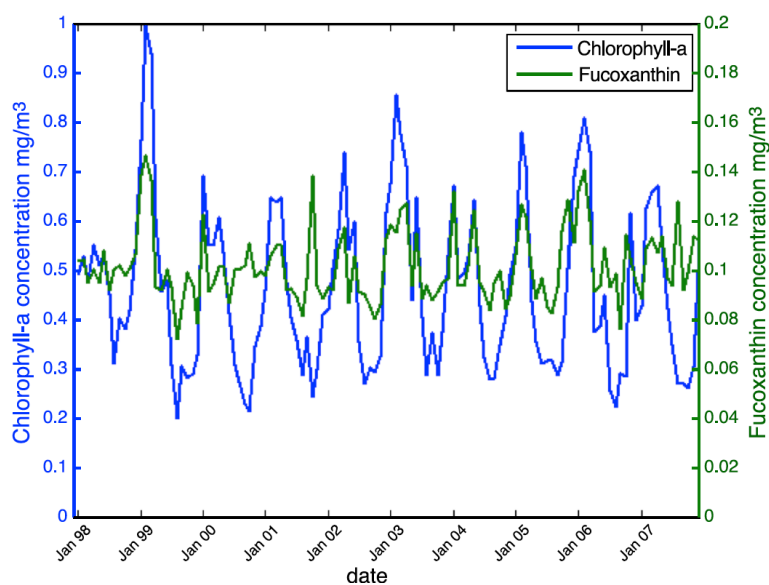
557  
558 Figure 14 : Monthly fucoxanthin concentration averaged for an 11- years (1998-2009) for December  
559 (A), March (B) and May (C).  
560

561  
562 In order to study the seasonal variability of the fucoxanthin concentration with some statistical  
563 confidence in the Senegalo-Mauritanian upwelling region, we constructed a monthly climatology for  
564 an 11-year period (1998–2009) of the SeaWiFS observations by summing the daily pixels of the  
565 month under study. The resulting climatology is presented in Figure 14 for December (Fig. 14a),  
566 March (Fig. 14b), and May (Fig. 14c). The fucoxanthin concentration, and consequently the  
567 associated diatoms, presents a well-marked seasonality. Fucoxanthin starts to develop in December  
568 North of 19°N, presents its maximum intensity in March when the upwelling intensity is maximum,  
569 extends up to the coast of Guinea (12°N) in April and begins to decrease in May where it is observed  
570 north of Cabo Verde peninsula (15°N) in agreement with the observations reported by *Farikou et al*,  
571 (2015) and *Demarcq and Faure*, (2000).





572 Figure 15 shows the fucoxanthin (in green) and the *chl-a* (in blue) concentrations computed from  
573 satellite observations for an 11-year period of SeaWiFS observations in the NSB region. There is a  
574 good correlation in phase between these two variables but not in amplitude (a good coincidence of  
575 peak occurrence but weak correlation in peak amplitude) showing that the relationship between  
576



577  
578

579 Figure 15 : . *chl-a* (in blue) and *fucoxanthin* (in green) concentrations for near-shore pixels (in the  
580 NSB region).

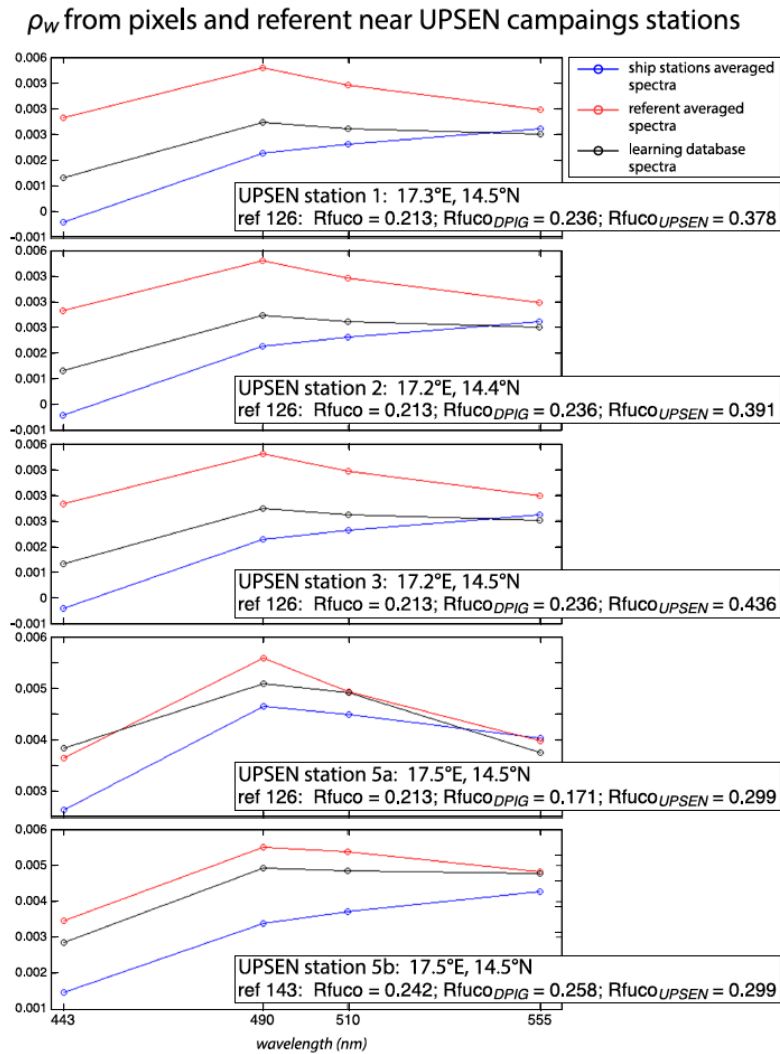
581  
582

583 fucoxanthin and *chl-a* is complex as mentioned by *Uitz et al.*, (2006). In particular, there is a weak  
584 peak in fucoxanthin in October 2001, which is not correlated with a *chl-a* peak.

585

### 586 5-2 Analysis of the UPSEN campaigns

587 Figure 16 shows, for the UPSEN stations 1, 2, 3, 5a and 5b (see figure 1 for their geographical  
588 position), the averaged in-situ spectrum (in blue), the 2S-SOM spectrum (in red) which is the  
589 spectrum of the 2S-SOM neuron captured by the collocated satellite VIIRS sensor observations, the  
590 spectrum (in black) of the learning database (DPIG) captured by that neuron that is the closest to the  
591 in-situ spectrum. These three spectra are close together showing the good functioning of the 2S-SOM.  
592



593

594

595

596 Figure 16 : For ship stations 1, 2, 3, 5a and 5b, the averaged spectrum of the in situ spectra of the  
 597 UPSEN station is shown in blue; the spectrum of the referent vector of the 2S-SOM neuron, which  
 598 has captured the satellite observations which are the closest to the UPSEN station is shown in red;  
 599 the spectrum of the learning database (DGIP) captured by the neuron that is the closest to the  
 600 averaged satellite spectra is shown in black.

601

602

603 Their shapes are close to these observed in the NSB region (Figure 12) but their intensity is lower  
 604 meaning that their waters are more absorbing than the NSB waters due to a higher pigment





605 concentration. In fact, the UPSEN stations were located near the coast (figure 1) in the Hann bight  
606 south off the Cap Verde peninsula, which is very rich in phytoplankton pigments. In table 3, we  
607 present the fucoxanthin ratios associated with the referent vectors ( $R_{\text{fuco}_{2S-SOM}}$ ), the closest DPIG  
608 fucoxanthin-ratios captured by the neuron of the referents and the fucoxanthin-ratios measured  
609 during the UPSEN campaign. We note that the fucoxanthin ratios of the in-situ measurements are in  
610 the range of the DPIG (see table 1), which allows a good functioning of the 2S-SOM. The pigment  
611 ratios obtained from ocean-color observations through 2S-SOM are close to pigment concentrations  
612 measured at the ship stations, which confirms the validity of the method we have developed. We  
613 remark that the best 2S-SOM estimate of fucoxanthin ratio with respect to the UPSEN in-situ  
614 measurement is given at station 5b which is the farthest off the coast. These results endorse the  
615 climatological study of the Senegalo-Mauritanian upwelling region we have done with the 2S-SOM  
616 (section 5.1).

617  
618

UPSEN STATION	REFERENT N°	RFUCO 2S-SOM	RFUCO DPIG	RFUCO 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 3 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

619  
620

621 Table 3 : For ship stations 1, 2, 3, 5a and 5b of the UPSEN campaigns, we show the referent  
622 captured by the VIIRS observations, the fucoxanthin-ratio associated with this referent ( $R_{\text{fuco}_{2S-SOM}}$ ) the fucoxanthin-ratio of the closest DPIG fucoxanthin-ratio captured by the neuron of the  
623 referent and the fucoxanthin-ratio measured in situ during the UPSEN campaign

624  
625  
626

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

633  
634



635 **6 - DISCUSSION**

636

637 Machine learning methods are powerful methods to invert satellite signals as soon as we have  
638 adequate database to support the calibration. Several technics have been used for retrieving  
639 biological information from ocean color satellite observations. First people employed MLPs, which  
640 are a class of neural networks suitable to model transfer function (*Thiria et al, 1993*). *Gross et al,*  
641 (2000, 2004) retrieved *chl-a* concentration from SeaWiFS, *Bricaud et al, (2006)* modeled the  
642 absorption spectrum with MLPs, *Raitos et al, 2008* and *Palacz et al, 2013* introduced additional  
643 environmental variables in their MLPs such as SST in the retrieval of PSC/PFT from SeaWiFS,  
644 which improved the skill of the inversion. Another suitable procedure was to embed NN in a  
645 variational inversion, which is very efficient way when a direct model exists (*Jamet et al, 2005;*  
646 *Brajard et al, 2006a,b; Badran et al, 2008*). Statistical analysis of absorption spectra of  
647 phytoplankton and of pigment concentrations were conducted by *Chazottes et al, (2006, 2007)*, by  
648 using a SOM.

649 In the present study, due to the fact that the learning dataset we used is quite small (515 elements)  
650 which makes MLPs and classical supervised learning methods unusable, we decided to use an  
651 unsupervised neural network classification method which is an extension of the SOM, method well  
652 adapted to deal with small database whose elements are very inhomogeneous; we cluster available  
653 satellite ocean-color reflectance at five wavelengths and their derived products, such as chlorophyll  
654 concentration, and the associated in situ pigment ratios.

655 The major points of this study are the following:

656 - The clustering was carried out by developing a new neural classifier, the so-called 2S-SOM, which  
657 presents several advantages with respect to the classical SOM. As in the SOM, we defined clusters  
658 that assemble vectors, which are close together in terms of a specified distance. This classifier was  
659 learned from a worldwide database (DPIG) whose vectors are ocean-color parameters observed by  
660 satellite multi-spectral sensors and associated pigment concentrations measured in situ. In the  
661 operational phase, SeaWiFS images are decoded allowing the estimation of the pigment  
662 concentration ratios. The major advantage of 2S-SOM with respect to the classical SOM is to  
663 cluster variables having similar physical significance in blocks having specific weights. The  
664 weights attributed to the four blocks are computed during the learning phase and vary with the  
665 quality of the variables and with respect to the location on the ocean (near the coast or in deep  
666 ocean). This permits to modulate the variable influence in the cost function, which makes the  
667 clustering more informative than this provided by the SOM. For offshore waters, the block  
668 decomposition allowed us to show that more influence is given to the reflectance ratios  $Ra(\lambda)$  and



669 less to the *chl-a* and pigment concentrations; on the contrary near the coast the weights indicate  
670 more active use of the pigment composition and the *chl-a* concentration. The resulting 2S-SOM  
671 clustering therefore takes into account at best the information that belongs to the specific water  
672 content.

673 - The 2S-SOM decomposes the DPIG into a large number of significant ocean-color classes allowing  
674 reproduction of the different possible situations encountered in the dataset we analyze. Besides, we  
675 assume that the relationship between the pigment concentration and the remote sensed ocean-color  
676 observations is independent of the location, which is justifiable since the relationship depends on  
677 the optical properties of ocean waters through well-defined physical laws which are region-  
678 independent. This also endorses the fact that we use a global database to retrieve pigments in a  
679 definite region. On the contrary the different phytoplankton species vary from one region to another  
680 making the relationship between pigment ratio and phytoplankton species strongly depending on  
681 the region, which justifies the fact we focused our study on the pigment retrieval rather than the  
682 PSC or PFT as mentioned above. Moreover, most of the recent phytoplankton in situ identifications  
683 have been made using pigment measurements with the HPLC method (*Hirata et al*, 2011). It is  
684 therefore more natural to retrieve the pigment concentrations, which is the quantity we measured,  
685 than the associated PSC or PFT, which are estimated from the pigment observations through  
686 complex non-linear and region-dependent algorithms (*Uitz et al*, 2006). Due to the characteristics  
687 of the DPIG, the method has the ability to retrieve pigment concentration patterns over a large  
688 range (0.02 – 2 mg m<sup>-3</sup>).

689 - We were able to analyze the pigment concentration in the Senegalo-Mauritanian region by  
690 processing satellite ocean color observations with the 2S-SOM. We found an important seasonal  
691 signal of fucoxanthin concentration with a maximum occurring in March. We evidenced a large  
692 offshore gradient of fucoxanthin concentrations, the near shore waters being richer than the  
693 offshore ones. We showed that the offshore region waters correspond to Case-1 waters while the  
694 near shore waters are close to Case-2 waters and are influenced by the variability of near shore  
695 process like turbidity, or presence of dissolved matters. The UPSEN measurements show that the  
696 pigment ratios of the Senegalo-Mauritanian region are in the range of the DPIG database used to  
697 calibrate the method, which justifies the use of the 2S-SOM algorithm to investigate this region.

698 - We used daily satellite observations to construct a monthly climatology of pigment concentrations  
699 of the Senegalo-Mauritanian upwelling region, which has been poorly surveyed by oceanic cruises.  
700 Due to the highly non-linear character of the algorithms for determining the pigment concentrations  
701 from satellite measurements, it is more rigorous mathematically to apply these algorithms to daily  
702 satellite data and to average this daily estimate for the climatology period under study than to



703 estimate them from the satellite data climatology as many authors have done (*Uitz et al., 2010*;  
704 *Hirata et al., 2011*). We found that Fucoxanthin starts to develop in December North of 19°N,  
705 presents its maximum intensity in March when the upwelling intensity is maximum, extends up to  
706 the coast of Guinea (12°N) in April and begins to decrease in May

707

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

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

717 - We then found that our fucoxanthin climatology is in agreement with in situ observations of  
718 phytoplankton reported in *Blasco et al. (1980)* in March to May 1974 off the coast of Senegal  
719 during the JOINT I experiment. These authors analyzed 740 water samples collected with Niskin  
720 bottles at 136 stations extending along a line at 21°40'N (in the northern part of the studied region)  
721 from 0 to 100 km offshore. The samples were taken at several depths (mostly at 100, 50, 30, 15, 5  
722 m). Phytoplankton cells were counted and identified by the Utermohl inverted microscope  
723 technique (*Blasco, 1977*). They found that diatoms reach their maximum concentration in April–  
724 May and are the most abundant group in that period, whereas the other cells predominate in  
725 March. Similar microscope observations have been reported in the ocean area south of Dakar by *A.*  
726 *Dia (1985)* during several ship surveys in February–March 1982–1983.

727 - Our method is also in agreement with the monthly eleven years climatology presented in *Farikou et*  
728 *al., (2015)* who used a modified PHYSAT method to retrieve the *PFT* in the Senegalo-Mauritanian  
729 region.

730 - The pigment concentrations provided by the 2S-SOM from the VIIRS sensor observations are in  
731 qualitative agreement with the in-situ measurements done at five stations during the two UPSEN  
732 campaigns in 2012 and 2013, showing that the method is able to function in waters where the  
733 pigment concentrations are quite high (fucoxanthin ratios of the order 0.4).

734

735

736



737

738

739 **7 - CONCLUSION**

740

741 We developed a new neural network clustering method, the so-called 2S-SOM algorithm to retrieve  
742 phytoplankton pigment concentration from satellite ocean color multi spectral sensors. The 2S-SOM  
743 algorithm is a SOM specifically designed to deal with a large number of heterogeneous components  
744 such as optical and chemical measurements. The major advantage of 2S-SOM with respect to the  
745 classical SOM is to cluster variables having similar significance in blocks having specific weights.  
746 The weights attributed to the blocks during the learning phase vary with the quality of the variables.  
747 This permits to modulate the variable influence in the cost function, which makes the clustering more  
748 informative than this provided by the SOM. Moreover, the 2S-SOM method is efficient and rapid as  
749 soon as the calibration is done since it uses elementary algebraic operations only. The 2S-SOM  
750 method is like a piecewise regression that takes advantage of the unsupervised classification of the  
751 SOM. We decomposed the DFIG database into a quite large number of partitions ( $9 \times 8 = 162$ ) when  
752 comparing our study to other studies (*Uitz et al*, 2006, 2012). The validity of the method has been  
753 controlled through a cross validation procedure and confirmed by three qualitative studies. Statistical  
754 parameters ( $R^2$  coefficients, RMSE and P-values) of the cross-validation between the DFIG in situ  
755 pigments and the pigments given by the 2S-SOM averaged for the 30 2S-SOM realizations which are  
756 presented in table 2, show the good performance of the method. It must be noticed that the  
757 performance mainly depends on the size of the learning set used to calibrate the 2S-SOM. This set  
758 must include all the situations encountered in the pigment retrieval. The larger the learning set, the  
759 better the method performs. Due to its generic character and its flexibility, the method could be used  
760 to determine a large variety of variables measure with satellite remote sensing observations.

761 The method was applied to study the seasonal variability of the fucoxanthin concentration in  
762 Senegalo-Mauritanian upwelling region. We showed a large offshore gradient of fucoxanthin, the  
763 higher concentration being situated near the shore. We were able to construct a monthly climatology  
764 for an 11-year period (1998–2009) of the SeaWiFS observations by summing the daily pixels of the  
765 month under study in a region which was poorly surveyed by oceanic cruises. The fucoxanthin  
766 concentration, and consequently the associated diatoms, presents a well-marked seasonality (Figure  
767 10). Fucoxanthin starts to develop in December North of  $19^\circ\text{N}$ , presents its maximum intensity in  
768 March when the upwelling intensity is maximum, extends up to the coast of Guinea ( $12^\circ\text{N}$ ) in April  
769 and begins to decrease in May where it is observed north of Cabo Verde peninsula ( $15^\circ\text{N}$ ) in  
770 agreement with the observations reported by *Farikou et al*, (2015) and *Demarcq and Faure*, (2000).



771 The UPSEN campaign results endorse the validity of the study of the Senegalo-Mauritanian  
772 upwelling region done with the 2S-SOM.

773

774 **Acknowledgments**

775 The study was supported by the projects CNES-TOSCA 2013-2014 and 2014-2015. The water-  
776 leaving reflectances were obtained from the SeaWiFS daily reflectances,  $\rho_{\text{obsTOA}}(\lambda)$ , provided by  
777 the NASA/GSFC/DAAC observed at the top of the atmosphere (TOA) and processed with the SOM-  
778 NV algorithm (Diouf et al., 2013) from 1998 to 2010. They are available at the web site:  
779 <http://poacc.locean-ipsl.upmc.fr/>. The DFIG data base was kindly provided by Dr. S. Alvain. We  
780 thank Dr. Alban Lazar and Dr. E. Machu for providing in situ data measured during the UPSEN  
781 experiments and stimulating discussions for their interpretation. We also thank Ray Griffiths for  
782 editing the manuscript.

783

784



785

786

## 787 **References**

788 Alvain S, Moulin C., Dandonneau Y. and Breon F. M. : Remote sensing of phytoplankton groups in  
789 case-1 waters from global SeaWiFS imagery. *Deep-Sea Res. Part1*, vol **52** (11), pp 1989-2004,  
790 2005.

791 Alvain S. Loisel H. and Dessailly D. : Theoretical analysis of ocean color radiances anomalies and  
792 implications for phytoplankton group detection. *Optics Express*, vol **20** (2), 2012.

793 Antoine D., André J. M. , Morel A. : Oceanic primary production : Estimation at global scale from  
794 satellite (Coastal Zone Color Scanner) chlorophyll. *Global Biogeochem Cy.* vol **10**, pp 57-69,  
795 1996.

796 Badran F., Berrada M. , Brajard J., Crepon M. , Sorrer C., Thiria S., Hermand J.P. , Meyer M.,  
797 Perichon L., Asch M. : Inversion of satellite ocean colour imagery and geoacoustic  
798 characterization of seabed properties : Variational data inversion using a semi-automatic adjoint  
799 approach *J. Marine Systems*, vol 69, pp 126-136, 2008

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

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

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

808 Blasco D. : Red tide in the upwelling region of Baja California. *Limnol. Oceanogr.* vol 22, pp 255-  
809 263, 1977

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

813 Bracher A., Bouman HA, Brewin RJW, Bricaud A, Brotas V, Ciotti AM, Clementson L, Devred E,  
814 Di Cicco A, Dutkiewicz S, Hardman-Mountford NJ, Hickman AE, Hieronymi M, Hirata T, Losa  
815 SN, Mouw CB, Organelli E, Raitzos DE, Uitz J, Vogt M and Wolanin A : Obtaining  
816 Phytoplankton Diversity from Ocean Color: A Scientific Roadmap for Future Development. *Front.*  
817 *Mar. Sci.* 4:55. doi: 10.3389/fmars.2017.00055, 2017





- 818 Brajard J., Jamet C., Moulin C. and Thiria S. : Atmospheric correction and oceanic constituents  
819 retrieval with a neuro-variational method. *Neural Networks*, Vol 19(2), p178-185, 2006
- 820 Brajard J., Jamet C., Moulin C. and Thiria S : Neurovariational inversion of ocean color images. *J.*  
821 *Atmos. Space Res.* Vol 38, n 2, pp 2169-2175, 2006
- 822 Brewin R. J. W., Hardman-Mountford N. J., Lavender S. J., Raitos D. E., Hirata T., Uitz J., et al. :  
823 An inter-comparison of bio-optical techniques for detecting dominant phytoplankton size class  
824 from satellite remote sensing. *Remote Sens. Environ.* 115, 325–339. doi:  
825 10.1016/j.rse.2010.09.004, 2011
- 826 Brewin R. J. W., Sathyendranath S., Hirata, T., Lavender, S.J., Barciela, R., Hardman-Montford, N.J :  
827 A three-component model of phytoplankton size class for the Atlantic Ocean. *Ecol. Model.* vol **22**,  
828 pp 1472-1483, 2010.
- 829 Bricaud A., Mejia C. , Blondeau Patissier D. , Claustre H., Crepon M. and Thiria S. : Retrieval of  
830 pigment concentrations and size structure of algal populations from absorption spectra using  
831 multilayered perceptrons. *Applied Optics Mars 2007* vol 46 n°8., 2006
- 832 Capet X., Estrade, P., Machu, E., Ndoye, S. et al. : On the Dynamics of the Southern Senegal  
833 Upwelling Center: Observed Variability from Synoptic to Superinertial Scales : *J. Phys.*  
834 *Oceanogr.* vol **47** (1), pp 155-180, 2017
- 835 Cavazos T. : Using Self-Organizing Maps to Investigate Extreme Climate Events: An Application to  
836 Wintertime Precipitation in the Balkans. *J. Climate*, vol **13**, 1718–1732, 2000.
- 837 Chazotte A., Crepon M., Bricaud A., Ras J. and Thiria S. : Statistical analysis of absorption spectra  
838 of phytoplankton and of pigment concentrations observed during three POMME cruises using a  
839 neural network clustering method. *Applied Optics*, 46 (18), 3790-3799, 2007
- 840 Chazottes A., Bricaud A., Crepon M. and Thiria S. : Statistical analysis of a data base of absorption  
841 spectra of phytoplankton and pigment concentrations using self-organizing maps. *Appl. Opt.* 45,  
842 8102-8115, 2006
- 843 Ciotti A. and Bricaud A. : Retrievals of a size parameter for phytoplankton and spectral light  
844 absorption by colored detrital matter from water-leaving radiances at SeaWiFS channels in a  
845 continental shelf region off Brazil. *Limnol. Oceanogr. Methods*, vol **4**, pp 237-253, 2006.
- 846 Demarcq H. and Faure V. : Coastal upwelling and associated retention indices from satellite SST.  
847 Application to *Octopus vulgaris* recruitment. *Oceanografica Acta*, vol **23**, pp 391-407, 2000.
- 848 Dia A. Biomasse et biologie du phytoplancton le long de la petite côte sénégalaise et relations avec  
849 l'hydrologie. Rapport interne N°44 du CRODT, Réf: 0C000798, 1981-1982. On line on the web  
850 site:<http://www.sist.sn/gsd/collect/publi/index/assoc/HASH2127.dir/doc.pdf>





- 851 Diouf D., Niang A., Brajard J., Crepon M. and Thiria S. : Retrieving aerosol characteristics and sea-  
852 surface chlorophyll from satellite ocean color multi-spectral sensors using a neural-variational  
853 method. *Remote Sens. Environ.* **vol 130**, pp 74-86, 2013.
- 854 Farikou O., Sawadogo S., Niang A., Brajard J., Mejia C., Crépon M. and Thiria S. : Multivariate  
855 analysis of the Sénégal-Mauritanian area by merging satellite remote sensing ocean color and SST  
856 observations. *J. Environ. Earth Sci.* **vol 5** (12), pp 756-768, 2013
- 857 Farikou O., Sawadogo S., Niang A., Diouf D., Brajard J., Mejia C., Dandonneau Y., Gasc G.,  
858 Crepon M., and Thiria S. : Inferring the seasonal evolution of phytoplankton groups in the  
859 Sénégal-Mauritanian upwelling region from satellite ocean-color spectral measurements, *J.*  
860 *Geophys. Res. Oceans*, **vol 120**, pp 6581-6601, 2015.
- 861 Gordon H. R. : Atmospheric correction of ocean color imagery in the Earth Observing System era. *J.*  
862 *Geophys. Res. Atmospheres*, **vol 102**(D14), pp 17081-17106, 1997.
- 863 Hewitson B.C. and Crane R. G. : Self organizing maps : application to synoptic climatology. *Climate*  
864 *research*, **vol 22**, pp 13-26, 2002
- 865 Friedrich T. and Oschlies A. : Basin-scale pCO<sub>2</sub> maps estimated from ARGO float data : A model  
866 study, *J. Geophys. Res.*, **vol 114**, C10012, doi: 10. 1029/2009JC005322, 2009.
- 867 Gross L., Frouin R., Dupouy C., Andre J. M. and Thiria S. : Reducing biological variability in the  
868 retrieval of chlorophyll\_a concentration from spectral marine reflectance. *Applied Optics*, **Vol. 43**  
869 **Issue 20** pp. 4041, 2004
- 870 Gross L., Thiria S., Frouin R., Mitchell B.G : Artificial neural networks for modeling transfer  
871 function between marine reflectance and phytoplankton pigment concentration *J. Geophys. Res.*  
872 **Vol 105**,no.C2, pp3483-3949, february 15, 2000
- 873 Jamet C., Thiria S., Moullin C., Crepon M. : Use of a neural inversion for retrieving Oceanic and  
874 Atmospheric constituents for Ocean Color imagery : a feasibility study.  
875 doi:10.1175/JTECH1688.1, *J. Atmos. Ocean. Techno. :/ Vol. 22*, No. 4, pp. 460–475, 2005
- 876 Hirata T. , Aiken J., Hardman-Mountford N., Smyth T. J. and Barlow R.G. : An absorption model to  
877 determine phytoplankton size classes from satellite ocean color, *Remote Sens. Environ.* **vol 112**,  
878 pp 3153-3159, 2008.
- 879 Hirata T. , Hardman-Mountford N.J., Brewin R.J.W., Aiken J., Barlow R., Suzuki K., Isada T.,  
880 Howell E., Hashioka T., Noguchi-Aita M. and Yamanaka Y. : Synoptic relationships between  
881 surface chlorophyll-*a* and diagnostic pigments specific to phytoplankton functional types.  
882 *Biogeosciences*, **vol 8** (2): pp 311-327, 2011.
- 883
- 884



- 885 Jeffreys S.W. and Vesk M. : Phytoplankton Pigment in Oceanography : Guidelines to Modern  
886 Methods, UNESCO, Paris, ed S. W. Jeffery, R.F.C. Mantoura and S. W. Wright, Introduction to  
887 marine phytoplankton and their pigment signatures, pp 33-84, 1997.
- 888 Jouini M., Lévy M. , Crépon M. and Thiria S. : Reconstruction of ocean color images under clouds  
889 using a neuronal classification method. Remote Sens. Environ. vol **131**, pp 232-246, 2013
- 890 Kohavi R. : A study of cross-validation and bootstrap for accuracy estimation and model selection.  
891 Proceedings of the Fourteenth International Joint Conference on Artificial Intelligence. San Mateo,  
892 CA: Morgan Kaufmann ed.. **2** (12): pp 1137–1143, 1995.
- 893 Kohonen T : Self-organizing maps (3<sup>rd</sup> ed.). Springer, Berlin Heidelberg New York. 2001
- 894 Kruizinga S. and Murphy A : Use of an analogue procedure to formulate objective probabilistic  
895 temperature forecasts in the Netherlands. Mon. Wea. Rev., vol **111**, pp 2244–2254, 1983.
- 896 Liu Y. and Weisberg R. H. : Patterns of ocean current variability on the West Florida Shelf using the  
897 self-organizing map, J. Geophys. Res., **110**, C06003, doi:10.1029/2004JC002786, 2005
- 898 Liu Y., Weisberg R. H., and He R. : Sea surface temperature patterns on the West Florida Shelf using  
899 growing hierarchical self-organizing maps, J. Atmos. Oceanic Technol., vol **23**(2), pp 325– 338,  
900 2006
- 901 Longhurst A. R., Sathyendranath S., Platt T., Caverhill C. : An estimation of global primary  
902 production in the ocean from satellite radiometer data. J. Plank. Res. vol **17**, pp 1245-1271, 1995
- 903 Lorenz E. N : Atmospheric predictability as revealed by naturally occurring analogs. J. Atmos. Sci.,  
904 vol 26, pp 639–646, 1969
- 905 Morel A. and Gentili G. : Diffuse reflectance of oceanic waters. III. Implication of bidirectionality  
906 for the remote-sensing problem. Appl. Opt. vol 35, pp 4850-4862, 1996.
- 907 Mouw C. B. and Yoder J. A. : Optical determination of phytoplankton size composition from global  
908 SeaWiFS imagery. J. Geophys. Res. vol **115**, C12018, doi:10.1029/2010JC006337, 2010.
- 909 Ndoye S. , Capet X., Estrade P., Sow B., Dagorne D., Lazar A., Gaye A. and Brehmer P. : SST  
910 patterns and dynamics of the southern Senegal-Gambia upwelling center. J. Geophys. Res. Oceans,  
911 vol 119, pp 8315–8335. 2014
- 912 Niang A., Badran F., Moulin C., Crépon M. and Thiria S. : Retrieval of aerosol type and optical  
913 thickness over the Mediterranean from SeaWiFS images using an automatic neural classification  
914 method. Remote Sens. Environ. vol 100, pp 82-94, 2006.
- 915 Palacz A. P., St. John, M. A., Brewin, R. J.W., Hirata, T., and Gregg, W.W. : Distribution of  
916 phytoplankton functional types in high-nitrate low-chlorophyll waters in a new diagnostic  
917 ecological indicator model. Biogeosciences 10, 7553–7574. doi: 10.5194/bg-10-7553, 2013.
- 918



919

920 Raitos D. E., Lavender, S. J., Maravelias, C. D., Haralambous, J., Richardson, A. J., and Reid, P.

921 C. : Identifying phytoplankton functional groups from space: an ecological approach. *Limnol.*

922 *Oceanogr.* 53, 605–613. doi: 10.4319/lo.2008.53.2.0605, 2008

923 Reusch D. B., Alley, R. B., and Hewitson, B. C : North Atlantic climate variability from a self-

924 organizing map perspective, *J. Geophys. Res.*, vol **112**, D02104, doi:10.1029/2006JD007460, 2007.

925 Sathyendranath S., Watts S., L., Devred E., Platt T., Caverhill C. M., and Maass H. : Discrimination

926 of diatom from other phytoplankton using ocean-colour data, *Mar. Ecol. Prog. Ser.*, vol 272, pp 59–

927 68, 2004.

928 Uitz J., Claustre H., Morel A. and Hooker S.B : Vertical distribution of phytoplankton communities

929 in open ocean: an assessment based on surface chlorophyll. *J. Geophys. Res.* **111**, C08005,

930 doi:10.1029/2005JC003207. 2006

931 Uitz J., Claustre H., Gentili B. and Stramski D. : Phytoplankton class-specific primary production in

932 the world's ocean: seasonal and interannual variability from satellite observations. *Global*

933 *Biogeochem. Cycles*, vol **24**, GB 3016, doi:10.1029/2009GB003680, 2010

934 Van den Dool H. : Searching for analogs, how long must we wait? *Tellus*, vol **46A**, pp 314–324,

935 1994.

936 Varma, S., Simon, R. : Bias in error estimation when using cross-validation for model selection; *BMC*

937 *Bioinformatics*. vol 7. PMC 1397873. PMID 16504092. doi:10.1186/1471-2105-7-91, 2006

938 Vidussi F., Claustre H., Manca B. B., Luchetta A. and Marty J. C. : Phytoplankton pigment distribution in

939 relation to upper thermocline circulation in the eastern Mediterranean sea during winter. *J. Geophys.*

940 *Res.*, vol 106, pp 19,939-19,956, 2001.

941 Westberry T., Behrenfeld M.J., Siegel D. A. and Boss E.: Carbon-based productivity modeling with

942 vertically resolved photoacclimatation. *Global Biogeochem. Cycles*, vol **22**, GB2024,

943 DOI:10.1029/2007GB003078, 2008

944 Zorita E. and Von Storch H. : The Analog Method as a Simple Statistical Downscaling Technique:

945 Comparison with More Complicated Methods. *Journal of Climate*, vol **12**, pp 2474-2489, 1999.

946