



# Bayesian Methodology for Ocean Color Remote Sensing

Robert Frouin, Bruno Pelletier

## ► To cite this version:

Robert Frouin, Bruno Pelletier. Bayesian Methodology for Ocean Color Remote Sensing. 66 pages. 2013. <hal-00822032>

HAL Id: hal-00822032

<https://hal.archives-ouvertes.fr/hal-00822032>

Submitted on 13 May 2013

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# BAYESIAN METHODOLOGY FOR OCEAN-COLOR REMOTE SENSING

Robert FROUIN <sup>a</sup> and Bruno PELLETIER <sup>b</sup>

<sup>a</sup> Scripps Institution of Oceanography  
University of California, San Diego  
9500 Gilman Drive, La Jolla, CA 92093-0224, USA  
[rfrouin@ucsd.edu](mailto:rfrouin@ucsd.edu)

<sup>b</sup> IRMAR, Department of Mathematics  
Université Rennes 2, CNRS, UEB  
Place du Recteur Henri Le Moal  
35043 Rennes Cedex, France  
[bruno.pelletier@univ-rennes2.fr](mailto:bruno.pelletier@univ-rennes2.fr)

## Abstract

The inverse ocean color problem, i.e., the retrieval of marine reflectance from top-of-atmosphere (TOA) reflectance, is examined in a Bayesian context. The solution is expressed as a probability distribution that measures the likelihood of encountering specific values of the marine reflectance given the observed TOA reflectance. This conditional distribution, the posterior distribution, allows the construction of reliable multi-dimensional confidence domains of the retrieved marine reflectance. The expectation and covariance of the posterior distribution are computed, which gives for each pixel an estimate of the marine reflectance and a measure of its uncertainty. Situations for which forward model and observation are incompatible are also identified. Prior distributions of the forward model parameters that are suitable for use at the global scale, as well as a noise model, are determined. Partition-based models are defined and implemented for SeaWiFS, to approximate numerically the expectation and covariance. The ill-posed nature of the inverse problem is illustrated, indicating that a large set of ocean and atmospheric states, or pre-images, may correspond to very close values of the satellite signal. Theoretical performance is good globally, i.e., on average over all the geometric and geophysical situations considered, with negligible biases and standard deviation decreasing from 0.004 at 412 nm to 0.001 at 670 nm. Errors are smaller for geometries that avoid Sun glint and minimize air mass and aerosol influence, and for small aerosol optical thickness and maritime aerosols. The estimated uncertainty is consistent with the inversion error. The theoretical concepts and inverse models are applied to actual SeaWiFS imagery, and comparisons are made with estimates from the SeaDAS standard atmospheric correction algorithm and in situ measurements. The Bayesian and SeaDAS marine reflectance fields exhibit resemblance in patterns of variability, but the Bayesian imagery is less noisy and characterized by different spatial de-correlation scales, with more realistic values in the presence of absorbing aerosols. Experimental errors obtained from match-up data are similar to the theoretical errors determined from simulated data. Regionalization of the inverse models is a natural development to improve retrieval accuracy, for example by including explicit knowledge of the space and time variability of atmospheric variables.

*Index Terms* — Remote sensing, ocean color, atmospheric correction, inverse problem, Bayesian statistics.

# Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
<b>2</b>	<b>Problem position and approach</b>	<b>5</b>
2.1	About ill-posed inverse problems . . . . .	6
2.2	Bayesian approach to atmospheric correction . . . . .	7
2.3	Practical implementation . . . . .	8
2.3.1	Forward modeling . . . . .	8
2.3.2	Inverse applications . . . . .	9
2.4	Smoothing in observation geometry . . . . .	11
<b>3</b>	<b>Modeling of the satellite signal</b>	<b>11</b>
3.1	Gaseous absorption . . . . .	11
3.2	Atmospheric transmittance and spherical albedo . . . . .	12
3.3	Sun glint reflectance . . . . .	13
3.4	Effective whitecap reflectance . . . . .	13
3.5	Water body reflectance . . . . .	13
3.6	Atmospheric reflectance . . . . .	14
<b>4</b>	<b>Approximation of the forward operator</b>	<b>16</b>
4.1	Atmospheric functions . . . . .	16
4.1.1	Definition of $\Theta_a$ . . . . .	16
4.1.2	Approximation of $\Phi_a$ . . . . .	17
4.2	Water body reflectance . . . . .	18
4.2.1	Definition of $\mathcal{X}_w$ . . . . .	18
4.2.2	The NOMAD data set . . . . .	19
4.2.3	The AERONET-OC data sets . . . . .	20
4.2.4	Assembling the NOMAD and AERONET-OC data sets . . . . .	20
<b>5</b>	<b>Implementation of the inverse models</b>	<b>22</b>
5.1	Noise distribution . . . . .	22
5.2	Prior distributions . . . . .	22
5.3	Inverse applications . . . . .	24
5.3.1	Construction of the partition . . . . .	24
5.3.2	Approximation of the model coefficients . . . . .	25
<b>6</b>	<b>Evaluation on simulated data</b>	<b>26</b>
6.1	About mean squared error . . . . .	26
6.2	Performance statistics . . . . .	27
6.2.1	Average errors . . . . .	28
6.2.2	Errors per observation geometry . . . . .	29
6.2.3	Errors per atmospheric parameter . . . . .	30
6.2.4	Detailed analysis for a typical geometry . . . . .	37
<b>7</b>	<b>Application to SeaWiFS imagery</b>	<b>37</b>
7.1	S1999045100113 image, South Africa . . . . .	37
7.2	Other application examples . . . . .	43
7.3	Comparison with in-situ data . . . . .	51

<b>8 Summary and Conclusions</b>	<b>56</b>
<b>Acknowledgments</b>	<b>61</b>
<b>Appendices</b>	<b>61</b>
<b>Appendix A Missing data inference</b>	<b>61</b>
<b>Appendix B Tree-based partition rules</b>	<b>62</b>

## 1 Introduction

The classic approach to ocean-color remote sensing from space (Antoine and Morel, 1999; Gordon et al., 1997; Wang, 2010) consists of (i) estimating the aerosol reflectance in the red and near infrared where the ocean can be considered black (i.e., totally absorbing), and (ii) extrapolating the estimated aerosol reflectance to shorter wavelengths. The water reflectance is then retrieved by subtraction. This process is referred to as atmospheric correction. Depending on the application context, the retrieved water reflectance may then related to chlorophyll-a concentration using a bio-optical model, semi-analytical or empirical (e.g., O'Reilly et al., 1998), or used in inverse schemes of varied complexity to estimate optical properties of suspended particles and dissolved organic matter (see Lee, 2006).

The process of atmospheric correction is inherently difficult to achieve with sufficient accuracy, since only a small fraction (10% or less) of the measured signal may originate from the water body. Furthermore, the surface and atmospheric constituents, especially aerosols, whose optical properties are influential, exhibit high space and time variability. However this two-step approach has been successful, and it is employed in the operational processing of imagery from most satellite ocean-color sensors. Variants and improvements to the classic atmospheric correction scheme have been made over the years, especially to deal with non-null reflectance in the red and near infrared, a general situation in estuaries and the coastal zone. The improvements in these regions consider spatial homogeneity for the spectral ratio of the aerosol and water reflectance in the red and near infrared (Ruddick et al., 2000) or for the aerosol type, defined in a nearby non-turbid area (Hu et al., 2000). They also use iteratively a bio-optical model (Bailey et al., 2010; Siegel et al., 2000; Stumpf et al., 2011), exploit differences in the spectral shape of the aerosol and marine reflectance Lavender et al. (2005), or make use of observations in the short-wave infrared, where the ocean is black, even in the most turbid situations (Oo et al., 2008; Wang et al., 2009, 2007).

Other empirical approaches to atmospheric correction have been proposed in the literature. In Frouin et al. (2006), the TOA reflectance is combined linearly, so that the atmosphere/surface effects are reduced substantially or practically eliminated. This algorithm assumes that the perturbing signal, smooth spectrally, can be modeled by a low-order polynomial, and the polynomial is selected so that the linear combination is sufficiently sensitive to chlorophyll-a concentration. In Steinmetz et al. (2011), the atmospheric reflectance is approximated by a polynomial with non-spectral and spectral terms that represent atmospheric scattering and surface reflection, including adjacency effects from clouds and white surfaces. The water reflectance is modeled as a function of chlorophyll concentration and a backscattering coefficient for non-algal particles, and spectral matching is applied to tune the atmospheric and oceanic parameters.

Another approach to satellite ocean-color inversion is to determine simultaneously the key properties of aerosols and water constituents by minimizing an error criterion between the measured reflectance and the output of a radiative transfer model (e.g., Chomko and Gordon, 1988; Kuchinke et al., 2009; Land and Haigh, 1996; Stamnes et al., 2007). This belongs to the family of deterministic solutions to inverse problems ; for a mathematical treatment of the subject, we refer the interested reader to Engl et al. (1996). Through systematic variation of candidate aerosol models, aerosol optical thickness, hydrosol backscattering

coefficient, yellow substance absorption, and chlorophyll-a concentration, or a subset of those parameters, a best fit to the spectral top-of-atmosphere reflectance (visible and near infrared) is obtained in an iterative manner. The advantage of this approach, compared with the standard, two-step approach, resides in its ability to handle both Case 1 and Case 2 waters. It also can handle both weakly and strongly absorbing aerosols, even if the vertical distribution of aerosols, an important variable in the presence of absorbing aerosols, is not varied in the optimization procedure. A main drawback is that convergence of the minimizing sequence may be slow in some cases, making it difficult to process large amounts of satellite data. To cope with this issue, a variant proposed in [Brajard et al. \(2006\)](#) and [Jamet et al. \(2005\)](#) consists of approximating the operator associated to the radiative transfer (RT) model by a function which is faster in execution than the RT code, e.g., by neural networks. Still, convergence speed of the minimization algorithm remains an issue. It may also not be easy to differentiate absorption by aerosols and water constituents like yellow substances, processes that tend to decrease the TOA signal in a similar way. As a result, the retrievals may not be robust to small perturbations on the TOA reflectance. This reflects the fact that atmospheric correction is an ill-posed inverse problem ; in particular, different values of the atmospheric and oceanic parameter can correspond to close values of the TOA reflectance. In the context of deterministic inverse problem, stability of the solution can be obtained by regularization (see [Engl et al., 1996](#)), but to the best of our knowledge, regularization strategies are not implemented in the approaches described above.

Another route is to cast atmospheric correction as a statistical inverse problem and to define a solution in a Bayesian context. In this setting, one group of approaches consist of estimating, based on simulations, a function performing a mapping from the TOA reflectance to the marine reflectance. In [Shroeder et al. \(2007\)](#), a neural network model is fitted to simulated data. A similar approach is studied in [Gross et al. \(2007a,b\)](#), where the (finite-dimensional) TOA signal, corrected for gaseous absorption and molecular scattering, is first represented in a basis which is such that the correlation between the ocean contribution and atmosphere contribution is, to some extent, minimized. This representation of the TOA reflectance makes the function approximation problem potentially easier to solve. In these studies, data are simulated for all the observation geometries. In [Frouin and Pelletier \(2007\)](#); [Pelletier and Frouin \(2004, 2005\)](#), the angular information is decoupled from the spectral reflectance, and atmospheric correction is considered as a collection of similar inverse problems indexed by the observation geometry. These methods can all be formalized in a Bayesian context ; see [Kaipio and Somersalo \(2004\)](#) and [Tarantola \(2005\)](#) for an introduction on the subject.

The Bayesian approach to inverse problem consists of first specifying a probability distribution, called the prior distribution, on the input parameters (atmospheric and oceanic) of the RT model. As the name implies, the prior distribution reflects prior knowledge that may be available before the measurement of the TOA reflectance. A probabilistic modeling of any perturbation of the TOA reflectance is also typically considered, in the form of an additive random noise. The solution to the inverse problem is then expressed as a probability distribution which, in the present context of atmospheric correction, measures the likelihood of encountering values of water reflectance given the TOA reflectance (i.e., after it has been observed). The posterior distribution is a very rich object, and its complete reconstruction and exploration can rapidly become prohibitive from the computational side. Instead, one may reduce the ambition to extracting useful quantities, like its expectation and covariance. In the present setting of atmospheric correction, the expectation provides an estimate of the water reflectance, while the covariance allows a quantification of uncertainty in the water reflectance estimate.

In this paper, we address ocean-color remote sensing in a Bayesian context. We make the following contributions. First of all, we formulate the atmospheric correction problem at a certain depth of physical modeling, and we use the angular decoupling as in [Frouin and Pelletier \(2007\)](#); [Pelletier and Frouin \(2004, 2005\)](#). Prior distributions suitable for use at a global scale, as well as a noise model, are determined. Second, we define and implement numerical approximations of the expectation and covariance of the posterior distribution (i.e., the complete Bayesian solution). The procedure is developed for the marine reflectance as well as for the atmospheric parameters, hence these quantities are retrieved simultaneously from the TOA

reflectance, and measures of uncertainties are provided along with the retrievals. The modeling choices in this work have been governed by keeping the execution time of the models small, and by having theoretical guarantees on the performance. Let us point out that it is a forward model which is inverted and that, as precise as the physical modeling can be, it is important to detect cases where the model is limited in view of the measured TOA reflectance. So as a final contribution, we define and implement a model, based on level sets to detect these situations where the retrievals become meaningless.

The paper is organized as follows. In Section 2, the inverse problem of atmospheric correction is defined, the Bayesian solution is formulated, and the inverse applications that will be implemented in practice are defined. Further satellite signal modeling operator approximation is exposed in Section 3 and 4. The implementation of the inverse applications is described in Section 5. Some technical details are gathered in Appendices A and B at the end of the paper. In Section 6, performance is evaluated on simulated data, and the ill posed-ness of the inverse problem is illustrated and discussed. In Section 7, the theoretical concepts and inverse models are applied to Sea-viewing Wide Field-of view Sensor (SeaWiFS) imagery, and comparisons are made with estimates from the standard atmospheric correction algorithm and in-situ measurements. In Section 8, conclusions are given about the Bayesian methodology in terms of performance, robustness, and generalization, as well as a perspective on future work. Regionalization of the inverse models is a natural development to improve retrieval accuracy, for example by including explicit knowledge of the space and time variability of atmospheric variables.

## 2 Problem position and approach

Let  $L_{toa}$  be the radiance measured by the satellite ocean-color sensor in a given spectral band. Express  $L_{toa}$  in terms of bidirectional reflectance  $\rho_{toa}$  as  $\rho_{toa} = \pi L_{toa}/(F_s \cos(\theta_s))$ , where  $F_s$  is the extraterrestrial solar irradiance (corrected for Earth-Sun distance), and where  $\theta_s$  is the Sun zenith angle. In clear sky conditions (i.e., a cloudless atmosphere), following lines devised in [Tanré et al. \(1979\)](#),  $\rho_{toa}$  may be modeled as

$$\rho_{toa} = T_g \left[ \rho_{mol} + \rho_{aer} + \rho_{mol-g} + \rho_{aer-g} + \rho_{mol-aer} + \rho_g t_a + \frac{T_a \rho_f}{1 - S_a \rho_f} + \frac{T_a \rho_w}{1 - S_a \rho_w} \right], \quad (2.1)$$

where  $T_g$  is the gaseous transmittance (accounts for absorption of photons by nitrous oxide, ozone, oxygen, and water vapor),  $\rho_{mol}$  and  $\rho_{aer}$  are the molecular and aerosol reflectance (account for multiple scattering of photons by molecules or aerosols only),  $\rho_{mol-g}$  and  $\rho_{aer-g}$  account for interactions between molecules or aerosols and photons reflected by a wavy surface,  $\rho_{mol-aer}$  accounts for the coupling between scattering by molecules and scattering and absorption by aerosols,  $\rho_g$  is the sun glint reflectance,  $t_a$  and  $T_a$  are the direct and total (direct plus diffuse) transmittance of the atmosphere along the path sun-to-surface and surface-to-sensor,  $\rho_f$  accounts for backscattering of photons by whitecaps,  $S_a$  is the spherical albedo of the atmosphere (accounts for successive photon interactions with the surface, the atmosphere, and the surface again), and  $\rho_w$  is the water reflectance (accounts for photons backscattered by the water body). In this decomposition, the perturbing signal from the atmosphere and surface is completely separated from the water body contribution.

In the presence of molecules only, the top-of-atmosphere signal from the atmosphere and surface ( $\rho_w = 0$ ) is reduced to:

$$\rho_{toa}^0 = T_g \left[ \rho_{mol} + \rho_{mol-g} + \rho_g t_{mol} + \frac{T_{mol} \rho_f}{1 - S_{mol} \rho_f} \right], \quad (2.2)$$

where  $t_{mol}$ ,  $T_{mol}$ , and  $S_{mol}$  are respectively the direct transmittance, the total transmittance, and the spherical albedo for molecules. Define the quantities  $\rho$  and  $\rho_a$  respectively by

$$\rho = \frac{1}{T_g} (\rho_{toa} - \rho_{toa}^0), \quad (2.3)$$

and

$$\rho_a = \rho_{aer} + \rho_{aer-g} + \rho_{mol-aer} + \rho_g t_a - \rho_g t_{mol} + \frac{T_a \rho_f}{1 - S_a \rho_f} - \frac{T_{mol} \rho_f}{1 - S_{mol} \rho_f}. \quad (2.4)$$

Then we have:

$$\rho = \rho_a + \frac{T_a \rho_w}{1 - S_a \rho_w}. \quad (2.5)$$

The inner term in the right-hand side of (2.2) can be accurately computed from atmospheric pressure and wind speed, while  $T_g$  can be well approximated in the spectral bands of ocean-color sensors given the absorber amounts along the optical path. Consequently,  $\rho_{toa}^0$  can be evaluated at the same time as the measurement of  $\rho_{toa}$ , which can thus be converted to  $\rho$  using (2.3). Note that transforming  $\rho_{toa}$  to  $\rho$  by (2.3) amounts at first correcting  $\rho_{toa}$  for gaseous absorption and next at subtracting for known effects due to molecules only. This pre-processing enhances the relative contribution of the ocean signal into  $\rho$  and, importantly, diminishes the influence of air pressure and wind speed.

In actuality, the observed (corrected) reflectance departs from the range of the theoretical model (2.5) due to measurement errors and modeling uncertainties. To represent these sources of variability, we shall add a random noise term  $\varepsilon$ , leading to the statistical model

$$\rho = \rho_a + \frac{T_a \rho_w}{1 - S_a \rho_w} + \varepsilon. \quad (2.6)$$

All the variables in (2.6) are functions of the wavelength  $\lambda$ . The observed data is finite-dimensional and is a vector  $y = (y_1, \dots, y_d)$  of measurements of  $\rho$  in spectral bands centered at wavelengths  $\lambda_1 < \dots < \lambda_d$ , i.e.,  $y_i = \rho(\lambda_i)$ , for  $i = 1, \dots, d$ .

With these notations, atmospheric correction refers to the process of estimating  $\rho_w$  from  $y$  and without knowledge of  $\rho_a$ ,  $T_a$  and  $S_a$  in (2.5). Note that since only  $y$  is observed, atmospheric correction is only one part of the complete inverse problem of estimating both the atmospheric and oceanic parameters from  $\rho$ , i.e.,  $(\rho_a, T_a, S_a)$  and  $\rho_w$ , even if in atmospheric correction, interest is only in  $\rho_w$ .

## 2.1 About ill-posed inverse problems

Atmospheric correction is an ill-posed inverse problem: even without noise, i.e., in model (2.5), different states of the atmosphere and of the ocean may correspond to very close values of the satellite signal. To see this using model (2.5), denote by  $\mathcal{X}_a$  and  $\mathcal{X}_w$  two sets of values for the triple  $(\rho_a, T_a, S_a)$  and the water reflectance  $\rho_w$ , respectively. Pick a point  $(\rho_a^*, T_a^*, S_a^*)$  in  $\mathcal{X}_a$  and a point  $\rho_w^*$  in  $\mathcal{X}_w$ , leading to a corrected reflectance  $\rho^*$  according to (2.5). Fix a threshold  $\delta > 0$ , which represents a noise level on  $\rho$ , and consider all the possible combinations of  $(\rho_a, T_a, S_a)$  and  $\rho_w$  such that  $\rho_a + T_a/(1 - S_a \rho_w)$  is at a distance no more than  $\delta$  from  $\rho^*$ . These combinations are the pre-images of a ball of radius  $\delta$  by the operator mapping  $((\rho_a, T_a, S_a), \rho_w)$  to  $\rho$  according to (2.5).

One example of pre-images is provided in Figure 1. In this example, the reflectance  $\rho$  is evaluated in 8 spectral bands from the visible to the near infra-red, and the Euclidean distance on  $\mathbb{R}^8$  is used. It is apparent in this case that quite different marine reflectance spectra are mapped to very close TOA reflectance spectra. As a consequence, in the presence of noise on the measurements, the uncertainty on the retrieved marine reflectance may be large. Naturally, the size of the set of pre-images depends on the choice of the threshold  $\delta$  and on the search spaces  $\mathcal{X}_a$  and  $\mathcal{X}_w$ . The noise level on  $\rho$  is taken as  $\delta = 0.001$ , as will be justified further in the paper. The set  $\mathcal{X}_a$  results from considering realistic aerosol models in varied proportions and load, and the set  $\mathcal{X}_w$  is composed of in-situ water reflectance spectra corresponding mainly to Case I waters. Similar results (not displayed here) are obtained for other cases. Since the variability of the marine reflectance component of the pre-images can be rather large, as illustrated in Figure 1, it is important to define and attach a measure of uncertainty to the retrieval of the marine reflectance from space.

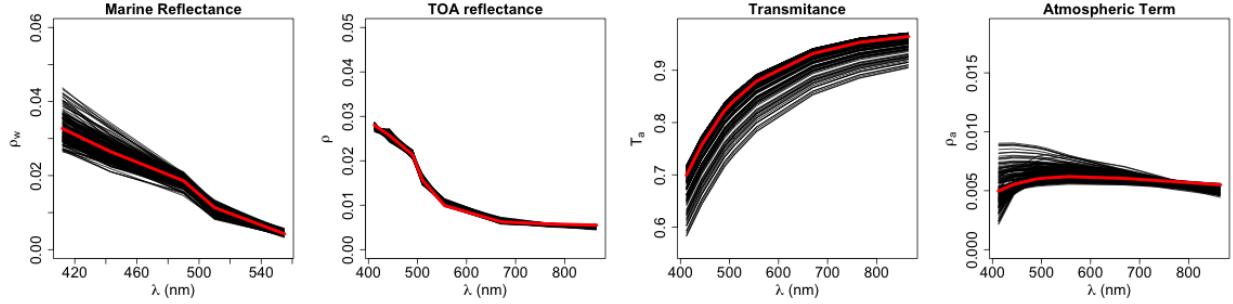


Figure 1: Example of pre-images. Actual values of  $\rho_w$ ,  $\rho$ ,  $T_a$  and  $\rho_a$  are displayed in red, and the pre-images at a distance no more than  $\delta$  (see text) are displayed in black.

## 2.2 Bayesian approach to atmospheric correction

We consider a finite-dimensional version of the statistical inverse problem (2.6) in which each quantity is evaluated at  $\lambda_i$ , for  $i = 1, \dots, d$ . To this end, let

$$x_{a,i} = (\rho_a(\lambda_i), T_a(\lambda_i), S_a(\lambda_i)) \quad \text{and} \quad x_{w,i} = \rho_w(\lambda_i) \quad \text{for all } i = 1, \dots, d,$$

and let

$$x_a = (x_{a,1}, \dots, x_{a,d}) \in \mathcal{X}_a \subset \mathbb{R}^{3d} \quad \text{and} \quad x_w = (x_{w,1}, \dots, x_{w,d}) \in \mathcal{X}_w \subset \mathbb{R}^d.$$

The subsets  $\mathcal{X}_a$  and  $\mathcal{X}_w$  in the above equations are constraint sets for the atmospheric parameters and the marine reflectance components respectively. The model reads as

$$y = \Phi(x_a, x_w) + \varepsilon, \tag{2.7}$$

where  $\varepsilon$  is a random vector in  $\mathbb{R}^d$ , and where  $\Phi : \mathcal{X}_a \times \mathcal{X}_w \rightarrow \mathbb{R}^d$  is the map with components  $(\Phi_i)_{1 \leq i \leq d}$  defined by  $\Phi_i(x_a, x_w) = \rho_a(\lambda_i) + T_a(\lambda_i)\rho_w(\lambda_i)/[1 - S_a(\lambda_i)\rho_w(\lambda_i)]$ .

In the Bayesian approach,  $x_a$ ,  $x_w$ , and  $y$  in (2.7) are treated as random variables. This defines a statistical model, and any vector of measurements  $y^{\text{obs}}$  is then considered as a realization of the random vector  $y$  in (2.7). To complete the definition of this model, there remains to specify a distribution  $\mathbb{P}_\varepsilon$  for the random noise  $\varepsilon$ , and a distribution for the pair  $(x_a, x_w)$ . The distribution of  $(x_a, x_w)$  is called the *prior distribution*. It describes, in a probabilistic manner, the prior knowledge one may have about  $x_a$  and  $x_w$  before the acquisition of the data  $y$ . Since there is no particular reason to expect that the atmospheric parameters and the marine reflectance should be correlated, such a distribution will be a product measure of the form  $\mathbb{P}_{x_a} \otimes \mathbb{P}_{x_w}$ , where  $\mathbb{P}_{x_a}$  and  $\mathbb{P}_{x_w}$  are probability measures on  $\mathbb{R}^{3d}$  and  $\mathbb{R}^d$ , respectively. The prior distribution allows one to incorporate known constraints in the model. For instance, that  $x_a$  and  $x_w$  must belong to the constraint sets  $\mathcal{X}_a$  and  $\mathcal{X}_w$  is specified in the model by considering prior distributions  $\mathbb{P}_{x_a}$  and  $\mathbb{P}_{x_w}$  with support in  $\mathcal{X}_a$  and  $\mathcal{X}_w$ . At last, the noise  $\varepsilon$  in model (2.7) will be considered as independent from  $(x_a, x_w)$ . Hence the prior distribution together with the noise distribution completely specifies the joint distribution of  $(x_a, x_w, y)$ .

The Bayesian solution to the inverse problem of retrieving  $(x_a, x_w)$  from  $y$  is defined as the conditional distribution of  $(x_a, x_w)$  given  $y$ . This distribution, further denoted by  $\mathbb{P}_{(x_a, x_w)|y}$ , is called the *posterior distribution*. Hence, given the observation  $y^{\text{obs}}$ , the solution is expressed as the probability measure  $\mathbb{P}_{(x_a, x_w)|y=y^{\text{obs}}}$ . From the computational side, though, the complete reconstruction of the posterior distribution is out of reach in many problems, despite the intense development of sampling techniques like Markov Chain Monte Carlo (MCMC). As an alternative, one generally restrict the objective to only estimating some

relevant characteristics of the posterior distribution, like its mean, its modes (i.e., points of local maximum), or its covariance matrix. In this work, we consider only the mean and the covariance matrix.

For the reader interested in materials on inverse problems, the book by [Kaipio and Somersalo \(2004\)](#) contains a comprehensive introduction to the Bayesian point of view, including computational aspects, as well as the book by [Tarantola \(2005\)](#), and the review article by [Stuart \(2010\)](#), while the text by [Engl et al. \(1996\)](#) provides a mathematical analysis of the method of deterministic regularization. Let us point out that, for the application that we have in mind, the inverse problem is finite-dimensional ; for approaches to infinite-dimensional inverse problems, see for example [Tarantola \(2005, Chapter 5\)](#), and the articles by [Stuart \(2010\)](#) and [Cotter et al. \(2009, 2010\)](#).

## 2.3 Practical implementation

We propose the following implementation, specific to atmospheric correction, of the general Bayesian approach to inverse problem defined above. First, we define prior distributions and a noise distribution that are suitable for an utilization at a global scale. Second, we define models to approximate the mean and covariance of the conditional distribution, which provide estimates for the marine reflectance, and to which we attach a measure of uncertainty. Finally, we introduce a quantity called a *p*-value which allows to detect those situations for which the observation  $y^{\text{obs}}$  is highly unlikely to have originated from model (2.7). This may occur for instance if the observation  $y^{\text{obs}}$  corresponds to one type of water (i.e., spectral dependance), not included in the support of the prior distribution on the marine reflectance, and in such a way that the observation  $y^{\text{obs}}$  lies far off the range of  $\Phi$  in (2.7). Such situations show evidence of a mismatch between the observed data and the model, and so this latter cannot be used reliably for inversion purposes.

### 2.3.1 Forward modeling

A prior distribution on the marine and atmospheric parameters describes, in terms of frequencies, the knowledge one may have about these parameters before the acquisition of the satellite data. For instance, if it is known in advance that, for the considered location and time, only marine spectra corresponding to open ocean waters will be encountered, then the prior distribution  $\mathbb{P}_{x_w}$  should reflect that. This can be achieved by requesting that its support  $\mathcal{X}_w$  be composed only of spectra of this type.

For the marine reflectance, we shall first define a compact set  $\mathcal{X}_w$  composed of realistic values of the marine reflectance encountered at a global scale. Next, in the absence of other information (e.g., relative proportion of water type), the prior distribution will be taken as the uniform measure on  $\mathcal{X}_w$ . Defining  $\mathcal{X}_w$  could be achieved by considering several models relating the marine reflectance with biological parameters, but we find it more appropriate to base our analysis on in-situ measurement of the marine reflectance. Indeed, modeling marine reflectance from measurements of inherent optical properties (e.g., [Garver and Siegel, 1997](#); [Gordon et al., 1988](#); [Morel and Maritorena, 2002](#); [Park and Ruddick, 2005](#)) does not fully take into account the natural correlations between controlling variables. For a given phytoplankton assemblage/population, absorption and backscattering coefficients are related in a unique way that is not captured in empirical bio-optical formulas determined from amalgamating measurements collected in diverse oceanic regions. Our analysis is therefore based on the NOMAD dataset ([Werdell and Bailey, 2005](#)), and on data acquired at several AERONET-OC sites ([Zibordi et al., 2010](#)).

To define the prior distribution  $\mathbb{P}_{x_a}$  on the atmospheric parameters,  $x_a$  is expressed as a function  $\Phi_a$  of other parameters  $\theta_a$ , taking values in some set  $\Theta_a$ . To keep the exposition simple at this point, we just mention that  $\theta_a$  is composed of physical variables with known range (pressure, wind speed), and of variables related to modeling of aerosol optical properties. As for the marine reflectance, uniform distributions will be considered for all these quantities, except for the aerosol optical thickness, for which we shall make use

of the study by Knobelspiesse et al. (2004) which shows that, over the oceans, this latter is approximately log-normally distributed. Details are provided in Section 5.2.

The noise term  $\varepsilon$  in (2.7) is intended to encapsulate all the sources of uncertainty leading to model (2.7). Hence its distribution  $\mathbb{P}_\varepsilon$  describes the frequencies of the magnitude by which the actual observations deviate from the physical forward modeling. Herein, we shall assume that  $\varepsilon$  is a gaussian random vector on  $\mathbb{R}^d$  with  $\mathbb{E}\varepsilon = 0$  and  $\text{Cov}(\varepsilon) = \sigma^2 \text{Id}_{\mathbb{R}^d}$ . Naturally these assumptions are restrictive, but since the actual ocean color sensors measure the spectrum in a limited number of spectral bands, more general models (see e.g. Bissantz et al., 2004) cannot be considered. Let us point out that the parameter  $\sigma$ , which governs the hypothesized global noise level, is akin to the regularization parameter in Tikhonov regularization scheme, and has to be properly selected. For this purpose, we consider a large number of TOA observations and we estimate  $\sigma$  by maximum likelihood.

### 2.3.2 Inverse applications

The prior distribution  $\mathbb{P}_{x_a} \otimes \mathbb{P}_{x_w}$  and the noise distribution  $\mathbb{P}_\varepsilon$  are now fixed as defined above. So in what follows,  $(x_a, x_w)$  denotes a random vector with distribution  $\mathbb{P}_{x_a} \otimes \mathbb{P}_{x_w}$ ,  $\varepsilon$  is a random vector independent from  $(x_a, x_w)$  and with distribution  $\mathbb{P}_\varepsilon$ , and  $y$  has distribution induced by (2.7). Expectations are taken according to these distributions.

Let  $r(y^{\text{obs}}) = \mathbb{E}[x_w|y = y^{\text{obs}}]$  and  $\Sigma(y^{\text{obs}}) = \text{Cov}(x_w|y = y^{\text{obs}})$  be the conditional mean and conditional covariance of the marine reflectance given the observation  $y^{\text{obs}}$ . A very large variety of numerical procedures and models may be employed to numerically approximate the applications  $m$  and  $\Sigma$ . To guide our choice, we considered the following objectives: (i) that the models be fast in execution, suitable for use on an operational basis, and (ii) that the theoretical applications  $m$  and  $\Sigma$  be approximated with a guaranteed accuracy and in a reasonable computer time. Based on these requirements, we define models based on a common partition of  $\mathbb{R}^d$  and which are either constant or linear above each element of the partition.

Formally, denote by  $A_1, \dots, A_M$  the elements, also called cells, of a partition of  $\mathbb{R}^d$ . We consider models which are linear over each cell (for  $r$ ) and constant over each cell (for  $\Sigma$ ), i.e., models of the form

$$\hat{r}(y^{\text{obs}}) = \sum_{m=1}^M (\alpha_m + B_m y^{\text{obs}}) \mathbf{1}_{A_m}(y^{\text{obs}}) \quad \text{and} \quad \hat{\Sigma}(y^{\text{obs}}) = \sum_{m=1}^M C_m \mathbf{1}_{A_m}(y^{\text{obs}}), \quad (2.8)$$

where  $\alpha_m \in \mathbb{R}^d$  and  $B_m \in \mathcal{M}_d(\mathbb{R})$  for all  $m = 1, \dots, M$ , and where  $C_m \in \mathcal{M}_d(\mathbb{R})$  is a covariance matrix, for all  $m = 1, \dots, M$ . From a numerical stand point, to process an observation  $y^{\text{obs}}$  using either of the models in (2.8), the first task is to determine which of the  $A_m$  among the  $M$  cells contain  $y^{\text{obs}}$ . To reduce the computational time of this operation, we shall use a partition induced by a perfect binary tree of depth  $K$ . This is a partition of  $\mathbb{R}^d$  into  $M = 2^K$  elements which is structured in such a way that determining cell membership requires only  $K$  evaluations of a simple rule (this is advantageous to an unstructured partition that could require  $2^K$  tests to determine cell membership).

The quality of the approximation of  $r$  by  $\hat{r}$  is measured by the  $L^2$  risk defined by  $\mathcal{R}(\hat{r}) = \mathbb{E}\|\hat{r}(y) - r(y)\|^2$ , where  $\|\cdot\|$  denotes the Euclidean norm on  $\mathbb{R}^d$ . It can be shown that  $\mathcal{R}(\hat{r}) = \mathbb{E}\|x_w - \hat{r}(y)\|^2 - \mathbb{E}\|x_w - r(y)\|^2$ , so that minimizing  $\mathcal{R}(\hat{r})$  with respect to  $\hat{r}$  is equivalent to minimizing  $\mathbb{E}\|x_w - \hat{r}(y)\|^2$ . Similarly, the quality of the approximation of  $\Sigma$  by  $\hat{\Sigma}$  is also measured by an  $L^2$  risk, where the norm used is induced by the embedding of  $\mathcal{M}_d(\mathbb{R})$  in  $\mathbb{R}^{d^2}$  endowed with the Euclidean norm.

Suppose that the partition  $A_1, \dots, A_M$  is fixed. Denote by  $\hat{\alpha}_m$ ,  $\hat{B}_m$ , and  $\hat{C}_m$ , for  $m = 1, \dots, M$ , the coefficients minimizing the risks for models  $\hat{r}$  and  $\hat{\Sigma}$  in (2.8). Observing that the risk decomposes into the sum of the errors over each  $A_m$ , the optimal coefficients are given by

$$(\hat{\alpha}_m, \hat{B}_m) \in \operatorname{argmin}_{\alpha, B} \mathbb{E} \left\{ \|\alpha + B y - r(y)\|^2 \mathbf{1}_{A_m}(y) \right\} \quad \text{and} \quad \hat{C}_m \in \operatorname{argmin}_C \mathbb{E} \left\{ \|C - \Sigma(y)\|^2 \mathbf{1}_{A_m}(y) \right\},$$

for all  $m = 1, \dots, M$ . Simple calculations allow to derive explicit expressions for the optimal coefficients which, in turn, can be approximated numerically to any arbitrary accuracy, since the joint distribution of  $(x_a, x_w, y)$  is known.

At this point, the only remaining free parameter of models (2.8) is the partition  $A_1, \dots, A_n$ . It can be shown that, under mild conditions, the  $L^2$  risk of these models will tend to zero as the partition is being refined. Basically, to ensure convergence, the number of cells must go to infinity while the cells have to shrink at an appropriate rate ; we refer the reader interested in these aspects to [Lugosi and Nobel \(1999\)](#); [Nobel \(1996\)](#) and [Gyorfi et al. \(2002, chapters 4 and 13\)](#). In the present work, the partition is grown from a perfect binary tree of depth  $K = 17$ , which yield a total of  $2^K = 131,072$  cells. The technical details of the construction are provided further in the paper. The main point here is that this partition gives a discretization of the set of possible values for the corrected reflectance (i.e.,  $\mathbb{R}^d$  for  $y$  in (2.7)).

It is essential to keep in mind that, as refined as the physical models considered in (2.1) and leading to (2.6) can be, the proposed methodology only offers an inversion of these models. Then given the observation  $y^{\text{obs}}$ , how much confidence can be placed in the retrievals, including in the proposed uncertainties ? In other words, how likely is  $y^{\text{obs}}$  to have originated from model (2.7) ? It is therefore critical to, at least, detect situations where there is reasonable evidence that the observation  $y^{\text{obs}}$  is incompatible with the model. For this purpose, we define a procedure based on level sets as follows.

Denote by  $f_y$  the density of  $y$  in model (2.7). For any  $t \geq 0$ , denote by  $\mathcal{L}(t)$  the upper level set of  $f_y$  at level  $t$  defined by  $\mathcal{L}(t) = \{y \in \mathbb{R}^d, : f_y(y) \geq t\}$ . Given the observation  $y^{\text{obs}}$ , let  $t(y^{\text{obs}})$  be the largest real number  $t > 0$  such that  $y^{\text{obs}}$  belongs to  $\mathcal{L}(t)$ , i.e.,  $t(y^{\text{obs}}) = \sup\{t > 0 : y^{\text{obs}} \in \mathcal{L}(t)\}$ . Then, we define the map  $p_V : \mathbb{R}^d \rightarrow (0; 1)$  by  $p_V(y^{\text{obs}}) = 1 - \mathbb{P}(y \in \mathcal{L}(t(y^{\text{obs}})))$ , i.e.,  $p_V(y^{\text{obs}})$  is the probability that a new observation  $y$  does not belong to the level set  $\mathcal{L}(t(y^{\text{obs}}))$ . We shall refer to  $p_V(y^{\text{obs}})$  as the *p*-value of  $y^{\text{obs}}$ .

To better understand this definition, fix a real number  $\alpha \in (0; 1)$ , and let  $t_\alpha > 0$  be the real number such that  $y$  belongs to  $\mathcal{L}(t_\alpha)$  with probability  $\alpha$ , i.e.  $\mathbb{P}(y \in \mathcal{L}(t_\alpha)) = \alpha$ . So, roughly speaking,  $\mathcal{L}(t_\alpha)$  is a set which contains a proportion  $\alpha$  of the data. Then, it can be proved that, for any measurable subset  $A$  of  $\mathbb{R}^d$  with  $\mathbb{P}(y \in A) = \alpha$ , we have  $\text{Vol}_d(\mathcal{L}(t_\alpha)) \leq \text{Vol}_d(A)$ , where  $\text{Vol}_d(\cdot)$  denotes the volume (Lebesgue measure) ; see e.g. [Polonik \(1995, 1997\)](#). In other words, for any  $0 < \alpha < 1$ ,  $\mathcal{L}(t_\alpha)$  is the smallest set (in terms of volume) which contains a proportion  $\alpha$  of the total mass and  $\mathcal{L}(t^{\text{obs}})$  is the smallest level set which contains  $y^{\text{obs}}$ . Then,  $p_V(y^{\text{obs}})$  can be interpreted as the probability that a new observation  $y$  be at least as extreme as the data  $y^{\text{obs}}$ . Hence a low value of  $p_V(y^{\text{obs}})$ , say lower than 1%, may indicate that the model and the observation are incompatible.

To construct an approximation of the *p*-value map  $p_V$ , we shall make use again of the partition  $A_1, \dots, A_M$ . For any  $1 \leq m \leq M$ , let  $p_m = \mathbb{P}(y \in A_m)$  and define the density estimate

$$\hat{f}_y(y) = \sum_{m=1}^M \hat{f}_m \mathbf{1}_{A_m}(y) \quad \text{with} \quad \hat{f}_m = \frac{p_m}{\text{Vol}_d(A_m)},$$

where it is implicitly assumed that  $\text{Vol}_d(A_m) > 0$  for all  $1 \leq m \leq M$ . Similarly, it can be shown that  $\hat{f}_y$  is a consistent estimate of  $f_y$  under appropriate conditions on the partition ([Lugosi and Nobel, 1999](#); [Nobel, 1996](#)). Then,  $p_V(y^{\text{obs}})$  can be approximated by

$$p_V(y^{\text{obs}}) = \sum_{m=1}^M p_m \mathbf{1}\{\hat{f}_m \leq \hat{f}_y(y^{\text{obs}})\}. \quad (2.9)$$

As for models (2.8), the coefficients  $p_m$ , for  $m = 1, \dots, M$ , can be approximated numerically to any arbitrary accuracy by simulating a large number of observations  $y$  and evaluating the proportion of them falling in each cell  $A_m$ .

## 2.4 Smoothing in observation geometry

In a satellite image, each pixel corresponds to a measurement of the incoming radiation flux in one direction, and atmospheric correction of the image is usually performed by processing each pixel independently of the others. This implies that, in this context, to process the whole image, one has to solve one inverse problem per pixel. More formally, the reflectance  $\rho$  in (2.6), as well as the operator  $\Phi$  in (2.7), depend on three angles  $t := (\theta_s, \theta_v, \Delta\varphi)$  characterizing the observation geometry, namely the sun zenith angle, the viewing zenith angle, and the azimuth difference angle (modulo  $\pi$ ). Following Pelletier and Frouin (2004, 2005) and Frouin and Pelletier (2007), we start with a discretization of the domain  $T$  of observation geometries into  $N$  points  $t_1, \dots, t_N$  forming a grid in  $T$ . Next, for each discretization point  $t_i$ , with  $1 \leq i \leq N$ , the models defined in Section 2.3 are constructed. The inverse models corresponding to an arbitrary observation geometry  $t$  are then defined by interpolation of the models over the grid.

The interpolation procedure is similar for each type of inverse models and is detailed next for the retrieval of the marine reflectance. For any  $t$  in  $T$ , denote by  $r_t(y)$  the condition expectation of  $x_w$  given  $y$  in model (2.7), where each quantity corresponds to the observation geometry  $t$ . Hence  $\{r_t : t \in T\}$  is the collection of functions that are the Bayesian solutions formed when  $t$  varies in  $T$ . For each  $1 \leq i \leq N$ , we first construct the model  $\hat{r}_i$  associated with observation geometry  $t_i$ . Then, given the observation  $y^{\text{obs}}$  corresponding to the observation geometry  $t$ , the Bayesian solution  $r_t(y^{\text{obs}})$  is approximated by

$$\hat{r}_t(y^{\text{obs}}) = \sum_{i=1}^N W_i(t) \hat{r}_i(y^{\text{obs}}). \quad (2.10)$$

In this equation, the coefficients  $W_i(t)$ , for  $i = 1, \dots, N$ , are defined in such a way that  $\hat{r}_t$  interpolates the  $N$  models  $\hat{r}_i$ ,  $i = 1, \dots, N$ . As per the quality of approximation of this final model, for each fixed  $t$  in  $T$ , the error of approximation  $\mathbb{E}\|\hat{r}_t(y) - r_t(y)\|^2$  can be made arbitrarily small since the joint distribution of  $(x_a, x_w, y)$  in (2.7) is known ; its value depends first on the quality of each  $\hat{r}_i$ ,  $i = 1, \dots, N$ , and second the number  $N$  of grid points for  $T$ , and the only limitations are the storage space for the model parameters and the execution time.

## 3 Modeling of the satellite signal

In this section, we provide details about the parameterizations, computations, and approximations used for the terms appearing in the decomposition of the satellite signal (equation 2.1), as well as a justification for this decomposition. In brief, the total signal is expressed as the sum of the marine signal (the signal of interest containing photons that have interacted with the water body) and the perturbing effects of the atmosphere and surface. Absorption by atmospheric gases is decoupled from scattering by molecules and aerosols and absorption by aerosols, and the water body is assumed to backscatter sunlight uniformly in all directions, i.e., the marine reflectance,  $\rho_w$ , is independent of viewing geometry. Also, the effect of reflectance contrast between the surface target and its environment, referred to as adjacency effect, is not taken into account, which is equivalent to considering that the pixel size is infinitely large or that spatial variations in surface reflectance are not important (large target formalism generally used in ocean color remote sensing). Note that the errors introduced by the various parameterizations and simplifications, as well as other uncertainties in the radiative transfer modeling (e.g., treatment of the air-sea interface and aerosols) are taken into account in the Bayesian inversion via the term  $\varepsilon$  of equation 2.6.

### 3.1 Gaseous absorption

In equation (2.1), the effects of scattering by molecules and aerosols are decoupled from the effect of absorption by atmospheric gases, and the gaseous transmittance is evaluated along the direct path from sun to

surface and surface to sensor. This is justified for ozone, which is located high in the atmosphere, where molecules and aerosols are rarified. Consequently, the incident and detected photons practically cross the ozone layer without being scattered. In the case of water vapor, the absorption bands occur where molecular scattering is weak, i.e., where aerosol scattering dominates. Since above 850 nm order 1 and 2 scattering events constitute the quasi-totality of the aerosol scattering signal, and since aerosols scatter mainly forward, the actual path followed by the photons does not differ much from the direct path sun-to-surface and surface-to-sensor. The same approximation is less justified for oxygen and nitrous oxide, especially when considering the completely atmospheric terms  $\rho_{mol}$ ,  $\rho_{aer}$ , and  $\rho_{mol-aer}$  (since photons do not interact with the surface or water body), but it is sufficiently accurate when gaseous absorption is weak (Deschamps et al. (1983)), which is the case in the spectral bands of ocean color sensors. Thus for measurements in spectral bands contaminated by gaseous absorption,  $T_g$  is expressed as:

$$T_g = \prod_i t_{g_i}(\theta_s, \theta_v, U_i) \quad (3.1)$$

where  $t_{g_i}$  and  $U_i$  are the transmittance and amount of gaseous absorber  $i$ , with  $t_{g_i}$  calculated for the direct path sun-to-surface and surface-to-sensor. Analytical expressions are obtained for  $t_{g_i}$  in each spectral band, by fitting random exponential band models (i.e., Goody (1964) for water vapor and Malkmus (1967) for the other gases) modified to take into account the variations of temperature and pressure along the atmospheric path (Buriez and Fouquart, 1980).

### 3.2 Atmospheric transmittance and spherical albedo

For a given zenith angle  $\theta$  (sun zenith angle  $\theta_s$  or view zenith angle  $\theta_v$ ), the direct transmittance, whether upward or downward (superscript  $u$  or  $d$ , respectively), is defined by  $t(\theta) = \exp[-(\tau_{mol} + \tau_{aer})/\cos(\theta)]$  where  $\tau_{mol}$  and  $\tau_{aer}$  are the molecular and aerosol optical thickness, and the total transmittance is given by  $T(\theta) = t(\theta) + E(\theta)$  where  $E(\theta)$  is the diffuse transmittance (molecules plus aerosols), i.e.,

$$t_a = t_a^d t_a^u = t(\theta_s) t(\theta_v) = \exp [-(\tau_{mol} + \tau_{aer})(1/\cos(\theta_s) + 1/\cos(\theta_v))] \quad (3.2)$$

$$\begin{aligned} T_a &= T_a^d T_a^u = T(\theta_s) T(\theta_v) = [t(\theta_s) + E(\theta_s)] [t(\theta_v) + E(\theta_v)] \\ &= \left[ \exp \left( -\frac{\tau_{mol} + \tau_{aer}}{\cos(\theta_s)} \right) + E(\theta_s) \right] \left[ \exp \left( -\frac{\tau_{mol} + \tau_{aer}}{\cos(\theta_v)} \right) + E(\theta_v) \right]. \end{aligned} \quad (3.3)$$

Note that  $T_a$ , unlike  $t_a$ , does not depend explicitly on the air mass  $m = 1/\cos(\theta_s) + 1/\cos(\theta_v)$ .

In these equations, the molecular optical thickness  $\tau_{mol}$  is computed as a function of wavelength  $\lambda$  using a fitting equation proposed by Hansen and Travis (1974), but modified to account for a depolarization factor of 0.0279 instead of 0.0310, i.e.,

$$\tau_{mol} \approx 0.00852 \lambda^{-4} (1 + 0.0113 \lambda^{-2} + 0.00013 \lambda^{-4}) P/P_0, \quad (3.4)$$

where  $P$  is the sea level atmospheric pressure,  $P_0$  the standard atmospheric pressure ( $P_0 = 1023.2$  hPa), and  $\lambda$  is expressed in  $\mu\text{m}$ . The molecules and aerosols vary with altitude according to:

$$\tau_a = \tau_{mol} \exp(-z/H_{mol}) + \tau_{aer} \exp(-z/H_{aer}), \quad (3.5)$$

where  $\tau_a$  is the total optical thickness of the atmosphere ( $\tau_a = \tau_{mol} + \tau_{aer}$ ),  $z$  is altitude, and  $H_{mol}$  and  $H_{aer}$  are the scale heights of molecules and aerosols ( $H_{mol} = 8$  km).

The spherical albedo of the atmosphere  $S_a$  has a small effect on the contributions from the water body and whitecaps, i.e., the last two terms of equation (2.1), because diffuse reflection by the surface (i.e.,  $\rho_f$  and  $\rho_w$ ) is small compared with atmospheric scattering. It is given by:

$$S_a = 2\pi \int_0^1 I^u(\mu) \mu d\mu \quad (3.6)$$

where  $\mu = \cos(\theta)$  and  $I^u$  is the diffuse irradiance of the atmosphere for a solar zenith angle  $\theta$ .

### 3.3 Sun glint reflectance

The sun glint reflectance  $\rho_g$  is considered spectrally flat (spectral dependence of the refractive index of water versus air is neglected) and computed as a function of wind speed,  $U$ , using the Cox and Munk (1954) model:

$$\rho_g = \pi G(U) R_{Fresnel}(\gamma) / [4 \cos^4(\beta) \cos^2(\theta_s)] \quad (3.7)$$

where  $R_{Fresnel}$  is the reflection coefficient at the air-sea interface for the scattering angle  $\gamma$ ,  $\beta$  is the wave inclination, and  $G$  is the Gaussian slope distribution (depends on wind speed  $U$ ) as given by Cox and Munk (1954). Dependence on wind direction is ignored. Wind speed is obtained from meteorological data. The above expression is not accurate for high values of  $\rho_g$  (i.e.,  $\rho_g > 0.01$ ) due to uncertainties on surface wind speed and intrinsic variability of Fresnel reflection on a rough surface.

### 3.4 Effective whitecap reflectance

The effective reflectance of whitecaps  $\rho_f$  is modeled as the product of the fraction of the surface contaminated by whitecaps,  $A$ , and the reflectance of whitecaps,  $\rho_{f0}$ . For  $A$ , the empirical relation between  $A$  and wind speed  $U$  obtained by Koepke (1984) is used, i.e.,  $A = 2.95 \cdot 10^{-6} U^{3.52}$  where  $U$  is expressed in m/s. For  $\rho_{f0}$ , the mean value of 0.22 measured by Koepke (1984) for combined patches and steaks (accounts for the thinning of whitecaps with time) is used with a spectral factor  $f(\lambda)$  based on in situ and aircraft measurements (Frouin et al., 1996; Nicolas et al., 2001). Thus  $\rho_f$  is computed as:

$$\begin{aligned} \rho_f &= A \rho_{f0} = 0.22 A f(\lambda) \\ &= \begin{cases} 0.65 \cdot 10^{-6} U^{3.52} \exp(-1.75(\lambda - 0.6)) & \text{if } \lambda \geq 0.6 \\ 0.65 \cdot 10^{-6} U^{3.52} & \text{if } \lambda < 0.6. \end{cases} \end{aligned} \quad (3.8)$$

Note that in equation (2.1) the effects of whitecaps and water reflectance are decoupled, which is convenient to isolate the signal backscattered by the water body. A more accurate modeling would be to replace the sum of the last two terms in equation 2.1 by  $(\rho_f + \rho_w)/[1 - S_a(\rho_f + \rho_w)]$ , but the decomposition is justified because  $\rho_f$ ,  $\rho_w$ , and  $S_a$  are small compared with unity.

### 3.5 Water body reflectance

The radiance backscattered by the water body, after transmission across the interface, is assumed independent of direction, i.e.,  $\rho_w$  is constant with respect to viewing geometry. In actuality,  $\rho_w$  exhibits some angular anisotropy, especially at large viewing and solar zenith angles (Morel and Gentili, 1993; Morel et al., 1995). The upward transmittance  $T_a^u$  is also modified due to the non-Lambertian character of the water body, by a few percent for the geometry conditions encountered in ocean-color remote sensing (Gordon et al., 2010;

Gordon and Franz, 2008). Our treatment is based on the premise that errors introduced by neglecting bidirectional effects on  $\rho_w$  are small compared with other atmospheric correction errors, i.e., those associated with the determination of the atmospheric reflectance  $\rho_a$ . Furthermore, taking into account bidirectional effects would have required knowing the scattering phase function of marine particles and their assemblages, but such knowledge is not comprehensive.

Spatial heterogeneity in the water reflectance is neglected in equation (2.1). This assumption, commonly made in ocean color remote sensing, may not be valid in the vicinity of land, clouds, and sea ice, or more generally in regions where the spatial water reflectance contrast is relatively high (Bélanger et al., 2007; Santer and Schmechtig, 2000). The adjacency effects in the TOA imagery can be estimated, however, from the spatial reflectance fields and atmospheric properties initially retrieved (i.e., in the large target formalism), then used to correct the TOA imagery before applying again the inversion scheme.

Water reflectance in the  $\rho$  simulations is taken directly from in-situ datasets acquired in a wide range of Case 1 and Case 2 situations (details are provided in Section 4). Using an average reflectance model for water reflectance would be restrictive (e.g., too rigid spectral constraints), and varying arbitrarily, even within observed boundaries, the parameters affecting water reflectance, i.e., backscattering and absorption coefficients, would not describe well natural variability, and might complicate the problem unnecessarily (some simulated cases may not be realistic).

### 3.6 Atmospheric reflectance

The various atmospheric functions and interaction terms in equation (2.1), except gaseous transmittance (expressed analytically), therefore  $\rho_a$  and  $\rho$  in equation (2.5), are computed using a successive-orders-of-scattering code (Deuzé et al., 1989; Lenoble et al., 2007). In the computations, unlike for water reflectance (see above), the aerosol characteristics are specified from models. This approach is justified for several reasons. On the one hand, the aerosol parameters of the forward model are generally not all available in existing archives (e.g., vertical distribution is often missing), and data, when they exist, are mostly from coastal or island sites, i.e., they do not represent well the open ocean. On the other hand, the objective is to retrieve water reflectance, not aerosol properties. In practice, after correction for  $T_g$ , the successive-orders-of-scattering code is run twice, once with aerosols and non-null water reflectance and once with only molecules and  $\rho_w = 0$ , and the second output is subtracted from the first to yield  $\rho$ .

Figure 2 displays, as an example, simulations of  $\rho_{toa} \cos(\theta_s)$  (no gaseous absorption) and  $\rho \cos(\theta_s)$  in the spectral range 400–900 nm for various aerosol conditions and angular geometries. The atmosphere contains molecules and aerosols and is bounded by a wavy surface. Sun zenith angle is 29.4 degrees, view zenith angle is 20.1 degrees, relative azimuth angle is 0 (forward scattering), 90, and 180 degrees (backscattering), wind speed is 5 m/s, and surface pressure is 1013 hPa. Maritime, continental, and urban aerosol models (WMO, 1983) are considered, and aerosol optical thickness is 0.1 and 0.5 at 865 nm. Aerosol scale height is 2 km. The water body is assumed black, i.e.,  $\rho_w = 0$ . The total signal, i.e.,  $\rho_{toa} \cos(\theta_s)$ , exhibits a strong spectral dependence for relative azimuth angles of 90 and 180 degrees, i.e., side and backscattering geometries, with values increasing rapidly with decreasing wavelength (Figures 2a, left and 2b, left). This is explained by the influence of molecular scattering. For a relative azimuth angle of 0 degree, i.e., forward scattering, the spectral variation of  $\rho_{toa} \cos(\theta_s)$  is much smaller, essentially due to the influence of sun glint (Figure 2c, left). The signal corrected for molecular effects, i.e.,  $\rho \cos(\theta_s)$ , is much smaller in magnitude, for example about twice less at 500 nm for the geometries considered (Figures 2a, right, 2b, right, and 2c, right). In the case of forward scattering,  $\rho$  is negative (Figure 2c, right), all the more as the aerosol optical thickness is large, because  $\rho_g$  is large and the term  $\rho_g(t_a - t_{mol})$  that appears in the expression of  $\rho_a$  is negative and  $t_a$  decreases as aerosol optical thickness increases.

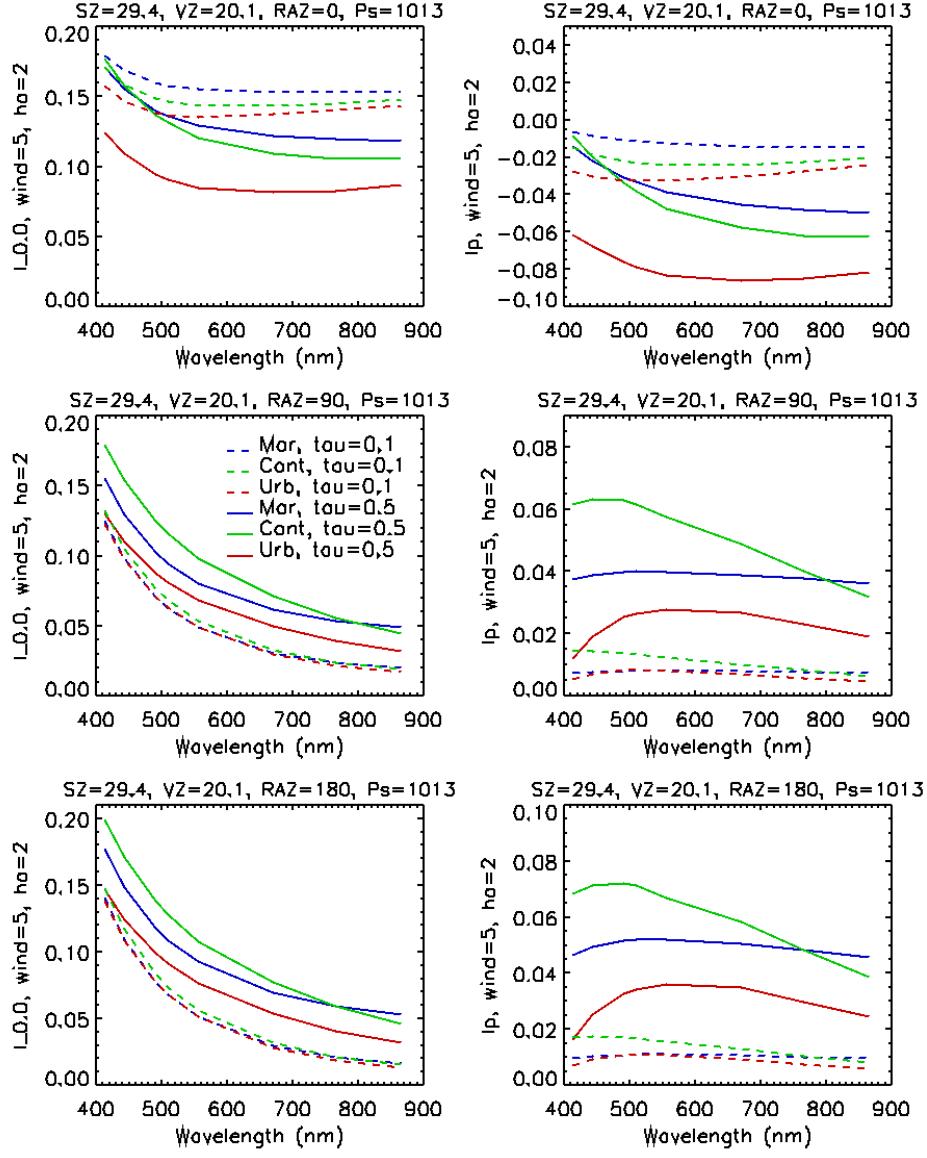


Figure 2: Simulations of the top-of-atmosphere normalized radiance ( $\pi L_{\text{toa}}/F_s$ ) by a vector radiation transfer code based on the successive-orders-of-scattering method. The atmosphere contains molecules and aerosols and is bounded by a wavy surface. Backscattering by the water body is null. Sun zenith angle is 29.4 deg., view zenith angle is 20.1 deg., relative azimuth angle is 0 (top panels), 90 deg. (middle panels), and 180 deg. (bottom panels), wind speed is 5 m/s, and surface pressure is 1013 hPa. Three types of aerosols are considered, i.e., maritime, continental, and urban, and aerosol optical thickness is 0.1 and 0.5 at 865 nm. Aerosol scale height is 2 km. The total signal is displayed in the left panels, and the signal after subtraction of the molecular signal (calculated assuming no aerosols) is displayed in the right panels.

## 4 Approximation of the forward operator

Constructing the inverse models requires multiple computations of the corrected TOA reflectance  $\rho$  in spectral bands of interest for our application, i.e., the SeaWiFS spectral bands centered at 412, 443, 490, 510, 555, 670, 765, and 865 nm. Of course, the construction scheme is general, and not restricted to a particular set of spectral bands. The angular grid, i.e., the Sun and view zenith angles,  $\theta_s$  and  $\theta_v$ , and the relative azimuth angle between Sun and view directions,  $\Delta\phi$ , are given in Table 4. The zenith angles are the Gauss angles in the successive-orders-of-scattering routine used to compute  $\rho$ . The range of angular conditions generally encountered in ocean-color remote sensing is covered, including observations in the Sun glint. As discussed later, the number of computations required to approximate the theoretical maps defined in Section 2 (i.e., conditional expectation, covariance, and  $p$ -value) will be of the order of the hundred of million points for each geometry. To keep the necessary computer time within reachable bounds, we first quantize the forward operator for each geometry, that is, we set up, once and for all, a data base stored on disk where  $\rho$  is evaluated at selected sampling points in the parameter space of atmospheric and oceanic variables. Next, to evaluate  $\rho$  at an arbitrary point, i.e., for an arbitrary state of the ocean and the atmosphere, the operator is approximated by interpolation. This yields a substantial gain, both in execution time and disk storage space, since the computational cost of running a radiative transfer code twice far exceeds the one of performing interpolation. Naturally there is a price to pay for this strategy, which results in an approximation error of the theoretical forward operator, i.e. the differences between actual computations with the radiative transfer code and the results from the interpolation procedure. That said, with properly chosen quantization points, located on a fine enough grid, the approximation error of the forward operator is intended to be small compared with the other sources of uncertainties and is accounted for via the noise term  $\varepsilon$  in (2.7).

$\theta_s, \theta_v$	0.11	1.43	3.28	5.14	7.00	8.87	10.73	12.59	14.46	16.32
	18.19	20.05	21.92	23.78	25.65	27.51	29.38	31.24	33.11	34.98
	36.84	38.71	40.57	42.44	44.30	46.17	48.03	49.90	51.76	53.63
	55.49	57.36	59.22	61.09	62.95	64.82	66.68	68.55	70.41	72.28
	74.15	76.01								
$\Delta\phi$	0.00	5.00	10.00	15.00	20.00	25.00	30.00	35.00	40.00	45.00
	50.00	55.00	60.00	65.00	70.00	75.00	80.00	85.00	90.00	95.00
	100.00	105.00	110.00	115.00	120.00	125.00	130.00	135.00	140.00	145.00
	150.00	155.00	160.00	165.00	170.00	175.00	180.00			

Table 1: Sun zenith angle, view zenith angle, and relative azimuth angle values of the grid.

Recall the notations from Section 2.2 where  $\mathcal{X}_a$  denotes the set of possible values for the vector of atmospheric parameters  $x_a = (x_{a,1}, \dots, x_{a,d})$  with  $x_{a,i} = (\rho_a(\lambda_i), T_a(\lambda_i), S_a(\lambda_i))$  at wavelength  $\lambda_i$ , for  $i = 1, \dots, d$ . Similarly,  $\mathcal{X}_w$  denotes the set of possible values for the vector of marine reflectance at wavelengths  $\lambda_1, \dots, \lambda_d$ . In the next section,  $x_a$  is expressed as a function  $\Phi_a$  of a vector of parameters  $\theta_a$  taking values in a set  $\Theta_a$ , so that we have  $\mathcal{X}_a = \Phi_a(\Theta_a)$ , and we define a discretization of  $\Theta_a$ . Next, we define a discretization of  $\mathcal{X}_w$  based on in-situ data.

### 4.1 Atmospheric functions

#### 4.1.1 Definition of $\Theta_a$

The components of the parameter vector  $\theta_a$  are the pressure, the wind speed, and the aerosol scale height, optical thickness, and type (i.e., model). The quantization points for these parameters are given in Table 2.

Parameter	Values	$N$
Pressure (hPa)	1003.0 ; 1013.25 ; 1023	3
Scale Height (km)	1.0 ; 2.0 ; 3.0	3
Aerosol models	15 WMO models ; see text.	15
Aerosol Optical Thickness	0.0 ; 0.05 ; 0.1 ; 0.15 ; 0.2 ; 0.3 ; 0.4 ; 0.5 ; 0.6	9
Wind speed (m/s)	1.0 ; 3.0 ; 5.0 ; 8.0 ; 12.0	5

Table 2: Discretization of the atmospheric parameters space: sampling values, and number  $N$  of different values.

The first five parameters take values in a closed interval of the real line. The last parameter, the aerosol model, refers to a convex combination of three basic types, namely Continental, Maritime, and Urban, whose optical properties are specified from [WMO \(1983\)](#). In this display, each aerosol mixture is parameterized by the mixing proportions of the three basic types, i.e., by a parameter  $\alpha := (\alpha_1, \alpha_2, \alpha_3)$  of the respective proportions. The proportions  $\alpha_i$ 's are comprised between 0 and 1 and sum to 1, so that the parameter space for  $\alpha$  is the two-dimensional unit simplex, further denoted by  $\mathcal{S}^2$ , of  $\mathbb{R}^3$ , i.e.,

$$\mathcal{S}^2 = \left\{ (\alpha_1, \alpha_2, \alpha_3) \in \mathbb{R}^3 : \sum_{i=1}^3 \alpha_i = 1 \text{ and } \alpha_i \geq 0 \text{ for all } 1 \leq i \leq 3 \right\}.$$

When looked at in  $\mathbb{R}^3$ , the set  $\mathcal{S}^2$  has the shape of a triangle which may then be drawn in the plane, as shown in Figure 3. Each corner of this triangle corresponds to one of the basic types, i.e. either 100% continental, 100% maritime, or 100% urban. In terms of the proportions  $\alpha$ , each corner corresponds to one of the  $\alpha_i$ 's equal to one and the others equal to zero. To save up storage space and computation time, the radiative transfer calculations have been conducted for mixtures of two basic types, with proportions varying by step of 20%. This leads to 15 aerosol mixtures, as stated in Table 2, with proportions parameter  $\alpha$  lying on the edges of the triangle, i.e., on the boundary of the simplex  $\mathcal{S}^2$ . These quantization points correspond to the 15 points on the edges of the larger triangle in Figure 3.

In conclusion, the resulting parameter space  $\Theta_a$  for the atmospheric parameters is the product of five intervals of the real line, with bounds given by the extremal values in Table 2, times the 2-simplex  $\mathcal{S}^2$ . We shall denote by  $\mathcal{R}_a$  the product of these five intervals, to arrive at the definition of  $\Theta_a$  as  $\Theta_a = \mathcal{R}_a \times \mathcal{S}^2$ , and we shall write any  $\theta_a$  in  $\Theta_a$  as  $\theta_a = (\zeta_a, \alpha)$  with  $\zeta_a$  in  $\mathcal{R}_a$  and  $\alpha$  in  $\mathcal{S}^2$ . The set  $\Theta_a$  is quantized into  $N = 6,075$  points (see Table 2), and we end up with the database  $\{(\theta_{a,i}, \Phi_a(\theta_{a,i})) : i = 1, \dots, N\}$  where the  $\theta_{a,i}$ 's denote the  $N$  quantization points. Note that these computations are performed for the  $d$  wavelengths  $\lambda_1, \dots, \lambda_d$  of interest.

#### 4.1.2 Approximation of $\Phi_a$

Write  $\theta_a = (\zeta_a, \alpha)$  with  $\zeta_a$  in  $\mathcal{R}_a$  and  $\alpha$  in  $\mathcal{S}^2$  as above. Denote by  $\{x_{\zeta,i} : i = 1, \dots, N_1\}$  the quantization points of  $\mathcal{R}_a$ , with  $N_1 = 405$ , and by  $\{\alpha_j : j = 1, \dots, N_2\}$  the initial quantization points of  $\mathcal{S}^2$ , with  $N_2 = 15$ . To approximate the value of  $\Phi_a$  at  $(\zeta_a, \alpha)$ , we proceed by successive linear interpolations, first along  $\zeta_a$ , and next along  $\alpha$ . As a remark, proceeding the other way around (i.e., first along  $\alpha$  and next along  $\zeta_a$ ) is equivalent for linear interpolation. Note first that  $\mathcal{R}_a$  is a product of closed intervals, and that the quantization points of  $\mathcal{R}_a$  are located on a regular grid. So in a first step, the value of  $\Phi_a((\zeta_a, \alpha_j))$  are approximated by multilinear interpolation over  $\mathcal{R}_a$  for all  $j = 1, \dots, N_2$ . There remains to interpolate along the proportions  $\alpha$  to approximate  $\Phi_a((\zeta_a, \alpha))$ . For this purpose, we first form a regular triangular mesh of  $\mathcal{S}^2$  as illustrated in Figure 3. This mesh contains 21 points forming the vertices of 25 triangles. Each of these

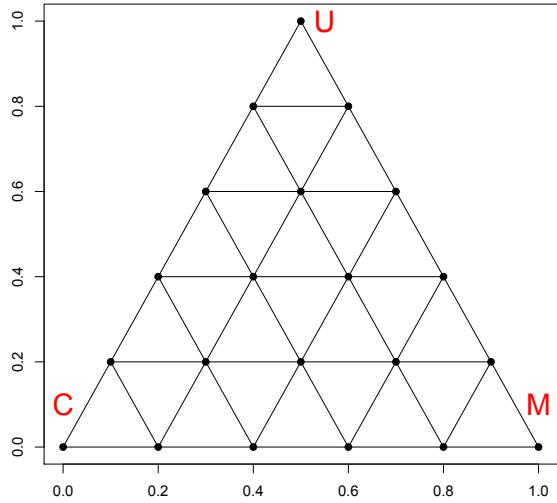


Figure 3: Parameter space for the proportions of the mixture of aerosol models.

21 points corresponds to a proportion parameter  $\alpha = (\alpha_1, \alpha_2, \alpha_3)$ , where each  $\alpha_i$  is an integer multiple of 20% ; recall that the initial computations have been performed for mixtures of two types of aerosol, i.e. for the 15 points on the boundary of the triangle of Figure 3 (for which at least one the  $\alpha_i$ 's is equal to 0). Successive linear interpolations between appropriate vertices of the mesh yields the approximated values of  $\Phi_a((\zeta_a, \alpha))$ , for any  $\alpha$  equal to one the 21 vertices of the mesh. Finally, given an arbitrary  $\alpha$  in  $S^2$ , there exists one triangle  $\mathcal{T}$  of the mesh containing  $\alpha$ , and the value of  $\Phi_a((\zeta_a, \alpha))$  is approximated by linear interpolation over  $\mathcal{T}$ , given the approximated values of  $\Phi_a$  at the three vertices of  $\mathcal{T}$  computed previously. Note that the triangle containing  $\alpha$  is unique, except if  $\alpha$  lies on the edges the mesh, in which case the triangle is chosen arbitrarily and the result of the interpolation does not depend on such a choice (the resulting interpolating function is continuous).

## 4.2 Water body reflectance

### 4.2.1 Definition of $\mathcal{X}_w$

Our goal here is to define an (approximately) uniform discretization of  $\mathcal{X}_w$ , that is, a finite collection of points of  $\mathcal{X}_w$  uniformly spread over  $\mathcal{X}_w$ . This could be achieved by considering several models relating the marine reflectance with biological parameters, but we find it more appropriate to base our analysis on in-situ measurement of the marine reflectance. As indicated in Section 2, modeling marine reflectance from measurements of inherent optical properties (e.g., [Garver and Siegel, 1997](#); [Gordon et al., 1988](#); [Morel and Maritorena, 2002](#); [Park and Ruddick, 2005](#)) does not fully take into account the natural correlations between controlling variables. For a given phytoplankton assemblage/population, absorption and backscattering coefficients are related in a unique way that is not captured in empirical bio-optical formulas determined from amalgamating measurements collected in diverse oceanic regions. To this aim our analysis is based on the NOMAD dataset ([Werdell and Bailey, 2005](#)), and on data acquired at several AERONET-OC sites ([Zibordi et al., 2010](#)). The NOMAD dataset includes coincident measurements of marine reflectance and chlorophyll concentration compiled by the NASA Ocean Biology Processing Group from tens of field experiments.

### Interpolation algorithm of $\Phi_a$

1. **Input:** The quantization points  $\{\zeta_{a,i} : i = 1, \dots, N_1\}$  and  $\{\alpha_j : j = 1, \dots, N_2\}$ , the corresponding values  $\Phi_a((\zeta_{a,i}, \alpha_j))$ , and the evaluation point  $\theta_a := (\zeta_a, \alpha)$ .
2. For all  $1 \leq j \leq N_2$ , approximate  $\Phi_a((\zeta_a, \alpha_j))$  by multilinear interpolation over  $\mathcal{R}_a$  using the values  $\Phi_a((\zeta_{a,i}, \alpha_j))$ , for  $i = 1, \dots, N_1$ , which gives the approximated values  $\hat{\Phi}_a((\zeta_a, \alpha_j))$ .
3. Form the regular triangular mesh of  $S^2$ , and approximate the value of  $\Phi_a$  by successive linear interpolation based on the values estimated at step 2.
4. Locate a triangle  $\mathcal{T}$  of the mesh containing  $\alpha$ , and linearly interpolate at  $\alpha$  using the values of step 3 located at the vertices of  $\mathcal{T}$ , yielding the approximated value  $\hat{\Phi}_a((\zeta_a, \alpha))$
5. **Output:** the approximated value  $\hat{\Phi}_a((\zeta_a, \alpha))$  of  $\Phi_a$  at  $\theta_a$ .

Table 3: Interpolation algorithm for the approximation of  $\Phi_a$ .

Above- and below-water radiometers are used in the measurement of marine reflectance. The major oceanic provinces are represented in NOMAD, but data are scarce in the southern Pacific and Indian Oceans, and in very clear waters. AERONET-OC provides marine reflectance and aerosol optical thickness through autonomous above-water radiometers operating on fixed platforms located near the coast. Situations of both Case 1 and Case 2 waters are sampled in NOMAD and AERONET-OC, but turbid waters containing sediments (e.g., from estuarine regions) are underrepresented in the ensembles.

#### 4.2.2 The NOMAD data set

The NOMAD data set contains above 4,000 measurements of the marine reflectance in 6 spectral bands from 412 nm to 670 nm (SeaWiFS spectral bands). Waters with chlorophyll concentrations in the range  $0.01 - 100 \text{ mg m}^{-3}$  are represented. Due to lack of ancillary information, the measurements made by above-water instruments, which constitute about 50% of the data set, were not corrected for anisotropy effects in the reflected light field, i.e., transformed to nadir values. The marine reflectance in the near infrared, i.e., at 765 nm and 865 nm, is not provided in NOMAD, but it is required for our analysis since we consider both Case I and Case II waters. While it is reasonable to assume a black ocean in the near infrared for Case I waters, this is not the case for Case II waters. In addition, the important variability in the observed Case II spectra makes it most difficult to reliably extrapolate the value of the marine reflectance at 865 nm, and even at 765 nm, based on measurements from 412 nm to 670 nm. Therefore from the NOMAD data set, the spectra identified as corresponding to Case I waters, and with at most one missing value, were extracted. This was accomplished using thresholds on the irradiance reflectance at 510 and 555 nm ([Morel and Belanger, 2006](#)). The number of extracted samples is equal to 2,651 and breaks down into 729 spectra without missing values, and 1,922 spectra with at most one missing value, as summarized in Table 4. The value of these spectra at 765 nm and 865 nm is set to 0 under the black ocean assumption in the near infrared. Each missing value is then inferred from the complete data. The procedure for inference of the missing values, described below, has been limited to only one missing value per spectrum. In spite of the apparent regularity of Case I marine reflectance spectra, inferring strictly more than one missing value could not be achieved with enough confidence for our analysis. Most of the missing values are located in the spectral band centered at 670 nm.

They can be estimated with reasonable confidence because of the relative regularity of the shape of Case I water spectra. That said, further inference, for more than one missing value, did not prove satisfactory for our analysis, with is the reason why Case I spectra with more than one missing value have been discarded from the NOMAD dataset.

Estimating the missing value in the marine reflectance spectra can be formulated as a regression problem, where the complete data is being used to fit a model predicting the value at the missing component. Numerous consistent regression estimation technique exist. Among these, we use  $k$ -nearest neighbour estimates. We do not comment further on these aspects at this point and we refer the interested reader to [Gyorfi et al. \(2002\)](#) for an exposition of the theory. Additional practical implementation details are gathered in Appendix A at the end of the paper. The estimation procedure is repeated for each of the 6 channels where one value is missing, and this results in a complete data set of 2,651 Case I marine reflectance spectra.

#### 4.2.3 The AERONET-OC data sets

Marine reflectance data from 9 AERONET-OC sites listed in Table 4, collected during 2002 to 2010, were extracted from the online web archive. The data, in the form of normalized water-leaving radiance in 8 spectral bands centered at 413, 440, 501, 530, 555, 674, 870, and 1019 nm, were transformed into marine reflectance and corrected for directional effects, i.e., converted to values at nadir (see [Zibordi et al., 2009](#)). Splines were then used to generate marine reflectance at the SeaWiFS wavelengths. The total number of marine reflectance spectra is equal to 12,397, but the sample size varies significantly from one site to another, ranging from 251 points (Wave CIS CSI 6 site) to 6,360 points (Venise site); see Table 4. In particular, about one half of the marine reflectance measurements come from the Venise site. The waters sampled have chlorophyll concentration in the approximate range  $0.1 - 20 \text{ mg/m}^3$ , which is much smaller than the NOMAD range. Many of the spectra exhibit values that correspond to Case II waters, with relatively high values at 560 and 510 nm, and non-zero values in the near infrared. Statistical differences also exist between the time series acquired at different locations. As an example, the estimated densities of the marine reflectance at 412 nm and 555 nm at each site are plotted in Figure 4, including those of the NOMAD data set. All the densities are essentially unimodal, but the location of the mode, as well as the spread, varies significantly with the geographical coordinates, i.e., water type or marine ecosystem. Narrow densities may be due, in some cases, to the limited amount of data at the stations. The larger spread is obtained with the NOMAD data set, presumably because it samples a wide range of chlorophyll concentrations and bio-geographic provinces.

#### 4.2.4 Assembling the NOMAD and AERONET-OC data sets

The measurements from the 9 AERONET-OC sites, combined with measurements extracted from the NOMAD data set, constitute a total of 15,048 in-situ marine reflectance spectra, sampled from 412 nm to 865 nm. Only 21% of the cases are from NOMAD. Because of sampling differences, the combined data points are not evenly distributed over the domain  $\mathcal{X}_w$ . Indeed, the extracted NOMAD data points are identified as corresponding to Case I waters, while the AERONET-OC data set covers both Case I and Case II waters. Moreover, from a statistical perspective, the NOMAD data come from separate field campaigns and can be considered as being roughly independent measurements. On the contrary, the AERONET-OC data come in the form of time series and are therefore naturally correlated. At last, the sample sizes of each group of data differ significantly. Note furthermore that the marine reflectance measurements were corrected for anisotropy in the light field, except the NOMAD above-water measurements (about 10% of the total ensemble). The forward model, on the other hand, assumes that the marine reflectance is isotropic (see Section 2). The solution of the inverse problem, therefore, will be a marine reflectance in the space defined by the NOMAD and AERONET-OC data sets (properly digitized, see below), i.e., a space of marine reflectance

Site/Dataset Name	Latitude(deg.)	Longitude(deg.)	Nb. Wavelengths	Nb. Measurements
NOMAD	Variable	Variable	5(6)	1722(729)
Venise	45.31	12.50	8	6360
Abu Al Bukhoosh	25.49	53.14	8	614
COVE SEAPRISM	36.90	-75.71	8	467
Gustav Dalen Tower	58.59	17.46	8	920
Helsinki Lighthouse	59.94	24.92	8	731
Lucinda	-18.51	146.38	8	303
MVCO	41.30	-70.55	8	2311
Palgrunden	58.75	13.15	8	440
Wave CIS site CSI 6	28.86	-90.48	8	251

Table 4: Summary of NOMAD and AERONET-OC datasets of water reflectance. For NOMAD, the number of complete sets (i.e., measurements at the 6 shorter SeaWiFS wavelengths) are given in parenthesis. The data ensemble is composed of 21% cases from NOMAD (Case I waters only) and 79% from AERONET-OC (Case I and Case II waters). The Venise and MVCO sites, located in the Adriatic Sea and the Mid-Atlantic Bight, contribute most of the AERONET-OC data, i.e., 50 and 19%, respectively.

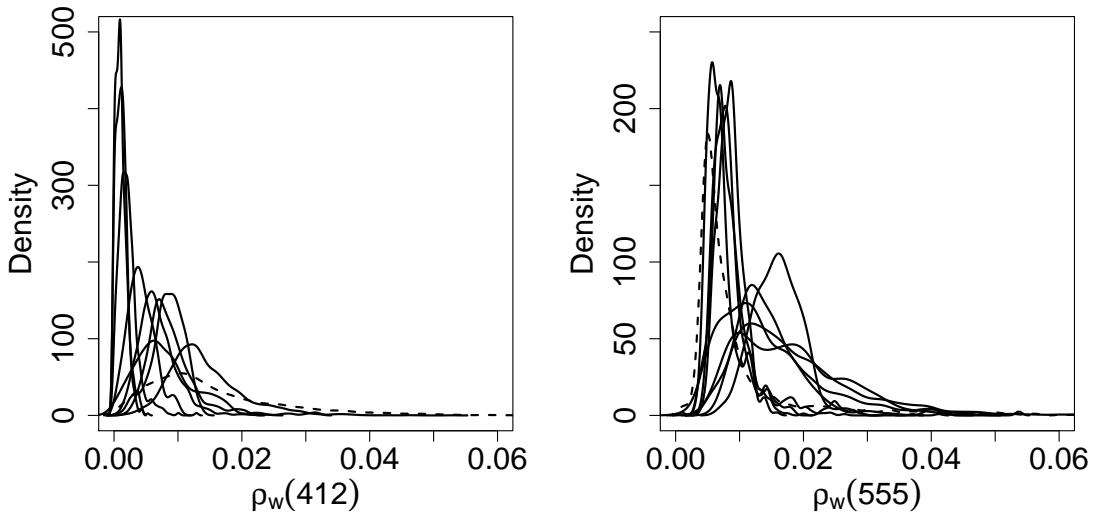


Figure 4: Estimated densities of the marine reflectance at 412 and 555 nm, left and right, respectively, acquired at each of the 9 AERONET sites (solid lines) and assembled in NOMAD (dashed line).

corrected for directional effects, under the assumption of isotropic marine reflectance. In actuality, i.e., using satellite observations, directional effects in the marine reflectance will introduce some errors in the retrievals.

To form a finite collection of marine spectra approximately uniformly distributed over  $\mathcal{X}_w$  we consider a standard accept/reject algorithm. More precisely, we first estimate  $\mathcal{X}_w$  by a union of balls centered at the in-situ data and of radius, say  $r$ . Next, we simulate an outcome of a random vector uniformly distributed over a rectangular domain containing all the balls. The simulated point is accepted if it belongs to one of the ball, and rejected otherwise, i.e., the simulated point is kept if it is at a distance no more than  $r$  from one of the in-situ data, and discarded otherwise. This procedure is repeated until a number, fixed in advanced, of simulated points have been accepted. The in-situ data being fixed, this algorithm yields a set of simulated points uniformly distributed over the union of balls. In turn, it may be shown that this union of balls converges to the unknown set  $\mathcal{X}_w$  under mild assumptions as the number of balls (here, the number of in-situ  $\rho_w$ ) goes to infinity ; see, e.g., [Biau et al. \(2008\)](#).

## 5 Implementation of the inverse models

### 5.1 Noise distribution

Let us recall model (2.7):

$$y = \Phi(x_a, x_w) + \varepsilon.$$

The additive term  $\varepsilon$  appearing in this model, referred to simply as the noise term, is intended to encapsulate all the sources of uncertainties in the forward modeling leading to (2.5). These include, in particular, measurement uncertainties, modeling uncertainties of the various radiative transfer processes, as well as the approximation error associated with the reconstruction of the forward operator by interpolation of the discrete (simulated) data. As stated in Section 2, we shall assume that  $\varepsilon$  is a gaussian random vector with  $\mathbb{E}[\varepsilon] = 0$  and  $\text{Cov}(\varepsilon) = \sigma^2 \text{Id}_{\mathbb{R}^d}$ . In this model, the parameter  $\sigma$  plays the role of the noise level in a deterministic setting.

Naturally this assumption is reductive. Indeed, modeling errors on the forward operator are likely to have a non-zero average (i.e., yielding a bias) and/or depend in magnitude on the input parameters (atmosphere and ocean states). Also, in a nonparametric setting, i.e. without assumption on the distribution of  $\varepsilon$ , this latter cannot be inferred since  $\varepsilon$  is never observed (and nor are the variables  $x_a$ ,  $x_w$ , and  $y$  measured simultaneously). It is therefore reasonable to impose a parametric assumption on the distribution of  $\varepsilon$ .

To determine a value of  $\sigma$ , we first extract  $N$  TOA reflectance  $y_1, \dots, y_N$  from various images acquired by the SeaWiFS sensor. Next, for each  $i = 1, \dots, N$ , let

$$\delta_i^2 = \inf \left\{ \|y_i - \Phi(x_a, x_w)\|^2 : (\rho_a, T_a, S_a) \in \mathcal{X}_a, \rho_w \in \mathcal{X}_w \right\},$$

i.e.,  $\delta_i$  is the distance from  $y_i$  to  $\Phi(\mathcal{X}_a, \mathcal{X}_w)$ . Then, we define a value  $\hat{\sigma}^2$  of  $\sigma^2$  by averaging the  $\delta_i^2$ 's, i.e., we set  $\hat{\sigma}^2 = (1/N) \sum_{i=1}^N \delta_i^2$ .

### 5.2 Prior distributions

By allowing to quantify the frequencies of occurrence of the parameters of interest before the acquisition of the data to be inverted, a prior distribution represents an a-priori information. When no such information is available to the user, one is frequently led to consider a uniform distribution over an appropriate set. In this study, we aim at implementing an atmospheric correction algorithm valid at a global scale. For this purpose, a suitable prior distribution would have to reflect the frequencies of occurrence for parameters like the marine reflectance, or the aerosol optical thickness, at a global scale, global being meant both in space

Parameter	Prior distribution
Pressure	Uniform over [1003, 1023]
Aerosol Scale Height	Uniform over [1.0, 3.0]
Aerosol model proportions	Uniform over the simplex $\mathcal{S}^2$ .
Aerosol Optical Thickness	Log-normal
Wind Speed	Uniform over [1, 12]
Marine Reflectance	Uniform over $\mathcal{X}_w$

Table 5: Prior distributions on the various parameters.

and time. In addition, since the prior distribution is one of the elements defining the Bayesian solution to the inverse problem, any information used to specify it must not originate from inversions of satellite observations, but instead from separate, independent field campaigns. As stated in Section 2, since there is no reason to expect, a-priori, that atmospheric and marine parameters be correlated, the prior distribution on the parameter space  $\mathcal{X}_a \times \mathcal{X}_w$  is of the product type  $\mathbb{P}_a \otimes \mathbb{P}_w$ .

The set  $\mathcal{X}_w$  has been defined in 4.2.1, and in the absence of other information,  $\mathbb{P}_w$  is taken as the uniform measure on  $\mathcal{X}_w$ . Recall that, in practice, we work with an (approximately) uniform discretization of  $\mathcal{X}_w$  in, say  $n$ , points  $x_{w,1}, \dots, x_{w,n}$ . So, in practice, we have  $\mathbb{P}_w = \frac{1}{n} \sum_{i=1}^n \delta_{x_{w,i}}$  and sampling from  $\mathbb{P}_w$  simply amounts at selecting one of the  $x_{w,i}$  uniformly at random and with replacement. Note that  $\mathbb{P}_w$  converges weakly to the uniform measure  $\mathcal{U}(\mathcal{X}_w)$  on  $\mathcal{X}_w$  as  $n$  goes to infinity ; this means that, whenever  $n$  is large, the two measures  $\mathbb{P}_w$  and  $\mathcal{U}(\mathcal{X}_w)$  are close in some sense.

Using the notations from Section 4.1, we have  $\mathcal{X}_a = \Phi_a(\Theta_a)$ , where  $\Theta_a = \mathcal{R}_a \times \mathcal{S}^2$ , with  $\mathcal{R}_a$  being the joint parameter space for the pressure, wind-speed, aerosol scale height and aerosol optical thickness, and  $\mathcal{S}^2$ , the 2-simplex, being the parameter space for the aerosol model. The study by Knobelispiesse et al. (2004) shows, from the analysis of 145 cruises in the Atlantic and Pacific oceans, and Asian seas (about 11,000 individual data points), that the aerosol optical thickness over the oceans is approximately distributed at the global scale as a log-normal distribution. Hence we shall use this distribution as a prior distribution on  $\tau_a$  with parameters taken as  $\mu = -2.5257$  and  $\sigma = 0.9854$  and extracted from Knobelispiesse et al. (2004). With these values for  $\mu$  and  $\sigma$ , the 95% quantile is equal to approximately 0.40, which means that if one simulates values of  $\tau_a$  from this prior distribution, then on average 95% of them will have a value lower than 0.40. For all the other parameters, to the best of our knowledge, no studies that are not based on satellite observations provide information on their relative frequencies at a global scale. Consequently, a uniform prior distribution of the domains of these parameters is retained. Formally, these considerations leads to defining a probability measure, say  $\mathbb{P}_{a,0}$ , on  $\mathcal{R}_a \times \mathcal{S}^2$  as the product measure of several uniform distributions and a log-normal distribution. Then, this yields the prior distribution  $\mathbb{P}_a$  as the image of  $\mathbb{P}_{a,0}$  through  $\Phi_a$ , i.e., we have  $\mathbb{P}_a = \mathbb{P}_{a,0} \circ \Phi_a^{-1}$ .

The prior distributions are summarized in Table 5. Let us point out that, for the marine reflectance, considering a uniform prior distribution on  $\mathcal{X}_w$  already reflects a prior information, namely that values of the marine reflectance not belonging to  $\mathcal{X}_w$  are not realistic. This has been possible by using the in-situ data from the NOMAD archive and the various AERONET-OC sites. Naturally, since at a global scale most of the oceans are Case I waters, while only a small proportion corresponds to more optically complex waters, it may seem desirable to favor the first type of waters or, in other words, that the prior distribution on  $\rho_w$  places more weight on Case I  $\rho_w$  than others. How to do so in an objective way is a non-trivial matter at a global scale, for this would require extensive field campaigns to estimate reliably the frequencies of the marine reflectance. On the other hand, this may certainly be envisioned at a regional scale, a perspective which is discussed at the end of this paper.

### 5.3 Inverse applications

The inverse models introduced in Section 2 for atmospheric correction are all based on the same partition  $A_1, \dots, A_M$  of  $\mathbb{R}^d$  and are defined by:

$$\hat{r}(y^{\text{obs}}) = \sum_{m=1}^M (\alpha_m + B_m y^{\text{obs}}) \mathbf{1}_{A_m}(y^{\text{obs}}), \quad (5.1)$$

$$\hat{\Sigma}(y^{\text{obs}}) = \sum_{m=1}^M C_m \mathbf{1}_{A_m}(y^{\text{obs}}), \quad (5.2)$$

$$p_V(y^{\text{obs}}) = \sum_{m=1}^M p_m \mathbf{1}\{\hat{f}_m \leq \hat{f}_y(y^{\text{obs}})\}, \quad (5.3)$$

where, for all  $y \in \mathbb{R}^d$ ,

$$\hat{f}_y(y) = \sum_{m=1}^M \hat{f}_m \mathbf{1}_{A_m}(y) \quad \text{with} \quad \hat{f}_m = \frac{p_m}{\text{Vol}_d(A_m)} \quad \text{and} \quad p_m = \mathbb{P}(y \in A_m). \quad (5.4)$$

The model (5.1) is an approximation of the conditional expectation  $\mathbb{E}[x_w|y = y^{\text{obs}}]$  given the observation  $y^{\text{obs}}$ ; it provides an estimation of the retrieved marine reflectance. The model (5.2) is an approximation of the conditional covariance  $\text{Cov}(x_w|y = y^{\text{obs}})$  given  $y^{\text{obs}}$ ; it provides a measure of uncertainty on the retrieved marine reflectance. The model (5.3) is an approximation of the probability that a new observation be at least as extreme as  $y^{\text{obs}}$ , in the sense of the definition given in Section 2; it can be used to detect cases where the forward modeling of the satellite signal and the data  $y^{\text{obs}}$  significantly differ to the extent that the results of the inversion are meaningless.

In addition, and as a by-product, we also construct similar inverse models to retrieve the following atmospheric parameters: the term  $\rho_a$  in (2.7), the aerosol optical thickness, the aerosol model type (i.e., the proportions of the mixture), with accompanying uncertainties. Since retrieving these parameters is not the primary focus of this study, we use a simpler model than (2.1) where the linear model in each cell is replaced by a constant, i.e., each model is of the following generic form

$$\hat{r}_a(y^{\text{obs}}) = \sum_{m=1}^M \alpha_m \mathbf{1}_{A_m}(y^{\text{obs}}), \quad (5.5)$$

with coefficients  $\alpha_1, \dots, \alpha_m$  in  $\mathbb{R}^p$ , where  $p$  denotes the dimension of  $x_a$ . The models for the conditional covariances are of the same type as (5.2).

#### 5.3.1 Construction of the partition

Common to all the inverse applications defined above is a partition  $A_1, \dots, A_M$  of  $\mathbb{R}^d$ . To invert a reflectance  $y^{\text{obs}}$ , the first operation to be performed is to determine which one of the  $A_m$ 's contains  $y^{\text{obs}}$ . The computational cost to determining cell membership of an arbitrary partition can be very large, but it may be significantly reduced when the partition has an appropriate structure. To keep the execution time of the inverse applications low, we consider in this work a partition based on a perfect binary tree.

A *tree* is a hierarchical structure formed by a collection of linked *nodes* together with associated rules. The top node is called the *root* of the tree. Each node may have several *children nodes* and at most one *parent node*. Nodes without children nodes are called *leaves*, and the *depth* of a node is the length of the

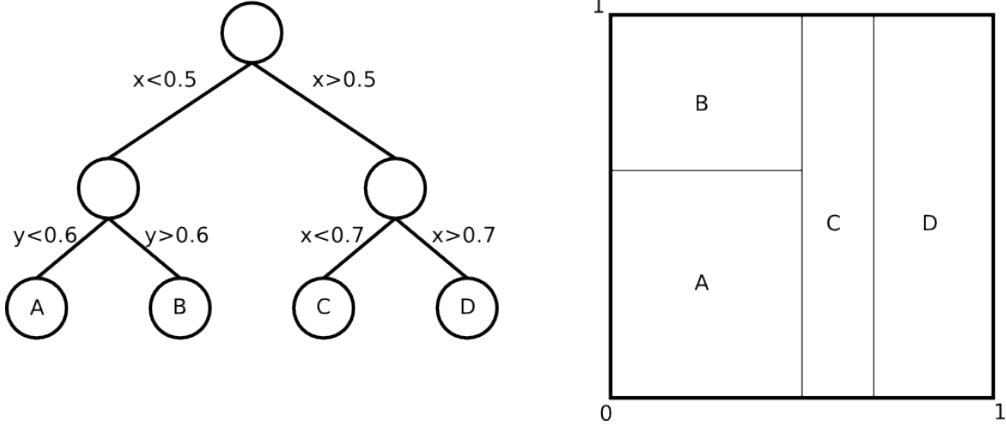


Figure 5: Illustration of a tree-structured partition of the unit square of  $\mathbb{R}^2$  based on a perfect binary tree with depth 2.

path from the root to the node. In a binary tree, each node has at most two children nodes. A binary tree is called *perfect* if every node other than the leaves has exactly two child nodes, and if all leaves are at the same depth. In this case, the two children of an internal nodes are commonly called the left and right children. The set of rules associated with the nodes other than the leaves induce a sequence defining paths from the root to the leaves. For our purposes, we consider binary comparison rules. Each rule is described by a pair  $(j, \delta)$ , where  $j$  is an integer and where  $\delta$  is a real number. To see how this effectively defines a partition, consider a vector  $y$  of  $\mathbb{R}^d$ , and denote by  $y_1, \dots, y_d$  its components. Starting from the root node, with rule  $(j, \delta)$ , the  $j^{\text{th}}$  component of  $y$  is compared with the threshold  $\delta$ . If  $y_j$  is smaller than  $\delta$ , then  $y$  is moved to the left child, otherwise,  $y$  is moved to the right child. The procedure is then repeated, either on the left child or on the right child, until the point  $y$  reaches a leaf. Hence this type of tree defines an axis-parallel partition of  $\mathbb{R}^d$  into hyper-rectangles: each cell of the induced partition is composed of the points  $y$  whose path in the tree attains the same leaf. So the resulting partition contains as many cells as leaves in the tree. In the case of a perfect binary tree, with depth  $K$ , the number of leaves is equal to  $M = 2^K$ . Note that, for an arbitrary point  $y$  of  $\mathbb{R}^d$ , the number of operations to perform to determine the cell to which  $y$  belongs is equal to  $K$  comparisons. We see then the strong computational interest in a tree-structured partition: the number of cells grows exponentially with the depth  $K$ , while the number of operations needed to determine cell memberships increases linearly with  $K$ . In addition, comparing two real numbers is computationally extremely fast. An illustrative example of a tree-structured partition of the plane  $\mathbb{R}^2$  is given in Figure 5.

In this work, we chose a depth  $K = 17$  which yields  $M = 2^K = 131,072$  cells. There remains to determine a set of splitting rules in such a way that convergence of the models is guaranteed. To this aim, we employ a fully data-driven procedure, i.e., the rules are determined from a large number of simulated data  $y$  in model (2.7). Details of the algorithm are postponed to Appendix B at the end of the paper. For additional materials on this subject, we refer the reader to Breiman et al. (1984) and Breiman (2001) as well as Lugosi and Nobel (1999); Nobel (1996).

### 5.3.2 Approximation of the model coefficients

From now on, the partition  $A_1, \dots, A_m$  is fixed and constructed as exposed above. The optimal values for the coefficients  $\alpha_m \in \mathbb{R}^d$ ,  $B_m \in \mathcal{M}_d(\mathbb{R})$ ,  $C_m \in \mathcal{M}_d(\mathbb{R})$  (a covariance matrix), and  $p_m \in \mathbb{R}$ , for  $m = 1, \dots, M$ , for the inverse applications (5.1), (5.2), and (5.3) are such that the  $L^2$  risk is being minimized. Simple calculations allow to derive their analytical expressions in terms of the joint distribution of

$(x_a, x_w, y)$ , which is known. In practice, their values are approximated numerically as follows.

Let  $(x_{a,1}, x_{w,1}, y_1), \dots, (x_{a,n}, x_{w,n}, y_n)$  be  $n$  data simulated according to model (2.7). Next, the data is split among each element of the partition. In the following, for all  $m = 1, \dots, M$ , we set  $\mathcal{I}_m = \{i = 1, \dots, n : y_i \in A_m\}$ , and the number of points in  $\mathcal{I}_m$  is denoted by  $\#\mathcal{I}_m$ .

**Marine reflectance models** For each cell  $A_m$ , the optimal coefficients  $\hat{\alpha}_m$  and  $\hat{B}_m$  are approximated by minimizing the empirical risk  $\sum_{i \in \mathcal{I}_m} \|x_{w,i} - \alpha - By_i\|^2$  over  $\alpha \in \mathbb{R}^d$  and  $B \in \mathcal{M}_d(\mathbb{R})$ . The solution (not displayed here) can be expressed analytically ; in practice, the optimal values are computed by fitting a linear model to the data  $(x_{a,i}, x_{w,i}, y_i)$  with  $i \in \mathcal{I}_m$ , for all  $m = 1, \dots, M$ . For each  $1 \leq m \leq M$ , the covariance matrix  $\hat{C}_m$  is approximated by the empirical covariance, i.e., by  $(\#\mathcal{I}_m)^{-1} \sum_{i \in \mathcal{I}_m} (x_{w,i} - \hat{\alpha}_m - \hat{B}_m y_i)(x_{w,i} - \hat{\alpha}_m - \hat{B}_m y_i)^t$ .

**Atmospheric parameters models** For the retrieval of the atmospheric parameters, the models defined in (5.5) are simpler: these are constant functions on each  $A_m$ ,  $m = 1, \dots, M$ . In this case, the optimal value  $\hat{\alpha}_m$  for  $\alpha_m$  is  $\mathbb{E}[x_a | y \in A_m]$  which can be approximated numerically by the average of the  $x_{a,i}$ 's with corresponding  $y_i$ 's belonging to  $A_m$ , i.e., by  $(\#\mathcal{I}_m)^{-1} \sum_{i \in \mathcal{I}_m} x_{a,i}$ . As for the marine reflectance models, the optimal covariance matrix  $\hat{C}_m$  is approximated by the empirical covariance  $(\#\mathcal{I}_m)^{-1} \sum_{i \in \mathcal{I}_m} (x_{a,i} - \hat{\alpha}_m)(x_{a,i} - \hat{\alpha}_m)^t$ .

**p-value model** Two type of quantities have to be calculated to implement the model giving the  $p$ -value for the observation  $y^{\text{obs}}$  to be inverted, namely the volume  $\text{Vol}_d(A_m)$  of the cell  $A_m$ , and the probability  $p_m = \mathbb{P}(y \in A_m)$  that a new observation  $y$  falls in  $A_m$ , for each cell  $A_m$ ,  $m = 1, \dots, M$ .

The volumes of the cells are computed at the time the partition is built. Note that, for the tree-structured partition, each cell  $A_m$  is the product of  $d$  intervals in  $\mathbb{R}^d$ . In theory, the cells can be unbounded. To cope with this issue, in practice, we first infer a hyper-rectangle  $\mathcal{B}$ , with sides parallel to the axes of  $\mathbb{R}^d$  which contains all the data points. Then we work with the intersections  $A_m \cap \mathcal{B}$  in place of the cells  $A_m$  for volume computations. This does not change the theoretical properties of models (5.3) and (5.4) so this is not made precise in the definitions of these models for clarity.

At last, the probability that a new observation falls in the  $m^{\text{th}}$  cell  $A_m$  is approximated by the average number of points falling in  $A_m$ , i.e., by  $\#\mathcal{I}_m/n$ . Note that, by construction of the partition, these numbers are almost all equal to  $1/(2^M)$ , i.e., each cell have approximately equal probability by construction (see Appendix B).

## 6 Evaluation on simulated data

### 6.1 About mean squared error

The quality of the inversion is measured by a quadratic criterion, defined as the average of squared errors of the retrievals, where the notion of average is relative to some distribution. So suppose that  $(x_a, x_w, y)$  has distribution  $\mathbb{Q}$ . The mean squared error (MSE)  $\mathcal{E}_{\mathbb{Q}}(\hat{r})$  associated with the retrieval of  $x_w$  from  $y$  by the model  $\hat{r}$  is defined by

$$\mathcal{E}_{\mathbb{Q}}(\hat{r}) = \mathbb{E}_{\mathbb{Q}} [\|x_w - \hat{r}(y)\|^2].$$

The subscript  $\mathbb{Q}$  in the equation above indicates that the expectation is with respect to  $\mathbb{Q}$  ; it is important to note that changing  $\mathbb{Q}$  typically changes the value of the criterion. Denote by  $r_{\mathbb{Q}}$  the regression function of  $x_w$  on  $y$ , i.e.,  $r_{\mathbb{Q}}(\cdot) = \mathbb{E}_{\mathbb{Q}}[x_w | y = \cdot]$ . Then, the MSE of any estimator  $\hat{r}$  decomposes into the sum

$$\mathcal{E}_{\mathbb{Q}}(\hat{r}) = \mathbb{E}_{\mathbb{Q}} [\|\hat{r}(y) - r_{\mathbb{Q}}(y)\|^2] + \mathbb{E}_{\mathbb{Q}} [\|r_{\mathbb{Q}}(y) - x_w\|^2]. \quad (6.1)$$

The first term in the right-hand side of (6.1) is the *approximation error*: its value reflects the error of approximating  $r_{\mathbb{Q}}$  by the model  $\hat{r}$ . The second term in the right-hand side of (6.1) does not depend on the model  $\hat{r}$ , but only on the distribution  $\mathbb{Q}$ . This term is therefore a lower bound on the MSE of the retrievals, under the assumption that the data has distribution  $\mathbb{Q}$ , i.e., for any estimator  $\hat{r}$  of the marine reflectance, we have

$$\mathcal{E}_{\mathbb{Q}}(\hat{r}) \geq \mathcal{E}_{\mathbb{Q}}(r_{\mathbb{Q}}) = \mathbb{E}_{\mathbb{Q}} [\|r_{\mathbb{Q}}(y) - x_w\|^2].$$

The minimal value of the MSE,  $\mathcal{E}_{\mathbb{Q}}(r_{\mathbb{Q}})$ , is typically different from 0. This is because, even without noise, the forward operator is not invertible. Still,  $\mathcal{E}_{\mathbb{Q}}(r_{\mathbb{Q}})$  depends non trivially on  $\mathbb{Q}$ . For instance, it is well known that atmospheric correction is better conditioned a problem in situations of low rather than large values of the aerosol optical thickness (AOT). Hence if  $\mathbb{Q}$  favors low value of AOT with respect to large values of AOT, the resulting MSE would tend to be lower than if  $\mathbb{Q}$  favored large versus low values of AOT. To the limit, if  $\mathbb{Q}$  is a distribution which charges one atmospheric state, say  $x_a^*$  (i.e.,  $x_a$  differs from  $x_a^*$  with probability 0), then, neglecting  $S_a$  in (2.7), we would have  $\mathcal{E}_{\mathbb{Q}}(r_{\mathbb{Q}})$  approximately equal to  $\mathbb{E}_{\mathbb{Q}}[\|T_a^{-1}\varepsilon\|^2]$ , i.e., the total variance of the noise modulated by the atmospheric transmittance  $T_a$ .

When the distribution  $\mathbb{Q}$  used to define the risk is taken equal to the joint distribution of  $(x_a, x_w, y)$  in model (2.7) with prior distributions and noise distribution as defined in Section 2, the MSE will be simply denoted by  $\mathcal{E}(\hat{r})$ , and the subscript  $\mathbb{Q}$  will be dropped from expectations. As above, we have the decomposition

$$\mathcal{E}(\hat{r}) = \mathbb{E} [\|\hat{r}(y) - r(y)\|^2] + \mathbb{E} [\|r(y) - x_w\|^2].$$

The approximation error  $\mathbb{E} [\|\hat{r}(y) - r(y)\|^2]$  is only due to the model  $\hat{r}$  and is expected to be low compared to the second term  $\mathbb{E} [\|r(y) - x_w\|^2]$ , since the number of cells in the partition of  $\mathbb{R}^d$  has been taken large. Also, we have  $\mathcal{E}(\hat{r}) = \sum_{i=1}^d \mathbb{E}[\|\hat{r}_i(y) - x_{w,i}\|^2]$ , and for all  $i = 1, \dots, d$ , each term can be decomposed as the sum of a squared bias and a variance, i.e., we have

$$\mathbb{E}[\|\hat{r}_i(y) - x_{w,i}\|^2] = (\mathbb{E}[\hat{r}_i(y) - x_{w,i}])^2 + \text{Var}(\hat{r}_i(y) - x_{w,i}).$$

But since  $r(.) = \mathbb{E}[x_w | y = .]$  by definition, we have  $\mathbb{E}[r(y)] = \mathbb{E}[x_w]$ , and since  $\hat{r}$  is close to  $r$ , the biases  $\mathbb{E}[\hat{r}_i(y) - x_{w,i}]$  are expected to be very low for all  $i = 1, \dots, d$ . Therefore, the MSE  $\mathcal{E}(\hat{r})$  is approximately equal to the sum of the variances  $\text{Var}(\hat{r}_i(y) - x_{w,i})$ , which in turn are close to the  $\text{Var}(r_i(y) - x_{w,i})$ 's. At last, note that  $\sum_{i=1}^d \text{Var}(r_i(y) - x_{w,i}) = \mathbb{E} [\|r(y) - x_w\|^2]$ .

## 6.2 Performance statistics

To quantify performance, we use the distribution  $\mathbb{Q}$  induced by the prior distributions and the noise distribution, and we consider the component-wise (i.e. per channel) biases, standard deviations, and mean squared errors defined respectively by

$$b_i = \mathbb{E}[\hat{r}_i(y) - x_{w,i}], \quad \sigma_i = \text{Var}(\hat{r}_i(y) - x_{w,i})^{\frac{1}{2}}, \text{ and } \mathcal{E}_i = \mathbb{E} [\|\hat{r}_i(y) - x_{w,i}\|^2], \quad (6.2)$$

for all  $i = 1, \dots, d$ . The dependence of these quantities on  $\hat{r}$  is omitted in the notations. In what follows, the component-wise root mean squared errors (RMSE) refer to the  $\mathcal{E}_i^{1/2}$ . Also, note that, as exposed above,  $\mathcal{E}_i = b_i^2 + \sigma_i^2$ .

These quantities can be evaluated, theoretically, from the joint distribution of  $(x_a, x_w, y)$  since this latter is known. In practice, their values are approximated numerically by simulating a large number  $n$  of data according to model (2.7) and by replacing expectations with empirical averages in (6.2).

This procedure has been accomplished for each observation geometry configuration in the discretization of the angular domain (see the angular grid in Table 4) with a number  $n$  of 1 million points. We shall also consider the same statistics (i) by bins of aerosol optical thickness at 865 nm, and (ii) by bins of

aerosol model proportions (in the 2-simplex). Formally, this amounts at replacing the expectations in 6.2 by conditional expectations given the event that the aerosol optical thickness [for (i)] or the aerosol model [for (ii)] belongs to a given bin. As before, these quantities are approximated numerically by the same expressions, but from the part of the simulated data falling in the bin under consideration.

For the aerosol optical thickness, the boundary values of the bins have been taken as the quantiles of the prior distribution of orders an integer multiple of 5% (see Table 6), i.e., each bin is of equal probability for the prior log-normal distribution. Hence on average, each bin contains about 50,000 points. For the aerosol proportions, the 2-simplex is partitioned into 66 regular equilateral triangles of equal area, yielding 66 bins of equal probability according to the prior distribution, which we recall is the uniform distribution on the simplex ; the resulting grid is a refined version of the one displayed in Figure 3. The corners of each triangle defining a bin correspond to a mixture of the three basic aerosol types with mixture coefficients an integer multiple of 10%.

Table 6: Boundary values of the bins in the aerosol optical thickness at 865nm. Each bin is of approximate equal probability according to the log-normal prior distribution. The total number of bin is 20.

0.0000	0.0158	0.0226	0.0288	0.0349	0.0412	0.0477
0.0547	0.0623	0.0707	0.0800	0.0905	0.1027	0.1170
0.1341	0.1555	0.1833	0.2222	0.2828	0.4046	0.6000

### 6.2.1 Average errors

In an attempt to provide a single measure of performance for the whole inversions, we defined a criterion by averaging biases and standard deviations over all the observation geometries ; see the angular grid points in Table 4. Note that the grid points are Gauss angles, so they are not uniformly distributed. Also, the error statistics depend on Sun and view directions ; some observation geometries are more favorable than others. Still, the sampling grid is finely defined, with well distributed grid points, so that the averaged statistics are a good measure of global accuracy.

Table 7: Geometry-averaged statistics: bias and standard deviation per channel, averaged over all the observation geometries.

Wavelength (nm)	412	443	490	510	555	670
Average Bias	1.81E-09	-8.06E-10	3.48E-09	3.06E-10	-1.08E-08	-4.66E-09
Average Standard Deviation	0.004321	0.003564	0.003220	0.002936	0.002652	0.001145

Table 7 displays the bias and standard deviation per spectral band, averaged over all the geometries. The retrievals of the marine reflectance ( $\rho_w$ ) are globally unbiased as expected ; see the theoretical justification in Section 6.1. The average standard deviations vary from about 0.004 at 412 nm to 0.001 at 670 nm. This is acceptable, in view of the  $\rho_w$  mode values in the NOMAD and AERONET-OC data sets (see Figure 4). The standard deviations would be smaller if unfavorable geometries (Sun glint, large air mass) were not included in the averages.

### 6.2.2 Errors per observation geometry.

Figure 6 and Figure 7 display histograms of the average biases and standard deviations over the geophysical conditions of the simulated data set (each error composing the histograms corresponds to an individual geometry). It can be noticed that bias and standard deviation depend on the angles. The distribution of the biases appears to be centered and unimodal, with 95% of the values having a magnitude less than about  $10^{-5}$  ( $< 3 \cdot 10^{-6}$  at 670 nm), i.e., the biases are small (Figure 6). For the standard deviations, 95% of the values are below 0.007, 0.006, 0.006, 0.005, 0.005, and 0.002 at 412, 443, 490, 510, 555, and 670 nm, respectively (Figure 7). Inversion errors, however, depend on the atmospheric parameters, as shown later in the analysis.

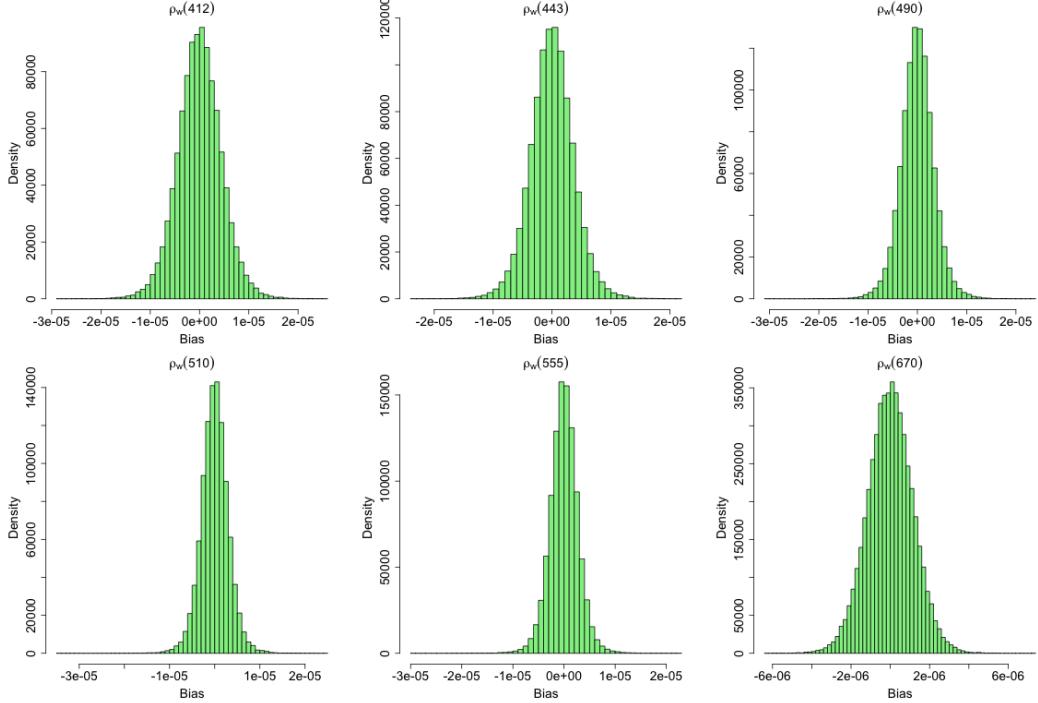


Figure 6: Global statistics: histograms of the bias per wavelength.

Figures 8, 9, and 10 display the standard deviations for each angular geometry at selected wavelengths, i.e., 412, 555, and 670 nm. Biases are not shown since close to 0 (see Section 6.2), and they do not exhibit a noticeable dependence with Sun and view angles. Consequently, the RMSE is practically equal to the standard deviation, and dependence with geometry is almost the same. Each panel correspond to a  $\theta_s$  value of the angular grid (see Table 4). For each fixed value of  $\theta_s$ , the (relative) viewing configuration is parameterized by a point in the upper half unit-disc of the plane. In polar coordinates, the radius corresponds to  $\sin \theta_v$ , while the angle corresponds to  $\cos \Delta\phi$ , with  $\Delta\phi$  expressed in the convention of the successive-orders-of-scattering code, i.e.,  $\Delta\phi = 0^\circ$  for forward scattering and  $\Delta\phi = 180^\circ$  for backscattering.

Error patterns are similar at all wavelengths, but values are higher at shorter wavelengths. Outside the Sun glint region, standard deviation generally increases with air mass (larger  $\rho_a$  signal to correct). Small values are observed at large scattering angles ( $> 120$  degrees), where the aerosol phase function, i.e., the influence of aerosols, is relatively small. Inside the Sun glint region, performance is degraded, but remains reasonable (e.g.,  $< 0.005$  at 412 nm) in many situations. At high Sun zenith angles (top three panels), the effect of Sun glint is less apparent, even absent, since it occurs at large view zenith angles, i.e., large air mass, for which the inversion is less accurate. Around the Sun glint directions, especially visible in the

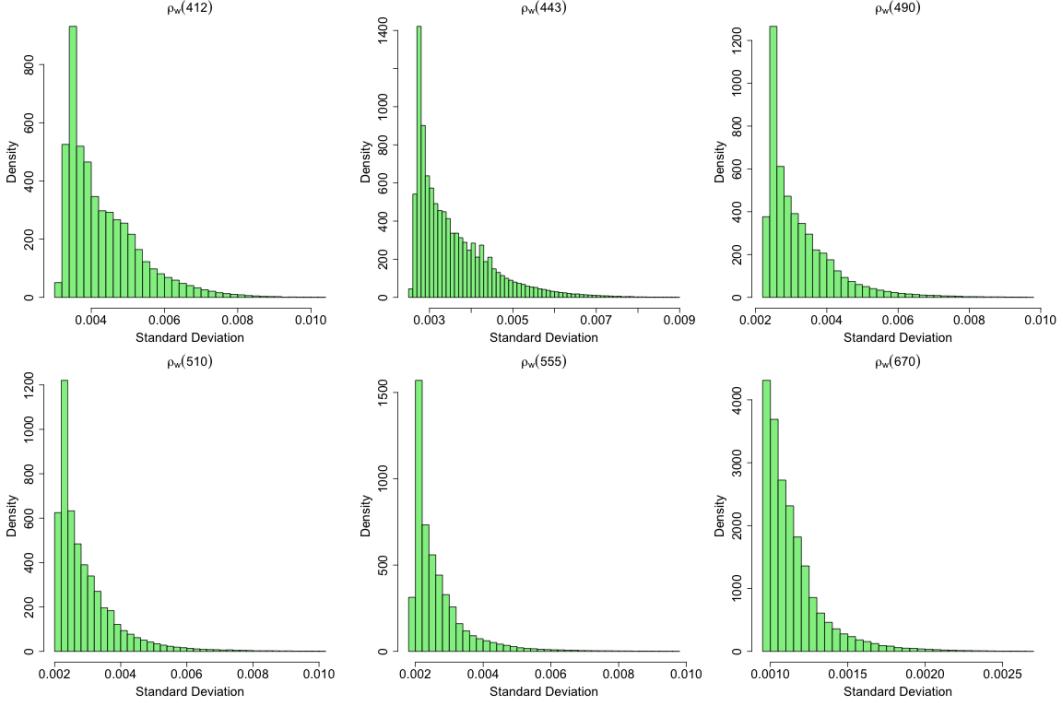


Figure 7: Global statistics: histograms of the standard deviation per wavelength.

three lower panels, the standard deviation exhibits some of the smallest values. This is explained by the combined effect of a large scattering angle in and around the Sun glint region and a relatively low air mass. The contrast between errors in the Sun glint and adjacent regions is sharp, because the presence of Sun glint is very sensitive to angular geometry (small changes in viewing direction may be associated with glint or no glint).

### 6.2.3 Errors per atmospheric parameter

**Errors per aerosol optical thickness.** The bias and standard deviation, averaged over all geometries, are displayed in Figures 11 and 12 for each spectral band as a function of aerosol optical thickness bin at 865 nm (defined in Table 6). The box plots show that the bias is generally less than 0.0005 in magnitude, almost negligible at 670 nm. For large  $\tau_a$  values (i.e., 0.4-0.6 bin), median absolute values reach 0.001 in the blue to yellow. The smallest biases are obtained for  $\tau_a$  values of 0.05 to 0.15 (typically encountered in the open ocean). Below and above this value, biases are negative and positive, respectively. Standard deviation is generally higher at the shorter wavelengths and increases fairly linearly with aerosol optical thickness. The median value in the 0.4-0.6  $\tau_a$  bin reaches approximately 0.009, 0.006, and 0.002 at 412, 555, and 670 nm, respectively. In this bin, the sample maximum is above 0.01 at 412 and 555 nm, and about 0.003 at 670 nm. Errors are much lower at  $\tau_a$  values generally encountered in the open ocean, for example 0.003, 0.002, and 0.001 at 412, 555, and 670 nm for  $\tau_a$  equal to 0.1. Note that the standard deviation does not go to zero as  $\tau_a$  decreases to zero because of the noise introduced in the data.

The influence of aerosol loading on the standard deviation varies with geometry. In fact, it can be modeled, at each wavelength, as a linear function of  $\tau_a$  (not shown here). The fit is very good for most of the geometries. Departures from linearity occur at large sun zenith angles. The slope and intercept follow patterns similar to those in Figures 8, 9, and 10, with higher slope and intercept in the Sun glint region and

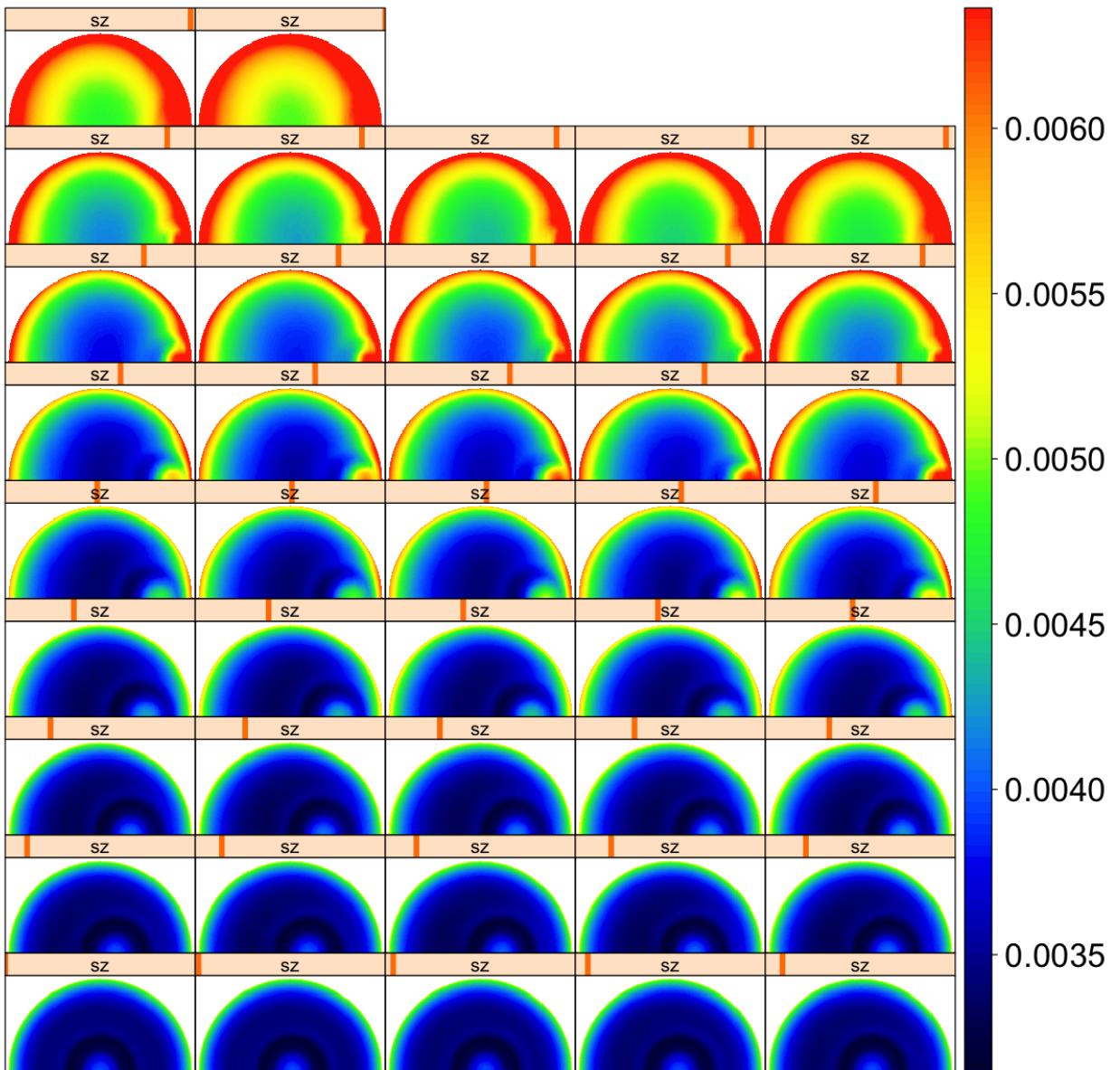


Figure 8: Standard deviation per geometry at 412 nm.

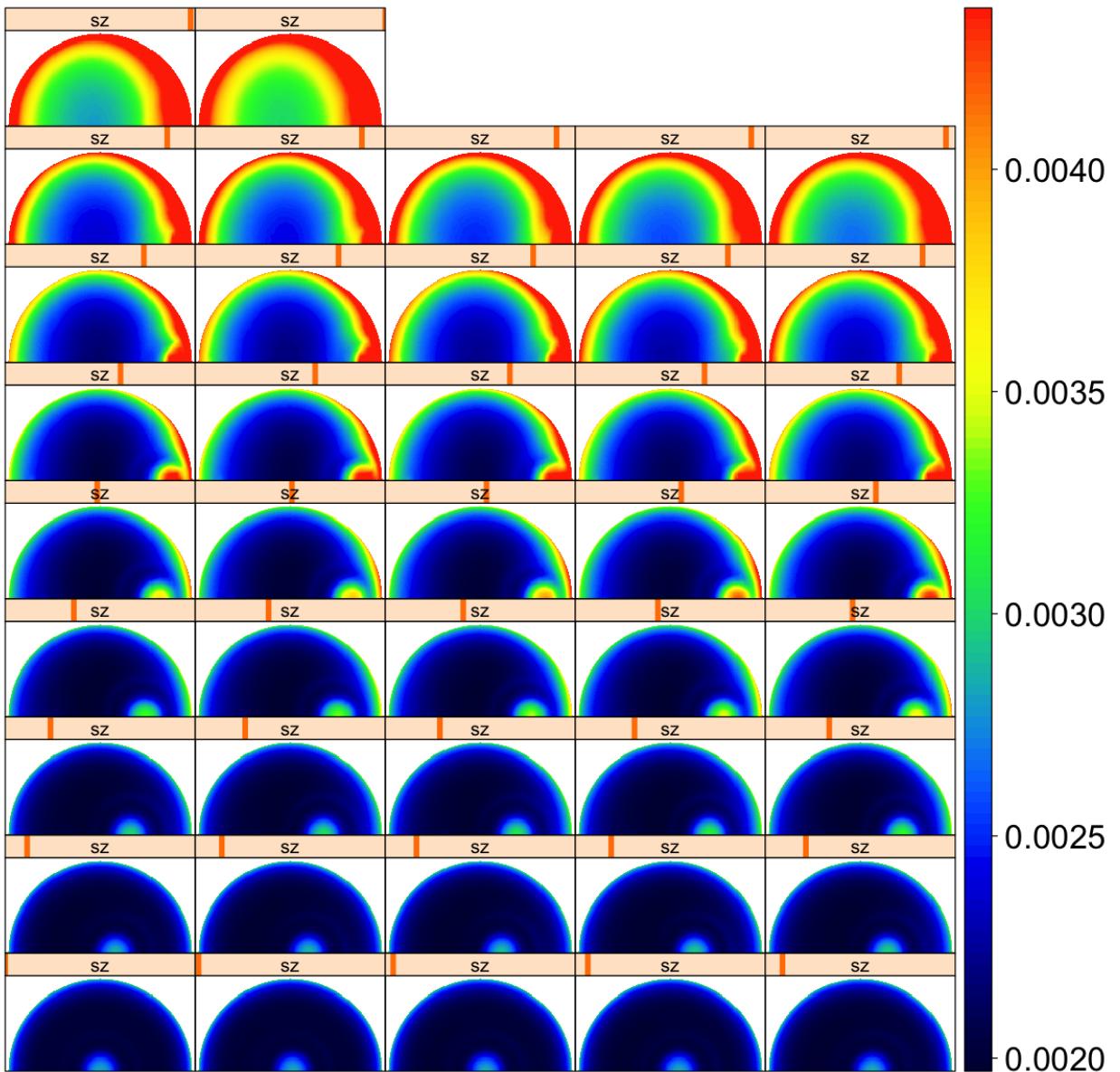


Figure 9: Standard deviation per geometry, at 555 nm.

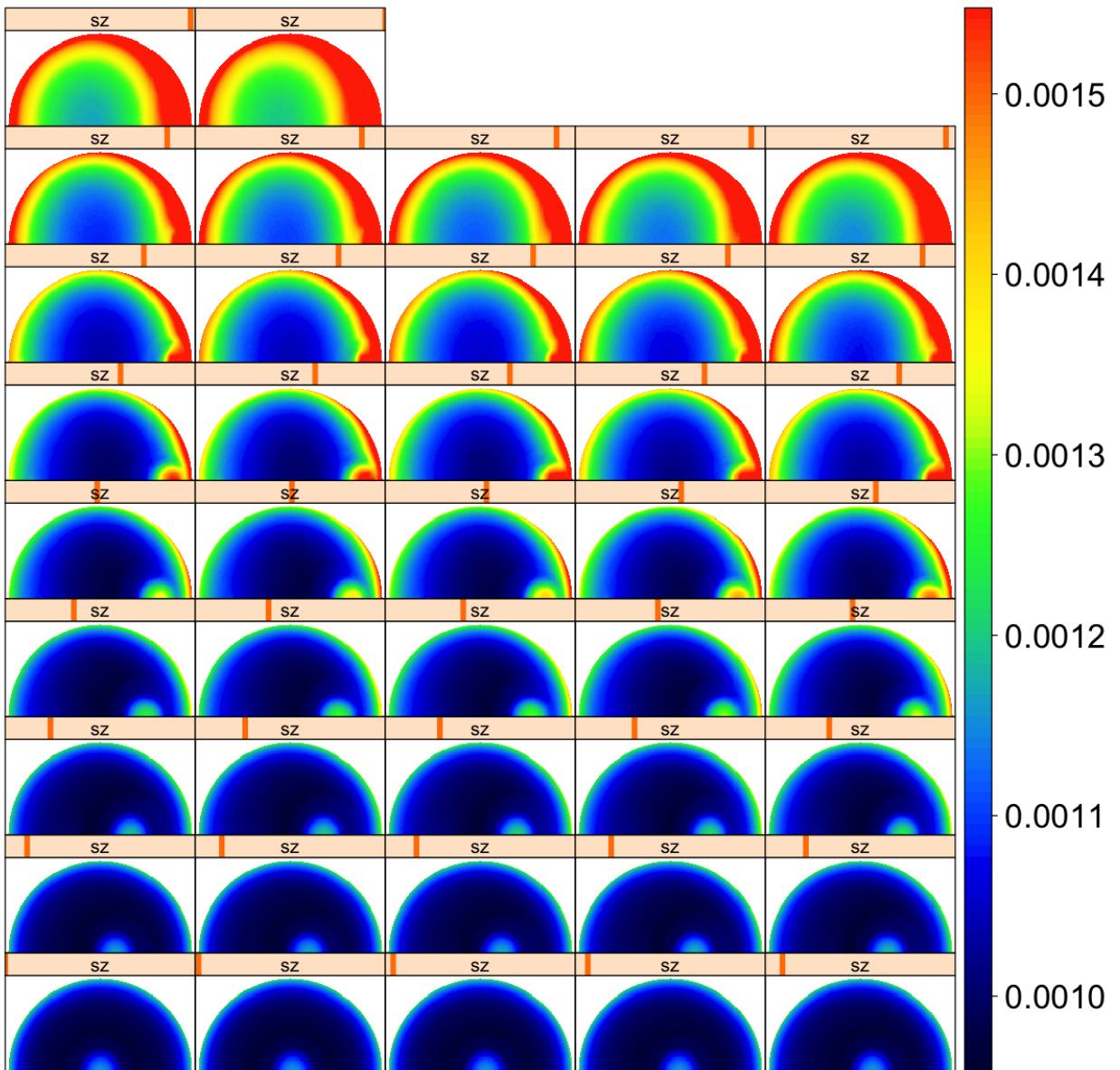


Figure 10: Standard deviation per geometry, at 670 nm.

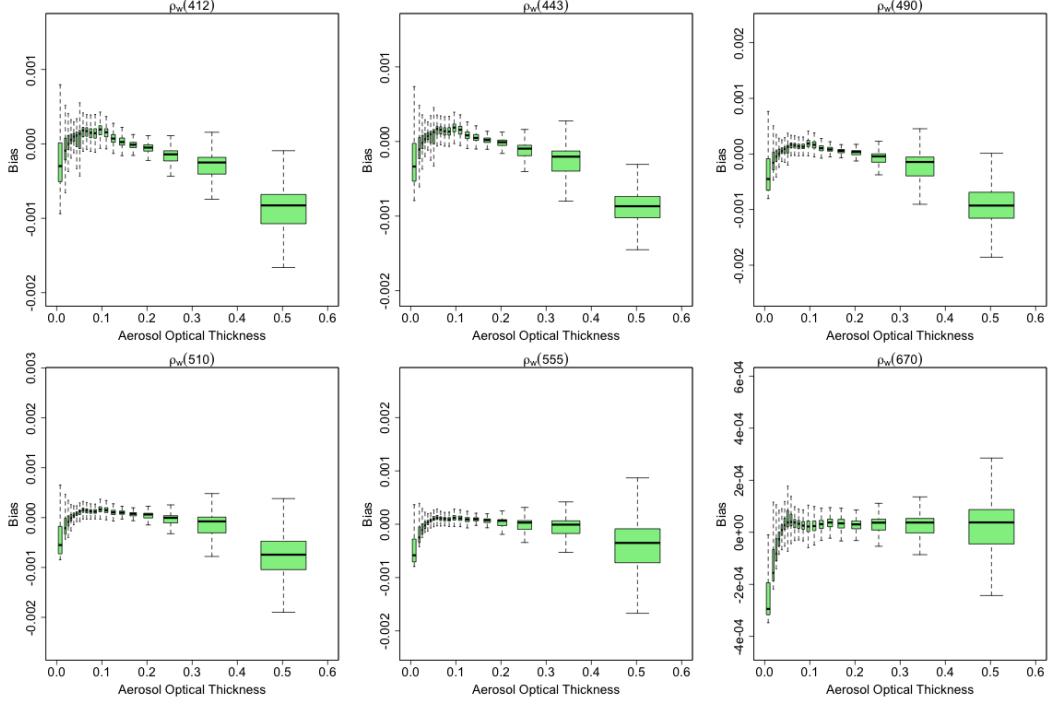


Figure 11: Bias per spectral band and aerosol optical thickness bin, with all the geometries.

at large air mass.

**Errors per aerosol type.** The bias and standard deviation, averaged over all geometries, are displayed in Figures 13 and 14 for each spectral band as a function of aerosol type. In the triangles, values at the right, top, and left corners correspond to purely maritime, urban, and continental aerosols, respectively, values at the edges to proportions of two basic types, and values inside the triangles to mixtures of the three basic types (see Section 3 and the beginning of Section 4.1).

The bias is negative at all wavelengths for aerosol mixtures dominated by the urban type, except at 670 nm, where values are slightly positive. The situation is quite reversed for mixtures dominated by the continental type, with positive bias at all wavelengths, except at 670 nm, where the bias is negligible. The mixtures dominated by the maritime type generally exhibit a small positive or negative bias. The standard deviation, on the other hand, is relatively small at all wavelengths when aerosols are mostly maritime, with minimum values occurring for mixtures of about 70% maritime, 15% urban, and 15% continental components.

The larger standard deviations are obtained for atmospheres with mostly urban or continental aerosols. The resulting RMSE (not shown) is also larger for those atmospheres, but due to the relative patterns of biases and standard deviations, the minimum values are shifted to mixtures containing 60% of maritime type. Note that the features in Figures 13 and 14 are general, with slight variations from one geometry to the other.

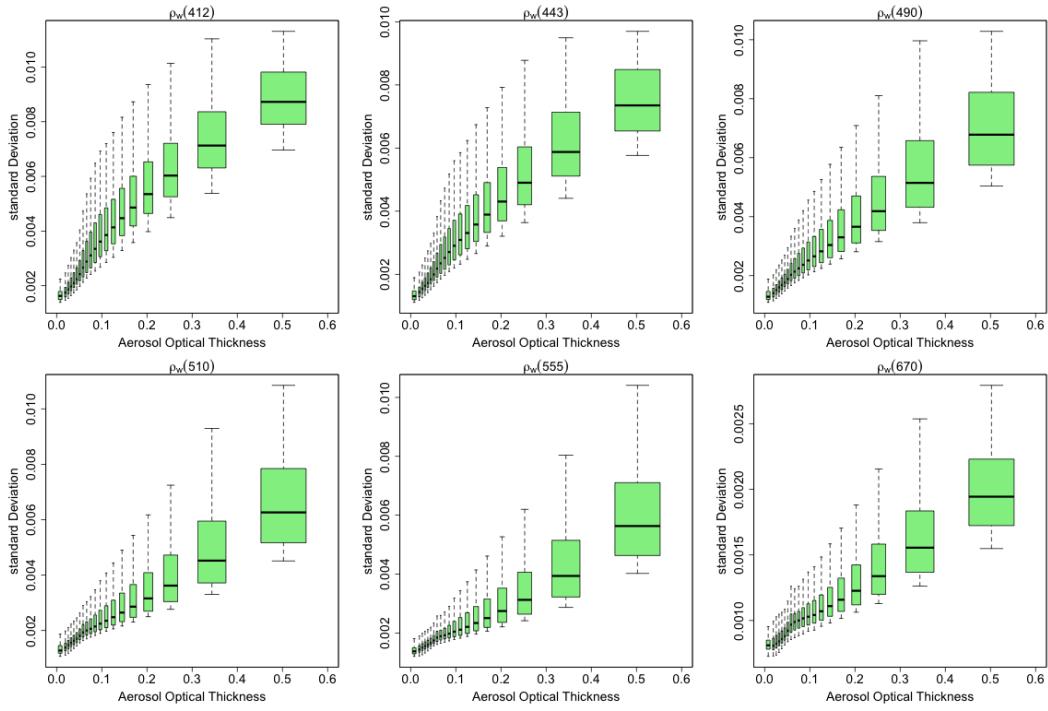


Figure 12: Standard deviation per spectral band and aerosol optical thickness bin, with all the geometries.

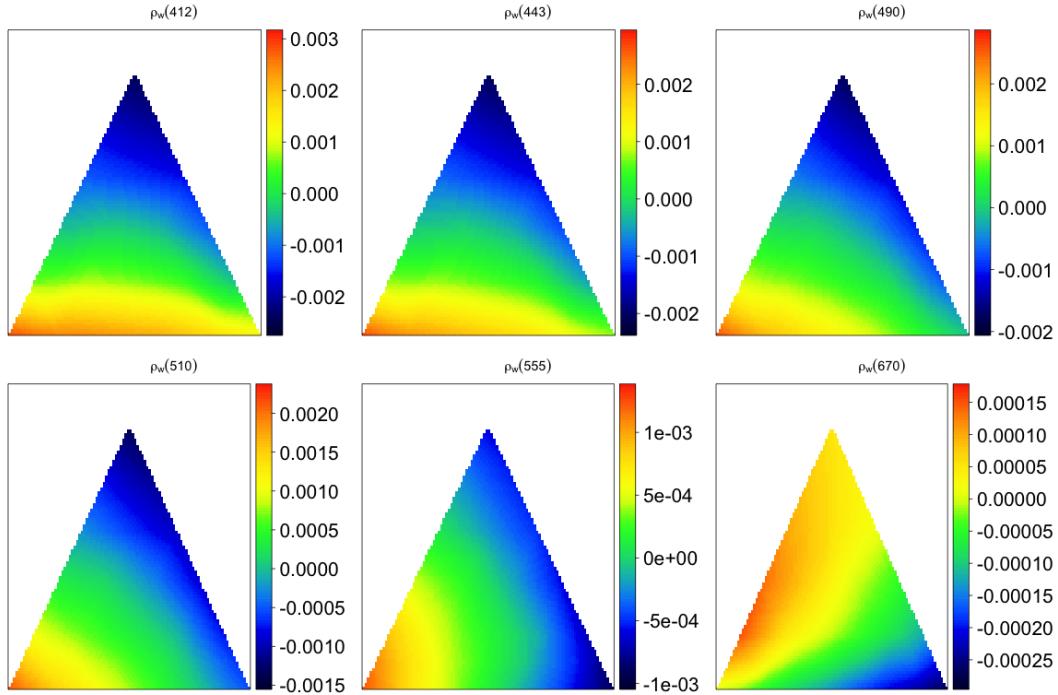


Figure 13: Bias per spectral band, averaged over all geometries, as a function of the aerosol type.

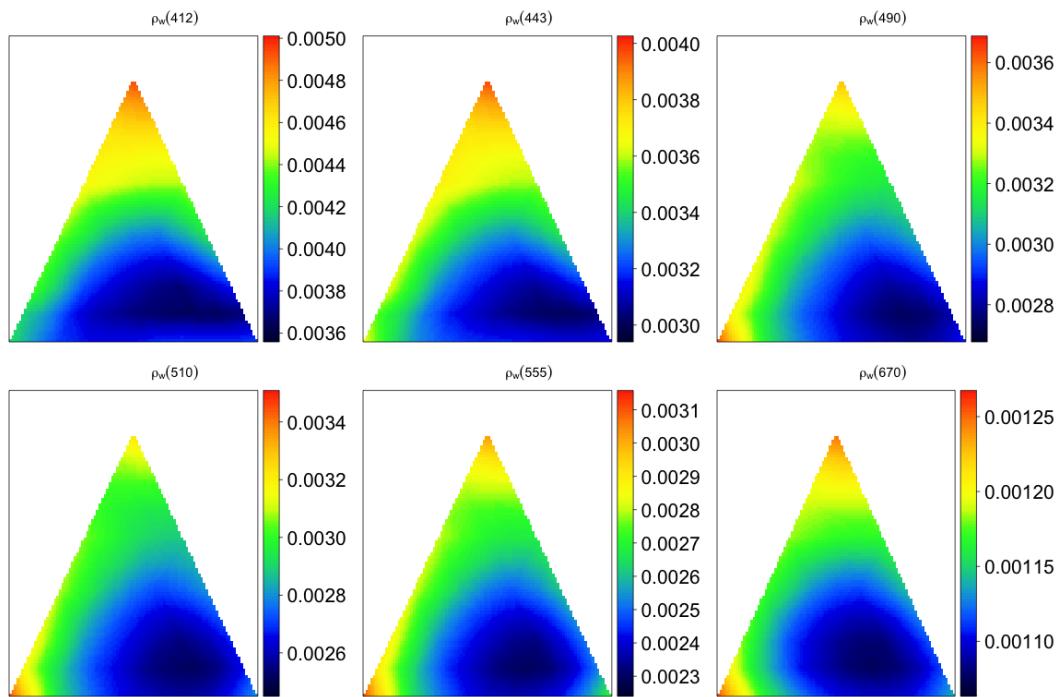


Figure 14: Standard deviation per spectral band, averaged over all geometries, as a function of the aerosol type.

#### 6.2.4 Detailed analysis for a typical geometry.

The inversion errors are examined for a typical geometry, i.e.,  $\theta_s = \theta_v = 30^\circ$ , and  $\Delta\phi = 120^\circ$ . Figure 15 displays estimated versus simulated (prescribed)  $\rho_w$ . The scatter plots indicate good performance at all wavelengths, with errors that do not depend significantly on the magnitude of  $\rho_w$  and sufficiently small to describe properly the variability of the simulated data set. In a few cases, however, the estimated values are inaccurate. This may probably be explained by the fact that these points lie close to the boundary of the support of the prior distribution of the marine reflectance. The biases are negligible and standard deviations are about 0.003, 0.002, 0.002, 0.001, 0.001, and  $< 0.001$  at 412, 443, 490, 510, 555, and 670 nm, respectively. These values are significantly lower than the globally averaged standard deviations reported in Table 7 (e.g., by 31% at 412 nm and 53% at 555 nm), for which all geometries are included. This illustrates the better performance for geometries typically encountered in ocean color remote sensing. Figure 16 displays for each wavelength the distribution of the total errors. About 70% of the cases have errors within the standard deviations indicated above. The estimated  $\rho_w$  standard deviation, a measure of uncertainty obtained from the posterior distribution, is plotted as a function of the  $\rho_w$  estimation error in Figure 17. The estimated uncertainty remains within the lines of slopes  $\pm 1/3$ , i.e., within three standard deviations of the inversion error, indicating consistency with the inversion error.

Similar information is displayed in Figure 18 displays for  $\tau_a$ , and in Figures 19, 20, and 21 for  $\rho_a$ , i.e., scatter plots of estimated versus simulated  $\tau_a$  or  $\rho_a$ , histograms of inversion errors, and plots of estimated uncertainty versus inversion error. The conclusions are similar to those for  $\rho_w$ , i.e., accurate retrievals over the range of geophysical conditions considered with negligible biases on average, allowing a good description of variability, and estimated uncertainty coherent and compatible with inversion errors.

## 7 Application to SeaWiFS imagery

The Bayesian inverse methodology has been applied to actual SeaWiFS data. Results are presented and discussed below for selected Local Area Coverage (LAC) images. These images were acquired over the oceans and seas around South Africa, Australia, Hawaii, East Asia, Southeast America, and the Mediterranean Sea. Table 8 lists the images, with name, date, time, size, and oceanic region. The marine reflectance retrievals obtained for the South Africa image are first analyzed in detail, including uncertainties,  $p$ -value, and spatial noise. Comparisons with estimates from the SeaDAS operational algorithm, are described. The other application examples are then examined, to illustrate robustness and generalization in contrasted oceanic regions.

### 7.1 S1999045100113 image, South Africa

Figure 22 displays SeaWiFS imagery acquired on 02 February 1999 around South Africa, in the region of the Agulhas and Benguela Currents, Agulhas retro-reflection, and South Atlantic Current. The RGB composite (Figure 22, left) reveals greener waters along the coast of South Africa due to coastal upwelling generated by the Agulhas and Benguela currents (offshore Ekman transport) and in the Agulhas retro-reflection region (turbulent zone of mixing). The TOA reflectance image at 865 nm (Figure 22, right) exhibits, outside of clouds, regions with relatively high values at the edge of the cloud system off the West coast of South Africa, where droplet concentration becomes low, and South of the Cape of Good Hope, where winds are strong (aerosols generated by wind action, whitecaps).

The marine reflectance retrieved by the Bayesian inversion scheme is displayed in Figure 23 for spectral bands centered on 412, 443, 490, 510, 555, and 670 nm. Clouds, masked using SeaDAS flags, are displayed in white. Near the coast, where upwelling brings nutrients to the surface, values are relatively low at 412, 443, and 490 nm and high at 555 and 670 nm near the coast (typically 0.005 at 443 nm and 555 nm),

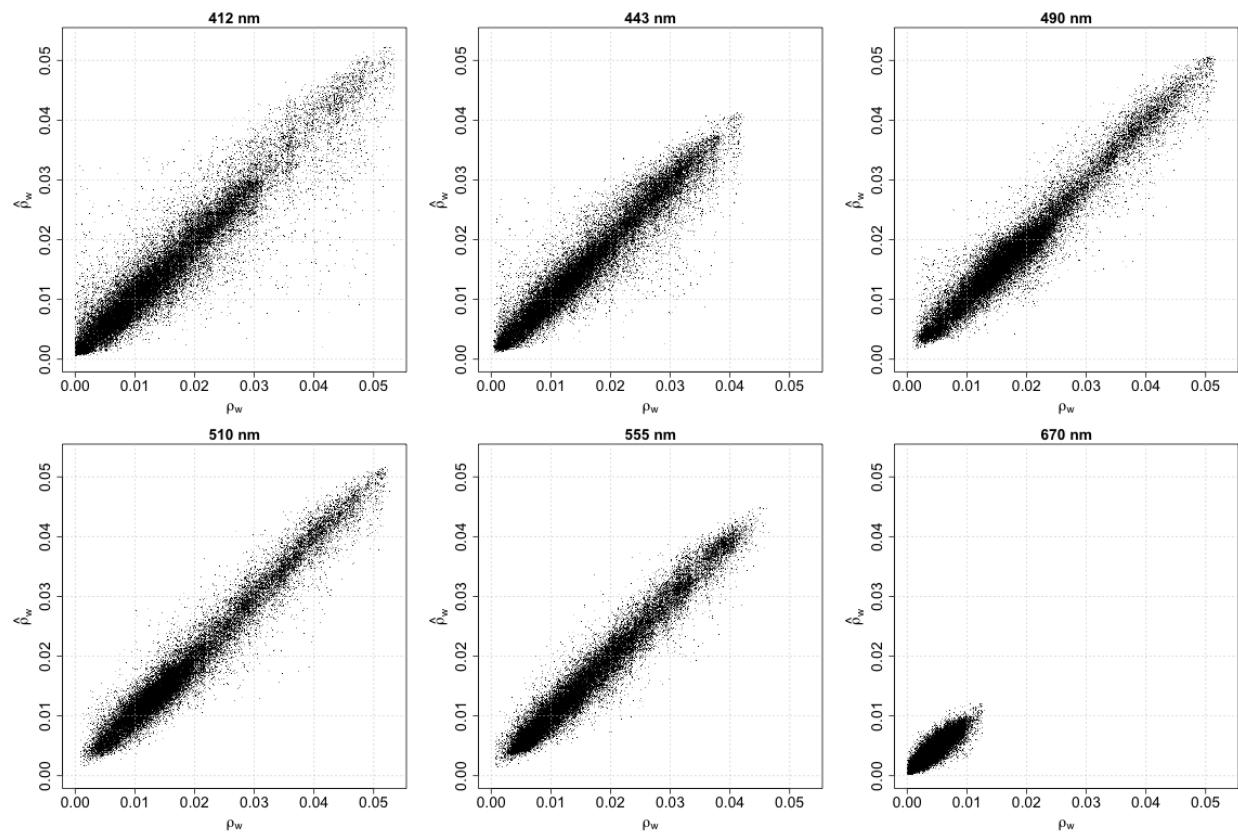


Figure 15: Estimated versus simulated  $\rho_w$ .

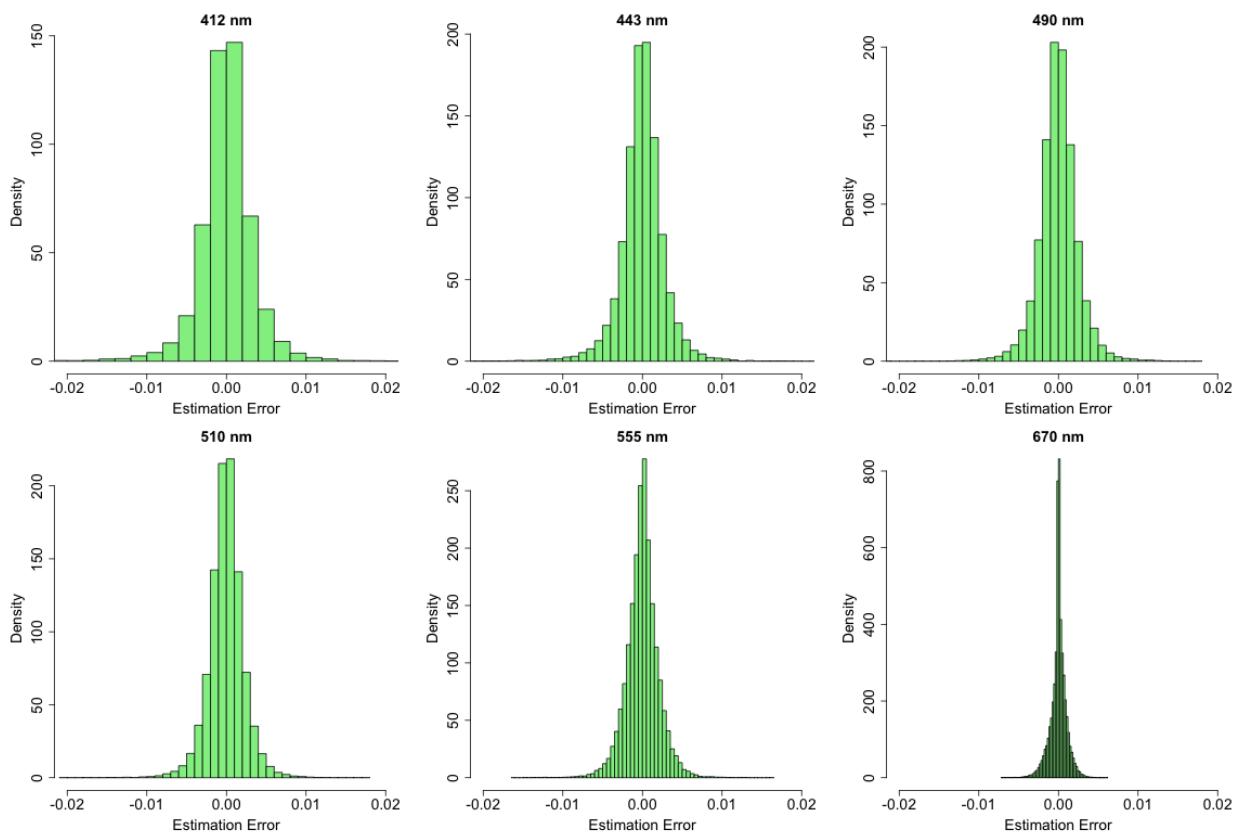


Figure 16: Histograms of  $\rho_w$  estimation error.

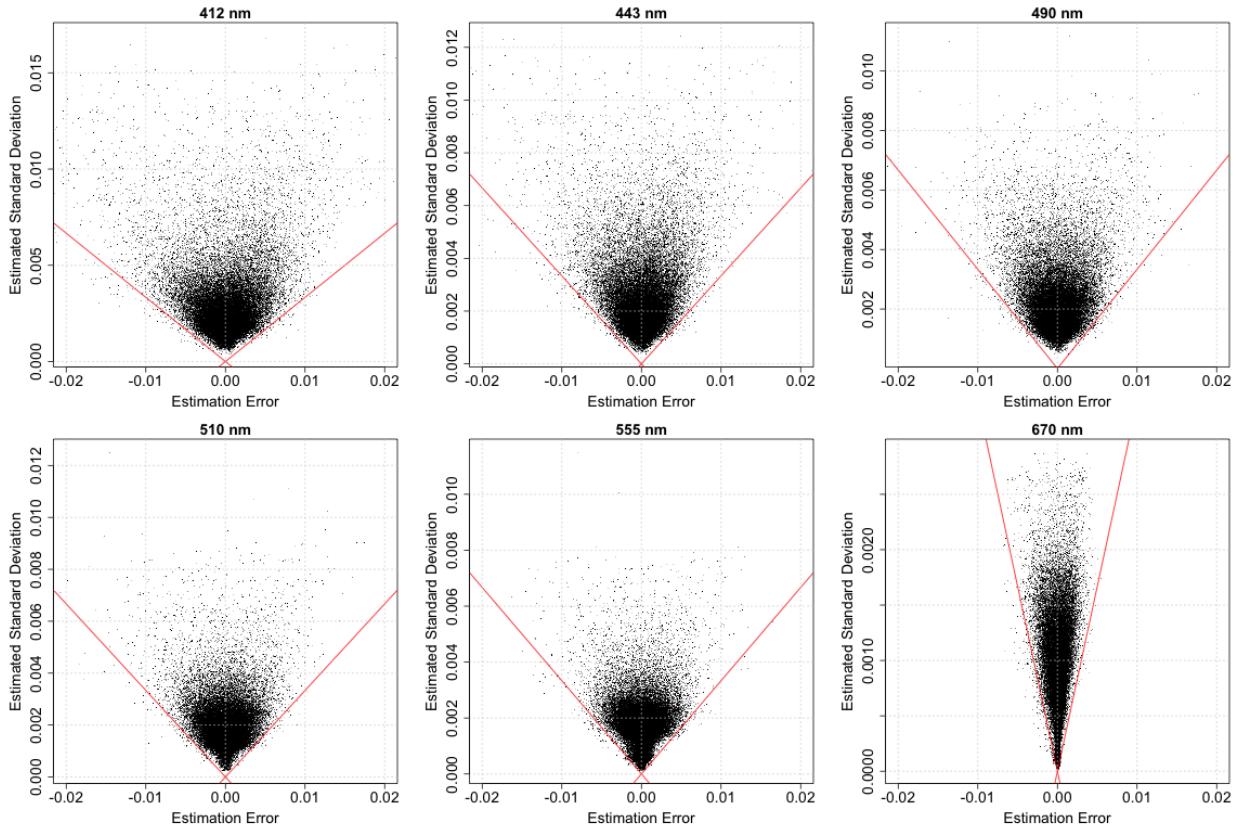


Figure 17: Estimated  $\rho_w$  Standard Deviation versus  $\rho_w$  estimation error, with lines of slopes  $\pm 1/3$ .

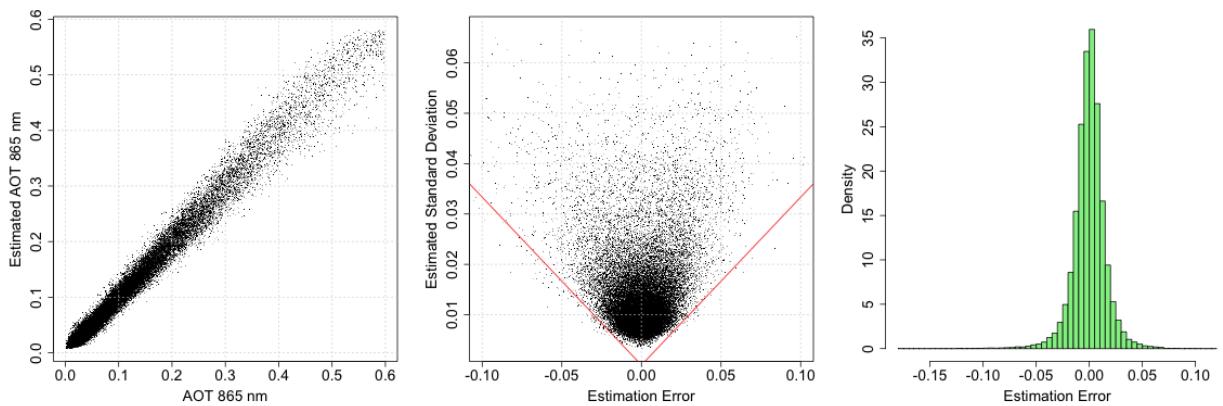


Figure 18: (Left) Estimated versus simulated  $\tau_a$ . (Center) Estimated standard  $\tau_a$  standard deviation versus  $\tau_a$  estimation error with lines of slopes  $\pm 1/3$ . (Right) histogram of  $\tau_a$  estimation errors.

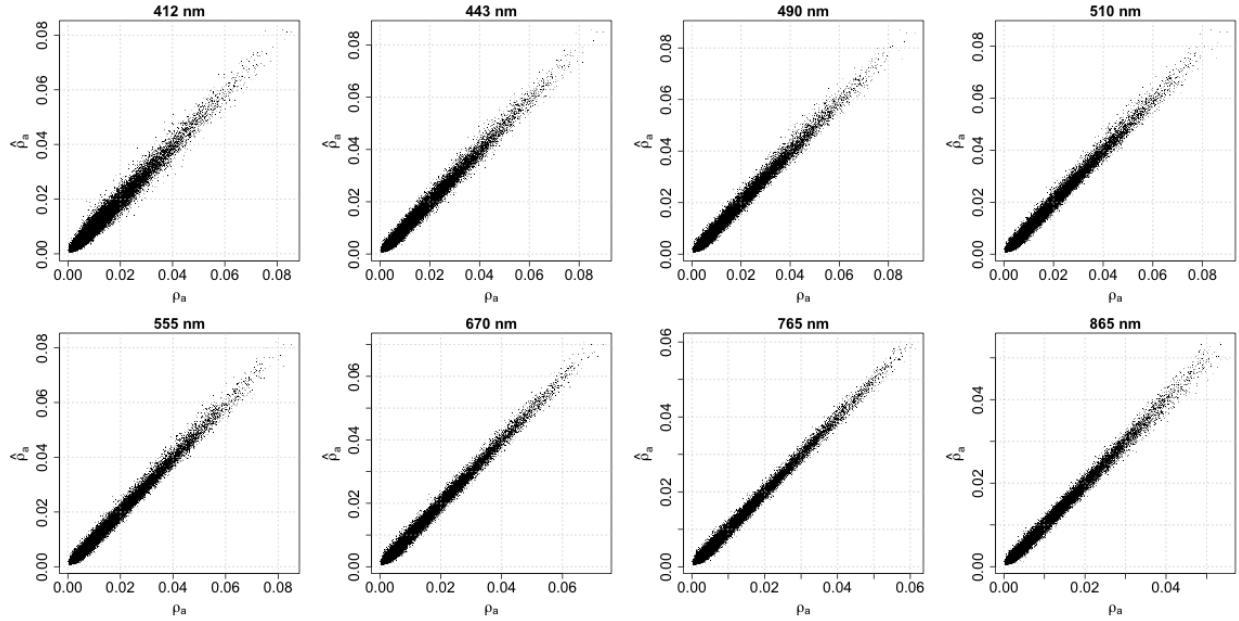


Figure 19: Estimated versus simulated  $\rho_a$ .

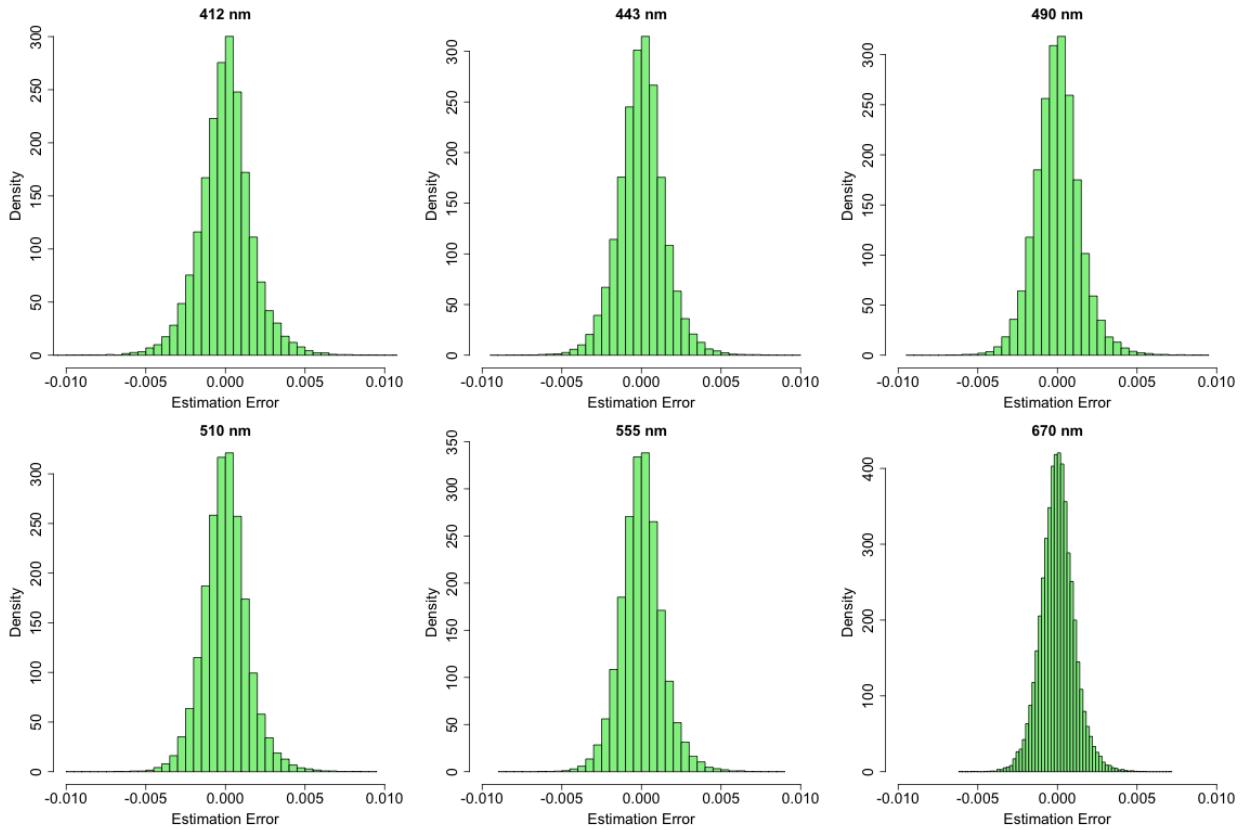


Figure 20: Histograms of  $\rho_a$  estimation error.

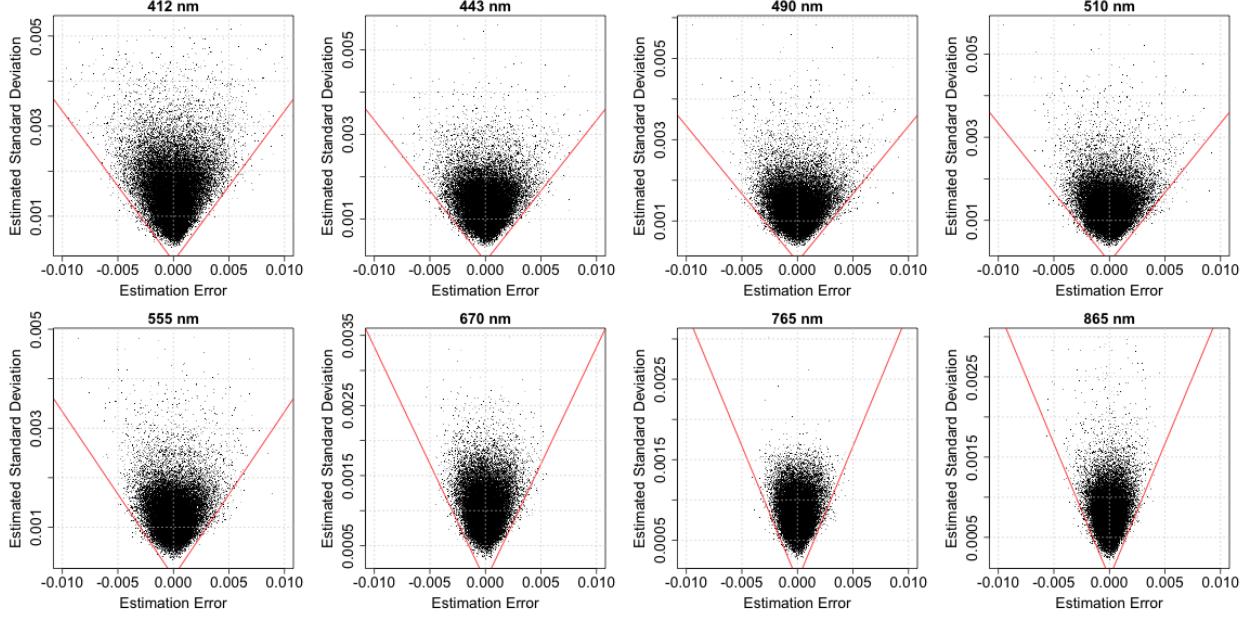


Figure 21: Estimated  $\rho_a$  standard deviation versus  $\rho_a$  estimation error, with lines of slopes  $\pm 1/3$ .

indicating productive waters. In the turbulent Agulhas retro-reflection zone (center of the image), marine reflectance in the blue is lower than in the surrounding regions of anti-cyclonic circulation associated with the Benguela drift (to the left) and the Agulhas return current (to the right), with typical values of 0.015 instead of 0.02 to 0.025. Near clouds, and in regions of broken cloudiness, the spatial features of marine reflectance exhibit continuity with respect to adjacent clear sky regions. No significant correlation exists between marine reflectance and wind speed, which exhibits a strong gradient from the center to the bottom of the image (not shown here), and between marine reflectance and TOA reflectance at 865 nm. At the edge of the cloud system on the left of the image, for example, mesoscale eddies of relatively higher marine reflectance are revealed, with spatial features apparently not affected by the large gradient of TOA reflectance (i.e., aerosol optical thickness) across the edge.

The absolute uncertainty associated with the marine reflectance estimates is provided in Figure 24 for each spectral band on a pixel-by-pixel basis. This uncertainty was calculated as the covariance of the conditional (posterior) distribution of  $\rho_w$  given  $\rho$  (see Section 2). Values remain within  $\pm 0.003$  in the blue,  $\pm 0.002$  in the green, and  $\pm 0.0005$  in the red for most pixels. Larger uncertainties, i.e.,  $\pm 0.005$  to  $\pm 0.01$  in the blue, are encountered near clouds (imperfect cloud masking or processes not accounted for in the modeling) and where the TOA reflectance at 865 nm is large. The amplitude of the uncertainty is generally larger when the marine reflectance is higher, and it represents a larger percentage of the marine reflectance when marine reflectance is lower, ranging between 40% (coastal regions) and 10% (outside the retro-reflection zone).

The marine reflectance fields obtained from SeaDAS, displayed in Figure 25, exhibit spatial features showing resemblance in characteristics and appearance with the corresponding Bayesian fields. The marine reflectance values are slightly smaller than the Bayesian estimates, but they vary within a similar range. More abnormal values are retrieved with SeaDAS near clouds, and the imagery is generally noisier. This is especially apparent in the marine reflectance field at 670 nm (Figure 25 compared with Figure 23, bottom right). The differences between the two types of estimates are also evidenced in the histograms of marine reflectance (Figure 28), which show distributions similar in shape, but shifted toward higher values in the

case of SeaDAS, by about 0.002 at 412, 443, 490 and 510 nm, 0.001 at 555 nm, and 0.0005 at 670 nm. At these last two wavelengths, the SeaDAS values are more broadly distributed, with peaks at 0.005 and 0.0007, respectively. Figure 26 displays variograms of marine reflectance, obtained for a  $128 \times 128$  pixel sub-area depicted in Figure 22. These variograms are defined as the variance of the difference between marine reflectance at two locations. They describe the degree of spatial dependence of the marine reflectance field. The SeaDAS values do not go to zero as spatial distance goes to zero, indicating more noisy retrievals. This may be attributed to the aerosol model selection in the SeaDAS algorithm (different models selected for neighboring pixels). Spatial correlation is lower at small scales with SeaDAS, and tends to decrease more slowly as distance increases. The differences have consequences on the determination of de-correlation scales and analysis of mesoscale ocean biological variability (see, e.g., Doney et al., 2003).

The Bayesian technique retrieves not only marine reflectance, but also aerosol optical thickness  $\tau_a$  and the atmospheric term  $\rho_a$  (see Equation 2.12) at all wavelengths. Figure 28 displays  $\tau_a$  and  $\rho_a$  at 865 nm and their associated uncertainties. The  $\tau_a$  values range from 0.02 to 0.2 and the  $\rho_a$  values from 0.001 to 0.015. High values are encountered near the edge to the cloud system off the West coast of South Africa and in the immediate vicinity of clouds, and they correspond to high TOA reflectance values at 865 nm (Figure 22). As expected the fields of  $\tau_a$  and  $\rho_a$  are well correlated, since  $\rho_a$  is essentially proportional to  $\tau_a$ . The uncertainty on  $\tau_a$  and  $\rho_a$  is generally within  $\pm 0.01$  and  $\pm 0.0005$ , with higher values in regions of relatively high aerosol optical thickness (e.g., upper left part of the image).

The  $p$ -value associated to each pixel of (or retrieval in) the image is also displayed in Figure 28. As explained in Section 2, this parameter quantifies how likely is the observation  $\rho$  with respect to the model. A low  $p$ -value (i.e.,  $< 0.01$  or 0.05) indicates that model and observation are incompatible. The  $p$ -value in Figure 28 is above 0.05 almost everywhere, except in the vicinity of clouds, where some processes may not be accounted for in the radiation-transfer modeling (e.g., adjacency effects due to the high reflectance contrast between clouds and the environment, large optical thickness). Low  $p$ -values are also encountered where winds are especially strong, i.e., in the region between the Cape of Good Hope and the Southeastern part of the cloud system in the upper left part of the image, where wind speed exceeds 12 m/s. Such conditions are outside the atmospheric parameter space used for approximating the forward operator  $\Phi_a$  (see Section 4.1).

## 7.2 Other application examples

Other SeaWiFS images, acquired over diverse oceanic regions, have been processed with the Bayesian algorithm. Figure 29 displays the retrieved marine reflectance imagery at 412 nm obtained for images over the Sea of Japan and the Northwest Atlantic on April 7, 2001 (S2001097031924, top left), the central Tropical Pacific around the Hawaiian Islands on March 17, 2001 (S2001076224946, top right), the Sea of Japan, Yellow Sea, and East China Sea on April 15, 2001 (S2001105040514, middle left), the Argentine Sea on March 10, 2002 (S2002069152729, middle right), the East Indian Ocean off Western Australia on August 11, 2002 (S2002233034805, bottom left), and the Mediterranean Sea on August 16, 2004 (S2004228121834, bottom right.). The images of the East Asian Seas (Figure 29, top left and middle left) were acquired during the Asian Pacific Regional Characterization Experiment (ACE-Asia), which took place off the coast of China, Japan, and Korea (Huebert et al., 2003). The atmosphere contained a variety of aerosols, including wind-blown dust from the China deserts and particles generated by human activity and industrial sources (e.g., Kahn et al., 2004), offering the opportunity to examine algorithm performance in the presence of absorbing aerosols.

The marine reflectance imagery at 412 nm exhibits low values (0.005 to 0.015) in the Sea of Japan, where waters are generally productive (Figure 29, top left and middle left), and relatively high values in the East China Sea, due sediments from the Yangtze river (Figure 29, middle left). Values are high ( $> 0.03$ ) in the nutrient-poor waters of the Tropical Pacific around Hawaii (Figure 29, top right). Alternate bands of low and

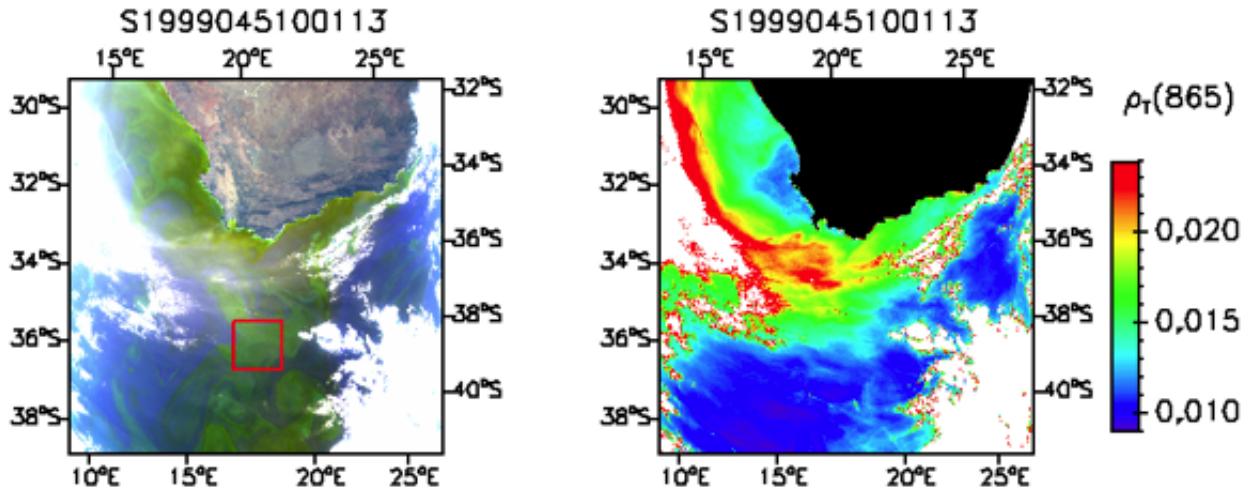


Figure 22: SeaWiFS imagery of the ocean off South Africa acquired on February 14, 1999. Left: True color image. Right: TOA reflectance at 865 nm. Clouds are masked in white according to SeaDAS.

high values are observed over the continental shelf off Patagonia and in the confluence zone of the Brazil and Malvinas currents, where turbulent eddies and swirls form, pulling up nutrients from the deep ocean (Figure 29 , middle right). Mesoscale features are revealed over the Northwest shelf of Australia, where poleward seasonal currents along the coast interact with the eastward circulation South of the Indonesian through flow (Figure 29 , bottom left). Contrasted situations are depicted in the Western Mediterranean Sea, with low values (0.015 to 0.025) in the Northern part, where phytoplankton concentration is relatively high due to runoff from continental margins (input of nutrients) and mixing from strong winds, and higher values (0.04) in the central, more oligotrophic part of the basin (Figure 29 , bottom right). Note that the band of high marine reflectance in the convergence zone off Patagonia (Figure 29 , middle right) is probably not associated with low chlorophyll concentration (i.e., low absorption in the blue), but rather to reflective phytoplankton species, presumably coccolithophores, as suggested by the marine reflectance in the green and red (not shown here), which is also high and not typical of oligotrophic waters.

The marine reflectance retrievals display good spatial continuity, but abnormal values are sometimes encountered in the vicinity of clouds, as evidenced in the images of the Argentine Sea and Eastern Indian Ocean. This may be due, as mentioned for the image of South Africa seas, to imperfect cloud masking and/or processes not accounted for in the forward modeling. Larger uncertainties are generally associated with abnormal marine reflectance values, as revealed in Figure 30. For example, uncertainties reach 0.007 in the patchy cloud region off Northwest Australia (Figure 30, bottom left) and at the edge of clouds in the Argentine Sea (Figure 30, middle right). In the clear sky region between Korea and Japan (Figure 30, middle left), uncertainties are also high ( $> 0.01$ ), which is due to absorbing aerosols (polluted air and dust from Yellow Sea, North China, and Northeast Mongolia, see Kahn et al. (2004)), and the likely inability of the forward model to represent local aerosol conditions, especially the effect of aerosol altitude on the aerosol reflectance. Note that the uncertainties are quite small (0.003-0.004) in the region influenced by the Yangtze river (Figure 30, middle left), even though the optically complex waters in that region are not represented in the forward model. Such situations may occur, for which part of the marine signal, generally not null in the near infrared, is interpreted as an aerosol signal, resulting in a posterior probability distribution with small conditional covariance. This suggests that the conditional covariance is not a sufficient measure of uncer-

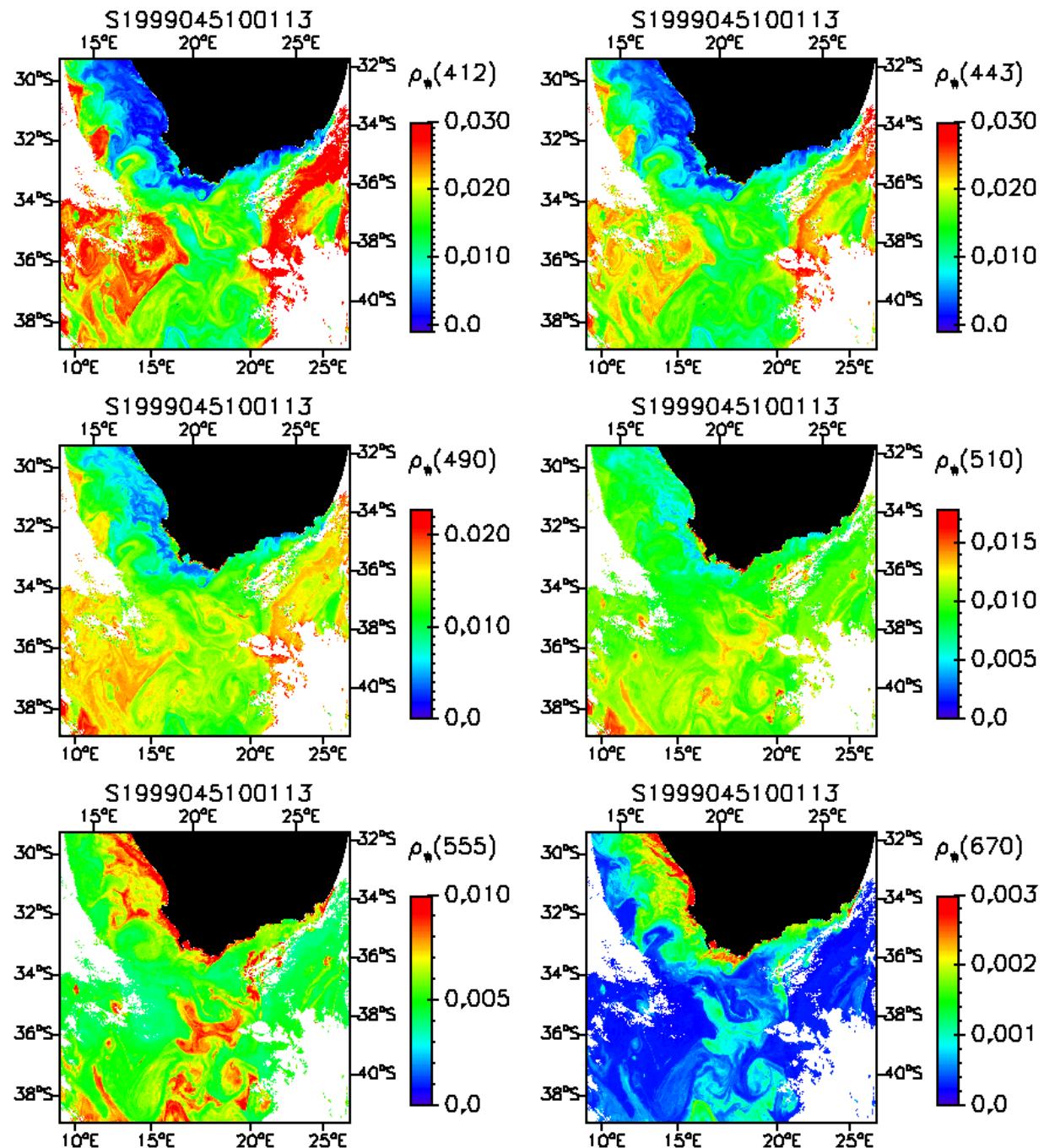


Figure 23: Estimated  $\rho_w$  by the Bayesian methodology.

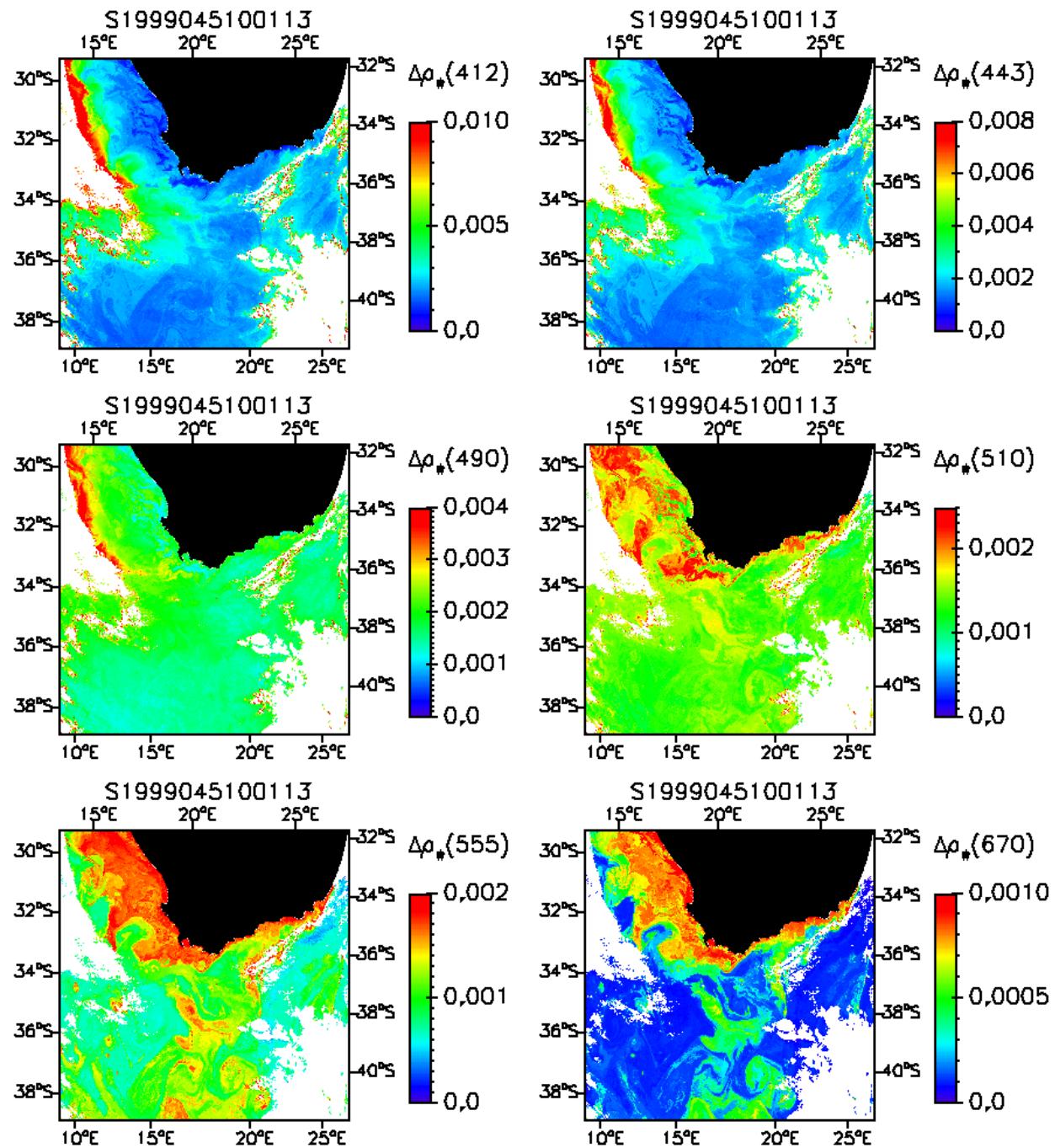


Figure 24: Estimated  $\rho_w$  standard deviation.

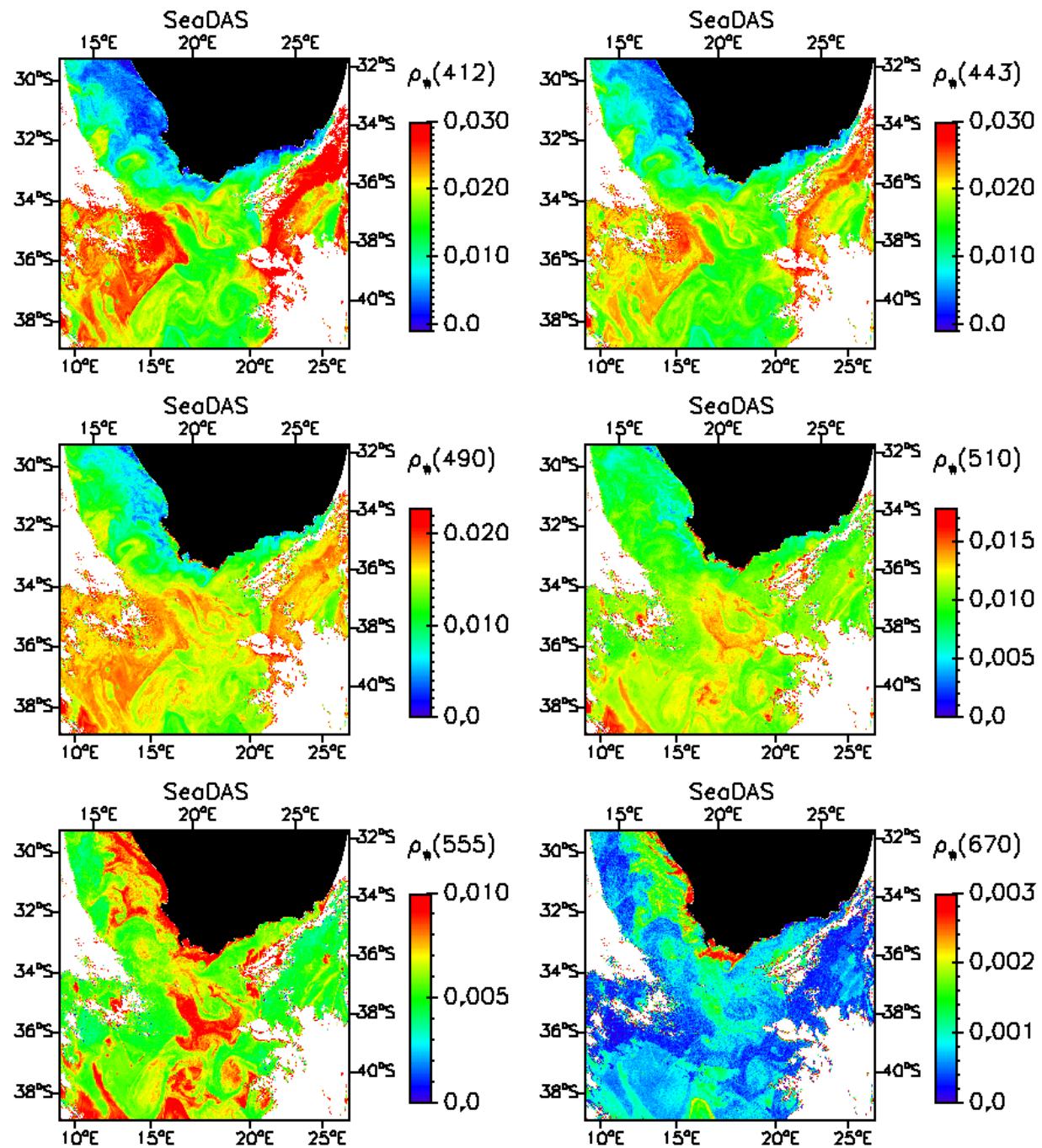


Figure 25: Top panels: Estimated  $\rho_w$  by the SeaDAS algorithm.

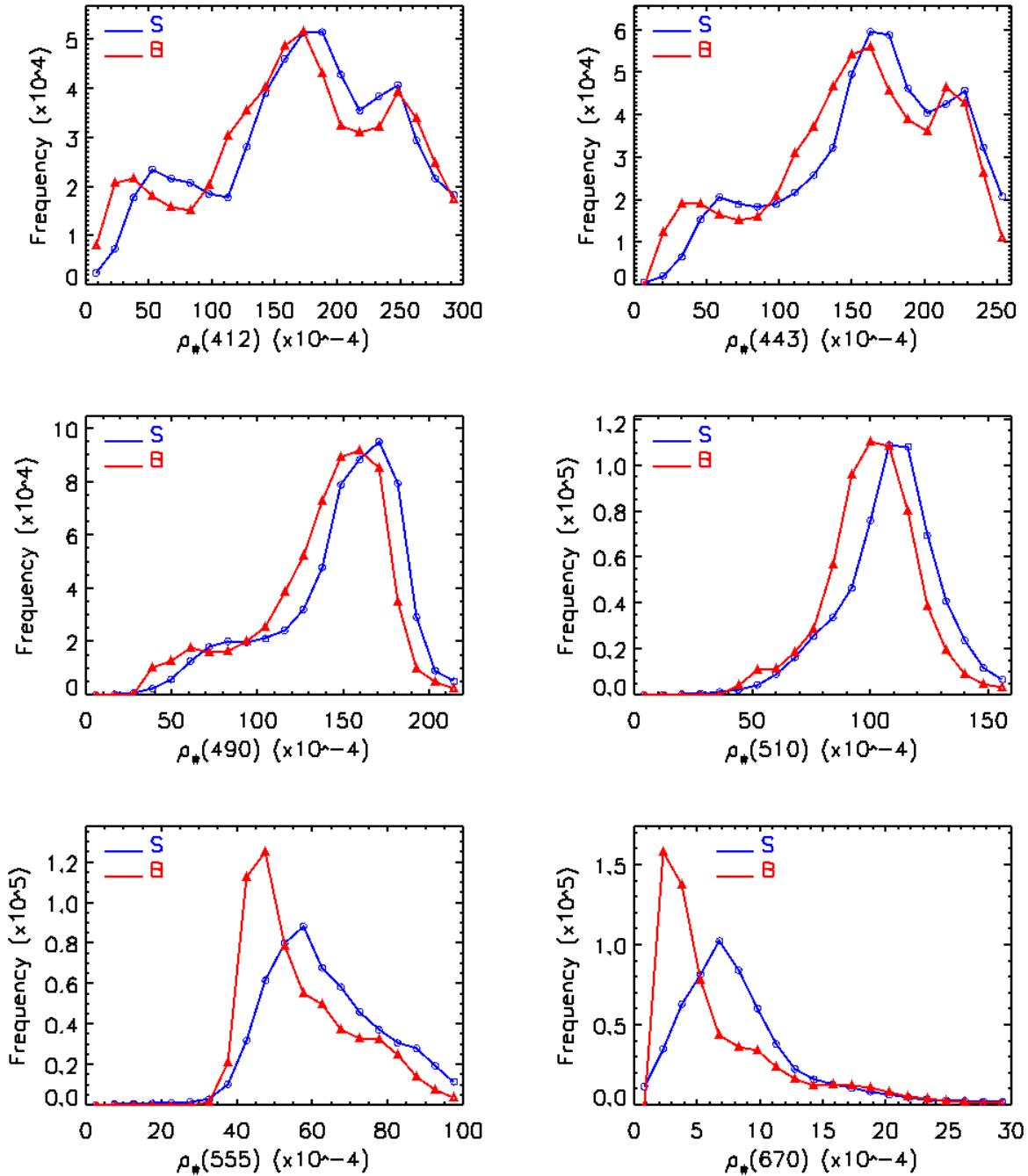


Figure 26: Histograms of valid  $\rho_w$  estimates obtained with the Bayesian technique (B) and the SeaDAS algorithm (S). The Bayesian estimates are shifted toward lower values.

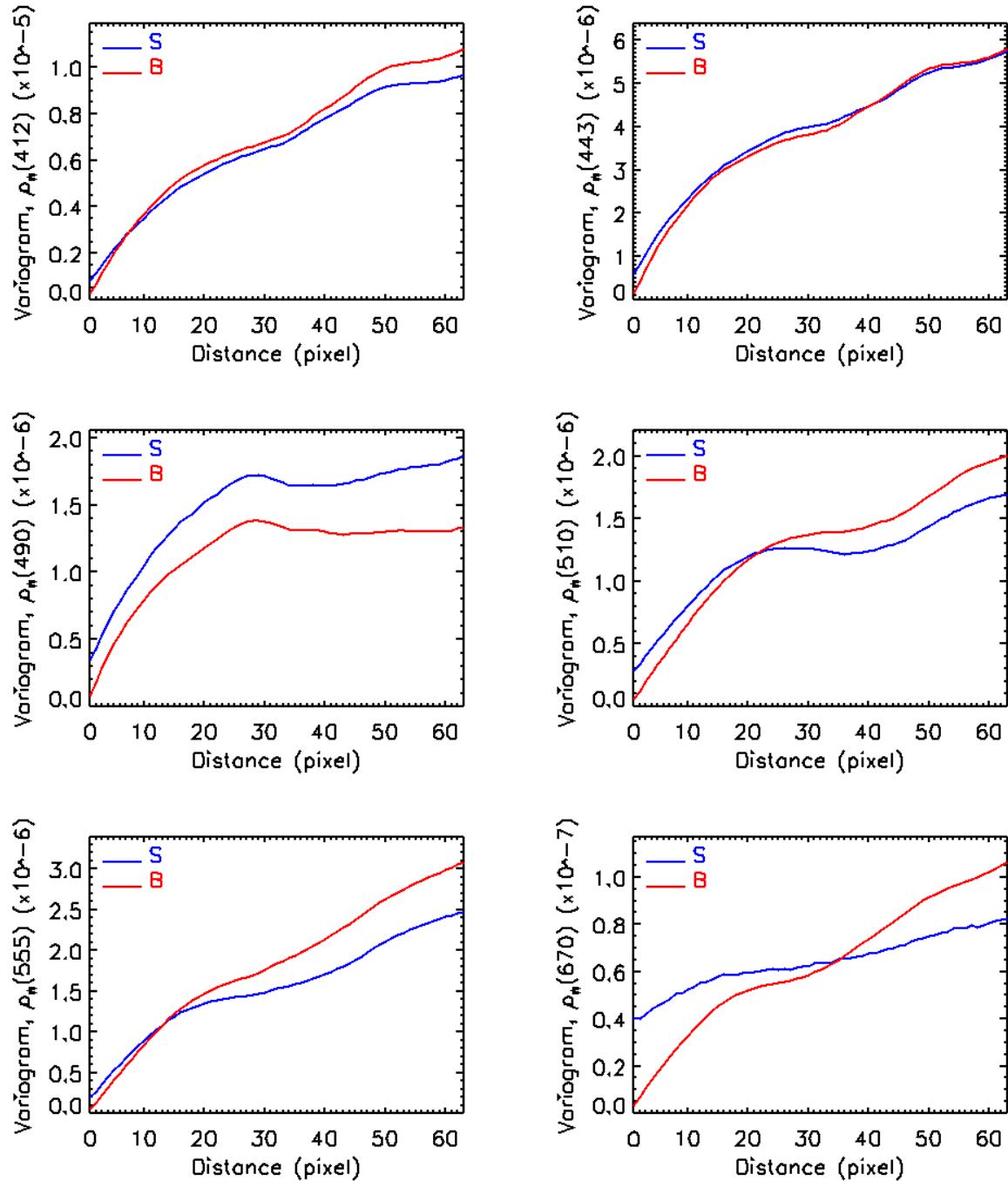


Figure 27: Variograms of valid  $\rho_w$  estimates obtained with the Bayesian technique (B) and the SeaDAS algorithm (S). Computations are made using data in a 128x128 pixel area South of the Cape of Good Hope. The variograms obtained with the statistical method indicate less noisy retrievals.

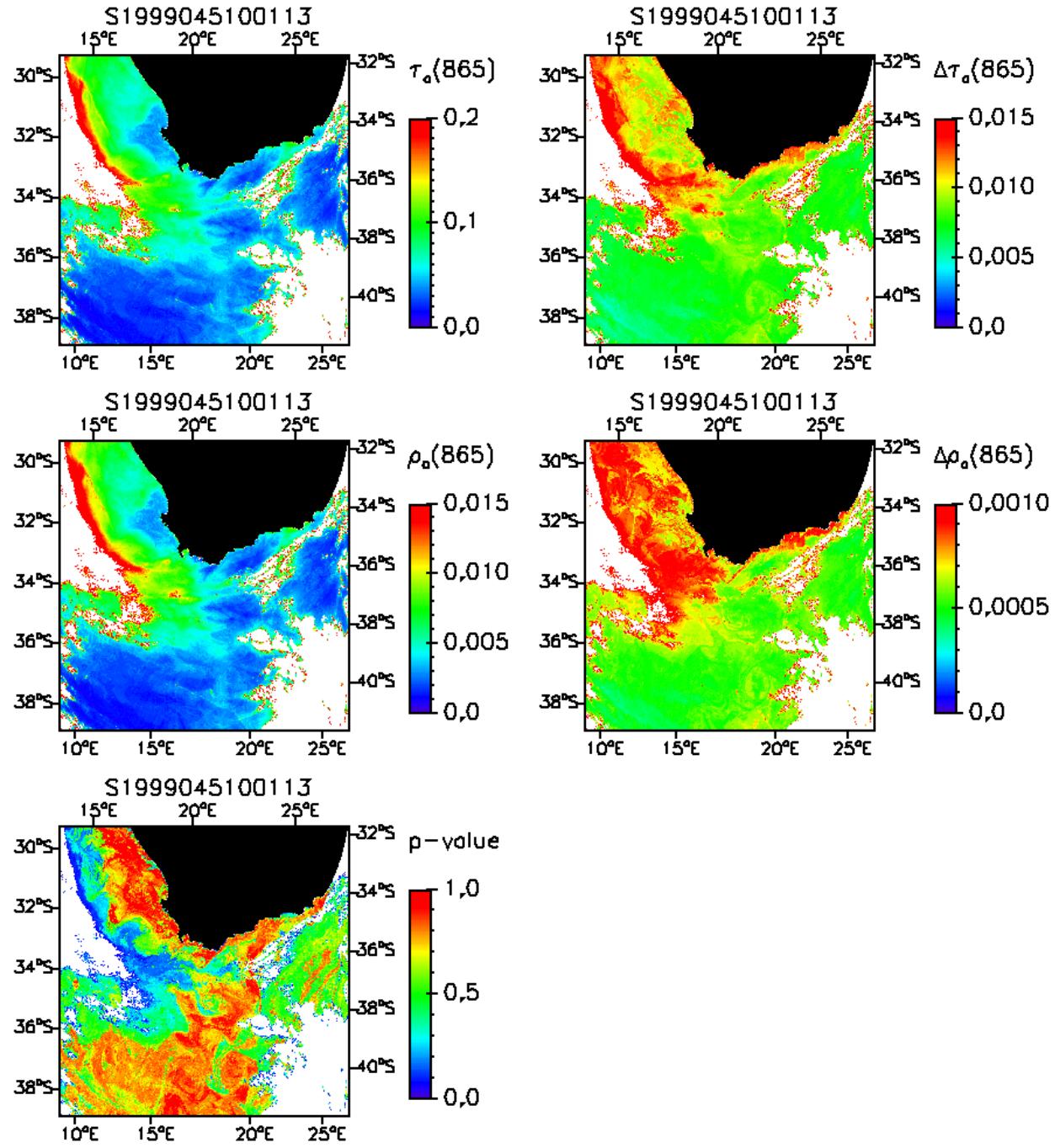


Figure 28: Top panels: Estimated  $\tau_a$  and  $\tau_a$  standard deviation. Middle panels: Estimated  $\rho_a$  and  $\rho_a$  standard deviation. Bottom panel:  $p$ -value.

tainty. One needs to take into consideration the adequacy of the forward model to fit the observations, or the probability that the model takes values as extreme as the observations, i.e., the *p*-value, in the assessment of performance.

Figure 31 gives the *p*-value associated to the marine reflectance images. In many parts of the images, the *p*-value is above 0.5, indicating good compatibility between model and observations. In the vicinity of clouds, or in non-cloudy regions affected by absorbing aerosols (Sea of Japan, Korea Strait), the *p*-values are relatively low, yet acceptable ( $> 0.1$ ). Unacceptable values ( $< 0.05$ ) are encountered in the East China Sea (Figure 31, middle left) and in the immediate vicinity of the Northwest Australia coast (Figure 31, lower left). In these regions, model and observations are inconsistent, and the retrievals should be rejected, even though the conditional covariance may be low (case of the East China Sea).

Compared with SeaDAS values, the Bayesian estimates do not exhibit systematic biases (Figure 32). Differences between the SeaDAS and Bayesian retrievals are generally negative in the Sea of Japan, the Korean Strait, the Yellow Sea, by as much as 0.01 in some regions, they are positive by up to 0.01 in the region influenced by the Yangtze river, they are negligible ( $< 0.001$  in magnitude) in the seas around Hawaii, they are mostly positive (0.001 to 0.004 except near clouds) in the Argentine sea, and they are either negative or positive with values between  $-0.01$  and 0.005 off Northwest Australia and in the Mediterranean Sea. In view of the absolute marine reflectance displayed in Figure 29, the SeaDAS-derived marine reflectance is often negative (and probably too low when positive) in the East Asian Seas, indicating the inability of the SeaDAS algorithm to deal with absorbing aerosols. This is not surprising, since in SeaDAS the aerosol reflectance is estimated from observations in the red and near infrared (765 and 865 nm bands), where the effect of aerosol absorption, essentially due to the coupling with molecular scattering, is negligible. The Bayesian technique, on the other hand, makes use of all the spectral information available, including observations in the blue that are sensitive to aerosol absorption effects.

### 7.3 Comparison with in-situ data

The performance of the Bayesian technique has been evaluated experimentally in comparisons with in situ measurements of marine reflectance. The measurements were taken from the MOBY and NOMAD data sets (Clark et al., 2003; Werdell and Bailey, 2005) and matched with the satellite data, within  $\pm 3$  hours of overpass. The closest  $3 \times 3$  pixel box was selected for processing by the inversion scheme, and the marine reflectance retrieved for each of the 9 pixels was interpolated to the geographic location of the in situ measurements. The cases for which some of the pixels in the box did not pass the SeaDAS cloud-screening flags were eliminated. This treatment, however, may not be sufficient when clouds are within a distance of 10 km, due to adjacency effects (Santer and Schmechtig, 2000). The TOA radiance was not corrected for vicarious calibration adjustment (only temporal calibration changes taken into account), allowing the MOBY data to be considered in the evaluation. Note that the NOMAD data can be included in the comparisons because they were only used in the model specification to define the support of the prior distribution on the marine reflectance. The match-up data set is not used in the construction of the models.

The marine reflectance match-up data sets covered the period from September 1997 to March 2004 (NOMAD) and December 1997 to March 2007 (MOBY), and consisted of 690 and 948 pairs of estimated and measured values, respectively. These included 132, 144, 144, 113, 129, and 28 pairs at 412, 443, 490, 510, 555, and 670 nm for NOMAD and 158 pairs at each of the 6 wavelengths for MOBY. Sun zenith angle varied from 3 to 58 degrees, view zenith angle from 22 to 58 degrees, and relative azimuth angle from 75 to 180 degrees, i.e., the match-up data encompassed a wide range of geometry conditions. Figure 33 gives the geographic location of the match-up data. Most of the points are located between 60S and 60N in the Atlantic and Pacific Oceans, and in coastal regions of the United States, but the Indian Ocean and the Mediterranean Sea are also sampled. Oligotrophic (e.g., Tropical Pacific) to productive (e.g., Patagonia shelf, Benguela current) biological provinces, Case 1 and Case 2 waters are represented in the match-up

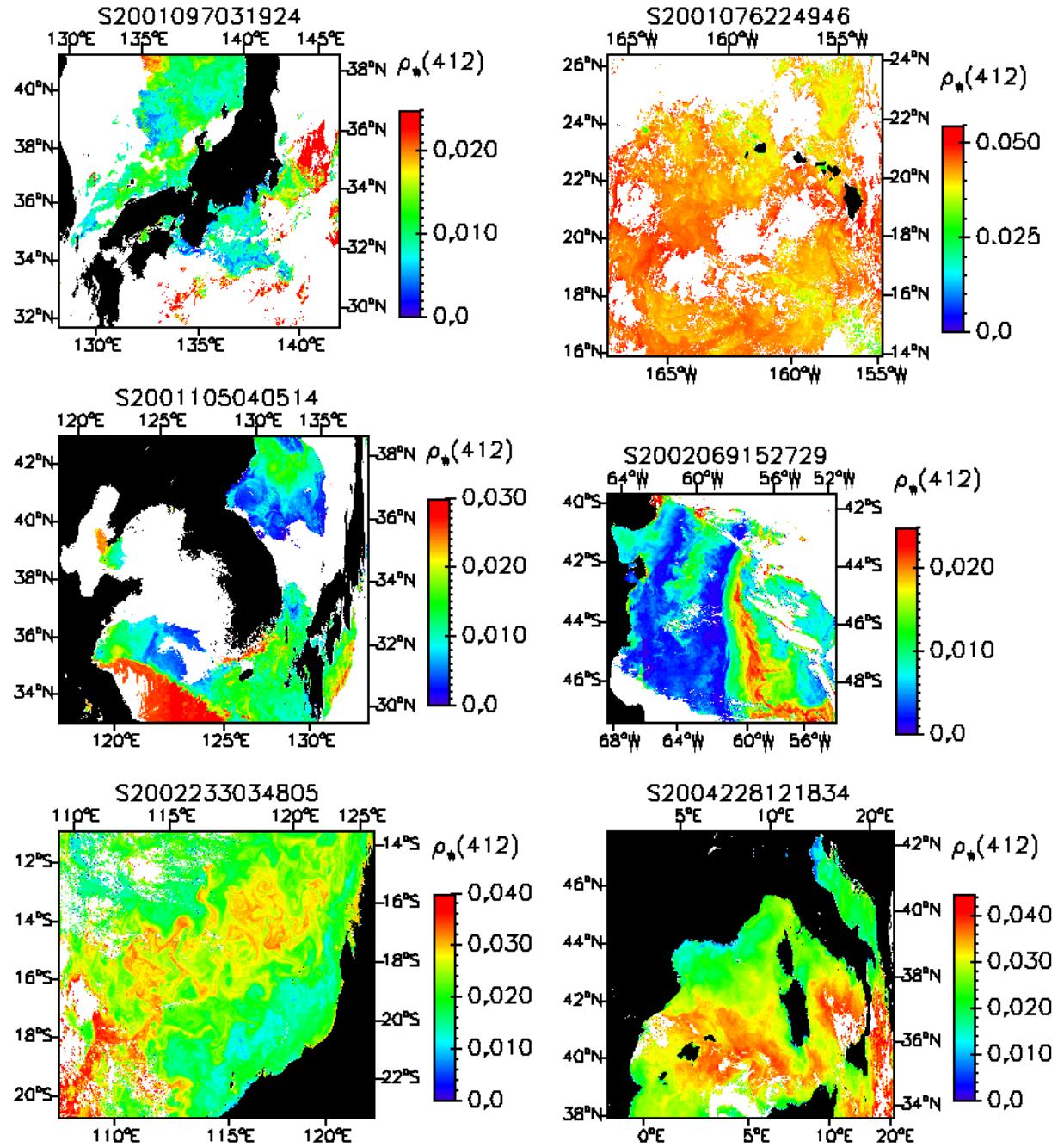


Figure 29: Marine reflectance imagery at 412 nm derived from SeaWiFS data using the Bayesian methodology. Top left: Sea of Japan and Northwest Atlantic. Top right: Pacific Ocean around Hawaii. Middle left: Sea of Japan and China sea. Middle right: Southwest Atlantic Ocean off Argentina. Bottom left: Indian Ocean off Northwest Australia. Bottom right: Central Mediterranean Sea.

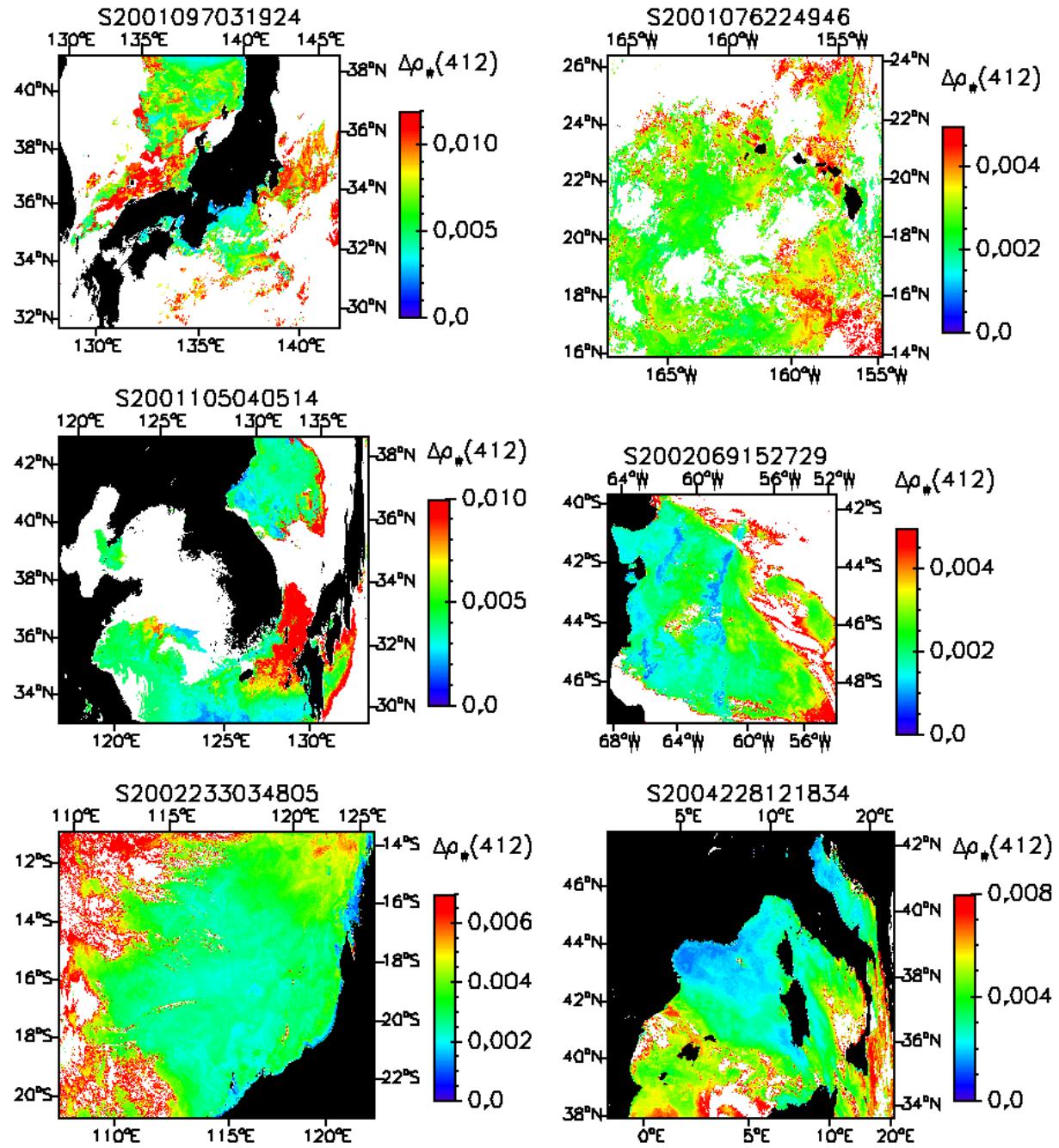


Figure 30: Uncertainty associated with marine reflectance retrievals at 412 nm for images of Figure 29.

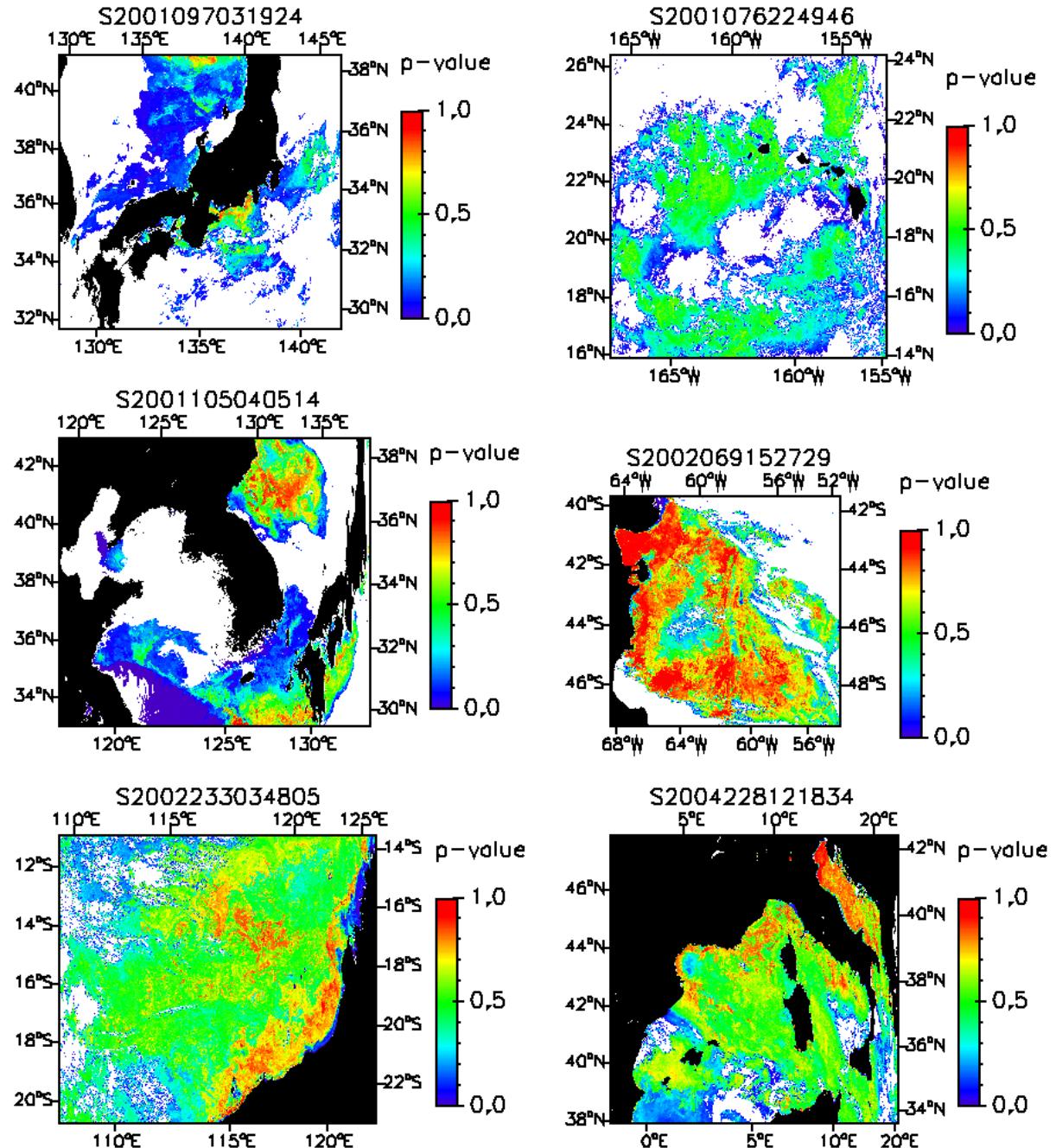


Figure 31: P-value (retrieval quality index) associated with marine reflectance estimates at 412 nm for images of Figure 29.

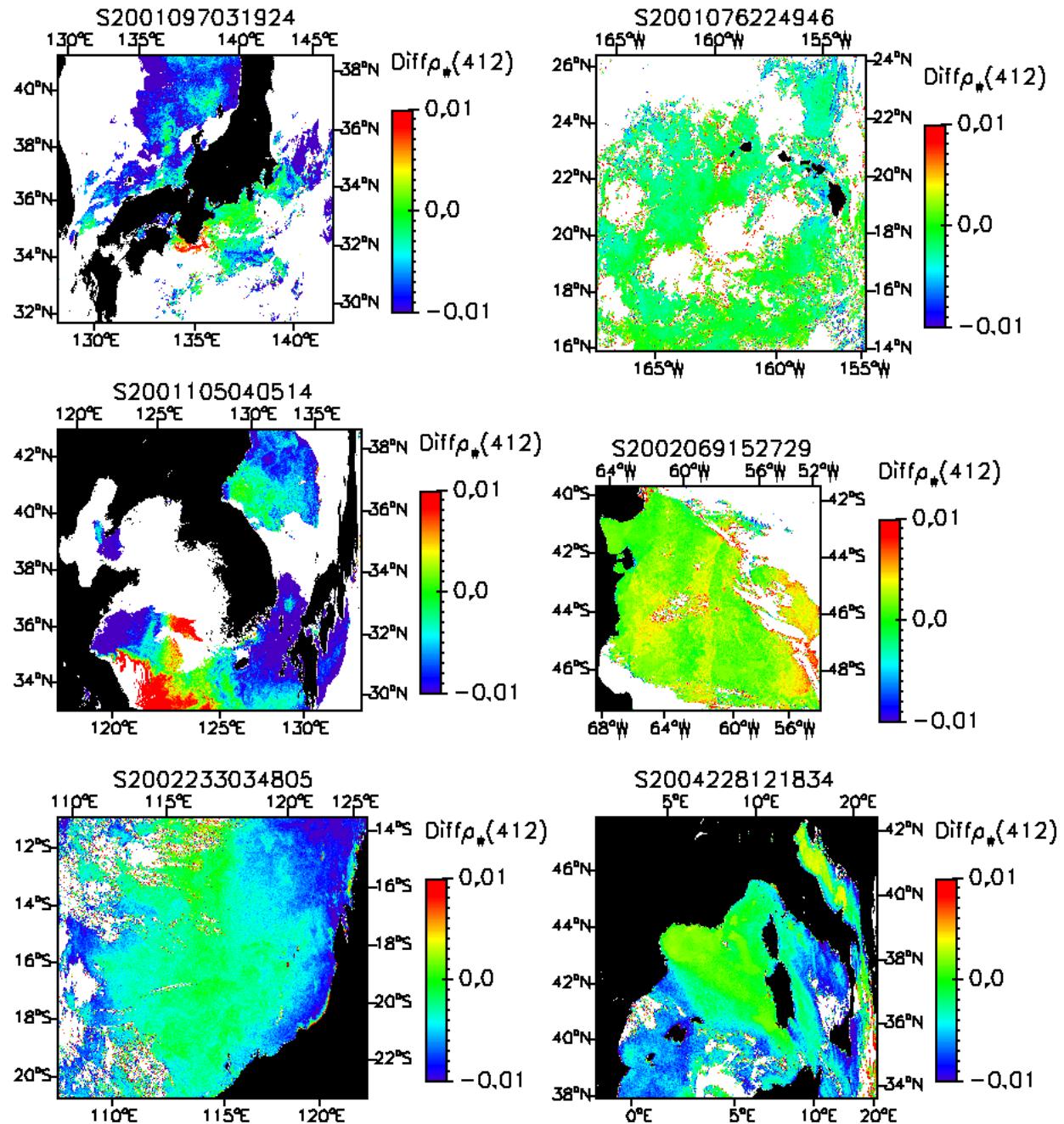


Figure 32: Difference between marine reflectance estimates at 412 nm from SeaDAS and the Bayesian methodology for images of Figure 29.

data, as well as various types of aerosols, (e.g., maritime in the open ocean, continental and pollution-type in coastal regions).

Figure 34 displays scatter plots of estimated versus measured marine reflectance for the MOBY and NOMAD data sets (top and bottom, respectively), and Tables 9 and 10 give the comparison statistics, in terms of coefficient of determination  $r^2$ , bias (estimated minus measured values) and RMS difference. In the scatter plots, the uncertainties associated with the marine reflectance retrievals are also displayed.

For the comparison using MOBY data (Figure 34, top), the scatter is relatively small, but the Bayesian estimates are biased high at 412, 443, and 490 nm, by 0.004, 0.002, and 0.001, respectively (Table 9). This may be explained in view of the theoretical performance (Section 7) and the aerosols prevailing at the MOBY site. According to Smirnov et al. (2003), at the Lanai AERONET site, located just a few kilometers from MOBY, the Angstrom exponent  $\alpha$  characterizing the spectral dependence of the aerosol optical thickness between 400 and 870 nm ( $\tau_a \approx \lambda^{-\alpha}$ ) exhibit average values of about 0.7. This would correspond in our modeling to a mixture of WMO maritime and continental aerosols since these aerosols have an Angstrom exponent of 0.1 and 1.2, respectively (urban aerosols are unlikely in Lanai). For such mixtures, Figure 13 indicates that positive biases, by a few 0.001, are expected theoretically. These biases may be reduced by algorithm regionalization, i.e., by taking into consideration a priori information on the aerosol properties prevailing in the region considered, which can be obtained from measurements or simulations from atmospheric transport models (see next section). Note, however, that the Bayesian estimates and the measurements generally agree within uncertainties, as determined from the posterior distribution. At 443 nm, the RMS difference is 0.003 or 9.2%, which falls a bit short of the requirements of 0.001 – 0.002 or 5% in clear waters for biological applications (Gordon et al., 1997). At 510, 555, and 670 nm, the  $r^2$  values are small, due to the lack of variability in the marine reflectance (Case 1, oligotrophic waters).

For the comparison using NOMAD data (Figure 34, bottom), the Bayesian estimates are less biased in magnitude than when using MOBY data, but the scatter is larger. The biases (higher Bayesian values) are 0.002 or 7.8% at 412 nm, 0.001 or 5% at 443 nm, and smaller at the other wavelengths (< 0.0001 at 670 nm) (Table 8-3). They represent a small component of the RMS errors, which decrease from 0.0059 at 412 nm to 0.0026 at 555 nm and to 0.0012 at 670 nm. These RMS errors are comparable with those obtained for the SeaDAS algorithm by the NASA Ocean Biology Processing Group and available from their web site using a much larger match-up data set sampling a wider range of conditions (4577 points), i.e., 0.0075 at 412 nm, 0.0045 at 555 nm, and 0.0019 at 670 nm. They are larger, however, than those obtained for the SeaDAS algorithm at the BOUSSOLE site, i.e., 0.0045 at 412 nm, 0.0012 at 555 nm, and 0.0003 at 670 nm (Antoine et al., 2008), but in this case, like for the MOBY match-up data set, the sampling is limited to a single site.

One cannot conclude, however, based on the analysis of such limited match-up data, whether or not the Bayesian technique performs better, in terms of accuracy, than the SeaDAS algorithm. In the previous discussion of imagery (Section 8), evidence was provided that, in general, agreement between marine reflectance estimated by the Bayesian technique and the SeaDAS algorithm may occur in one part of an image, but not everywhere in the same image. The lack of comprehensive evaluation data set emphasizes the importance of generalization in developing inversion schemes for global application, i.e., in our Bayesian approach, proper approximation of the forward operator.

## 8 Summary and Conclusions

The inverse ocean-color problem, i.e., the retrieval of spectral marine reflectance from spectral TOA reflectance measurements, has been examined in a Bayesian context. This is motivated by the ill-posed nature of the problem (many possible antecedents), which stems first from the complexity of the forward operator that relates the variable to recover to the set of measurements and second from the noise in the measure-

NOMAD blue, MOBY red

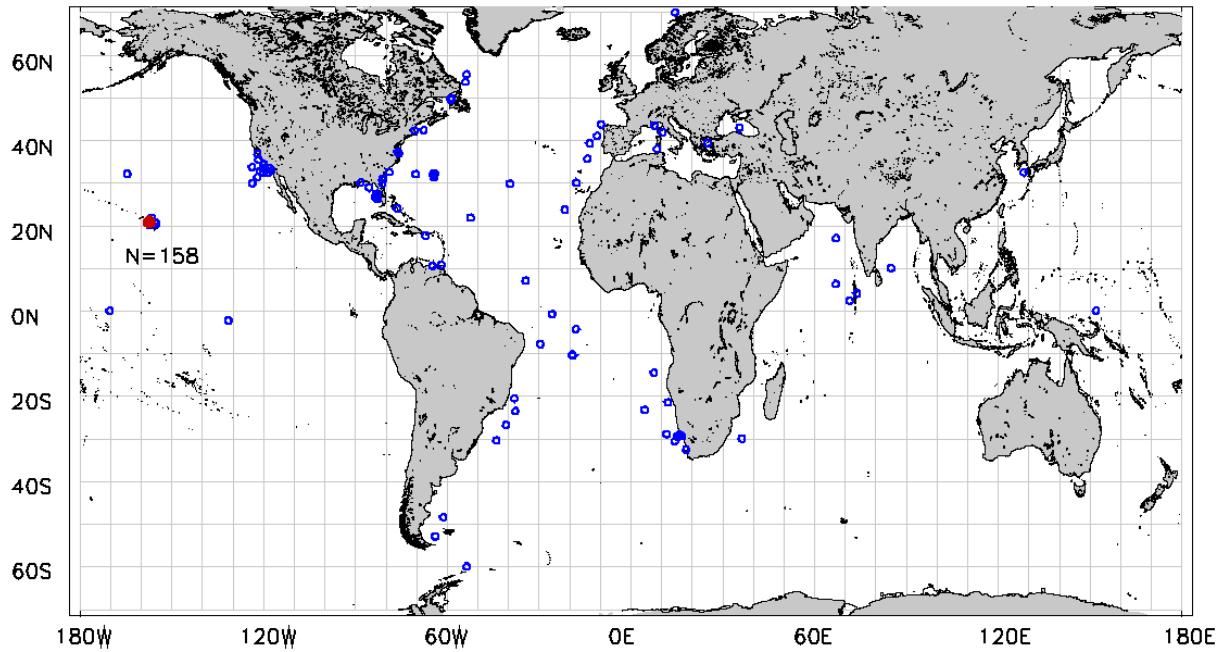


Figure 33: Geographic location of the NOMAD and MOBY match-up data sets.

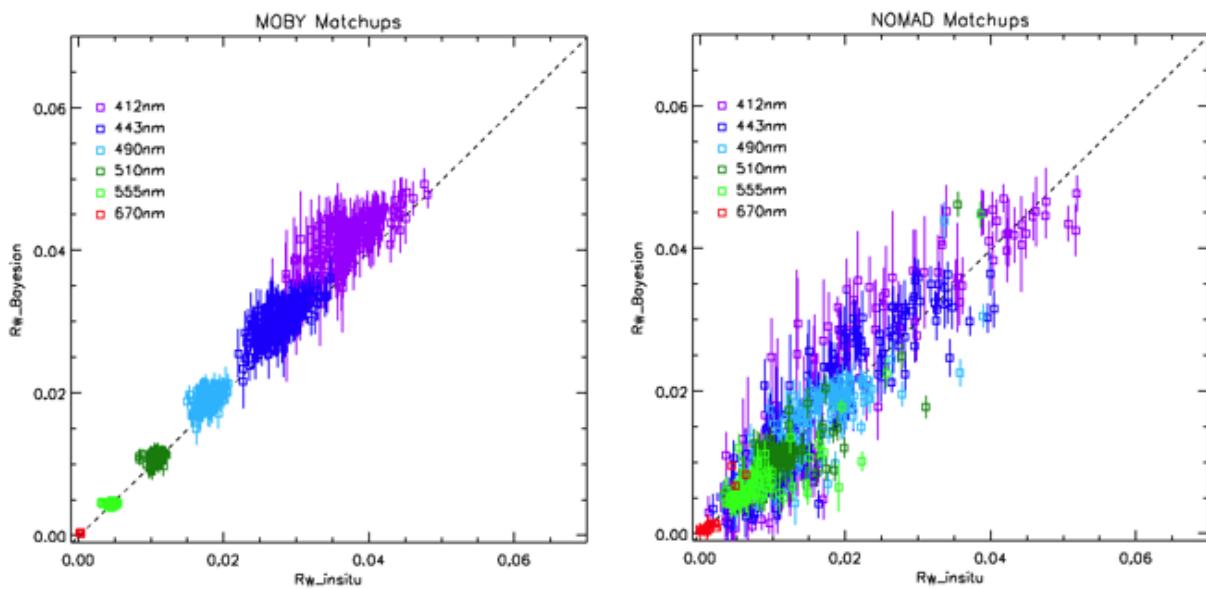


Figure 34: Estimated versus measured marine reflectance for MOBY and NOMAD match-up data sets.

Table 8: List of SeaWiFS LAC images processed and analyzed in the study.

Image Name	Date (mm/dd/yy)	Time (GMT)	Image Size (pixels*)	Oceanic Region
S1999045100113	02/14/99	10:01	1000 × 1000	Agulhas, Benguela, and South Atlantic Currents
S2001076224946	03/17/01	22:50	1050 × 960	Central Tropical Pacific North Equatorial Current
S2001097031924	04/07/01	03:19	950 × 980	Sea of Japan, Northwest Atlantic
S2001105040514	04/15/01	04:05	1000 × 1080	Sea of Japan, Yellow Sea, East China Sea
S2002069152729	03/10/02	15:28	800 × 800	Argentine Sea, Brazil Current
S2002233034805	08/11/02	03:48	1000 × 1100	East Indian Ocean, South Equatorial Current
S2004228121834	08/16/04	12:18	1000 × 1000	Western Mediterranean Sea

\*Rows (across-track) × columns (along-track).

Table 9: Comparison statistics of marine reflectance estimated by the Bayesian technique and measured in situ (MOBY match-up data set).

$\lambda$ (nm)	Average $\rho_w$	$r^2$	Bias	RMS Difference	No. Points
412	0.03687	0.599	0.00407	0.00491	132
443	0.02778	0.485	0.00236	0.00306	144
490	0.01801	0.201	0.00112	0.00159	144
510	0.01071	0.034	0.00010	0.00072	113
555	0.00464	0.001	-0.00018	0.00040	129
670	0.00034	0.003	-0.00001	0.00005	28
All	0.01639	0.987	0.00124	0.00247	690

Table 10: Same as Table 9, but NOMAD data set.

$\lambda$ (nm)	Average $\rho_w$	$r^2$	Bias	RMS Difference	No. Points
412	0.02049	0.838	0.00164	0.00593	158
443	0.01799	0.806	0.00088	0.00449	158
490	0.01545	0.670	-0.00025	0.00346	158
510	0.01155	0.587	-0.00049	0.00301	158
555	0.00712	0.722	-0.00074	0.00258	158
670	0.00133	0.820	-0.00008	0.00121	158
All	0.01418	0.852	0.00023	0.00403	948

ments. By defining the general solution of the inverse problem as a probability distribution (the posterior distribution), the Bayesian paradigm allows one to quantify the likelihood of encountering specific values of the input variable (marine reflectance) given the observed output variable (TOA reflectance or, in our modeling, the TOA reflectance corrected for molecular scattering effects).

The Bayesian approach makes it possible to incorporate known constraints of the marine reflectance (i.e., correlation between components) and to account for the varied sources of uncertainty (i.e., measurement noise, radiation transfer modeling errors). Importantly, it also permits the construction of reliable multi-dimensional confidence domains of the retrieved marine reflectance. These confidence domains are specific to each sample and can be constructed for any probability value. Specifically, the mean and covariance of the posterior distribution are computed. These quantities provide, for each pixel, an estimate of the marine reflectance and a measure of its uncertainty. The  $p$ -value, which quantifies how the TOA observation fits the forward model value, is also computed, allowing one to identify situations for which observation and model are incompatible. Thus the methodology offers the means to analyze and interpret ocean-color imagery in view of confidence limits and model adequacy, on a pixel-by-pixel basis. This is a definite advantage compared with standard atmospheric correction techniques, which rely essentially on evaluation against too few in situ measurements for accuracy assessments.

The definitions of the Bayesian solution and of the inverse applications require a prior distribution on the oceanic and atmospheric parameters, and a noise distribution. The prior distribution on the marine reflectance has been defined by using the NOMAD and AERONET-OC data sets for marine reflectance. The data sets include situations of Case1 and Case 2 waters, but not situations of very turbid waters containing sediments (e.g., estuarine waters). Our choice of using data instead of simulations for the marine reflectance was deliberate, dictated by the fact that models do not take fully into account the natural correlations between the intervening optical parameters (it is desirable not to introduce assumptions on the variable to retrieve). This makes it more difficult, however, to discretize properly the marine reflectance space (data not uniformly distributed, missing values). Since relatively little independent information (i.e., not resulting from inversion of satellite observations) is available about the global distribution of the atmospheric parameters and marine reflectance over the ocean, the prior distributions were considered uniform over the data sets, except for the aerosol optical thickness, which, based on a comprehensive in situ data set, was assumed to be distributed log-normally. For the noise, a normal distribution with zero mean was used, and its diagonal covariance matrix was determined by comparing a large and diverse ensemble of TOA reflectance extracted from actual satellite imagery with the output from the theoretical forward operator.

The theoretical inverse applications were approximated numerically based on extensive simulated data. The inverse applications models are all defined as piecewise constant, or piecewise linear, functions over the cells of a common partition of the space of TOA reflectances. The choice of models for the numerical approximations was based on several considerations, i.e., models fast in execution, convenience to approximate the conditional covariance (a second order quantity), and detection of abnormal cases (limitation of the forward model) by the  $p$ -value. Partition-based models are suitable for these purposes ; the partition is hierarchical, with hierarchy induced by a perfect binary tree, which drastically reduces the computational cost of determining cell membership.

The inverse models were constructed for application to SeaWiFS imagery. Theoretical performance for this ocean-color sensor is good globally, i.e., on average over all the geophysical conditions and geometries considered, with negligible biases and standard deviation decreasing from 0.004 at 412 nm to 0.001 at 670 nm. Errors are smaller, however, for geometries that avoid Sun glint contamination and minimize air mass and aerosol influence. For example the standard deviation is reduced to 0.002 – 0.003 in the blue when the Sun and view zenith angles are 30 degrees and the relative azimuth angle is 120 degrees. Performance is degraded in the presence of Sun glint, but remains acceptable (< 0.005 at 412 nm) in many situations. Errors increase with increasing optical thickness, reaching 0.009, 0.006, and 0.002 at 412, 555, and 670 nm, respectively, when the aerosol optical thickness at 550 nm is in the range 0.4 – 0.6. With respect to

aerosol type, the larger standard deviations are obtained for atmospheres with continental or urban aerosols, and the smaller for atmospheres with maritime aerosols (the case over most of the open ocean). Biases are negative for mixtures dominated by the urban type, and positive for mixtures dominated by continental aerosols, by a few 0.001 in magnitude in the blue and smaller at longer wavelengths, and they are small for mixtures dominated by the maritime type. Importantly, the estimated uncertainty (conditional covariance) is consistent with the inversion error, i.e., it provides a good measure of uncertainty.

Application to actual SeaWiFS imagery yielded marine reflectance fields with realistic features and patterns in view of the current knowledge of ocean circulation and biogeochemistry. The retrieved fields exhibited good continuity near clouds, and they did not exhibit significant correlation with the corresponding fields of atmospheric variables. Uncertainty generally remained within  $\pm 0.003$  in the blue, except near clouds and where the aerosol optical thickness was large. The  $p$ -value was mostly above 0.05, but often above 0.5, indicating good compatibility between forward model and observation. Unacceptable  $p$ -values ( $< 0.01$ ) were encountered in certain regions, in particular coastal regions influenced by river outflow (optically complex waters containing sediments). For those regions, the retrievals should be discarded. Compared with the marine reflectance fields obtained from the SeaDAS algorithm, the Bayesian fields do not exhibit systematic biases, but they are less noisy, and they have different de-correlation scales. Unlike the Bayesian estimates, the SeaDAS values were too low, sometimes negative, in the presence of absorbing aerosols. Note that standard, yet very conservative flags were applied to the selected SeaWiFS imagery to eliminate observations with clouds, Sun glint, and too high aerosol reflectance, i.e., the Bayesian methodology was not tested for such situations. In view of the theoretical performance, however, reasonable estimates are expected in the presence of Sun glint and fairly high aerosol optical thickness, and possibly in the presence of thin clouds since they may be interpreted as aerosols. Those situations will be examined in a future study.

Compared with marine reflectance measurements, the Bayesian estimates exhibit RMS differences of 0.005, 0.0004, and 0.00005 at 412, 555, and 670 nm (MOBY data set) and 0.006, 0.003, and 0.001, respectively (NOMAD data set). These values are comparable with those expected from the theoretical analysis of performance and the values obtained with the SeaDAS algorithm during various evaluation activities. The Bayesian estimates, however, are biased high at the MOBY site, which was plausibly explained by the aerosols prevailing at the site and the algorithm performance for those aerosols. One cannot conclude, from examining such a limited match-up data set, whether the Bayesian methodology is more accurate, but there is evidence, from analyzing the marine reflectance imagery, of better performance in the presence of absorbing aerosols.

The performance of the Bayesian methodology depends critically on the characteristics of the prior distributions (i.e., how they are specified). Due to lack of information, the distributions were taken as uniform over the space of the various variables, except for the aerosol optical thickness. But one could have taken into account that Case 2 waters and continental and pollution aerosols are more likely to be encountered in coastal regions. Such information may help to constrain the Bayesian solution and, therefore, improve retrieval accuracy. It may originate from various (independent) sources, in particular simulations by global numerical models of the atmosphere and ocean. These models predict, regionally, the temporal variability of key variables in the forward modeling, for example the likelihood of encountering a certain aerosol type and vertical profile or a certain chlorophyll concentration (from which one may deduce some information about marine reflectance variability). Due to the ill-posed nature of the inverse problem, this “regionalization” aspect is key to improving performance in situations difficult to deal with, such as absorbing aerosols and optically complex waters. Since these situations occur in biologically productive regions, in general the coastal zone, the expected gain in accuracy for biogeochemistry studies would be significant. “Regionalization” of the Bayesian methodology, as well as extending the methodology to other satellite sensors, and further evaluation, will be addressed in future work.

## Acknowledgments

Funding for this work has been provided by NASA under various grants. The technical support of Mr. John McPherson from the Scripps Institution of Oceanography, University of California at San Diego, is gratefully acknowledged. The authors also thank the NASA OBPG, the NOAA MOBY Project, and the NASA AERONET-OC Project for making available the satellite and in-situ data sets used in the study.

## Appendix A Missing data inference

To illustrate the methodology, consider the problem of estimating the (missing) value of  $\rho^{670}$  based on measurements of  $\rho^{412}, \dots, \rho^{555}$ , where  $\rho^\lambda$  denotes the marine reflectance in spectral band  $\lambda$ . Suppose that at our disposal is a random sample of size  $n$  of complete observations  $(\rho_i^{412}, \dots, \rho_i^{670})$ , for  $i = 1, \dots, n$ . The procedure consists in estimating the conditional expectation  $\mathbb{E} [\rho^{670} | \rho^{412}, \dots, \rho^{670}]$  from the  $n$  complete data. The rationale behind this estimate is that the conditional expectation is the best approximation of  $\rho^{670}$  that can be constructed based solely on the information conveyed by  $\rho^{412}, \dots, \rho^{555}$ , where optimality is understood in the sense of the average quadratic loss criterion. The different steps of the

Estimating the conditional expectation from the data is a standard nonparametric regression estimation problem. Numerous techniques exist for this purpose (see, e.g., Gyorfi et al., 2002) and in this work, we considered a  $k$ -nearest neighbour regression estimate. The regression estimate is constructed using the  $n$  complete data (here  $n = 729$ ) and is next applied to each spectrum containing one missing value. The algorithm is detailed in Table 11 for the estimation of the missing value of  $\rho^{670}$  from  $\rho^{412}, \dots, \rho^{555}$ . The procedure is repeated for each of the 6 channels where one value is missing. This results in a complete data set of 2,651 Case I marine reflectance spectra.

Table 11:  $k$ -nearest neighbor ( $k$ -NN) algorithm for the estimation of the missing values in the marine reflectance spectra. Illustration for the estimation of the marine reflectance at 670 nm based on the measurements at smaller wavelengths.

***k*-NN Algorithm for Missing Value Estimation: Estimation of  $\rho^{670}$  from  $\rho^{412}, \dots, \rho^{555}$**

1. **Input:**  $\rho := (\rho^{412}, \dots, \rho^{555})$  and the complete data  $(\rho_i^{412}, \dots, \rho_i^{670})$  for  $i = 1, \dots, n$ .

2. For each complete data  $(\rho_i^{412}, \dots, \rho_i^{670})$ , compute its distance  $d_i$  to  $\rho$  according to:

$$d_i = \left[ (\rho_i^{412} - \rho^{412})^2 + \dots + (\rho_i^{555} - \rho^{555})^2 \right]^{\frac{1}{2}}.$$

3. Sort the distances in increasing order:  $d_{(1)} \leq \dots \leq d_{(n)}$ .

4. Select the  $k$  observations  $\rho_{(1)}^{670}, \dots, \rho_{(k)}^{670}$  corresponding to the  $k$  smallest distances of step 3.

5. **Output:** the estimate  $\hat{\rho}^{670}$  as the average of the observations selected at step 4, i.e.,

$$\hat{\rho}^{670} = \frac{1}{k} \sum_{j=1}^k \rho_{(j)}^{670}.$$

## Appendix B Tree-based partition rules

To guarantee convergence of a density or regression model, for instance, based on a partition, typically the number of cells must go to infinity while the cells must shrink at an appropriate rate ; see, e.g., [Lugosi and Nobel \(1999\)](#); [Nobel \(1996\)](#). Several techniques have been developed to try to infer an optimal set of splitting rules (or even an optimal tree) from the data, especially to prevent over-fitting ; see [Breiman et al. \(1984\)](#) and the literature on random forests ([Breiman, 2001](#)). As exposed above, to keep the execution time of the models low, the partition is induced by a perfect binary tree. To determine suitable splitting rules, we first simulate a number  $n$  of TOA reflectances  $y_1, \dots, y_n$  according to model (2.7). The splitting rules are then defined recursively as follows. For each node to split, we select the axis  $j$  such that the  $j^{\text{th}}$  component of the simulated data has maximal variance. Next, the split threshold  $\delta$  is set as the median of the component of the data corresponding to this axis. The procedure is then repeated successively, starting from the root node to the leaves of the tree, until all the splitting rules are computed. The whole algorithm is summarized in Table 12. Note that since the split threshold is taken as the median along a certain axis, at each node, a data point has an equal probability of being moved to the left child as to the right child. Therefore all the cells of the partition have equal probability content, i.e., the probability that  $y$  falls in a given cell of the partition is constant and equal to  $1/2^K$  for a tree of depth  $K$ , since it contains  $2^K$  leaves.

In the present work, we have taken  $K = 17$  which yields a partition with  $2^K = 131,072$  cells. The number  $n$  of simulated points used to determine the splitting rules is 140 millions. As a final comment, the algorithm is not run on the canonical basis of  $\mathbb{R}^d$ , as described herein for simplicity, but on the basis of the eigenvectors of the covariance matrix of  $y$ .

Table 12: Construction algorithm of the tree structured partition associated with a perfect binary tree of depth  $K$ .

### Construction algorithm of the tree-structured partition

1. **Input:** Data  $y_1, \dots, y_n$  and depth  $K$ .
2. Associate the  $n$  data to the root node.
3. Compute the variances on each axes.
4. Select the axis  $j$  corresponding to the maximal variance.
5. Sort the  $j^{\text{th}}$  component of the data in ascending order:  $y_{(1)}^j \leq y_{(2)}^j \leq \dots \leq y_{(n)}^j$ , and set the split value as the median of the (univariate) ordered sample, e.g., as  $\delta = \frac{1}{2} (y_{(n/2)}^j + y_{(n/2+1)}^j)$  if  $n$  is odd.
6. Move the data whose  $j^{\text{th}}$  component is lower than  $\delta$  to the left child, and otherwise to the right child.
7. Repeat the above procedure on the left and right children until all the nodes up to depth  $K$  are constructed.
8. **Output:** A perfect binary tree with depth  $K$  and the splitting rules of the form  $(j, \delta)$  for each internal node.

## References

- Antoine, D., F. d'Ortenzio, S. B. Hooker, G. Bcu, B. Gentilli, D. Tailliez, and A. J. Scott (2008). Assessment of uncertainty in the ocean reflectance determined by three satellite ocean color sensors (meris, seawifs, and modis-a) at an offshore site in the mediterranean sea (boussole project). *J. Geophys. Res.*. doi:10.1029/2007JC004472.
- Antoine, D. and A. Morel (1999). A multiple scattering algorithm for atmospheric correction of remotely sensed ocean colour (meris instrument): principle and implementation for atmospheres carrying various aerosols including absorbing ones. *Int. J. Remote Sen.* 20, 1875–1916.
- Bailey, S., B. Franz, and P. Werdell (2010). Estimation of near-infrared water-leaving reflectance for satellite ocean color data processing. *Opt. Express* 18, 7521–7527.
- Bélanger, S., J. Ehn, and M. Babin (2007). Impact of sea ice on the retrieval of water-leaving reflectance, chlorophyll a concentration and inherent optical properties from satellite ocean color data. *Rem. Sen. Environ.* 111, 51–68.
- Biau, G., B. Cadre, and B. Pelletier (2008). Exact rates in density support estimation. *Journal of Multivariate Analysis* 99, 2185–2207.
- Bissantz, N., T. Hohage, and A. Munk (2004). Consistency and rates of convergence of nonlinear tikhonov regularization with random noise. *Inverse Problems* 20, 1773–1789.
- Brajard, J., C. Jamet, C. Moulin, and S. Thiria (2006). Use of a neuro-variational inversion for retrieving oceanic and atmospheric constituents from satellite ocean color sensor: Application to absorbing aerosols. *Neural Networks* 19, 178–185.
- Breiman, L. (2001). Random forests. *Machine Learning* 45, 5–32.
- Breiman, L., J. Friedman, R. Olshen, and C. Stone (1984). *Classification and Regression Trees*. Chapman and Hall.
- Buriez, J. and Y. Fouquart (1980). Generalization of the curtis-godson approximation to inhomogeneous scattering atmospheres. *J. Quant. Spectrosc. Radiat. Transfer* 24, 407–419.
- Chomko, R. and H. Gordon (1988). Atmospheric correction of ocean color imagery: use of the junge power-law aerosol size distribution with variable refractive index to handle aerosol absorption. *Appl. Opt.* 37, 5560–5572.
- Clark, D. K., M. A. Yarbrough, M. Feinholz, S. Flora, W. Broenkow, Y. S. Kim, B. C. Johnson, S. W. Brown, M. Yuen, and J. L. Mueller (2003). Moby, a radiometric buoy for performance monitoring and vicarious calibration of satellite ocean color sensors: measurement and data analysis protocols. NASA Tech. Memo. 2003 211621/Rev4, vol. VI, edited by J. L. Mueller, G. S. Fargion, and C. R. McClain, 139 pp., NASA GSFC, Greenbelt, Md.
- Cotter, S., M. Dashti, J. Robinson, and A. Stuart (2009). Bayesian inverse problems for functions and applications to fluid mechanics. *Inverse Problems* 25, 115008.
- Cotter, S., M. Dashti, and A. Stuart (2010). Approximation of bayesian inverse problems. *SIAM Journal of Numerical Analysis* 48, 322–345.
- Cox, C. and W. Munk (1954). Measurement of the roughness of the sea surface from photographs of the suns glitter. *J. Opt. Soc. Am.* 44, 838–850.
- Deschamps, P., M. Herman, and D. Tanré (1983). Modélisation du rayonnement réfléchi par latmosphère et la terre, entre 0.35 et 4  $\mu\text{m}$ . *ESA Contract Report No. 4393/8C/F/DD(SC)*. 156 pp.
- Deuzé, J., M. Herman, and R. Santer (1989). Fourier series expansion of the transfer equation in the atmosphere-ocean system. *J. Quant. Spectrosc. Radiat. Transfer* 41, 483–494.
- Doney, S. C., D. M. Glover, S. J. McCue, and M. Fuentes (2003). Mesoscale variability of sea-viewing wide field-of-view sensor (seawifs) satellite ocean color: Global patterns and spatial scales. *J. Geophys. Res.* 108. doi:10.1029/2001JC000843.
- Engl, H., K. Hanke, and A. Neubauer (1996). *Regularization of Inverse Problems*. Kluwer.

- Frouin, R., P.-Y. Deschamps, L. Gross-Colzy, H. Murakami, and T. Nakajima (2006). Retrieval of chlorophyll-a concentration via linear combination of global imager data. *J. Oceanogr.* 62, 331–337.
- Frouin, R. and B. Pelletier (2007). Fields of non-linear regression models for atmospheric correction of satellite ocean-color imagery. *Rem. Sen. Environ.* 111, 450–465.
- Frouin, R., M. Schwindling, and P.-Y. Deschamps (1996). Spectral reflectance of sea foam in the visible and near infrared: In-situ measurements and remote sensing implications. *J. Geophys. Res.* 101, 14361–14371.
- Garver, S. and D. Siegel (1997). Inherent optical property inversion of ocean color spectra and its biogeochemical interpretation. *J. Geophys. Res.* 102, 18607–18625.
- Goody, R. (1964). *Atmospheric radiation, 1. Theoretical basis*. Oxford University Press.
- Gordon, H., O. Brown, R. Evans, J. Brown, R. Smith, and K. Baker (1988). A semi-analytical radiance model of ocean color. *J. Geophys. Res.* 93, 10909–10924.
- Gordon, H., P.-Y. Deschamps, M. Wang, and R. Frouin (2010). Atmospheric diffuse transmittance. Technical report, IOCCG report No. 10 “Atmospheric correction for remotely-sensed ocean color products”.
- Gordon, H., T. Du, and T. Zhang (1997). Remote sensing of ocean color and aerosol properties: resolving the issue of aerosol absorption. *Appl. Opt.* 36, 8670–8684.
- Gordon, H. and B. Franz (2008). Remote sensing of ocean color: Assessment of the water-leaving radiance bidirectional effects on the atmospheric diffuse transmittance for seawifs and modis intercomparisons. *Rem. Sen. Environ.* 112, 2677–2685.
- Gross, L., S. Colzy, R. Frouin, and P. Genry (2007a). A general ocean color atmospheric correction scheme based on principal components analysis part i: Performance on case 1 and case 2 waters. In R. Frouin and Z. Lee (Eds.), *Proc. SPIE*, Volume 6680. doi:10/12.738508.
- Gross, L., S. Colzy, R. Frouin, and P. Genry (2007b). A general ocean color atmospheric correction scheme based on principal components analysis part ii: Level 4 merging capabilities. In R. Frouin and Z. Lee (Eds.), *Proc. SPIE*, Volume 6680. doi:10/12.738514.
- Gyorfi, L., M. Kohler, A. Krzyzak, and H. Walk (2002). *A Distribution-Free Theory of Nonparametric Regression*. Springer: New-York.
- Hansen, J. and L. Travis (1974). Light scattering in planetary atmospheres. *Space Sci. Rev.* 16, 527–610.
- Hu, C., K. Carder, and F. Muller-Karger (2000). Atmospheric correction of seawifs imagery over turbid coastal waters: A practical method. *Rem. Sens Environ.* 74, 195–206.
- Huebert, B., T. Bates, P. B. Russell, G. Shi, Y. J. Kim, K. Kawamura, G. Carmichael, and T. Nakajima (2003). An overview of ace-asia: Strategies for quantifying the relationships between asian aerosols and their climate impacts. *J. Geophys. Res.* 108. doi:10.1029/2003JD003550.
- Jamet, C., S. Thiria, C. Moulin, and M. Crépon (2005). Use of a neurovariational inversion for retrieving oceanic and atmospheric constituents from ocean color imagery: A feasibility study. *J. Atmos. Ocean Technol.* 22, 460–475.
- Kahn, R., J. Anderson, T. L. Anderson, T. Bates, F. Brechtel, C. M. Carrico, A. Clarke, S. J. Doherty, E. Dutton, R. Flagan, R. Frouin, H. Fukushima, B. Holben, S. Howell, B. Huebert, A. Jefferson, H. Jonsson, O. Kalashnikova, J. Kim, S.-W. Kim, P. Kus, W.-H. Li, J. Livingston, C. McNaughton, J. Merrill, S. Mukai, T. Murayama, T. Nakajima, P. Quinn, J. Redemann, M. Rood, P. Russell, I. Sano, B. Schmid, J. Seinfeld, N. Sugimoto, J. Wang, E. J. Welton, J.-G. Won, and S.-C. Yoon (2004). Environmental snapshots from ace-asia. *J. Geophys. Res.* 109. doi: 10.1029/2003JD004339.
- Kaipio, J. and E. Somersalo (2004). *Statistical and Computational Inverse Problems*. Springer: Berlin.
- Knobelgesse, K. D., C. Pietras, G. Fargion, M. Wang, R. Frouin, M. Miller, A. Subramaniam, and W. Balch (2004). Maritime aerosol optical thickness measured by handheld sunphotometers. *Rem. Sen. Environ.* 93, 87–106.
- Koepke, P. (1984). Effective reflectance of oceanic whitecaps. *Appl. Opt.* 23, 1816–1824.
- Kuchinke, C., H. Gordon, L. Harding, and K. Voss (2009). Spectral optimization for constituent retrieval in

- case 2 waters i: Implementation and performance. *Rem. Sen. Environ.* 113, 571–587.
- Land, P. and J. Haigh (1996). Atmospheric correction over case 2 waters using an iterative fitting algorithm. *Appl. Opt.* 35, 5443–5451.
- Lavender, S., M. Pinkerton, G. Moore, J. Aiken, and D. Blondeau-Patissier (2005). Modification of the atmospheric correction of seawifs ocean colour images over turbid waters. *Cont. Shelf Res.* 25, 539–555.
- Lee, Z. (2006). Remote sensing of inherent optical properties: Fundamentals, tests of algorithms and applications. Technical Report 5, IOCCG.
- Lenoble, J., M. Herman, J.-L. Deuzé, B. Lafrance, and D. Tanré (2007). A successive orders of scattering code for solving the vector equation of transfer in the earths atmosphere with aerosols. *J. Quant. Spec. Rad. Trans.* 107, 479–507.
- Lugosi, G. and A. Nobel (1999). Consistency of data-driven histogram methods for density estimation and classification. *Annals of Statistics* 27, 1830–1864.
- Malkmus, W. (1967). Random lorentz band model with exponential-tailed s-1 line intensity distribution function. *J. Opt. Soc. Am.* 53, 323–329.
- Morel, A. and S. Belanger (2006). Improved detection of turbid waters from ocean color sensors information. *Rem. Sen. Environ.* 102, 237–249.
- Morel, A. and B. Gentili (1993). Diffuse reflectance of oceanic waters, 2. bidirectional aspects. *Appl. Opt.* 32, 6864–6879.
- Morel, A. and S. Maritorena (2002). Bio-optical properties of oceanic waters: A reappraisal. *J. Geophys. Res.* 106, 7163–7180.
- Morel, A., K. Voss, and B. Gentili (1995). Bidirectional reflectance of oceanic waters: A comparison of modeled and measured upward radiance fields. *J. Geophys., Res.* 100, 13143–13150.
- Nicolas, J.-M., P.-Y. Deschamps, and R. Frouin (2001). Spectral reflectance of oceanic whitecaps in the visible and near infrared: Aircraft measurements over open ocean. *Geophys. Res. Lett.* 28, 4445–4448.
- Nobel, A. (1996). Histogram regression estimation using data dependent partitions. *Annals of Statistics* 24, 1084–1105.
- Oo, M., M. Vargas, A. Gilerson, B. Gross, F. Moshary, and S. Ahmed (2008). Improving atmospheric correction for highly productive coastal waters using the short wave infrared retrieval algorithm with water-leaving reflectance constraints at 412 nm. *Appl. Opt.* 47, 3846–3859.
- O'Reilly, J., S. Maritorena, B. Mitchell, D. Siegel, K. Carder, S. Garver, M. Kahru, and C. McClain (1998). Ocean colour chlorophyll algorithms for seawifs. *J. Geophys. Res.* 103, 24937–24953.
- Park, Y. and K. Ruddick (2005). Model of remote sensing reflectance including bidirectional effects for case 1 and case 2 waters. *Appl. Opt.* 44, 1236–1249.
- Pelletier, B. and R. Frouin (2004). Fields of nonlinear regression models for inversion of satellite data. *Geophysical Research Letters* 31. L16304, doi 10.1029/2004GL019840.
- Pelletier, B. and R. Frouin (2005). Remote sensing of phytoplankton chlorophyll-a concentration by use of ridge function feilds. *Applied Optics* 45(4), 784–798.
- Polonik, W. (1995). Measuring mass concentrations and estimating density contour clusters – an excess mass approach. *Annals of Statistics* 23, 855–881.
- Polonik, W. (1997). Minimum volume sets and generalized quantile processes. *Stochastic Processes and their Applications* 69, 1–24.
- Ruddick, K., F. Ovidio, and M. Rijkeboer (2000). Atmospheric correction of seawifs imagery for turbid and inland waters. *Appl. Opt.* 39, 897–912.
- Santer, R. and C. Schmechtig (2000). Adjacency effects on water surfaces: Primary scattering approximation and sensitivity study. *Appl. Opt.* 39, 361–375.
- Shroeder, T., I. Behnert, M. Schaale, J. Fischer, and R. Doerffer (2007). Atmospheric correction algorithm for meris above case-2 waters. *Int. J. Remote Sens* 28, 1469–1486.
- Siegel, D., M. Wang, S. Maritorena, and W. Robinson (2000). Atmospheric correction of satellite ocean

- color imagery: the black pixel assumption. *Appl. Opt.* 39, 3582–3591.
- Smirnov, A., B. N. Holben, O. Dubovik, R. Frouin, and I. Slutsker (2003). Maritime component in aerosol optical models derived from aeronet (aerosol robotic network) data. *J. Geophys. Res.* 108. doi:10.1029/2002JD002701.
- Stamnes, K., W. Li, B. Yan, H. Elde, A. Barnard, S. Pegau, and J. Stamnes (2007). Accurate and self-consistent ocean color algorithm: simultaneous retrieval of aerosol optical properties and chlorophyll concentrations. *Appl. Opt.* 42, 939–951.
- Steinmetz, F., P. Deschamps, and D. Ramon (2011). Atmospheric correction on the presence of sun glint: application to meris. *Optics Express* 19, 9783–9800.
- Stuart, A. (2010). Inverse problems: A bayesian perspective. *Acta Numerica* 19, 451–559.
- Stumpf, R., R. Arnone, R. Gould, P. Martilonich, and V. Ransibrahmanakul (2011). A partially coupled ocean-atmosphere model for retrieval of water-leaving radiance from seawifs in coastal water. In *Algorithm Updates for the Fourth SeaWiFS Data Reprocessing*, Volume 22 of *SeaWiFS Post-Launch Technical Series*.
- Tanré, D., M. Herman, P.-Y. Deschamps, and A. De Leffe (1979). Atmospheric modeling for space measurements of ground reflectances, including bidirectional properties. *Appl. Opt.* 18, 3587–3594.
- Tarantola, A. (2005). *Inverse Problem Theory*. SIAM.
- Wang, M. (2010). Atmospheric correction for remotely-sensed ocean-colour products. Technical Report 10, IOCCG.
- Wang, M., S. Son, and W. Shi (2009). Evaluation of modis swir and nir-swir atmospheric correction algorithm using seabass data. *Remote Sen. Environ.* 113, 635–644.
- Wang, M., T. Tang, and W. Shi (2007). Modis-derived ocean color products along the china east coastal region. *Geophys. Res. Lett.* L06611, doi:06610.01029/ 02006GL02859.
- Werdell, P. and S. Bailey (2005). An improved in-situ bio-optical data set for ocean color algorithm development and satellite data product validation. *Remote Sen. Environ.* 98, 122–140.
- WMO (1983). Radiation commission of iamap meeting of experts on aerosol and their climatic effects. WCP55, Williamsburg VA, 28-30.
- Zibordi, G., B. Holben, F. Mélin, D. DALimonte, J. Berthon, I. Slutsker, and D. Giles (2010). Aeronet-oc: an overview. *Can. J. Rem. Sens.* doi: 10.5589/m10-073.
- Zibordi, G., B. Holben, I. Slutsker, D. Giles, D. DALimonte, F. Mélin, J. Berthon, D. Vandemark, H. Feng, G. Schuster, B. Fabbri, S. Kaitala, and J. Seppala (2009). Aeronet-oc: A network for the validation of ocean color primary products. *Atmos. Oceanic Technol.* 26, 16341651.