



Cloud mask via cumulative discriminant analysis applied to satellite infrared observations: scientific basis and initial evaluation

U. Amato¹, L. Lavanant², G. Liuzzi³, G. Masiello^{3,4}, C. Serio^{3,4}, R. Stuhlmann⁵, and S. A. Tjemkes⁵

¹IAC/CNR, Napoli, Italy

²Meteo-France, DP, Centre de Meteorologie Spatiale BP 50747 22307 Lannion, France

³School of Engineering, University of Basilicata, Potenza, Italy

⁴CNISM, Research Unit of Potenza, University of Basilicata, Potenza, Italy

⁵EUMETSAT, Darmstadt, Germany

Correspondence to: C. Serio (carmine.serio@unibas.it)

Received: 20 April 2014 – Published in Atmos. Meas. Tech. Discuss.: 6 June 2014

Revised: 5 September 2014 – Accepted: 6 September 2014 – Published: 7 October 2014

Abstract. We introduce a classification method (cumulative discriminant analysis) of the discriminant analysis type to discriminate between cloudy and clear-sky satellite observations in the thermal infrared. The tool is intended for the high-spectral-resolution infrared sounder (IRS) planned for the geostationary METEOSAT (Meteorological Satellite) Third Generation platform and uses IASI (Infrared Atmospheric Sounding Interferometer) data as a proxy. The cumulative discriminant analysis does not introduce biases intrinsic with the approximation of the probability density functions and is flexible enough to adapt to different strategies to optimize the cloud mask. The methodology is based on nine statistics computed from IASI spectral radiances, which exploit the high spectral resolution of the instrument and which effectively summarize information contained within the IASI spectrum. A principal component analysis prior step is also introduced, which makes the problem more consistent with the statistical assumptions of the methodology. An initial assessment of the scheme is performed based on global and regional IASI real data sets and cloud masks obtained from AVHRR (Advanced Very High Resolution Radiometer) and SEVIRI (Spinning Enhanced Visible and Infrared Imager) imagers. The agreement with these independent cloud masks is generally well above 80 %, except at high latitudes in the winter seasons.

1 Introduction

Modern meteorological satellites carry infrared sensors onboard to sense the Earth emission spectrum at very high spectral resolution. These include, for example, AIRS (Atmospheric Infrared Sounder), IASI (Infrared Atmospheric Sounding Interferometer) and CrIS (Cross-track Infrared Sounder). All of these spectrometers are characterized by broadband spectral coverage (3.7–15.5 μm) and a spectral sampling rate in the range 0.25–2 cm^{-1} .

EUMETSAT (European Centre for the Exploitation of Meteorological Satellites) is also preparing for METEOSAT (Meteorological Satellite) Third Generation (MTG), which will carry the infrared sounder (IRS) at a hyperspectral resolution of 0.625 cm^{-1} wave numbers.

For the MTG programme, a twin configuration – comprising the MTG imaging satellite, called MTG-I, and the MTG sounding satellite, called MTG-S – has been selected and consolidated as a baseline. Therefore the MTG-S satellite will not carry an imager in the visible onboard; as a consequence the cloud-screening process of MTG-IRS spectral radiances has to rely on a stand-alone system. One could argue that MTG-IRS has itself imaging capability; however its spectral coverage is limited below $\approx 2200 \text{ cm}^{-1}$ and, therefore, it will miss the near-infrared and visible portion of the spectrum. The need to develop a stand-alone scene analysis for MTG-IRS has mostly motivated this work, in which the IASI instrument will be used as a proxy for MTG-IRS.

MTG-IRS is a hyperspectral sounder orbiting onboard a geostationary platform, and as such it is expected to further improve numerical weather prediction (NWP) forecast performance compared with that already reached with hyperspectral sensors on polar satellites. The high spectral resolution of new advanced infrared (polar) sensors has resulted in better coverage and significantly improved temperature and moisture soundings capabilities compared with past sounding instrumentation (e.g. Hilton et al., 2012). Infrared data, however, are frequently affected by clouds. Thus, observations must be processed for operational data assimilation and inversion for geophysical parameters, either by screening to remove cloud-contaminated soundings or by a so-called process of cloud clearing. It should be stressed that operational numerical weather prediction centres currently use cloud-screened or cloud-cleared data.

As said, IASI will be used as a proxy for MTG-IRS; therefore the methodology we present in this paper will be exemplified for this instrument. IASI has been developed in France by the Centre National d'Etudes Spatiales (CNES) and is onboard the Metop (Meteorological Operational Satellite) platform, a series of three satellites belonging to the EUMETSAT European Polar System (EPS). The instrument has a spectral coverage extending from 645 to 2760 cm^{-1} , which, with a sampling interval $\Delta\sigma = 0.25 \text{ cm}^{-1}$, gives 8461 data points or channels for each single spectrum. Data samples are taken at intervals of 25 km along and across track, each sample having a minimum diameter of about 12 km. Further details on IASI and its mission objectives can be found in Hilton et al. (2012).

Most cloud detection methods are based on the definition of some statistics and related statistical tests that are a measure of some attribute of the whole spectrum, of a suitable spectral interval, of radiance at two or very few wavelengths or even of a single radiance. For example it is widely acknowledged that clouds have a higher reflectance and are generally colder than the underlying surface, which motivated the use of visible and infrared regions for the discrimination. Earlier statistics have compared radiance observations with those calculated assuming clear sky conditions (based on radiative transfer (RT) models), as in the ISCCP algorithm (Rossow, 1989); essentially they label a pixel as cloudy if its measured radiance is less than the calculated clear-sky value by a certain amount that takes into account the variability of the latter. Such algorithms suffer from the drawback of a very large variability of clouds and of an underlying land surface, which makes estimating the status of the pixel inaccurate. Therefore it is nowadays preferred to develop statistics that involve two or more wavelengths and that are least sensitive to the underlying surface, especially on land, thereby being able to discard the estimate of radiance in clear sky conditions. On some occasions ancillary information is used coming from NWP models. It is useful to adopt several statistics to discriminate cloudy from clear sky

in specific conditions (e.g. cirrus or phase of the cloud) and to estimate a sort of probability map of the cloudy condition.

Most statistics of all present operational imagers onboard satellites are based on this principle, and they differ in the choice of the involved wavelengths or spectral interval. For example the CLAVR-x (Cloud Advanced Very High Resolution Radiometer Extended) cloud mask product for AVHRR (Advanced Very High Resolution Radiometer) (Heidinger, 2004), which supersedes CLAVR-1 (Stowe et al., 1999), includes a series of 12 tests based on the six AVHRR spectral radiances depending on the type of surface and, in some cases, on the time of the day. The MAIA cloud mask (Lavanant et al., 2007) developed for the AVHRR onboard Metop uses a similar approach. For MODIS (Moderate-Resolution Imaging Spectroradiometer) (Ackerman et al., 1998, 2008) nine statistics and corresponding tests are introduced, each of them involving two or three wavelengths at most; in the case of SEVIRI (Spinning Enhanced Visible and Infrared Imager) (EUMETSAT, 2009) about 10 statistics and corresponding tests are defined depending on the time of the day (daytime, sun glint, twilight, night-time) and on the surface type (land or sea). Often tests based on spatial coherence statistics are also introduced that rely on the different spatial correlation in images in clear and cloudy conditions (Sandhya et al., 2004). Besides statistics strictly based on radiance (thresholding, comparison of radiances or their ratios), many others have been developed aimed at picking particular features that are different in clear and cloudy conditions. For example the h_s index is introduced in Serio et al. (2000) based on correlation and cross-correlation of spectral radiances in the atmospheric infrared window. Many other methodologies for cloud detection have been developed in the last decade in the frameworks of statistics and machine learning. It is not the purpose of the present paper to discuss them here (we defer to Tapakis and Charalambides (2012) for a recent review) because they will not be used for the reasons that are going to be discussed. In addition, we will focus mostly on the problem of discriminating the presence of cloud contamination in the radiances, but not on cloud classification in terms of phase and top pressure and temperature.

One of the most appropriate approaches for cloud detection within the footprint or field of view of a sounding instrument is its co-location with a suitable imager radiometer, or even better to design and develop the sounder with a built-in imager. Both approaches have been developed for IASI. However, the built-in imager has only two channels and is not appropriate for cloud detection. Co-location of IASI footprint with AVHRR imagery has been used, for example, by Lavanant et al. (2007). However, technological constraints, such as those envisaged for MTG-IRS, or simply the need of real-time processing, as required, for example, in NWP applications, can severely limit the synergy between sounders and imagers. Thus, in some circumstances a stand-alone cloud detection strategy for a given sounder is not a choice but a constraint we are forced to deal with. In this

case, most limitations are related to the reduced spectral coverage, which for high-spectral-resolution infrared sounders can miss, for example, the near-infrared or visible portion of the spectrum. The fact that the spectral regions of interest for cloud detection are limited motivated this paper to use the very well consolidated physical arguments already available for present sensors to IASI, suitably adapted to its design features. Therefore in this paper some new or adapted from past sensor statistics will be introduced.

In addition a new cloud detection methodology (cumulative discriminant analysis, CDA) will be introduced that uses several arguments from the statistics framework; it can be fully considered as a classification methodology where a training data set is needed to set the thresholds for discriminating clouds from clear sky. CDA relies on the nonparametric empirical cumulative distribution function of each statistic in clear and cloudy conditions that has optimal properties from the statistical point of view. To choose the thresholds of the tests, a cost function is to be introduced to be minimized. Starting from the estimate of the type I and type II errors, CDA is flexible enough to adapt to the requirements of the user in terms of preferring the type I or type II error or a mixture of the two, also putting some constraints on the minimum acceptable rate for one of them. When the cost function is the sum of type I and type II errors, the methodology reverts to the classical discriminant analysis for a loss function of 0–1. It can also naturally handle receiver operating characteristics (ROC) graphs relying on sensitivity and specificity (see Fawcett, 2004) and other measures of a test's accuracy (e.g. F measure, G measure). The CDA is extended in this paper to more statistics assuming their independence. To partially overcome the approximation of independence, a principal component analysis prior transform has been applied as in Amato et al. (2008), giving rise to a very fast and accurate methodology. By its very construction, the methodology naturally provides a quality indicator of the retrieved status for each pixel (clear or cloudy).

The present paper is mostly concerned with the scientific basis of CDA. An initial evaluation of the cloud mask is performed with global and regional IASI data complemented with cloud masks derived from AVHRR and SEVIRI.

The paper is organized as follows. Section 2 deals with the data used for developing statistics, training and validating the cloud detection scheme. The series of statistics defined and used in the paper are described in Sect. 3. Section 4 includes the full statistical development of CDA. Experiments are shown in Sect. 5, while conclusions are drawn in Sect. 6.

2 Data sets

We have individuated and developed two data sets of IASI spectra to train, cross-check and validate the cloud detection scheme. These are referred to as IASI1 and IASI2, and their characteristics are listed below:

IASI1. It comprises 888 380 IASI spectra corresponding to a 12 h global acquisition on 17–18 November 2009 (Lavanant et al., 2011). The set is qualified for sky type (clear/cloudy) through the CMS (Centre de Meteorologie Spatiale, Lannion, France) cloud mask (Lavanant et al., 2011), which is based on the co-location of the IASI footprint with AVHRR imagery.

IASI2. It comprises 1 072 050 spectra corresponding to a 15 h global acquisition on 22–23 July 2007. Also in this case the set is qualified for sky type (clear/cloudy) through the CMS cloud mask.

Both sets have been endowed with skin temperature fields from the ECMWF (European Centre for Medium-range Weather Forecasts) analysis, co-located in space and time with the IASI footprints.

To take into account the possible dependence of the test statistics and of ancillary data (reference spectra and ECMWF skin temperature) on season and air mass type, we have defined 10 climatic zones. To increase statistics homogeneity, some models are also split into submodels according to the hemisphere (Northern Hemisphere, NH, or Southern Hemisphere, SH) or to the time of the day. The climate zones are listed in Table 1.

It is important to stress that the definition of the climatic zones 1 to 5 does not include sea ice or snow/ice in the case of land. For these models we assume the skin temperature to be above the freezing point of 273 K. The cases of sea ice and land ice/snow are also split with latitude. Also note that we have included Antarctica as a special climate zone because this region is covered with permanent ice.

Only for the tropical zone do the two data sets play a complementary role (one for training and the other one for validation). In the other climatologies some data sets are populated with only a few IASI spectra, which makes it impossible to use them as training and/or validation data sets; therefore we had to merge them with homogeneous zones in order to reach a statistically significant size. In some cases, we had no data at all and the cloud mask thresholds were not produced.

This is indeed a limitation for operational uses. However, as already mentioned the main goal of this paper is to show the scientific basis of the methodology and to outline its initial evaluation.

2.1 Silver standard

Choice of the reference training data set is an important step in classification, because properties of the statistics in clear and cloudy sky conditions are estimated from it. Due to its importance, this data set is also called “gold” standard, for which the class label of each sample of the data set is assumed to be exactly known. Often, this is obtained by an expert who manually trains each sample of the data set and assigns the class. In practice in many applications this is not possible, or it is possible only for a limited sample; as

Table 1. Definition of climate zones as a function of latitude and season.

Model	Climatic zone	Latitude (degrees)	Season (months)	Data set
1	Tropical	–35 to 35	Jan–Dec	Both
2 ₁	Mid-lat summer NH	35 to 60	Apr–Oct	IASI2
2 ₂	Mid-lat summer SH	–60 to –35	Nov–Mar	IASI1
3 ₁	Mid-lat winter NH	35 to 60	Nov–Mar	IASI1
3 ₂	Mid-lat winter SH	–60 to –35	Apr–Oct	IASI2
4 ₁	High-lat summer NH	60 to 90	Apr–Oct	IASI2
4 ₂	High-lat summer SH	–90 to –60	Nov–Mar	IASI1
5 ₁	High-lat winter NH	60 to 90	Nov–Mar	IASI1
5 ₂	High-lat winter SH	–90 to –60	Apr–Oct	IASI2
6 ₁	Sea ice NH	0 to 90	Jan–Dec	Both
6 ₂	Sea ice SH	–90 to 0	Jan–Dec	Both
7 ₁	Land below 1 km with snow/ice NH	0 to 90	Jan–Dec	IASI1
7 ₂	Land below 1 km with snow/ice SH	–60 to 0	Jan–Dec	IASI2
8 ₁	Land above 1 km with snow/ice NH	0 to 90	Jan–Dec	IASI2
8 ₂	Land above 1 km with snow/ice SH	–60 to 0	Jan–Dec	IASI1
9	Antarctica below 1 km		Jan–Dec	IASI1
10	Antarctica above 1 km		Jan–Dec	IASI1

a consequence the sample is not fully representative of the full population and the statistical properties estimated from the data set are not accurate. In the case of cloud detection there is a further problem coming from the fact that it can be difficult to recognize clouds from satellite imagery even for a well-trained expert in particular conditions like night-time or cirrus. When a gold standard is not available or not fully representative, it is usual to choose a training data set starting from the results of another classification algorithm that is proven to be very reliable. In this case this data set is called “silver” standard. The silver standard used in the present paper relies on the CMS cloud mask which has been used as the IASI reference cloud mask.

2.1.1 The CMS reference cloud mask

The IASI reference cloud mask we have used in our analysis has been developed at CMS (Lavanant et al., 2011); in turn, it is based on the cloud mask derived from the AVHRR data at full resolution (Lavanant et al., 2007).

AVHRR pixels are co-located with IASI footprints. Within each IASI footprint, the cloud fraction, C_f , is determined according to the ratio

$$C_f = 100 \frac{N_{\text{cloudy}}}{N_{\text{total}}} \times 100,$$

with N_{cloudy} and N_{total} being the number of AVHRR cloudy and total pixels, respectively. C_f ranges from 0 to 100, with 0 corresponding to a totally clear IASI observation. For the present analysis, the training of the IASI cloud mask has been performed considering a threshold $C_f = 5$; that is, the IASI footprint is considered clear if the AVHRR-based cloud fraction is below or equal to 5 %.

The CMS-AVHRR cloud detection scheme is based on a series of tests, which use the AVHRR observations in its six channels and take advantage of a priori or background information for the surface and atmospheric parameters. The a priori state vector is based on the ECMWF 12–18 h forecast for the total water vapour content and land surface temperature and on climatologies for the other parameters (Lavanant et al., 2007, 2011). In addition to the tests based on the visible channels, the AVHRR cloud detection scheme also uses a series of thermal infrared brightness temperature tests and a series of local uniformity tests, which are designed to detect cloud edges, thin cirrus and small cumulus, by using their high spatial variations in the visible, near infrared or infrared channels.

Comparisons of CMS with CLAVR-x show an agreement which is close to 90 % (Lavanant et al., 2007). To have an idea of the quality and accuracy of AVHRR-based cloud masks, we quote the validation of CLAVR-x (Heidinger et al., 2012) on the basis of a data set that coincides with the training one. The validation indicates that the cloud mask is estimated with an accuracy ranging from 78 to 94 % according to the surface type (best for deep ocean, worst for Antarctica). Because of the good agreement between CLAVR-x and CMS, we can conclude that these figures are representative of the accuracy of the CMS cloud mask as well.

For the specific case of IASI, the CMS cloud mask has been intercompared with other IASI cloud detection schemes in Lavanant et al. (2011). The best agreement is achieved with schemes relying themselves on the AVHRR imager. The agreement ranges from 78 to 85 %.

2.1.2 Validation data set based on the SEVIRI cloud mask

To have an independent validation data set, we have also used the SEVIRI cloud mask (Derrien and Le Gléau, 2005). A portion of the data set IASI2, covering the African continent, has been co-located with SEVIRI imagery, and the SEVIRI cloud mask has been used to identify clear and cloudy IASI pixels.

A second, smaller, independent validation data set makes use of observations collected over Europe/Africa from 25 September to 4 October 2012 in the framework of the inter-calibration IASI/SEVIRI experiment, performed within the activities of EUMETSAT GSICS (Global Space-based Inter-Calibration System).

Validation of the SEVIRI cloud mask is provided in Derrien and Le Gléau (2005) and then updated in Derrien (2012) according to the latest release of the product. It is based on a subset of case studies (366 298 pixels over the European and African areas along the period from 10 December 2010 to 21 March 2011) where corresponding in situ observations were available in terms of octal in the framework of the World Meteorological Organization synoptic code (surface synoptic observation, SYNOP). Results show that the SEVIRI cloud mask agrees with SYNOP for 96.5 % of pixels, with percentages reaching 98.9 % during daytime and 95.7 % during night-time.

These two SEVIRI validation data sets have to be considered silver standards, since they rely on the operational SEVIRI cloud mask.

3 Statistics

This section summarizes the statistics computed from the IASI radiances that will be used to discriminate between cloudy and clear pixels. A detailed account about their definition and capability to discriminate between cloudy and clear scenes can be found in Serio et al. (2013). For the sake of brevity, in this section we limit ourselves to showing the basic aspects of the statistics, insisting more on the ones which are less common or rely on an implementation specifically developed for this work.

First we provide some notations that will be used in the paper. Brightness temperature, BT; spectrum, $T(\sigma)$; and spectral radiance, $R(\sigma)$, at wave number σ are related through the Planck function

$$R(\sigma) = \frac{C_1 \sigma^3}{\exp(C_2 \sigma / T(\sigma)) - 1},$$

where wave number σ is in units of cm^{-1} , $R(\sigma)$ is in units of $\text{W m}^{-2} (\text{cm}^{-1})^{-1} \text{sr}^{-1}$, $C_1 = 1.1911 \times 10^{-8} \text{ W m}^{-2} (\text{cm}^{-1})^{-4} \text{sr}^{-1}$ and $C_2 = 1.4388 \text{ K} \cdot (\text{cm}^{-1})^{-1}$. As a consequence $T(\sigma)$ is in units of K. We shall denote by R_σ and

Table 2. Definition of the brightness temperatures involved in the statistics.

BT (K)	Spectral range [σ_1, σ_2] (cm^{-1})
$T_{790.5}$	790.5
$T_{791.75}$	791.75
\bar{T}_{832}	[830, 834]
\bar{T}_{874}	[872.5, 875.5]
\bar{T}_{900}	[899.5, 900.5]
\bar{T}_{1168}	[1167.5, 1168.5]
\bar{T}_{2003}	[2001, 2005]
\bar{T}_{2700}	[2650, 2750]

T_σ the observed radiance and the corresponding BT at wave number σ , respectively.

Many of the statistics are defined as a function of suitable BT in the atmospheric window (see Table 2 for the definition of the BT).

The two BTs $T_{790.5}$ and $T_{791.75}$ are defined at specific wave numbers. The other temperatures are computed as averages over the corresponding spectral ranges listed in Table 2:

$$\bar{T}_\sigma = \frac{1}{N_\sigma} \sum_{\tau=\sigma_1}^{\sigma_2} T_\tau,$$

where σ is the central wave number of the temperatures (832, 874, 900, 1168, 2003 and 2700 cm^{-1} , orderly), σ_1 and σ_2 are the range of the spectral interval for each central wave number (see Table 2) and N_σ is the number of spectral ordinates within the spectral range.

Table 3 summarizes the nine statistics computed from the IASI radiances. It should be stressed that at the present stage of defining statistics we do not discuss the estimate of the thresholds for discriminating clear and cloudy pixels, because this task will be accomplished with the CDA methodology (Sect. 4). CDA objectively determines the appropriate threshold for each statistic.

The four statistics W_1 , W_2 , W_3 and W_4 are Inoue-like window slope statistics (Inoue, 1985; Inoue and Ackerman, 2002). W_1 is the classical Inoue slope test, which is highly sensitive to cirrus clouds. W_2 is a variant of W_1 and is mostly effective in the case of surface features rapidly changing with the wave number, such as desert sand. The statistic W_4 has been widely used in AVHRR-based cloud detection (e.g. Lavanant et al., 2007) and is motivated by the fact that the cirrus and stratus cloud types have a reflectance at 3.7 μm which is higher than that of most surface features. The channel at 3.7 μm has been used both during night-time and daytime. Among the four W statistics, W_3 is the most original one and has been defined to be sensitive to low and thin water clouds. The water droplet mode radius of most cloud types is 5 μm (e.g. Liou, 1992), which means that at 5 μm (2000 cm^{-1}) scattering effects dominate over absorption. The reverse happens at 12 μm (833 cm^{-1}), where the absorption dominates

Table 3. Statistics derived from the IASI radiances. Temperatures are defined in Table 2.

Number	Statistic	Method
1	h_s	Based on shape similarity between a couple of spectra (observation, reference)
2	χ_s^2	Based on a χ^2 -like variable defined on a couple of skin temperature values (\hat{T}_s, T_s^R), with \hat{T}_s estimated from the spectrum and T_s^R a suitable reference
3	T_0	Based on the BT \bar{T}_{832}
4	ΔT_{CO_2}	Split window test based on the CO ₂ Q-branch at 791 cm ⁻¹ , $\Delta T_{CO_2} = T_{790.5} - T_{791.75}$
5	W_1	Based on the difference $\bar{T}_{900} - \bar{T}_{832}$
6	W_2	Based on the difference $\bar{T}_{900} - \bar{T}_{1168}$
7	W_3	Based on the difference $\bar{T}_{832} - \bar{T}_{2003}$
8	W_4	Based on the difference $\bar{T}_{832} - \bar{T}_{2700}$
9	sh	Spatial homogeneity statistic based on the standard deviation of T_0 corresponding to a cluster of $n \times n$ nearby pixels. For IASI we consider the cluster of 2×2 pixels within a given field of view.

over scattering. Thus, in the presence of a semi-transparent water cloud a strong contrast is expected between the BT at 5 μm and that at 12 μm . W_3 is expected to play a significant role over land in detecting warm clouds during night-time. Over land, because of emissivity, the two equivalent, long-wave statistics W_1 and W_2 can have a variability in clear sky, which is much larger than that expected for cloudy sky. In contrast, we expect that W_3 can assume only positive values in clear sky and negative ones in the case of cloudy scenes.

A slope window test would be almost useless without suitable driver (non ho capito driver) thermal contrast tests which capture the physical evidence that, normally, the land is warmer than the cloud top. We have basically two thermal contrast statistics. The first one is just the brightness temperature at 833 cm⁻¹, that is, T_0 . Sea surface emissivity has a peak at 833 cm⁻¹, and at this wave number natural land features have the smallest emissivity variability. In clear sky conditions, independently of the kind of surface, the channel at 833 cm⁻¹ is the brightest point in the spectrum.

The spatial homogeneity statistic, sh, is similar to those used with AVHRR-based cloud detection. It uses the standard deviation of T_0 corresponding to a cluster of 2×2 IASI pixels. This statistic is normally very well suited to detect cloud edges; however for IASI it has not proved to be effective because of the discontinuous IASI scan pattern geometry and the relatively large field of view. Furthermore, in contrast it is expected to be very effective for MTG-IRS because of the imaging capability of the instrument and the smaller field of view.

The CO₂ in-band-out-band statistic, ΔT_{CO_2} (Masiello et al., 2003), exploits the strong thermal contrast which is present in clear sky conditions between the lower and upper troposphere. It is based on the absorption feature of the weak CO₂ Q-branch at 791 cm⁻¹. CO₂ absorption yields a very well defined and sharp spectral feature centred at 791.75 cm⁻¹ in between a window region with weak water vapour absorption.

The χ_s^2 statistic makes use of information provided by NWP forecasts and/or analysis for sea and surface temperature, T_s . This temperature is contrasted with that directly estimated by IASI time-space co-located observations.

For the case of sea surface, the IASI estimate of T_s relies on the classical split window algorithm for the estimation of skin temperature. It is based on the two temperatures \bar{T}_{874} and \bar{T}_{900} . The algorithm we have developed for IASI is dependent on the field-of-view angle and on the type of air mass. For the calculation of the regression coefficients involved in the algorithm we have assumed the Northern and Southern Hemisphere to be climatically equivalent. For practical computations, one needs a suitable database of atmospheric and surface parameters. To this end, we have used the ECMWF Chevalier database (Chevalier, 2001). The accuracy of the split window technique ranges between 0.3 and 1 K, depending on the columnar load of water vapour.

The statistic χ_s^2 for sea surface compares \hat{T}_s , estimated by the IASI split window technique, against a suitable reference, $T_{s,R}$. We have

$$\chi_s^2 = \frac{(\hat{T}_s - T_{s,R})^2}{v_I^2 + v_R^2},$$

where v_I^2 and v_R^2 are the squared uncertainties (variances) of \hat{T}_s and $T_{s,R}$, respectively.

If appropriately built up, and in absence of biases affecting estimated and reference T_s , the statistic χ_s^2 is distributed according to a χ^2 density function with one degree of freedom. As far as v_R^2 is concerned, its value depends on the quality of the reference. In the present implementation, $T_{s,R}$ is obtained from the ECMWF analysis, which is generally recognized to be accurate within 1 K. Again, a detailed account of the IASI split window and χ^2 test for sea surface can be found in Serio et al. (2013).

For sea surface, the χ^2 test is most powerful; the quality of the split window and the overall reliability of the statistic have been variously assessed through IASI spectra co-located in time and space with ECMWF analysis. Figure 1

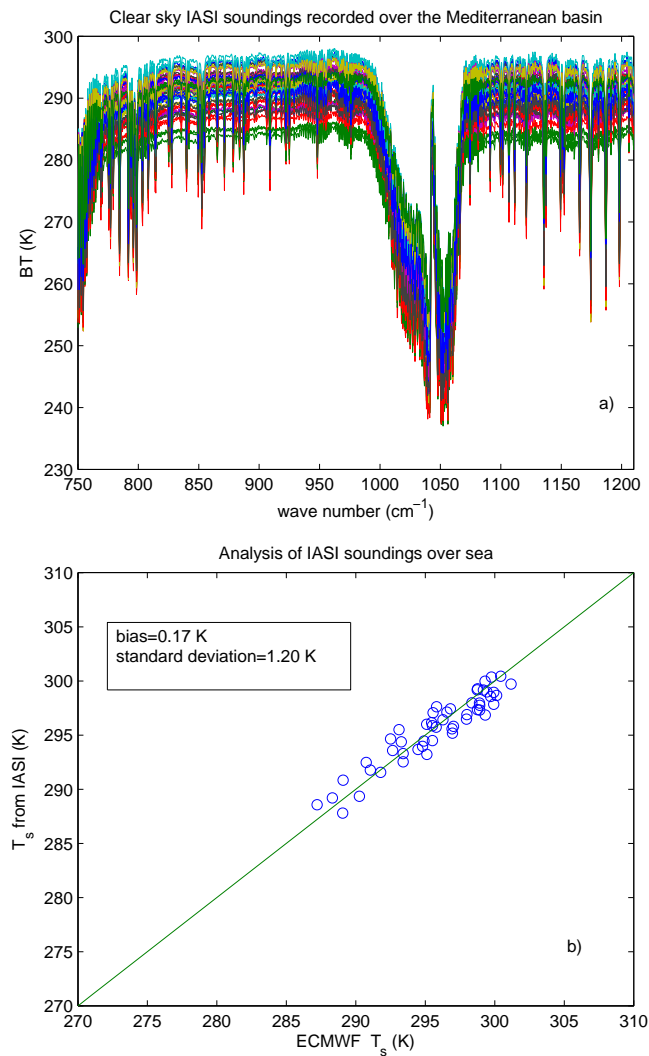


Figure 1. Example of the performance of the split window technique for the estimation of the skin temperature. Panel (b) compares the IASI T_s obtained from the spectra shown in (a) with the time-space co-located ECMWF analysis.

shows an example obtained over the Mediterranean area with a sample of spectra covering the period March–October 2010 and qualified for clear sky according to the methodology developed in this work.

For land, mostly motivated by the lack of a robust and fast method to estimate the surface temperature, the χ^2 test is substituted by a statistic which considers the difference

$$|\bar{T}_{832} - T_{s,R}|,$$

where again the reference temperature is obtained from the ECMWF analysis. In other words, \bar{T}_{832} is considered as a proxy of the surface temperature. The fact that \bar{T}_{832} could be affected by water vapour is not of any concern here, because our aim is not to estimate the surface temperature. The

statistic is used to capture the high thermal contrast expected among clear and cloudy conditions.

Unlike the statistics so far described, spatial homogeneity statistic h_s (Maseillo et al., 2002; Serio et al., 2000; Masiello et al., 2003) fully exploits the hyperspectral capabilities of IASI with regard to MTG-IRS.

Basically, it is designed to exploit the unique spectral signature of sea/land surface in the atmospheric window region 800–950 cm^{-1} . It uses the observed spectrum and a clear-sky reference spectrum, T_σ^O and T_σ^R , respectively, both converted to BT, according to the formula

$$h_s = \frac{\sum_{j=1}^{N_L} |r_j^O - r_j^R|}{\sum_{j=1}^{N_L} |r_j^{OR}|},$$

where r_j^O , r_j^R and r_j^{OR} are correlation and cross-correlation coefficients of (T_σ^O, T_σ^R) and N_L is the number of considered wave number lags (e.g. Maseillo et al., 2002; Serio et al., 2000). In the present paper the reference spectra for sea surface are obtained by σ -IASI (Amato et al., 2002) with a set of atmospheric profiles derived from the Chevalier data set (Chevalier, 2001) using Masuda’s emissivity model for sea surface (Masuda et al., 1988). In the case of land surface, uncertainty associated with surface emissivity could make the method less effective, mostly due to troubles in defining a suitable reference spectrum. Therefore we have developed an approach that is as independent of surface emissivity as possible: instead of the radiance spectrum, we consider its first difference

$$\Delta R_\sigma = R_{\sigma+\Delta\sigma} - R_\sigma,$$

which is an approximation of the differential of the spectrum with respect to the wave number, σ . This is a high-pass filter, which removes the relatively smooth, and hence slow, component introduced from surface emissivity.

Figure 2 shows a demonstration of the procedure, comparing R_σ and ΔR_σ within the spectral range [800, 950] cm^{-1} for very different surface emissivities. It is seen that in the most transparent regions the filter largely suppresses the surface emission component and leaves the atmospheric line component unaffected. This atmospheric component is largely affected by the presence of clouds.

With regard to MTG-IRS the use of computed reference spectra to calculate h_s is envisaged at the very beginning of its operational life. Once we have collected enough observations, it is desirable to use directly MTG-IRS observations, suitably screened for clear sky.

To summarize, among the nine statistics we consider for the problem of cloud detection, W_1 , W_2 , W_4 , T_0 , χ_s^2 and sh are heritage of AVHRR/MODIS cloud detection algorithms; W_3 is quite new since the channel at 2003 cm^{-1} is not available with current satellite radiometer imagers; the thermal contrast statistic ΔT_{CO_2} relies on the high spectral resolution

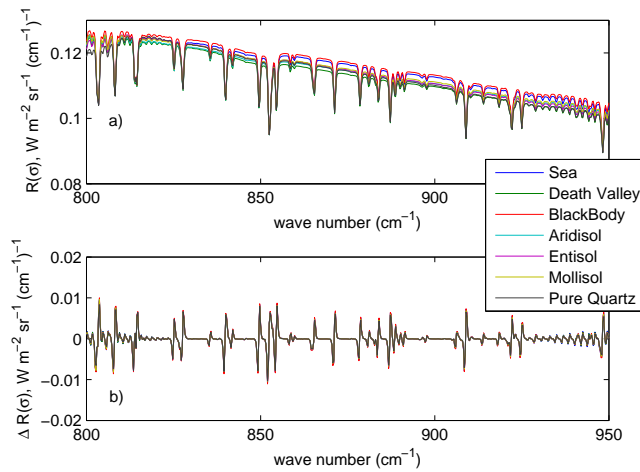


Figure 2. Top panel: radiance spectra in the atmospheric window segment $[800\text{--}950]\text{cm}^{-1}$ for very different surface emissivities. Bottom panel: the same radiance spectra after first-order differentiation. The synthetic spectra are based on the same set of atmospheric parameters; only emissivity is changed according to the surface features listed in the legend.

of IASI and, in perspective, of MTG-IRS and cannot be designed, for example, for imagers like AVHRR. Finally, *hs* is specialized for hyperspectral instrumentation and heavily relies on the concept of spectrum and/or spectral radiance, as opposed to the coarse spectral resolution of radiance which up to now has characterized coarse- and moderate-spectral-resolution satellite imagers, such as AVHRR, MODIS and SEVIRI.

4 Cumulative discriminant analysis

Estimate of suitable thresholds for each statistic is required to produce a cloud mask. This is accomplished by (a) choosing the most effective statistics among the considered ones and by (b) actually determining the corresponding thresholds. Towards this objective, this section introduces the cumulative discriminant analysis, which is the engine that will drive the discrimination/classification methodology for producing the final cloud mask. It is important to stress that we look at discrimination/classification in a probabilistic fashion. Also the CDA engine has to be run only once to discriminate the most effective statistics and to generate thresholds. This is done on the basis of a suitable training data set. Once we get the thresholds, we do not need to run CDA again to apply the scheme to a given IASI spectrum. Of course, the full procedure is needed in case the training data set is changed.

4.1 Notations

We begin with some notations to make clear conventions that will be used in the paper. We assume clear as the target sky

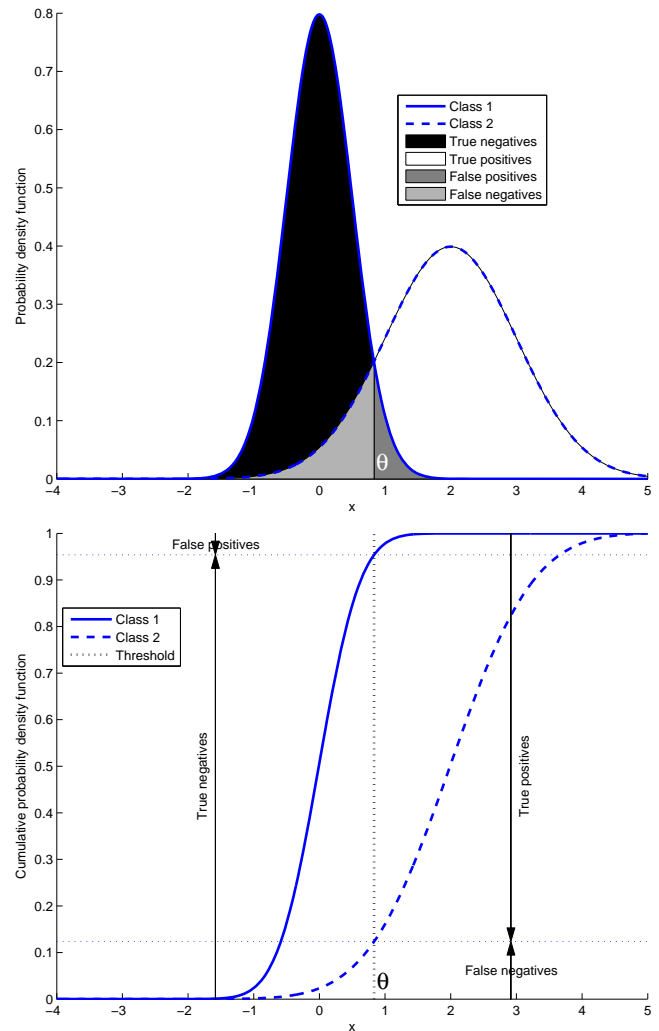


Figure 3. Visual representation of type I and II errors by density (top panel) and cumulative distribution (bottom panel) functions.

condition. Then type I error is defined as the fraction of pixels being clear and classified as cloudy; this fraction is also called false positives, FP, or miss. Analogously we define as type II error the fraction of pixels being cloudy and classified as clear (also called false negatives, FN, or false alarm). The fraction of pixels exactly classified as clear or cloudy will be denoted by true positives, TP, and true negatives, TN, respectively. The performance of the classification scheme can be summarized through the fraction, S , of successful classifications,

$$S = \frac{TPN_{\text{clear}} + TNN_{\text{cloudy}}}{N_{\text{clear}} + N_{\text{cloudy}}}, \quad (1)$$

with N_{clear} and N_{cloudy} being the number of true clear and cloudy pixels, respectively.

As a graphical representation (see Fig. 3, top) we plot the probability density function of two classes, with x being a generic statistic. We also assume that the classifier acts by

means of a threshold θ that gives back class 1 if $x \leq \theta$ and class 2 if $x > \theta$. Therefore in the plot true negatives are given by the area below the higher peaked function (continuous line) up to the threshold θ (the black and light-grey areas); true positives are given by the area below the lower peaked function (dashed line) from the threshold upwards (dark-grey and white areas). Analogously, false positives (type I error) are given by the dark-grey area and false negatives (type II error) by the light-grey area. An alternative graphical representation can be given in terms of the cumulative distribution functions (Fig. 3, bottom).

Intercept of the threshold line with the two cumulative distribution functions for the two classes locates the true negatives (below the intercept) and the false positives (above) from the curve of the first class and the true positives (above the intercept) and false negatives (below) from the curve of the second class.

4.2 Nonparametric estimates

Many classification methodologies require estimate of the probability density function of the classes. While this can be accomplished relying on parametric methods, most density functions coming from real applications are not well approximated by such methods for their poor ability to fit actual distributions. For this reason nonparametric density functions are resorted where the shape of the density functions is not assigned from the beginning through a family of functions but is the result of the approximation. The most widespread of such methods is kernel density estimation (Silverman, 1986). Despite its popularity, however, nonparametric density estimation requires tuning of one or more parameters (choice of the kernel function, bandwidth). In particular, bandwidth heavily affects accuracy of the nonparametric density estimate through its smoothness: the larger the bandwidth (i.e. the width of the kernel function), the smoother the density estimate. Some criteria, both based on asymptotic arguments and on the actual data, have been developed (Wand and Jones, 1995) for its estimation; however choice of the bandwidth continues to be the crux of the methodology.

On the other hand if we consider the cumulative distribution function, an estimator exists that enjoys many optimal theoretical properties and is very attractive from the computational point of view: the empirical cumulative distribution function $F(\vartheta)$. Given the sample $x_i, i = 1, \dots, N$, it is easily defined as

$$F(\vartheta) = \frac{\text{Number of elements } \leq \vartheta}{N} = \frac{1}{N} \sum_{i=1}^N \mathbf{1}\{x_i \leq \vartheta\}.$$

The most important result is that the estimator almost surely converges to the true cumulative distribution function asymptotically; therefore the estimator is consistent (van der Vaart, 1998). In addition convergence holds uniformly over ϑ . From the computational point of view the estimator is simply computed by sorting the sample $x_i, i = 1, \dots, N$, in

$O(N \log N)$ operations. We observe that no parameter exists to be tuned.

4.3 Cumulative discriminant analysis

Basing on the theoretical properties of the cumulative distribution functions, a CDA methodology is proposed for the classification.

We first consider the univariate case corresponding to one statistic. We suppose that the classifier is of the discriminant analysis type where the decision rule $\Gamma(x, \mathbf{x})$ is based on a threshold ϑ as

$$\Gamma(x; \mathbf{x}) = \begin{cases} 1 \text{ (Clear)} & \text{if } x \leq \vartheta; \\ 2 \text{ (Cloudy)} & \text{otherwise,} \end{cases} \quad (2)$$

where $\mathbf{x} \equiv \{x_1, \dots, x_N\}$ is the sample statistic of size N that defines the training data set and x is the actual statistic to be classified basing on the training data set.

Classical discriminant analysis would yield a threshold rule similar to Eq. (2) under the assumption of Gaussianity of the density function of the statistic for both clear and cloudy conditions, which does not occur in practice. By relying on the CDA we refrain from assuming this hypothesis and instead use an optimal estimate of the cumulative distribution function.

Now an estimate of the threshold has to be given in order that the method is fully operative. For this purpose we define a cost function $\mathcal{C}(x, \vartheta)$ whose minimization gives back the estimate of the optimal threshold, $\hat{\vartheta}$:

$$\hat{\vartheta} = \operatorname{argmin}_{\vartheta} \mathcal{C}(x, \vartheta).$$

An example of cost function could be the following one. Let E^I be the type I error and E^{II} the type II one. Clearly it is

$$E^I = 1 - F^{\text{Clear}}(\vartheta) \text{ and } E^{II} = F^{\text{Cloudy}}(\vartheta). \quad (3)$$

Then we define the cost function as

$$\mathcal{C}(x, \vartheta) = E^I + E^{II} = 1 - F^{\text{Clear}}(\vartheta) + F^{\text{Cloudy}}(\vartheta).$$

Such a formula is of the linear discriminant analysis type without the assumption of Gaussianity of the density function. Indeed it aims at minimizing the number of misclassifications.

A second possibility is given by the following:

$$\mathcal{C}(x, \vartheta) = \max(E^I, E^{II}) = \max(1 - F^{\text{Clear}}(\vartheta), F^{\text{Cloudy}}(\vartheta)). \quad (4)$$

The rationale behind choice (4) is that we want to simultaneously minimize both type I and II errors, and our objective is a proper balancing of the error between the two classes. Actually, a classical discriminant analysis could be misleading in the case of the training data set not being balanced

between the two classes, cloudy and clear. Then, since the cost function minimizes the overall error, the threshold naturally outweighs the most populated class that will be better classified at the possible detriment of the smaller class. If the actual scene to be classified is poorly balanced in favour of the other one (less populated in the training data set), then the overall misclassifications will increase. This problem is very common in cloud detection scenes because the relative fraction of clear and cloudy pixels depends on seasonal and meteorological arguments and therefore is very variable with the actual scene.

In the present project we use the cost function (4) since the objective is to estimate the rate of misclassifications for both clear and cloudy conditions. This choice is the least non-committal with respect to the average conditions of the sky, cloudy or clear. We observe that it corresponds to assume that clear and cloudy classes have equal size, which has a counterpart with a uniform prior of discriminant analysis. If one were to take account of real or different relative size of the clear and cloudy sample, a weighted CDA could be devised by introducing proper weights to the type I and II errors, in analogy with prior probabilities in classical discriminant analysis.

Note that at the minimum of the cost function (4) it is $E^I = E^{II}$; that is, type I and II errors coincide when no weights are introduced in Eq. (4); otherwise the minimum of the cost function is given by the ratio of the weights.

Finally we mention that, in minimizing the cost function with respect to the threshold ϑ , we have to include in the decision rule also the sign direction of the rule, that is, $x \leq \vartheta$ or $x \geq \vartheta$ for clear conditions. In the latter case the decision rule becomes

$$\Gamma(x; \mathbf{x}) = \begin{cases} 1 \text{ (Clear)} & \text{if } x \geq \vartheta \\ 2 \text{ (Cloudy)} & \text{otherwise,} \end{cases} \quad (5)$$

and the type I and II errors are written as

$$E^I = F^{\text{Clear}}(\vartheta), \quad E^{II} = 1 - F^{\text{Cloudy}}(\vartheta).$$

4.4 Multivariate CDA

Formulas of the cost function can be generalized to the case of D statistics giving rise to a procedure that can be implemented effectively.

Generalization of the decision rule (2) results in

$$\Gamma(x^1, x^2, \dots, x^D; \mathbf{x}) = \begin{cases} 1 \text{ (Clear)} & \text{if } x^d \leq \vartheta_d, 1 \leq d \leq D; \\ 2 \text{ (Cloudy)} & \text{otherwise;} \end{cases} \quad (6)$$

analogously, generalization of the rule (5) becomes

$$\Gamma(x^1, x^2, \dots, x^D; \mathbf{x}) = \begin{cases} 1 \text{ (Clear)} & \text{if } x^d \geq \vartheta_d, 1 \leq d \leq D; \\ 2 \text{ (Cloudy)} & \text{otherwise.} \end{cases} \quad (7)$$

In a similar way all 2^D decision rules with mixed inequality signs can be built.

We assume that the theory underlying the computation of the type I and II errors falls in the framework of the probability of the union of two or more events that are not mutually exclusive. For such events in the case of decision rule (7) we have the generalization of the type I error to the two-dimensional case as

$$E^I = F_1^{\text{Clear}}(\vartheta_1) + F_2^{\text{Clear}}(\vartheta_2) - F_1^{\text{Clear}}(\vartheta_1)F_2^{\text{Clear}}(\vartheta_2).$$

Generalization to the three-dimensional case is

$$\begin{aligned} E^I = & F_1^{\text{Clear}}(\vartheta_1) + F_2^{\text{Clear}}(\vartheta_2) + F_3^{\text{Clear}}(\vartheta_3) \\ & - F_1^{\text{Clear}}(\vartheta_1)F_2^{\text{Clear}}(\vartheta_2) - F_1^{\text{Clear}}(\vartheta_1)F_3^{\text{Clear}}(\vartheta_3) \\ & - F_2^{\text{Clear}}(\vartheta_2)F_3^{\text{Clear}}(\vartheta_3) \\ & + F_1^{\text{Clear}}(\vartheta_1)F_2^{\text{Clear}}(\vartheta_2)F_3^{\text{Clear}}(\vartheta_3) \end{aligned} \quad (8)$$

and so on for higher D .

Similar formulas, not written here for the sake of brevity, can be devised for E^{II} and other decision rules. For example for the decision rule (6) and $D = 3$, the type I and II errors of (3) generalize to

$$\begin{aligned} E^I = & 1 - F_1^{\text{Clear}}(\vartheta_1)F_2^{\text{Clear}}(\vartheta_2)F_3^{\text{Clear}}(\vartheta_3) \\ E^{II} = & F_1^{\text{Cloudy}}(\vartheta_1)F_2^{\text{Cloudy}}(\vartheta_2)F_3^{\text{Cloudy}}(\vartheta_3). \end{aligned} \quad (9)$$

Furthermore, Eq. (9) can be obtained from Eq. (8) and vice-versa with abuse of notation via

$$\begin{aligned} F_1^{\text{Clear}}(\vartheta_1) &= 1 - F_1^{\text{Cloudy}}(\vartheta_1), \quad F_1^{\text{Cloudy}}(\vartheta_1) = 1 - F_1^{\text{Clear}}(\vartheta_1); \\ F_2^{\text{Clear}}(\vartheta_2) &= 1 - F_2^{\text{Cloudy}}(\vartheta_2), \quad F_2^{\text{Cloudy}}(\vartheta_2) = 1 - F_2^{\text{Clear}}(\vartheta_2); \\ F_3^{\text{Clear}}(\vartheta_3) &= 1 - F_3^{\text{Cloudy}}(\vartheta_3), \quad F_3^{\text{Cloudy}}(\vartheta_3) = 1 - F_3^{\text{Clear}}(\vartheta_3). \end{aligned}$$

4.5 Limitations of CDA

CDA is based on the assumption that the decision rule involves one threshold for each statistic. If we recall how a discriminant analysis based on the probability density function works, this means that the threshold has to split the domain of the statistic in two parts, and in each of them one of the density functions in turn has to be higher than the other one in general. Visually, this occurs when the two density functions look to be translated between each other. From the cumulative distribution function point of view, we require that the two functions do not cross (or cross at very far tails). This condition strictly occurs with parametric discriminant analysis (LDA, QDA), due to the hypothesis of Gaussianity of the density functions. For the IASI data this assumption is verified for most statistics. Actually CDA could be similarly generalized to eventually two and more thresholds.

We already mentioned the assumption of independence of statistics; for the general D -dimensional case this means factorization of the multivariate density functions into univariate functions. In addition, the multidimensional decision rule is chosen as the cartesian product of unidimensional rules based on thresholds.

We noticed that in the case of multiple statistics a pixel is labelled as clear if all corresponding tests are satisfied with the right directions (e.g. Eq. 6). This procedure is common to several operational cloud mask algorithms. However when the number of statistics and corresponding tests increases, the probability that a pixel is labelled as cloudy just due to the intrinsic randomness of the statistic, while the pixel was actually clear, increases. As a consequence there is an increase of the type I error. To overcome this effect it is advisable to adopt a strategy of Bonferroni or false discovery rate (FDR) (Benjamini et al., 1995) correction.

Finally, CDA also assumes a binary (2-class) problem; generalization to more classes could be accomplished, for example, with the same techniques that apply to support vector machines (SVM) (e.g. dealing with all couples of bands and then adopting a voting strategy).

4.6 Computational aspects

The code for the computation of the multidimensional CDA has been developed in Matlab. In particular, an optimized one-dimensional version has been implemented where a direct minimization of the cost function is performed by means of a direct exhaustive search over the step-wise estimated cumulative distribution function.

Its generalization to the fully general D -dimensional case is not computationally possible due to the exponential increase of the computational cost with D . Therefore the hypothesis of independence of statistics is resorted to, and minimization of the cost function is accomplished with an optimization method that is derivative-free (Nelder–Mead, based on simplex).

We observe by simple arguments that increase of the number of involved statistics can not degrade the cost function. In fact, by setting the threshold of the new statistic to $+\infty$ or $-\infty$ according to the sign of the inequality, \leq or \geq , respectively, makes the decision rule not depend on the new threshold. As a consequence the cost function is the same, and therefore it can only improve with the number of dimensions by properly choosing all thresholds.

We also observe that of course, when the number of statistics increases, all the thresholds will vary as well in general, even though we should not expect dramatic variations. For this reason the initialization of the Nelder–Mead optimization is performed starting from D -independent unidimensional optimizations of the thresholds (indeed also of the corresponding inequality signs). Then the threshold of the statistic that, singularly taken, yields the minimum value of the cost function is used as initialization; for the remaining thresholds an intermediate value between the optimal one-dimensional case and $\pm\infty$ (according to the estimated sign of the inequality) is set.

Finally we do not perform a full optimization of the sign of threshold inequalities in the D -dimensional case because the number of cases to be considered would explode

exponentially with D as 2^D . We assume that the sign of the threshold inequality is preserved when moving from one to more dimensions.

Starting from the formulation (9) the algorithm for the computation of the cost function for D statistics goes through the following steps:

- compute unidimensional cumulative distribution functions $F_d^{\text{Clear}}, F_d^{\text{Cloudy}}, d = 1, \dots, D$, for each statistic d starting from the training set;
- compute directional cumulative distribution functions, $\tilde{F}_d^{\text{Clear}}, \tilde{F}_d^{\text{Cloudy}}, d = 1, \dots, D$, according to the selected direction of threshold inequalities:

$$\tilde{F}_d^{\text{Clear}}(x^d) = \begin{cases} F_d^{\text{Clear}}(x^d), & \text{if } x^d \leq \vartheta_d \Rightarrow \text{Clear} \\ 1 - F_d^{\text{Clear}}(x^d), & \text{if } x^d \geq \vartheta_d \Rightarrow \text{Clear} \end{cases}$$

$$\tilde{F}_d^{\text{Cloudy}}(x^d) = \begin{cases} F_d^{\text{Cloudy}}(x^d), & \text{if } x^d \leq \vartheta_d \Rightarrow \text{Clear} \\ 1 - F_d^{\text{Cloudy}}(x^d), & \text{if } x^d \geq \vartheta_d \Rightarrow \text{Clear}; \end{cases}$$

- compute type I and II errors as

$$E^{\text{I}} = 1 - \prod_{d=1}^D \tilde{F}_f^{\text{Clear}}(x^d)$$

$$E^{\text{II}} = \prod_{d=1}^D \tilde{F}_f^{\text{Cloudy}}(x^d);$$

- compute the cost function as

$$\mathcal{C} = \max(E^{\text{I}}, E^{\text{II}}).$$

4.7 Principal component analysis

In the case of strong departures from the assumed hypothesis of independence of statistics, Eq. (8) and similar ones are not strictly valid anymore for $D > 1$. On the other side type I and II errors can be easily generalized; for example in the case of decision rule (6) we get

$$E^{\text{I}} = 1 - F^{\text{Clear}}(\vartheta_1, \vartheta_2, \dots, \vartheta_D)$$

$$E^{\text{II}} = F^{\text{Cloudy}}(\vartheta_1, \vartheta_2, \dots, \vartheta_D).$$

However a direct efficient generalization of the algorithm described in Sect. 4.6 is not possible because it would require a direct estimation of the full D -dimensional cumulative distribution functions and therefore an exponential growth with D of the computational cost and of the memory requirement. For this reason a specific full algorithm can be developed only for the case $D = 2$.

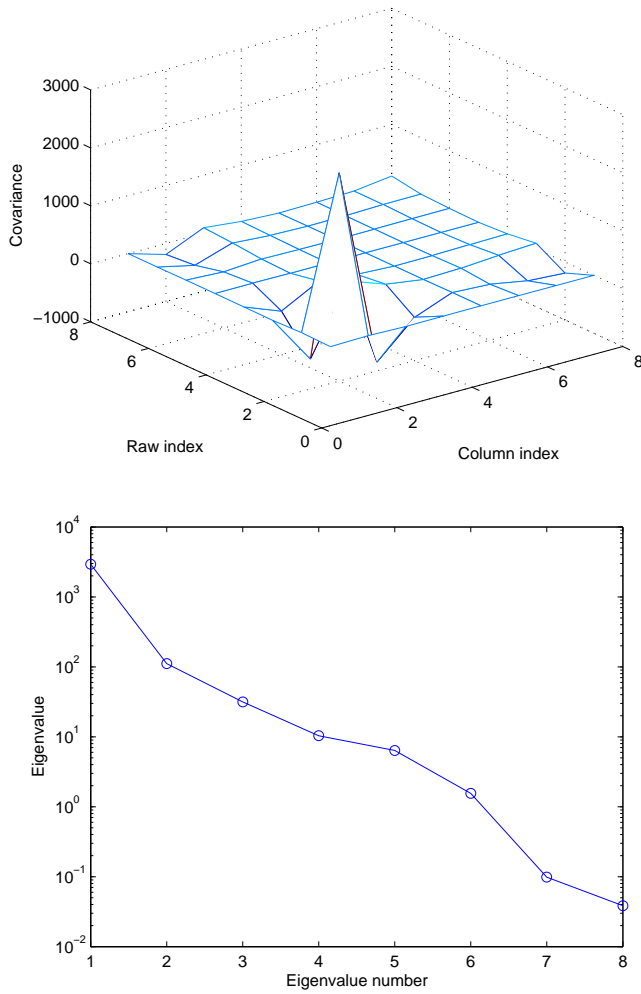


Figure 4. Example of covariance matrix of the vector x for the case of the tropical climate zone and sea surface; top panel, covariance matrix (numbering of rows and columns is consistent with the statistic numbering shown in Table 3); bottom panel, eigenvalues.

One way to partially overcome the independence assumption is to perform a previous transform of statistics to nearly independent ones. This can be accomplished by a classical principal component analysis, PCA (that gives independence under the assumption of Gaussian density functions of the statistics), or independent component analysis, ICA, that does not need the hypothesis of Gaussianity. PCA also allows one to easily reduce dimensionality of the problem by retaining only a small number of principal components, PC (for ICA this is more questionable because according to theoretical arguments independent components are not sorted and no agreed consolidated methodology exists for their selection). In the case that one or two principal components are sufficient to approximate the statistics, then the approximation of independence is completely avoided.

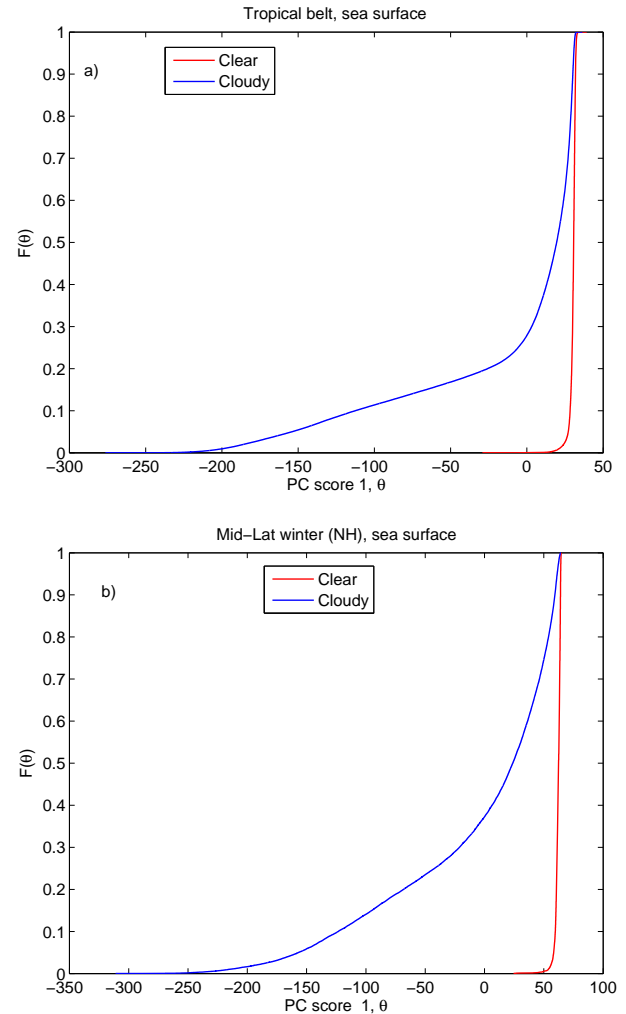


Figure 5. Example of cumulative distribution function, F , of the first PC score of the statistic vector x (Eq. 10) over sea surface. (a) tropical climatic zone; (b) mid-latitude winter (NH) climatic zone.

Use of principal components is easily accomplished by performing a transform of the original radiance data (x) into principal components (c) by PCA, then retaining the first $\bar{D} \leq D$ principal components $c^1, \dots, c^{\bar{D}}$, and finally straightforwardly applying the algorithms described in the present section to them.

5 Results

This section will show the main results of the numerical experiments worked out on the databases IASI1 and IASI2 and with the additional validation data set qualified for clear sky with the SEVIRI cloud mask. Finally, a comparison with the IASI level 2 cloud mask implemented at EUMETSAT will be shown.

Table 4. CDA results based on coinciding training and test data sets over sea surface and ice. The percentage of hit or success (TP) and of correct non-events (TN) can be obtained by $100 - \text{FP}$ and $100 - \text{FN}$, respectively.

Climate zone	PC	FN (miss)	FP (false alarm)	Merit function
Tropical	1	17.4	17.4	82.6
Mid-lat summer NH	1	15.2	15.2	84.8
Mid-lat summer SH	1	11.2	11.2	88.8
Mid-lat winter NH	1	9.4	9.4	90.6
Mid-lat winter SH	1	16.4	16.4	83.6
High-lat summer NH	1,2	7.1	7.1	92.9
High-lat summer SH	1	7.1	7.1	92.9
High-lat winter NH	1	7.8	7.8	92.2
High-lat winter SH	1,2	11.6	11.6	88.4
Sea ice NH day	1	23.4	9.6	85.0
Sea ice SH day	1,2	21.4	23.7	77.8
Sea ice NH night	1,2	44	15	84.0
Sea ice SH night	1	31	1.8	97.4

The present work is not intended to yield a fully operative cloud detection algorithm, which requires extensive and independent training and validation data sets covering all climatologies. Rather, we work out some sample test cases to prove feasibility of the full methodology.

We will present the results obtained by a transform to PCs. We will not show the optimal thresholds because they make sense only if endowed with the prior PC transform that, as is well known, is unique only apart from rotations. All thresholds and transform matrices are available upon request. For the present work, the statistic *sh* (see Table 3) was not considered, mostly because it added a poor value or no value to the remaining eight statistics.

Let us define the vector of statistics, \mathbf{x} , as

$$\mathbf{x} = \left(h_s, \chi_s^2, T_0, \Delta T_{\text{CO}_2}, W_1, W_2, W_3, W_4, \right)^T, \quad (10)$$

where the superscript T stands for transpose. Note that the statistics are ordered in the vector according to their list number as defined in Table 3. For the case of tropical climate zone and sea surface, Fig. 4 (left) shows the covariance matrix of \mathbf{x} used to build up the PC basis along with the eigenvectors. The covariance matrix was obtained on the basis of the data set IASI2. It is seen that the matrix is not diagonal, which supports the transformation of \mathbf{x} to principal components. The plot of eigenvalues, shown in the same figure (right) indicates that the first component exhibits most variance of the statistics. It is also noteworthy that the vector \mathbf{x} is made up of unlike quantities, some dimensionless and others in units of K. In principle we could standardize it before the transformation; however we did not get better results. We also have to keep in mind that if the units are changed (e.g. Kelvin to centigrade degrees for temperature), then the thresholds have to be re-estimated.

5.1 Validation – coinciding training and validation data sets

To summarize the performance of the cloud mask we introduce the merit function

$$\mathcal{M} = 100(1 - C).$$

\mathcal{M} is, indeed, the fraction of successful classifications S (see Eq. 1).

The CDA methodology has been applied to select the thresholds for the cloud detection scheme. In doing so, the statistics introduced in Sect. 3 have been transformed to the PC space as in Sect. 4.7. In the training phase the thresholds have been estimated on the basis of IASI1 or IASI2 data set depending on which population has the larger size for the climate zones at hand. In the present Section validation is made for the same training set; later on (Sect. 5.2) validation will be made on the other (smaller) data set. For example we used IASI2 as a training data set and will use IASI1 as a validation one in Sect. 5.2 for the mid-latitude climatic zone in winter (96 522 and 40 941 samples, respectively); we made the reverse choice in summer (81 729 samples for IASI1 and 40 177 for IASI2). Both data sets are quite equivalent for the tropical climatology (186 787 samples for IASI1 and 230 267 for IASI2). Note that we could have jointly used both data sets in the training phase when training and validation data sets coincide; however we kept their distinction to have the same training data set in both experiments of coinciding (this section) and different (Sect. 5.2) training and validation data sets.

Thresholds have been estimated for all the climatic zones of Table 1.

Table 5. CDA results based on coinciding training and validation data sets for daytime soundings over land. The percentage of hit or success (TP) and of correct non-events (TN) can be obtained by $100 - \text{FP}$ and $100 - \text{FN}$, respectively.

Climate zone	PC	FN (miss)	FP (false alarm)	Merit function
Tropical	1	8.4	8.4	91.6
Mid-lat summer NH	1	17.6	17.6	82.4
Mid-lat summer SH	1	6.9	6.9	93.1
Mid-lat winter NH	1	11.1	11.1	89.9
Mid-lat winter SH	1	10.4	10.4	89.6
High-lat summer NH	1	11.8	11.8	88.2
High-lat winter NH	1	10.3	10.3	89.7
Land below 1 km with snow/ice NH	1	17.6	23	77.3
Land below 1 km with snow/ice SH	1	16.3	14.1	85.3
Land above 1 km with snow/ice NH	1	7.5	7.5	92.5
Land above 1 km with snow/ice SH	1	16.2	12.4	86.4
Antarctica below 1 km with snow/ice	1	7.4	7.4	82.6
Antarctica above 1 km with snow/ice	1	34.9	34.9	65.1

Table 6. CDA results based on coinciding training and validation data sets for night-time soundings over land. The percentage of hit or success (TP) and of correct non-events (TN) can be obtained by $100 - \text{FP}$ and $100 - \text{FN}$, respectively.

Climate zone	PC	FN (miss)	FP (false alarm)	Merit function
Tropical	1	17.8	17.8	82.2
Mid-lat summer NH	1	12.2	12.2	87.8
Mid-lat summer SH	–	–	–	–
Mid-lat winter NH	1	10.1	10.1	89.9
Mid-lat winter SH	1	12	12	88.0
High-lat summer NH	1	8.5	8.5	91.5
High-lat winter NH	1	13.7	13.7	86.3
Land below 1 km with snow/ice NH	1	35.7	10	81.2
Land below 1 km with snow/ice SH	1	19.5	18.5	81.3
Land above 1 km with snow/ice NH	2	12	53	69.0
Land above 1 km with snow/ice SH	1	10	45	65.0
Antarctica below 1 km with snow/ice	2	35.2	35.2	65.0
Antarctica above 1 km with snow/ice	2	26.3	26.3	73.7

5.1.1 Sea surface

The performance of the CDA methodology for the relevant climatic zones is summarized in Table 4. We found that CDA always selects only the first principal component to optimally discriminate clear from cloudy sky, but high-lat summer NH, high-lat winter SH and two sea ice climatologies seem to yield a better score with the first and second principal components.

Figure 5 exemplifies the cumulative distribution function, F , for clear and cloudy sky as estimated from the training data set. Figure 5a applies to tropics and it is seen that the right tail of F^{Cloudy} has a relatively wide overlap with F^{clear} . This overlap is likely due to warm clouds, which tend to be radiometrically equivalent to the colder clear-sky scenes. Another source of possible misclassification is due to IASI

sub-pixel clouds, that is, broken clouds which do not uniformly fill the IASI field of view. At nadir, the IASI field of view is 12 km (compare with the 1 km pixel size of AVHRR). Finally, it should be stressed that the overlap could also be the result of a non-perfect AVHRR reference/training cloud mask. However, whatever its origin may be, because of this overlap the performance for the tropics is 82.6 % with error of first and second type of 17.4 % (see Table 4).

Figure 5b exemplifies the case of the mid-lat winter (NH) climate zone. It is now possible to see that the overlap zone is smaller than that observed for the previous case. As a result, we have a total performance above 90 %, while the error of first and second type is smaller than 10 % (see Table 4).

For sea ice, it was not possible to find an optimal solution with both the fraction of miss and false alarm below 20 %. This fraction was considered to be a sort of limiting

condition for an acceptable cloud scene analysis. When this was not possible, we chose to keep, whenever possible, the false alarm fraction below 20 %, and the thresholds were individuated by trial and error in order to meet this condition.

The lack of optimal thresholds (according to the cost function of Eq. 4) for sea ice is mostly due to the loss of thermal contrast between clouds and surface, which becomes very severe during polar nights at high latitudes. Once again, the performance for sea ice is also a result of the difficulty to develop a good training data set for these surface features.

5.1.2 Land surface

We recall that for land surface we distinguish between day and night for each climatic zone of Table 1. Tables 5 and 6 summarize the results of CDA for day and night, respectively. Missing values mean that there were not enough data available to perform the CDA analysis. We see that the performance of the cloud detection, according to the thresholds selected by CDA, is quite close to 90 %, with the type II error close to 10 % for most cases.

For land, surface temperature is strongly driven by the daily sun cycle and yields a very high thermal contrast during daytime. This contrast makes it quite easier to discriminate clouds from clear sky at the tropics. Its effect is reflected in the shape of the cumulative distribution function, F^{Clear} , as it is possible to see from Fig. 6a, which applies to the tropical climatic zone during daytime. It is seen that the right far tail F^{Clear} has no overlap with F^{Cloudy} . Because of the daytime thermal contrast between clouds and surface, the cloud detection performance is quite close to 90 %, apart from a few cases corresponding to snow/ice.

Conversely, for night-time the performance tends to decrease because of loss of thermal contrast, which is more severe in the case of snow/ice. For snow/ice the cloud detection becomes much more difficult because F^{Clear} and F^{Cloudy} tend to overlap (see, e.g., Fig. 6b). When this occurs the thresholds have been tuned ad hoc in some cases to compensate for the driving effect of the cloudy population. When it was not possible to keep the fraction of false alarms below 20 %, we used the alternate rule to have a total performance higher than 65 %. These rules are somewhat arbitrary and just reflect the difficulty to have a good cloud detection over snow/ice.

5.2 Validation – different training and validation data sets

Apart from the tropics, the data sets IASI1 and IASI2 cover different climatic zones. Therefore we present results only for the tropical climatic zone.

For the case of sea surface we have used the data set IASI2 for the training and the data set IASI1 for validation. Performance of the methodology is shown in Table 7.

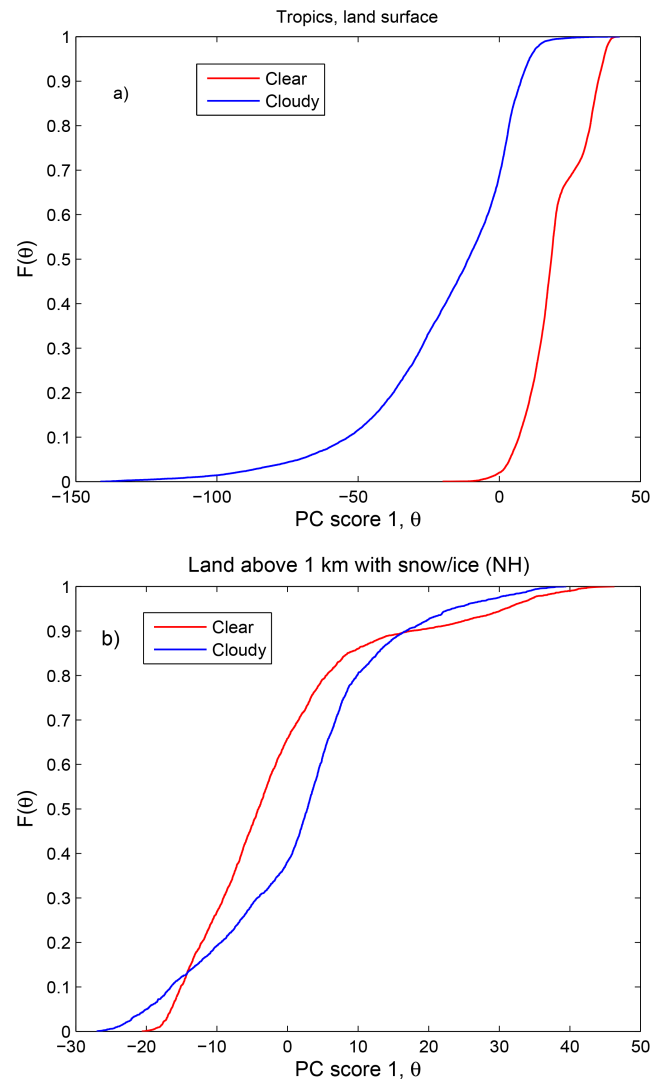


Figure 6. Example of cumulative distribution function, F , of the first PC score of the statistic vector x (Eq. 10) over land surface. **(a)** Tropical climatic zone (daytime); **(b)** land above 1 km with snow/ice (NH, night-time).

The same experiment has been carried out for land, where we distinguish between day and night (results also shown in Table 7). During daytime the data set used for training is IASI1, whereas validation is performed on the data set IASI2. The reverse occurs during night-time because of the better coverage of the latter.

Comparing Table 7 with Table 4, it is seen that the performance obtained on the validation data set is largely consistent with the case of coinciding training and validation data sets. This occurs even when the two data sets are largely nonhomogeneous, showing the robustness of the methodology.

Table 7. Validation results using data sets IASI1 and IASI2. The total score (S) or equivalently merit function, false alarms and miss are shown in percentage.

Area	Training set size	FN (miss)	FP (false alarm)	Merit function
Tropics, sea surface	186 787	15.2	16.7	83.3
Tropics, land daytime	172 216	15.5	6.7	93.3
Tropics, land night-time	85 882	5.4	20.0	80.0

Table 8. Validation results based on the co-location of IASI/SEVIRI. The first row refers to the IASI2 data set and the second row to the GSICS data set.

Area	Training set size	FN (miss)	FP (false alarm)	Merit function
IASI2 set (Europe/Africa)	47 160	7.9	5.7	96.7
GSICS set (Europe/Africa)	4277	11.3	16.7	86.7

5.3 Validation with SEVIRI

Validation has been also performed with a cloud mask independent of CMS, on which the training data set relies. For southern Europe and the African continent three consecutive, daytime, IASI orbits belonging to the set IASI2 have been co-located with SEVIRI imagery.

The IASI footprints considered in this experiment are shown in Fig. 7. The first IASI footprint observation took place on 22 July 2007 at 06:42:12 UTC, and the last observations at 10:14:18 UTC. To have an idea of the cloud coverage and type over the target area, Fig. 7 shows an RGB composite SEVIRI full-disk image starting at 09:27:43 UTC of the same day.

Since the orbits are in the daytime, the SEVIRI cloud mask also benefits from visible channels. In contrast only information in the infrared spectral region is used with IASI. The IASI footprints have been co-located with the SEVIRI operational cloud mask and flagged as clear sky if 100 % of SEVIRI pixels falling within the IASI field of view are themselves clear sky. In this way we have a SEVIRI-based cloud mask for IASI which can be compared with that obtained with our procedure.

The results for this experiment are shown in Table 8. It is seen that the agreement with the SEVIRI-based cloud mask is excellent, which further testifies the very good performance of the cloud detection over land in the case of the tropical climatology.

The SEVIRI cloud mask has been also used to classify a data set independent of IASI1 and IASI2. We have used a set of IASI data co-located with SEVIRI imagery from 25 September to 4 October 2010. It was acquired over Europe/Africa in the framework of the IASI/SEVIRI inter-calibration experiment, performed within the activities of EUMETSAT GSICS. We have a total of 4277 IASI spectra, whose footprints are shown in Fig. 8.

As for the previous exercise, the IASI footprints have been co-located with the SEVIRI operational cloud mask and flagged as clear sky if 100 % SEVIRI pixels falling within the IASI field of view are themselves clear sky. The results are shown in the second row of Table 8.

Again we have a very good comparison with a total score of 86.7 %. In this case we have an excess of false alarms with respect to the Africa case study also shown in Table 8. However, we have ascertained that the problem is likely due to the IASI/SEVIRI co-location used in the GSICS system.

6 Conclusions

With regard to developing a cloud mask for MTG-IRS, we have designed and implemented a stand-alone cloud detection tool expected to run only using MTG-IRS spectral radiances.

The strategy has been implemented and tested using IASI as a proxy for MTG-IRS. The scheme relies on a set of nine statistics. Six of them are heritage of the AVHRR-based cloud classification and discrimination, while three more have been specifically designed to take into account the hyperspectral characteristics of MTG-IRS.

Furthermore, a methodology (cumulative discriminant analysis) has been developed in the context of classification of binary variables, which, using all the statistics, optimally computes thresholds needed to classify a given pixel as clear or cloudy. The system necessitates a training data set as well as a corresponding reference cloud mask. Once the training has been performed, the algorithm only needs the thresholds and therefore is very fast from the computational point of view. The probabilistic approach allows us to quantify the type I (miss) and type II (false alarm) errors. Within the methodology, the two types of error can have the role of design parameters, so that the cloud detection scheme can have

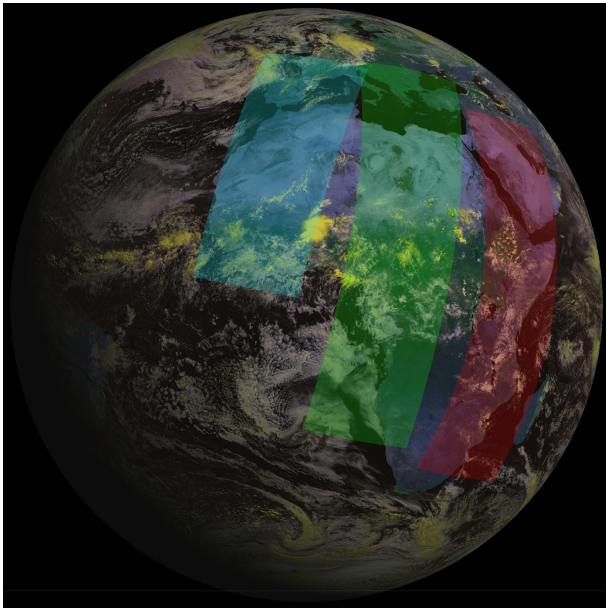


Figure 7. SEVIRI full-disk imagery starting at 09:27:43 UTC for the day 22 July 2007. The IASI scan patterns used for the validation exercise are the three coloured rectangles.

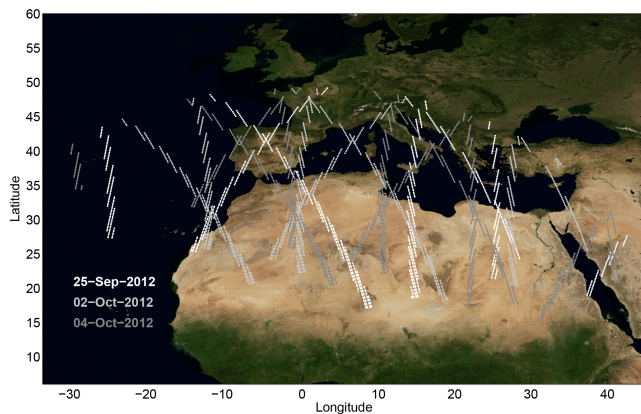


Figure 8. Footprint positions of IASI observations for the GSICS IASI/SEVIRI inter-calibration experiment.

a prescribed rate of false alarms. Alternately, we can design and develop a cloud detection scheme in which both types of error are simultaneously minimized.

The methodology we have developed is much more than a classification tool. In fact, the nine statistics are intended to summarize the information content of the IASI spectra into a limited number of variables as an input to the classification methodology. Once a transform of the statistics into principal components has been performed, CDA discriminates how many and which principal components are needed for an effective cloud detection. Nevertheless, CDA is general enough to be coupled with different other statistics.

Apart from a few situations, we have found that the first principal component of the statistics is enough to build up an effective cloud screening. This largely simplifies the scene analysis algorithm because we need to compute thresholds just for one principal component; moreover it circumvents assumption of independence of statistics; finally the overall algorithm is very efficient.

Basically, the methodology relies on a training data set and a reference cloud mask. It is recommended that the training data set be made up of real observations, having a global coverage and a full coverage of the possible real clear and cloudy scenes.

The reference cloud mask is needed to extract from the training data set the two classes of spectra, one for clear sky and the second one for cloudy sky. At present, we rely on the AVHRR cloud mask, which is operationally produced by CMS.

The algorithm has been tested against independent data sets as well as cloud masks. As a result, we have the agreement being normally well above 80 %, apart from regions covered by sea ice or land snow/ice, where also the reference cloud masks tend to have a relatively high rate of failure. We think that, once applied to real MTG-IRS scenes, the cloud detection should even improve performance because MTG-IRS is an imager and its pixel size is expected to be 4 times smaller than that of IASI.

The paper has been mostly dealing with the scientific basis of CDA and an initial evaluation of the scheme. Although the methodology yields results which are equivalent or better than those of concurrent schemes for IASI (e.g. Lavanant et al., 2011), we think that the scheme still deserves a more comprehensive training/validation. In this respect, we now have new evidence (Lavanant and Roquet, 2013) that all weather types (hence cloud types) are present on the full disk whatever the day considered. This is particularly true over sea, whereas over land the situation is a bit more complicated. In fact, because of orography all weather types are not all found in both hemispheres simultaneously, and in general the amount of clouds is smaller during the NH summer than during the NH winter. These findings support our choice of the 2-day global data set, with 1 day in the NH Summer and the other in the NH Autumn, which we have used for training. However, at least 2 more days would be needed for a better and comprehensive validation. Towards this objective and with regard to future use of the methodology for operational applications, we are currently developing an appropriate seasonal and global IASI/AVHRR/SEVIRI data set along the lines suggested by Lavanant and Roquet (2013).

Acknowledgements. Work was performed under the EUMETSAT EUM/CO/12/4600001033/SAT contract.

Edited by: B. Kahn

References

- Ackerman, S., Strabala, K., Menzel, W., Frey, R., Moeller, and C., Gumley, L.: Discriminating clear-sky from clouds with MODIS, *J. Geophys. Res.*, 103, 32–141, 1998.
- Ackerman, S., Holz, R., Frey, R., Eloranta, E., and McGill, B. M.: Cloud detection with MODIS. Part II: Validation, *J. Atmos. Oceanic Technol.*, 25, 1073–1086, 2008.
- Amato, U., Masiello, G., Serio, and C., Viggiano, M.: The σ -IASI code for the calculation of infrared atmospheric radiance and its derivatives, *Environ. Model. Softw.*, 17, 651–667, doi:10.1016/S1364-8152(02)00027-0, 2002.
- Amato, U., Antoniadis, A., Cuomo, V., Cutillo, L., Franzese, M., Murino, L., and Serio, C.: Statistical cloud detection from SEVIRI multispectral images, *Remote Sens. Environ.*, 112, 750–766, doi:10.1016/j.rse.2007.06.004, 2008.
- Benjamini, Y., and Hochberg, Y.: Controlling the false discovery rate: a practical and powerful approach to multiple testing, *J. R. Stat. Soc. Ser. B*, 57, 289–300, 1995.
- Chevalier F.: Sampled Database of 60 Levels Atmospheric Profiles from the ECMWF Analysis, Technical Report ECMWF EUMETSAT SAF programme Research, 4, ECMWF, Reading, 2001.
- Derrien, M. and Le Gléau, H.: MSG/SEVIRI cloud mask and type from SAFNWC, *Int. J. Remote Sens.*, 26, 4707–4732, 2005.
- Derrien, M.: Validation Report for “Cloud Products” (CMa-PGE01 v3.2, CT-PGE02 v2.2 andCTTH-PGE03 v2.2), EUMETSAT, Darmstadt, 2012.
- EUMETSAT, CM-SAF Product CM-02, CM-08 and CM-14 Cloud Fraction, Cloud Type and Cloud Top Parameter Retrieval from SEVIRI, EUMETSAT, Darmstadt, 2009.
- Fawcett, T.: ROC graphs: Notes and practical considerations for researchers, *ReCALL*, 31, HPL-2003-4, 1–38, 2004.
- Heidinger, A. K.: CLAVR-x Cloud Mask Algorithm Theoretical Basis Document, NOAA/NESDIS/Office of Research and Applications, Washington, D.C., 2004.
- Heidinger, A. K., Evan, A. T., Foster, M. J., and Walther, A.: A Naive Bayesian Cloud-Detection Scheme Derived from CALIPSO and Applied within PATMOS-x, *J. Appl. Meteor. Climatol.*, 51, 1129–1144, 2012.
- Hilton, F., Armante, R., August, T., Barnet, C., Bouchard, A., Camy-Peyret, C., Capelle, V., Clarisse, L., Clerbaux, C., Coheur, P. F., Collard, A., Crevoisier, C., Dufour, G., Edwards, D., Faijan, F., Fourrié, N., Gambacorta, A., Goldberg, M., Guidard, V., Hurtmans, D., Illingworth, S., Jacquinet-Husson, N., Kerzenmacher, T., Klaes, D., Lavanant, L., Masiello, G., Matricardi, M., McNally, A., Newman, S., Pavelin, E., Payan, S., Péquignot, E., Peyridieu, S., Phulpin, T., Remedios, J., Schlüssel, P., Serio, C., Strow, L., Stubenrauch, C., Taylor, J., Tobin, D., Wolf, W., and Zhou, D.: Hyperspectral Earth Observation from IASI: four years of accomplishments, *B. Am. Meteorol. Soc.*, 93, 347–370, doi:10.1175/BAMS-D-11-00027.1, 2012.
- Inoue, T.: On the temperature and effective emissivity determination of semi-transparent cirrus clouds by bi-spectral measurements in the 10 mm window region, *J. Meteor. Soc. Jpn.*, 63, 88–99, 1985.
- Inoue, T. and Ackerman, S.A.: Radiative effects of various cloud types as classified by the split window technique over the Eastern Sub-tropical Pacific derived from co-located ERBE and AVHRR data, *J. Meteor. Soc. Jpn.*, 80, 1383–1394, 2002.
- Lavanant, L., Marguinaud, P., Harang, L., Lelay, J., Péré, S., and Philippe, S.: Operational cloud masking for the OSI SAF global METOP/AVHRR SST product, Joint 2007 EUMETSAT Meteorological Satellite Conference and the 15th Satellite Meteorology & Oceanography Conference of the American Meteorological Society, EUMETSAT P. 50, available at: <http://www.eumetsat.int> (last access: 4 October 2014), 2007.
- Lavanant, L., Fourrie, N., Gambacorta, A., Grieco, G., Heilliette, S., Hilton, F. I., Kim, M.-J., McNally, A. P., Nishihata, H., Pavelin, E. G., and Rabier, F.: Comparison of cloud products within IASI footprints for the assimilation of cloudy radiances, *Quart. J. Royal Meteorol. Soc.*, 137, 1988–2003, 2011.
- Lavanant, L. and Roquet P.: Method to generate a reference cloud mask for MTG-IRS, Technical Memorandum, Météo-France/DP/CMS/R&D, Lannion, France, 56 pp., 2013.
- Liou, K. N.: Radiation and Cloud Processes in the Atmosphere, Oxford University Press, Oxford, 1992.
- Masiello, G., Matricardi, M., Rizzi, R., and Serio, C.: Homomorphism between cloudy and clear spectral radiance in the 800–900 CM-1 window region, *Appl. Opt.*, 41/6, 965–973, doi:10.1364/AO.41.000965, 2002.
- Masiello, G., Shimoda, H., and Serio, C.: Qualifying IMG Tropical Spectra for Clear Sky, *J. Quant. Spectrosc. Rad. Transf.*, 77/2, 131–148, doi:10.1016/S0022-4073(02)00083-3, 2003.
- Masuda, K., Takashima, T., and Takayma, Y.: Emisivity of pure and sea waters for the model sea surface in the infrared window regions, *Remote Sens. Environ.*, 24, 313–329, 1988.
- Rossow, W. B.: Measuring cloud properties from space: A review, *J. Climate*, 2, 201–213, 1989.
- Sandhya, K. N., Rajeev, K., and Parameswaran, K.: Cloud screening in IRS-P4 OCM satellite data: potential of spatial coherence method in the absence of thermal channel information, *Remote Sens. Environ.*, 90, 259–267, 2004.
- Serio, C., Lubrano, A. M., Romano, F., Shimoda, H.: Cloud detection over sea-surface using auto-correlation functions of upwelling infrared spectra in the 800–900 cm^{-1} region, *Appl. Opt.*, 29, 3565–3572, doi:10.1364/AO.29.003565, 2000.
- Serio, C., Masiello, G., Venafra, S., and Amato, U.: Application Data for MTG-IRS Cloud Detection Method, EUMETSAT, Darmstadt, 2013.
- Silverman, B. W.: Density estimation for statistics and data analysis, Chapman & Hall, London, 1986.
- Stowe, L. L., Davis, P. A., and McClain, E. P.: Scientific basis and initial evaluation of the CLAVR-1 global clear/cloud classification algorithm for the advanced very high resolution radiometer, *J. Atmos. Oceanic Technol.*, 16, 656–681, 1999.
- Tapakis, R. and Charalambides, A. G.: Equipment and methodologies for cloud detection and classification: A review, *Solar Energy*, 95, 392–430, 2012.
- van der Vaart, A. W.: Asymptotic statistics, Cambridge University Press, Cambridge, 1998.
- Wand, M. P. and Jones, M. C.: Kernel smoothing, Chapman & Hall, London, 1995.