



Retrieval of cloud fraction using machine learning algorithms based on FY-4A AGRI observations

Jinyi Xia and Li Guan

China Meteorological Administration Aerosol-Cloud and Precipitation Key Laboratory,
Nanjing University of Information Science and Technology, Nanjing 210044, China

Correspondence: Li Guan (liguan@nuist.edu.cn)

Received: 1 April 2024 – Discussion started: 10 June 2024

Revised: 10 October 2024 – Accepted: 11 October 2024 – Published: 25 November 2024

Abstract. Cloud fraction as a vital component of meteorological satellite products plays an essential role in environmental monitoring, disaster detection, climate analysis and other research areas. Random forest (RF) and multilayer perceptron (MLP) algorithms were used in this paper to retrieve the cloud fraction of AGRI (Advanced Geosynchronous Radiation Imager) on board the Fengyun-4A (FY-4A) satellite based on its full-disk level-1 radiance observation. Corrections have been made subsequently to the retrieved cloud fraction in areas where solar glint occurs using a correction curve fitted with sunglint angle as weight. The algorithm includes two steps: the cloud detection is conducted firstly for each AGRI field of view to identify whether it is clear sky, partly cloudy or overcast within the observation field. Then, the cloud fraction is retrieved for the scene identified as partly cloudy. The 2B-CLDCLASS-lidar cloud fraction product from the CloudSat and CALIPSO active remote sensing satellite is employed as the truth to assess the accuracy of the retrieval algorithm. Comparison with the operational AGRI level-2 cloud fraction product is also conducted at the same time. The results indicate that both the RF and MLP cloud detection models achieved high accuracy, surpassing that of operational products. However, both algorithms demonstrated weaker discrimination capabilities for partly cloudy conditions compared to clear-sky and overcast situations. Specifically, they tended to misclassify fields of view with low cloud fractions (e.g., cloud fraction = 0.16) as clear sky and those with higher cloud fractions (e.g., cloud fraction = 0.83) as overcast. Between the two models, RF exhibited higher overall accuracy. Both RF and MLP models performed well in cloud fraction retrieval, showing lower mean error (ME), mean absolute error (MAE) and root mean

square error (RMSE) compared to operational products. The ME for both RF and MLP cloud fraction retrieval models was close to zero, while RF had slightly lower MAE and RMSE than MLP. During daytime, the high reflectance in sunglint areas led to larger retrieval errors for both RF and MLP algorithms. However, after correction, the retrieval accuracy in these regions improved significantly. At night, the absence of visible light observations from the AGRI instrument resulted in lower classification accuracy compared to daytime, leading to higher cloud fraction retrieval errors during nighttime.

1 Introduction

Clouds occupy a significant proportion within satellite remote sensing data acquired for Earth observation. According to the statistics from the International Satellite Cloud Climatology Project (ISCCP), the annual average global cloud coverage within satellite remote sensing data is around 66 %, with even higher cloud coverage in specific regions (such as the tropics) (Zhang et al., 2004). The impact of clouds on the radiation balance of the Earth's atmospheric system is influenced by the optical properties of clouds. Cloud detection, as a vital component of remote sensing image data processing, is considered a critical step for the subsequent identification, analysis and interpretation of remote sensing images. Therefore, accurately determining cloud coverage is essential in various research domains, such as environmental monitoring, disaster surveillance and climate analysis.

Fengyun-4A (FY-4A) is a comprehensive atmospheric observation satellite launched by China in 2016. The uploaded AGRI (Advanced Geosynchronous Radiation Imager) has 14

channels and captures full-disk observation every 15 min. In addition to observing clouds, water vapor, vegetation and the Earth's surface, it also possesses the capability to capture aerosols and snow. Moreover, it can clearly distinguish different phases and particle size of clouds and obtain high- to mid-level water vapor content. It is particularly suitable for cloud detection due to its simultaneous use of visible, near-infrared and longwave infrared channels for observation with 4 km spatial resolution.

Numerous cloud detection algorithms have been provided based on observations from satellite-borne imagers. The threshold method has been widely employed by researchers, including the early ISCCP (International Satellite Cloud Climatology Project) method (Rossow and Leonid, 1993) and the proposed threshold methods based on different spectral features or underlying surfaces (Kegelmeyer, 1994; Solvsteen, 1995; Baum and Trepte, 1996). However, there is a significant subjectivity in selection of thresholds as to whether it is the single and fixed threshold in the early days, multiple thresholds, dynamic thresholds, or adaptive thresholds. The selection of thresholds is influenced by season and climate. Surface reflectance varies significantly between different seasons, such as increased reflectance from snow in winter and vegetation flourishing in summer affecting reflectance. As a result, changes in surface features during different seasons lead to variations in the distribution of grayscale values in images, requiring adjustments to thresholds based on seasonal characteristics. Climate conditions like cloud cover and atmospheric humidity impact the distinguishability of clouds and other features. For instance, in humid or cloudy climates, the reflectance of the surface and clouds may be similar, necessitating stricter thresholds for differentiation. Therefore, climate conditions also influence threshold selection.

The other category of cloud detection algorithms is based on statistical probability theory. For example the principal component discriminant analysis and quadratic discriminant analysis methods were used for SEVIRI (Spinning Enhanced Visible and Infrared Imager) cloud detection (Amato et al., 2008). The cloud detection algorithm for the Thermal Infrared (TIR) sensor was based on the Bayesian theory of total probability (Merchant et al., 2010) and the naive Bayes algorithm for AGRI (Yan et al., 2022). The unsupervised clustering cloud detection algorithms for MERIS (Medium Resolution Imaging Spectrometer) (Gomez-Chova et al., 2006) and the combining *k*-means clustering and Otsu's method for MODIS (Xiang, 2018) all have achieved high accuracy in cloud detection.

More and more machine learning algorithms are being utilized by researchers in cloud detection studies with the development of machine learning. For instance, probabilistic neural networks, especially radial basis function networks, were used for AVHRR cloud detection (Zhang et al., 2001). The utilization of convolutional neural network methods (Chai et al., 2024) offers important perspectives for cloud detection research.

Currently, there is limited research literature on cloud detection and cloud fraction retrieval algorithms for FY-4A/4B AGRI. The operational cloud fraction product of FY-4A AGRI utilized a threshold method with 4 km spatial resolution. Differences in climatic and environmental factors lead to varying albedo and brightness temperature observations for the instrument at different times and locations. Therefore, the choice of thresholds is easily influenced by factors such as season, latitude and land surface type (Gao and Jing, 2019). Using multiple sets of thresholds for discrimination would significantly slow down the cloud detection process. Moreover, most algorithms focus solely on cloud detection, which classified the observed scenes as cloud or clear sky without providing the specific cloud fraction information for the scenes. The use of active remote sensing instruments carried by Cloudsat and CALYPSO is not influenced by thresholds when retrieving cloud fraction, enabling a more accurate cloud fraction retrieval. However, due to Cloudsat and CALYPSO being polar-orbiting satellites, the cloud fraction over the full disk cannot be obtained. Utilizing the Cloudsat and CALYPSO level-2 product 2B-CLDCLASS-lidar as the reference truth, a random forest model trained based on FY-4A AGRI full disk radiation data can address the shortcomings of threshold methods and achieve a high accuracy of cloud fraction over the full disk.

In summary, this paper established cloud detection and cloud fraction retrieval models using multilayer perceptron (MLP) and random forest (RF) algorithms, based on FY-4A AGRI full-disk level-1 observed radiance data. The cloud fraction from the CloudSat and CALYPSO level-2 product 2B-CLDCLASS-lidar was used as the label. The results were compared with the 2B-CLDCLASS-lidar product and the official AGRI operational products for validation.

2 Research data and preprocessing

2.1 FY-4A data

FY-4A was successfully launched on 11 December 2016. Starting from 25 May 2017, FY-4A drifted to a position near the main business location of the Fengyun geostationary satellite at 104.7° east longitude on the Equator. Its successful launch marked the beginning of a new era for China's next-generation geostationary meteorological satellites as an advanced comprehensive atmospheric observation satellite. The Advanced Geosynchronous Radiation Imager (AGRI), one of the main payloads of the Fengyun-4 series geostationary meteorological satellites, can perform large-disk scans and rapid regional scans at a minute level. It has 14 observation channels in total with the main task of acquiring cloud images. The channel parameters and main uses of AGRI are detailed in Table 1 (<https://www.nsmc.org.cn/nsmc/cn/instrument/AGRI.html>, last access: 20 November 2024). The first six visible light channels have no values at

Table 1. FY-4A AGRI channel parameters.

Channel number	Band range (μm)	Central wavelength (μm)	Spatial resolution (km)	Main applications
1	0.45–0.49	0.47	1	clouds, dust, aerosols
2	0.55–0.75	0.65	0.5	clouds, sand dust, snow
3	0.75–0.90	0.825	1	vegetation
4	1.36–1.39	1.375	2	cirrus
5	1.58–1.64	1.61	2	clouds, snow
6	2.10–2.35	2.225	2	cirrus, aerosols
7	3.50–4.00	3.75 H	2	fire point, the intense solar reflection signal
8	3.50–4.00	3.75 L	4	low clouds, fog
9	5.80–6.70	6.25	4	upper-level water vapor
10	6.90–7.30	7.1	4	mid-level water vapor
11	8.00–9.00	8.5	4	subsurface water vapor
12	10.30–11.30	10.8	4	surface and cloud-top temperatures
e	11.5 0–12.50	12.0	4	surface and cloud-top temperatures
14	13.2–13.8	13.5	4	cloud-top height

night, meaning that channels with a central wavelength less than or equal to 2.225 μm are unavailable during nighttime. FY-4A AGRI data were downloaded from the official website of the China National Satellite Meteorological Center (<http://satellite.nsmc.org.cn>, last access: 20 November 2024), including level-1 full disk radiation observation data preprocessed through quality control, geolocation and radiation calibration as well as the level-2 cloud fraction (CFR) product. The spatial resolution of all these data is 4 km at nadir, and the temporal resolution is 15 min.

2.2 CloudSat and CALIPSO cloud product

CALIPSO (Cloud-Aerosol Lidar and Infrared Pathfinder Satellite Observations) is a satellite jointly launched by NASA and CNES (the French National Center for Space Studies) in 2006. It is a member of the A-Train satellite observation system. CALIPSO is equipped with three payloads, among which CALIOP (the Cloud and Aerosol Lidar with Orthogonal Polarization) is a primary observational instrument. Observing with dual wavelengths (532 and 1064 nm) CALIOP can provide high-resolution vertical profiles of clouds and aerosols with 30 m vertical resolution. As the first satellite designed to observe global cloud characteristics in a sun-synchronous orbit, CloudSat is also among NASA’s A-Train series satellites. The CPR (Cloud Profiling Radar) installed on it operates at 94 GHz millimeter wavelength and is capable of detecting the vertical structure of clouds and providing vertical profiles of cloud parameters. The scanning wavelengths of CPR and CALIOP are different. CALIOP is capable of observing the top of mid-level to high-level clouds, whereas CPR can penetrate optically thick clouds. Combining the strengths of these two instruments enables the acquisition of precise and detailed information on cloud layers and cloud fraction.

The joint level-2 product 2B-CLDCLASS-lidar is mainly utilized in this study. It provides the cloud fraction at different heights with horizontal resolution 2.5 km (along-track) × 1.4 km (cross-track) through combining the observations from CPR and CALIOP. Since the two instruments have a different spatial domain such as vertical resolution, spatial resolution and spatial frequency, the spatial domain of the output products is defined in terms of the spatial grid of the CPR. In the algorithm, the cloud fraction is calculated using a weighted scheme based on the spatial probability of overlap between the radar and lidar observations. The calculation of the lidar cloud fraction within a radar footprint is represented by Eq. (1) (Mace, 2014):

$$C_1 = \frac{\sum_{i=1}^{\text{no. of lidar obs}} w_i \delta_i}{\sum_{i=1}^{\text{no. of lidar obs}} w_i}, \tag{1}$$

where C_1 represents the lidar cloud fraction within a radar footprint; w_i is the spatial probability of overlap for a particular lidar observation; δ_i indicates the lidar hydrometeor occurrence, where a value of 1 signifies the presence of hydrometeor and 0 indicates the absence; and i counts the lidar profile in a specific radar observational domain.

This calculation considers the contributions of multiple lidar observations within a radar resolution volume to determine the cloud fraction within that volume. The CloudSat product manual (Wang, 2019) can be referred to for more detailed information on 2B-CLDCLASS-lidar. The data used are available to download from the ICARE Data and Services Center (<https://www.icare.univ-lille.fr/data-access/data-archive-access/>, last access: 20 November 2024).

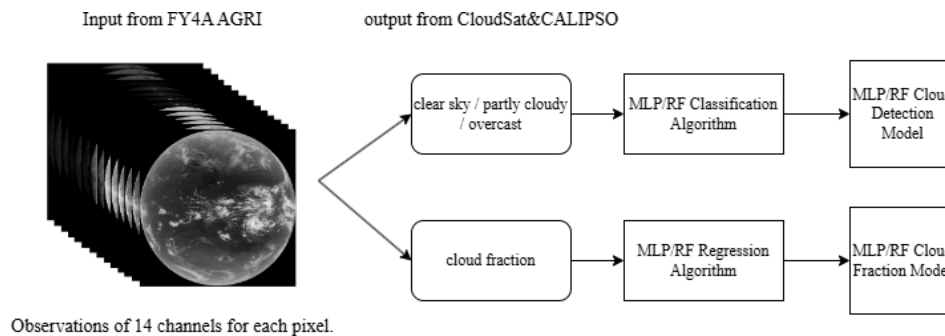


Figure 1. Method workflow. The input consists of 14 channel observation values for each pixel from FY-4A AGRI, and the ground truth labels or outputs are sourced from the CloudSat and CALIPSO cloud fraction products. The cloud detection classification model and the cloud fraction retrieval model are established separately.

2.3 Establishment of training data

The crucial aspect of establishing a training data in machine learning algorithms is how to obtain the cloud fraction values (ground truth) as labels. The error in cloud fraction retrieved solely from passive remote sensing instruments is significant. Using active remote sensing data can provide more accurate cloud fraction information in the vertical direction. Therefore, the spatiotemporally matched 2B-CLDCLASS-lidar cloud fractions are utilized as output labels in this paper.

The FY-4A AGRI and 2B-CLDCLASS-lidar data with a spatial difference between fields of view within 1.5 km and a time difference within 15 minutes are spatiotemporally matched. To make the 2B-CLDCLASS-lidar cloud fraction data collocated within AGRI pixels more effective, at least two 2B-CLDCLASS-lidar pixels are required within each AGRI field of view. The cloud fraction average of these pixels is used as the cloud fraction for that AGRI pixel. However, the errors in the matched dataset are unavoidable. The AGRI scanning method operates from left to right and top to bottom. Each complete scan of the full disk takes 15 min and generates a dataset. It is impossible to determine the exact moment of a specific point within the full disk. This limits the time range for matching datasets to within 15 minutes. However, in areas with higher wind speeds, clouds can move a significant distance within that 15 min window. Therefore, errors arising from timing issues cannot be avoided.

Cloud detection and cloud fraction label generation for 2B-CLDCLASS-lidar are as follows. There may be multiple layers of clouds in each field of view. If there is at least one layer cloud with cloud fraction of 1 in the 2B-CLDCLASS-lidar profile, then the scene is labeled as overcast with a cloud fraction of 1. If all layers in the profile are cloud-free, the scene is labeled as clear sky. The scene between the above two situations is labeled as partly cloudy, and the cloud fraction is the average of cloud fractions at different layers.

The algorithm includes two steps: the cloud detection is conducted firstly for each AGRI field of view to identify whether it is clear sky, partly cloudy or overcast within the

observation field. Then, the cloud fraction is retrieved for the scene identified as partly cloudy. So the training data include dataset A used for cloud detection and dataset B for cloud fraction retrieval. The input variables in dataset A are the FY-4A AGRI level-1 radiative observations from 14 channels, and the output variable is the temporally and spatially matched 2B-CLDCLASS-lidar cloud detection label. The output is categorized into three types: overcast, partly cloudy and clear sky, with values 1, 2 and 3 respectively. The cloud fraction product from 2B-CLDCLASS-lidar consists of discrete values: 0, 0.16, 0.33, 0.50, 0.66, 0.83 and 1. According to the result statistics, the cloud fractions of 2B-CLDCLASS-lidar pixels within the AGRI field of view are mostly the same. After averaging, the proportions of cloud fractions of [0.16, 0.33, 0.5, 0.67, 0.83] are extremely high. Therefore, other cloud fraction situations with extremely small proportions can be ignored. Doing so can also better balance the training samples. Here, 0 indicates clear sky, values from 0 to 1 represent varying cloud fractions for partly cloudy conditions and 1 signifies overcast. To ensure the balance and representativeness of the samples, the proportions of different cloud fraction samples in dataset A are set at 5 : 1 : 1 : 1 : 1 : 1 : 5. Regarding the samples for partly cloudy type in dataset A, the collocated 2B-CLDCLASS-lidar cloud fraction products serve as output labels for cloud fraction retrieval model B. The input of training dataset B remains the FY-4A AGRI level-1 radiative observations.

Due to the instrument's limited lifespan, only 2B-CLDCLASS-lidar data up to August 2019 can be obtained. The sample time range used in this paper is from August 2018 to July 2019. A period of 5 d was randomly selected each month as daytime samples and 5 d as nighttime samples. A total of 120 d of time- and space-matched FY-4A AGRI full-disk observations and 2B-CLDCLASS-lidar data were used as training and testing samples. Among them, 80 % of the data were used for training, and 20 % were used for testing. The total number of daytime samples in dataset A is 91 073, while dataset B contains 30 358 samples. The to-

Table 2. Recall rate (POD) and false alarm rate (FAR) of operational cloud detection products and multiple models.

	Sky classification	Daytime product	Nighttime product	Daytime RF	Nighttime RF	Daytime MLP	Nighttime MLP
POD	Clear sky	0.6359	0.5781	0.964	0.919	0.959	0.905
	Partly cloudy	0.7174	0.7449	0.914	0.845	0.895	0.808
	Overcast	0.7736	0.7384	0.959	0.919	0.957	0.920
FAR	Clear sky	0.1778	0.0934	0.047	0.102	0.064	0.131
	Partly cloudy	0.1819	0.2117	0.078	0.153	0.085	0.172
	Overcast	0.2499	0.2683	0.038	0.061	0.039	0.063

Table 3. Errors of cloud fraction.

	Daytime product	Nighttime product	Daytime RF	Daytime MLP	Nighttime RF	Nighttime MLP
ME	0.1987	0.2121	0.0006	-0.0009	-0.0028	-0.0032
MAE	0.2279	0.2441	0.1011	0.1053	0.1221	0.1322
RMSE	0.2776	0.2938	0.1285	0.1332	0.1510	0.1623

tal number of nighttime samples in dataset A is 95 493, and dataset B includes 31 831 samples.

Although the model was trained and tested using data from 2018 to 2019, to test the universality of the algorithm, it was applied to real-time observations from FY-4A and FY-4B AGRI in 2023.

3 Algorithms

Our preliminary experiments involved multiple algorithms, including LIBSVM, MLP, BP neural network and random forest. These experiments highlighted that, among the baselines, random forest and MLP achieved the highest overall accuracy. For this reason, we selected them to perform additional experiments. Using RF and MLP algorithms to train the model with the established sample set, the overall process is shown in Fig. 1.

3.1 Random forest (RF)

This algorithm integrates multiple trees based on the bagging idea of ensemble learning, with the basic element being the decision tree (Breiman, 1999). When building a decision tree, N sets of independent and dependent variables are randomly sampled with replacement from the original training samples to create a new training sample set; m variables are randomly sampled without replacement from all independent variables, the dependent variable data are split into two parts using the selected variables, and the purity of the subsets is calculated for each split method. The variable utilized by the split method with the highest purity is used to partition the data, completing the decision at that node. This process of binary splitting continues to grow the decision tree un-

til stopping criteria are met, completing the construction of a single decision tree. These steps are repeated N_{tree} times to build a random forest model consisting of N_{tree} decision trees (Breiman, 2001). Random forest adopts ensemble algorithms, with the advantage of high accuracy. It can handle both discrete and continuous data, without the need for normalization, making it more efficient compared to other algorithms.

3.2 Multilayer perceptron (MLP)

This algorithm consists of a fully connected artificial neural network (Duda, et al., 2007). The classifier/regressor takes feature vectors or tensors as input. The input is mapped through multiple fully connected hidden layers containing hidden weights, which produce classifications/regressions at the output layer. A nonlinear activation function (such as sigmoid or rectified linear unit (ReLU)) is applied in each hidden layer to facilitate a nonlinear model. For classifiers, the output of the final hidden layer is combined and passed through a softmax function to generate class predictions. For the loss function, the cloud detection model is cross-entropy, and the cloud fraction model is the mean square error (MSE). The model’s weights are trained in a supervised manner using backpropagation.

3.3 Hyperparameters

In this paper, a total of eight models were established, including daytime/nighttime random forest classification/regression models and daytime/nighttime MLP classification/regression models. For the random forest, we first conducted experiments using the following hyperparameter ranges: trees – [200, 300, 400, 500, 600, 700], minleaf – [1, 2, 5,

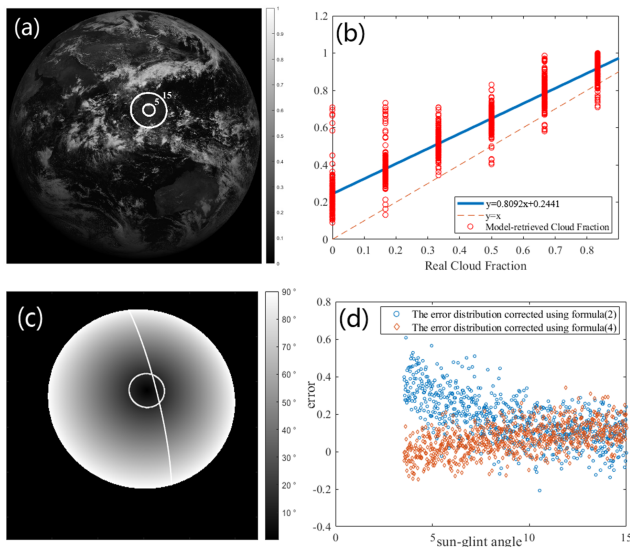


Figure 2. (a) Albedo image of the $0.67\ \mu\text{m}$ channel (the circles are the contours of the sunglint angle). (b) Scatter plot of cloud fraction in the sunglint region (The blue line represents the curve (namely Eq. 2) fitted by the least-squares method between the retrievals and the truths.). (c) Distribution of SunGlintAngle and satellite flight track of CloudSat and CALYPSO at 04:00 on 5 June 2019. (d) Distribution of cloud fraction retrieval error with sunglint angle.

10] and criterion – [Gini, entropy]. Ultimately, the best selections were as follows: daytime RF classification model – trees = 500, nighttime RF classification model – trees = 600, daytime RF regression model – trees = 400 and nighttime RF regression model – trees = 500. All four models have min-leaf = 1 and criterion = Gini.

For the MLP, experiments were conducted using the following hyperparameter ranges: number of hidden layers – [2, 3, 4, 5, 6, 7, 8, 9], hidden layer size – [8, 16, 32, 64, 128], epochs – [30, 50, 100], solver hyperparameter – [lbfgs, sgd, adam]. The optimal parameters found are as follows: (1) MLP classification model for daytime – number of hidden layers = 5, (2) MLP classification model for nighttime – number of hidden layers = 5, (3) MLP regression model for daytime – number of hidden layers = 4 and (4) MLP regression model for nighttime – number of hidden layers = 6. All four models have hidden layer size = 64, epochs = 50, solver = adam, BatchSize = 1500, initial learning rate = 0.01, learning rate schedule = piecewise, factor for dropping the learning rate = 0.1 and number of epochs for dropping the learning rate = 10.

4 Results and analysis

To assess the accuracy and stability of the retrieval model, two types of validation methods are utilized. One way involves a direct comparison from images, qualitatively comparing the model’s retrieval results and official cloud fraction

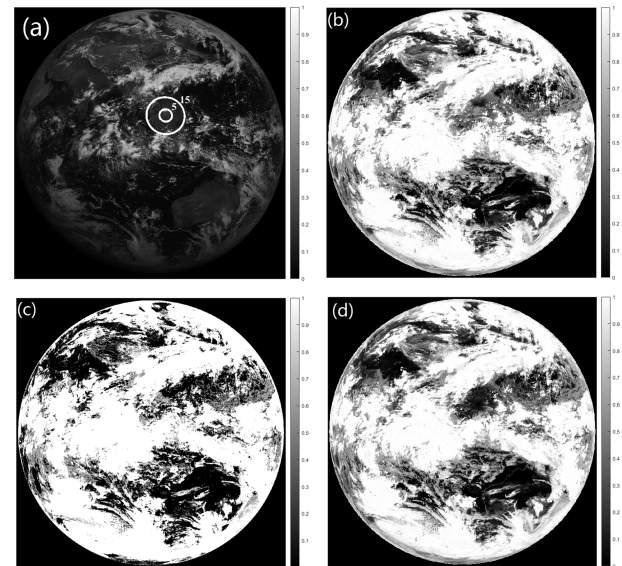


Figure 3. FY-4A AGRI at 04:00 on 1 June 2023. (a) Albedo image of the $0.67\ \mu\text{m}$ channel (The circles are the contours of the sunglint angle.). (b) Random forest cloud fraction retrieval without sunglint correction. (c) Operational cloud fraction product. (d) Random forest cloud fraction retrieval with sunglint correction.

products with AGRI-observed cloud images. Another approach uses 2B-CLDCLASS-lidar as the ground truth and introduces five parameters for quantitative comparison: recall, false alarm rate (FAR), mean error (ME), mean absolute error (MAE) and root mean square error (RMSE). To evaluate the ability of operational products, RF and MLP cloud detection models to distinguish overcast, partly cloudy and clear sky, the recall (probability of detection, POD) is calculated using the formula $\text{POD} = \text{TP} / (\text{TP} + \text{FN})$, and the false alarm rate is calculated using the formula $\text{FAR} = \text{FP} / (\text{TP} + \text{FP})$. Taking the overcast scene as an example, TP represents the number of correctly identified overcast conditions, FN represents the number of overcast conditions misidentified as partly cloudy or clear sky, and FP represents the number of clear sky or partly cloudy conditions misidentified as overcast. When assessing the accuracy of operational products and cloud fraction models for the cloud fraction retrieval results of partly cloudy scenes, mean error (ME), mean absolute error (MAE) and root mean square error (RMSE) are used.

4.1 Objective analysis of cloud fraction retrievals

First, using the 2B-CLDCLASS-lidar cloud fraction product as the ground truth, we calculated the accuracy of the operational cloud detection products. The results are shown in columns 3–4 of Table 2. The samples used for this statistic are the same as those for testing the model below (20 % of dataset A).

Table 4. POD and FAR of cloud detection in the sunglint area.

	Sky classification	Operational product	RF	RF after correction
POD	Clear sky	0.4120	0.0987	0.9023
	Partly cloudy	0.7371	0.9663	0.9587
	Overcast	0.8856	0.9845	0.9845
FAR	Clear sky	0.1229	0.1633	0.0938
	Partly cloudy	0.3332	0.7943	0.0276
	Overcast	0.2983	0.1321	0.1321

Based on the cloud detection model trained above, cloud detection experiments were conducted using the test samples from dataset A. The time- and space-matched 2B CLDCLASS-lidar cloud fraction product served as the ground truth to assess the accuracy of cloud detection. The results are shown in columns 5–8 of Table 2. During the day, the random forest model achieved an overall accuracy of 94.2 %, while the MLP model had an overall accuracy of 93.7 %. The random forest model exhibited slightly higher recall rates for clear skies, partly cloudy and overcast conditions compared to the MLP model, and its FAR was lower as well. Both models performed poorly in recognizing partly cloudy conditions, as the models tended to classify true cloud fractions of 0.16 as clear skies and those of 0.83 as overcast. At night, the random forest model achieved an overall accuracy of 89.4 %, while the MLP model had an accuracy of 88.7 %. The random forest model had higher recall rates for clear skies and partly cloudy conditions compared to the MLP, while the recall rates for overcast conditions were similar for both models. The FAR for the random forest model was lower than that of the MLP. Overall, both the random forest and MLP models showed higher classification accuracy for clear skies, partly cloudy and overcast conditions compared to operational products, with the random forest model performing better.

Based on the previous model’s assessment of the field of view as partly cloudy, the cloud fraction in this AGRI field of view is retrieved using the cloud fraction model established earlier. For model evaluation, both the operational product and the 2B-CLDCLASS-lidar cloud fraction product are classified as partly cloudy, with the matched 2B-CLDCLASS-lidar cloud fraction product considered the ground truth. The average error, mean absolute error and root mean square error for both daytime and nighttime operational products and cloud fraction model retrieval (Table 3) are calculated. It can be observed that the average errors of both models are close to 0 during both daytime and nighttime. The errors are smaller during the day than at night, with the RF model exhibiting lower errors than the MLP model. In summary, the errors of both models are smaller than those of the operational products, and the RF model performs better in the cloud fraction retrieval task.

Table 5. Cloud fraction errors in the sunglint area.

	Operational product	RF retrievals	RF after correction
ME	0.2354	0.1741	0.0670
MAE	0.2511	0.1820	0.0849
RMSE	0.2771	0.2166	0.1041

Based on the experiments mentioned above, the performance of RF in cloud detection and cloud fraction retrieval slightly outperforms that of MLP. Therefore, subsequent experiments will utilize the RF algorithm.

4.2 Cloud fraction correction in sunglint regions

Sunglint refers to the bright areas created by the reflection of sunlight to the sensors of observation systems (satellites or aircraft). This phenomenon usually occurs on extensive water surfaces, such as oceans, lakes or rivers. This specular reflection of sunlight will cause an increase in the reflected solar radiation received by onboard sensors, manifested as an enhancement of white brightness in visible images. The increase in visible channel observation albedo will affect various subsequent applications of data, including cloud detection and cloud cover retrieval.

The position of sunglint area can be determined using the SunGlintAngle value in the FY-4A GEO file. SunGlintAngle is defined as the angle between the satellite observation direction or reflected radiation direction and the mirror reflection direction on a calm surface (horizontal plane). It is generally accepted that the range of SunGlintAngle < 15° is easily affected by sunglint (Kay et al., 2009). The positions of the SunGlintAngle contour lines at 5 and 15° are marked in Fig. 1a. It can be observed that the edge of sunglint in Fig. 1a essentially overlaps with the position of SunGlintAngle = 15°. Thus, the region where SunGlintAngle < 15° is defined as the sunglint range in this paper, and only the cloud fraction within this range will be adjusted in the subsequent correction.

To correct the cloud fraction in the sunglint areas, we first identified the fields of view (FOVs) where sunglint occurred during FY-4A AGRI observations from August 2018 to July 2019, totaling 1476 FOVs. When matching the sample set of the sunglint area, two issues need to be explained. (1) Cloud fraction is the average of cloud fractions of different layers: among the matched pixels, only one-layer cloud and two-layer cloud appear. When there are two layers of cloud, there is always one layer with a cloud fraction of 1. According to the previous description, when there is one layer with a cloud fraction of 1, this pixel should be regarded as fully cloudy. (2) The average cloud fraction of at least two CloudSat and CALIPSO pixels is taken as the cloud fraction of the AGRI pixel: due to the very small area of the sunglint area, the matching is very difficult. If at least two

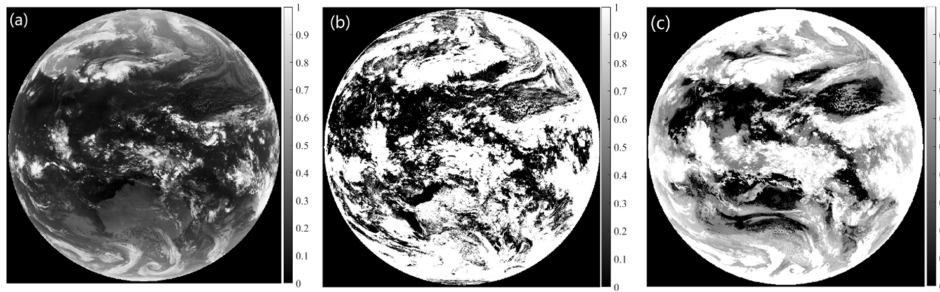


Figure 4. FY-4B AGRI at 17:00 on 18 April 2023. (a) Brightness temperature of the 10.8 μm channel, (b) operational cloud fraction product and (c) random forest cloud fraction retrieval.

CloudSat and CALIPSO pixels within an AGRI pixel are required, this will make the available sample size very small. Therefore, when making the sample set of the sunglint area, only one CloudSat and CALIPSO pixel within an AGRI pixel is required. For the above two reasons, the true cloud fraction in the sample is a discrete value. Subsequently, a direct least-squares fitting method was conducted between the retrieved cloud fraction and the collocated 2B-CLDCLASS-lidar cloud fraction (ground truth). The scatter plot is illustrated in Fig. 2b, where the x axis is the 2B-CLDCLASS-lidar cloud fraction and the y axis is the model-retrieved cloud fraction. The blue line represents the curve (namely Eq. 2) fitted by the least-squares method between the retrievals and the truths. The dashed thin line is the $x = y$ line. It is evident that the retrieved cloud fraction is generally slightly overestimated.

Taking observations at 04:00 UTC (all times in the paper correspond to UTC) on 5 June 2019 as an example, Fig. 2c presents the distribution of SunGlintAngle and the flight trajectory of the Cloudsat and CALYPSO satellite. White circles denote the sunglint region with SunGlintAngle $< 15^\circ$, and the white line represents the satellite flight track. As depicted in the figure, the majority of Cloudsat and CALYPSO flight trajectories do not pass through the central position of sunglint area but instead traverse locations with larger SunGlintAngle values. The intensity of sunglint effect decreases with the increase in SunGlintAngle. This suggests that the true values for spatial and temporal matching mostly do not fall within the strongest sunglint region. From Fig. 2d, it can be seen that the impact of sunglint becomes stronger as SunGlintAngle decreases, which results in a higher observation albedo. This further leads to the overestimated cloud fraction values in the retrieval. It is evident that the cloud fraction error is related to the value of SunGlintAngle, and this influence is not considered in Eq. (2). Directly applying Eq. (2) to correct the cloud fraction retrievals would result in too small a correction intensity for the FOVs near the center of sunglint and an excessively large correction intensity for the FOVs in the sunglint edge region (even erroneous clear sky may appear). Considering this, a correction formula (3)–(4) using SunGlintAngle as the weight is introduced, where

W_i represents the angle weight for a certain pixel i in the sunglint region, n is the number of pixels within the SunGlintAngle $< 15^\circ$ range, y_i is the initial model retrieval of cloud cover for the field of view i and x_i is the final corrected cloud fraction.

$$x = (y - 0.2441)/0.8092 \quad (2)$$

$$W_i = \frac{\text{glint angle}_i}{\frac{1}{n} \sum_{i=0}^n \text{glint angle}_i} \quad (3)$$

$$x_i = W_i \left(\frac{y_i - 0.2441}{0.8092} \right) \quad (4)$$

Figure 2d shows the distribution of errors with respect to SunGlintAngle, where the blue dots represent the error distribution corrected using formula (2), and the orange dots represent the error distribution corrected using formula (4). It can be seen from Fig. 2d that after correction by formula (4), the errors in the smaller range of SunGlintAngle are significantly reduced.

4.3 Algorithm universal applicability testing

Although the retrieval model in this article was built based on data from May 2019 due to the limited lifespan of the instrument, how effective is it in real-time FY-4A AGRI observations and even subsequent FY-4B AGRI applications? The algorithm's universal applicability was tested using real-time observations from FY-4A and FY-4B AGRI in 2023.

Taking the full-disk observation of FY-4A AGRI at 04:00 on 1 June 2023 as an example, the radiance observations from 14 channels are initially fed into the random forest cloud detection model to determine the sky classification (overcast, partly cloudy or clear sky) in each AGRI field. The random forest cloud fraction retrieval model is utilized to retrieve the cloud fraction in scenes identified as partly cloudy. Figure 3a is the observed albedo at 0.67 μm , where the circles represent the contours of the sunglint angle; (b) is the cloud fraction retrievals from random forest algorithm; (c) is

the official operational cloud fraction product; and (d) is random forest cloud fraction retrievals with sunglint correction. It can be seen from Fig. 3 that many clear-sky scenes are erroneously identified as cloudy by the operational product, and the cloud fraction is generally overestimated, with many scenes having a cloud fraction of 1. The random forest algorithm identifies more regions as clear skies or partly cloudy than the operational products, matching better with the observations in the $0.67\ \mu\text{m}$ albedo image. Brighter regions in the visible image correspond to cloud cover areas, and darker areas represent clear-sky conditions. The sunglint region in the central South China Sea (the circled area in Fig. 3a) is depicted in Fig. 3b, where the clear-sky scenes over the ocean are misidentified as partly cloudy by the random forest algorithm due to the increase in observed albedo. Although the operational product in this area also suffers from the impact of unremoved sunglint, it identifies more clear-sky scenes, and the cloud fraction is relatively low. Thus, it is evident that the random forest algorithm exhibits significant cloud detection and cloud fraction errors in these sunglint regions. Correction is necessary for the cloud fraction retrievals in the sunglint region.

Figure 3d shows the cloud fraction distribution after correction using Eq. (9) in the sunglint region. The correction eliminates the influence of sunglint comparing to the cloud fraction in sunglint area before correction in Fig. 3b. The scenes misjudged as partly cloudy are corrected to clear sky and match well with the actual albedo observations in Fig. 3a, which accurately restores the true cloud coverage over the South China Sea.

Statistical analysis was conducted on the correction effect using samples with sunglint in the training data. The POD and FAR in the sunglint area are listed in Table 5, and the error is in Table 6. It can be seen that after correcting for the cloud fraction, the POD for clear skies increased from 0.0987 to 0.9023. The FAR for partly cloudy decreased from 0.7943 to 0.0276. ME, MAE and RMSE show significant reductions, and the results after correction outperform operational products.

FY-4B launched in 2021 has a total of 15 channels with an additional low-level water vapor channel at $7.42\ \mu\text{m}$ compared to FY-4A. Taking the full-disk observation of FY-4B AGRI at 17:00 on 18 April 2023, as an example, the radiance observation data of the remaining eight channels (near-infrared and infrared channels) except for the $7.42\ \mu\text{m}$ channel and the visible light channels were input into the random forest cloud detection model. Figure 4a shows the brightness temperature distribution observed in the $10.8\ \mu\text{m}$ channel of FY-4B AGRI, (b) represents the operational cloud fraction product for FY-4B AGRI and (c) shows the cloud fraction retrieved by this algorithm. Figure 4 illustrates that the random forest algorithm identifies more regions as clear skies or partly cloudy than the operational products, aligning better with the brightness temperature observations in the $10.8\ \mu\text{m}$ channel. Especially in high-latitude regions of the Southern

Hemisphere and areas with strong convection near the Equator, the cloud cover provided by operational products is too high and even misjudged. It can be seen that the random forest algorithm is also suitable for cloud fraction retrieval of FY-4B AGRI.

5 Conclusions

This paper used random forest and multilayer perceptron (MLP) algorithms to retrieve cloud fraction from FY-4A AGRI full-disk level-1 radiance observation data and verified the accuracy of the algorithms using the Cloudsat and CALYPSO active remote sensing satellite's 2B CLDCLASS-lidar cloud fraction product. The following conclusions were drawn:

1. The random forest and MLP algorithms performed well in cloud detection and cloud fraction retrieval tasks, and their accuracy was higher than that of operational products. The accuracy of cloud detection can reach over 93 %, and the error of cloud fraction retrieval is close to zero. Compared with the MLP algorithm, the RF algorithm has a slightly higher accuracy in cloud detection and a slightly lower error in cloud fraction retrieval, showing better performance.
2. At night, the classification accuracy is lower than during the day due to the lack of observations in the visible channel of AGRI, resulting in higher cloud fraction errors at night.
3. The accuracy of identifying partly cloudy scenes is lower than that of identifying clear sky and overcast scenes for both RF and MLP algorithms. Scenes with a very low cloud fraction (0.16) are often misclassified as clear sky, while scenes with a high cloud fraction (0.83) are often misclassified as overcast.
4. The sunglint area cloud fraction correction curve, fitted with SunGlintAngle as the weight, greatly improves the accuracy of cloud fraction retrieval and reduces the misclassification rate of clear-sky scenes as partly cloudy or partly cloudy scenes as overcast due to increased reflectance.

Data availability. FY-4A AGRI data are available at <http://satellite.nsmc.org.cn> (National Satellite Meteorological Center, 2024) and the 2B-CLDCLASS-lidar data at <https://www.icare.univ-lille.fr/data-access/data-archive-access/> (Université de Lille, 2024).

Supplement. The supplement related to this article is available online at: <https://doi.org/10.5194/amt-17-6697-2024-supplement>.

Author contributions. JX: formal analysis, methodology, software, visualization and writing (original draft preparation). LG: conceptualization, data curation, funding acquisition, supervision, validation and writing (review and editing).

Competing interests. The contact author has declared that neither of the authors has any competing interests.

Disclaimer. Publisher's note: Copernicus Publications remains neutral with regard to jurisdictional claims made in the text, published maps, institutional affiliations, or any other geographical representation in this paper. While Copernicus Publications makes every effort to include appropriate place names, the final responsibility lies with the authors.

Acknowledgements. We acknowledge the High Performance Computing Center of the Nanjing University of Information Science and Technology for their support of this work.

Financial support. This research has been supported by the National Natural Science Foundation of China (grant no. 41975028).

Review statement. This paper was edited by Jian Xu and reviewed by three anonymous referees.

References

- Amato, U., Antoniadis, A., Cuomo, V., Cuttillo, L., Franzese, M., Murino, L. and Serio, C.: Statistical cloud detection from SEVIRI multispectral images, *Remote Sens. Environ.*, 112, 750–766, <https://doi.org/10.1016/j.rse.2007.06.004>, 2008.
- Baum, B. and Trepte Q.: A Grouped Threshold Approach for Scene Identification in AVHRR Imagery, *J. Atmos. Ocean. Technol.*, 16, 793–800, [https://doi.org/10.1175/1520-0426\(1999\)016<0793:AGTAFS>2.0.CO;2](https://doi.org/10.1175/1520-0426(1999)016<0793:AGTAFS>2.0.CO;2), 1999.
- Breiman L.: Random Forests-Random Features [J], *Machine Learn.*, 45, 5–32, 1999.
- Breiman, L.: Random Forests, *Machine Learn.*, 45, 5–32, <https://doi.org/10.1023/A:1010933404324>, 2001.
- Chai, D., Huang, J., Wu, M., Yang, X., and Wang, R.: Remote sensing image cloud detection using a shallow convolutional neural network[J], *ISPRS J. Photogramm.*, 2024, 20966–20984, <https://doi.org/10.1016/j.isprsjprs.2024.01.026>, 2024.
- Duda, R. O., Hart, P. E., and Stork, D. G.: *Pattern Classification*, New York: John Wiley & Sons, 2001, xx + 654 pp., ISBN: 0-471-05669-3, *J. Classif.*, 24, 305–307, <https://doi.org/10.1007/s00357-007-0015-9>, 2007.
- Gao, J. and Jing, Y.: Satellite Remote Sensing Cloud Detection Method Based on Fully Convolutional Neural Network, *Infrared Technology*, 41, 607–615, 2019.
- Gomez-Chova, L., Camps-Valls, G., Amoros-Lopez, J., Guanter, L., Alonso, L., Calpe, J., and Moreno, J.: New Cloud Detection Algorithm for Multispectral and Hyperspectral Images: Application to ENVISAT/MERIS and PROBA/CHRIS Sensors, *IEEE International Symposium on Geoscience and Remote Sensing*, 2757–2760, <https://doi.org/10.1109/igarss.2006.709>, 2006.
- Kay, S., Hedley, J., and Lavender, S.: Sun Glint Correction of High and Low Spatial Resolution Images of Aquatic Scenes: a Review of Methods for Visible and Near-Infrared Wavelengths, *Remote Sens.*, 1, 697–730, <https://doi.org/10.3390/rs1040697>, 2009.
- Kegelmeyer, W. P. J.: *Extraction of cloud statistics from whole sky imaging cameras*, March 1994, Livermore, California, University of North Texas Libraries, UNT Digital Library, <https://doi.org/10.2172/10141846>, 1994.
- Mace, G. G. and Zhang, Q.: The CloudSat radar-lidar geometrical profile product (RL-GeoProf): Updates, improvements, and selected results, *J. Geophys. Res.*, 119, 9441–9462, <https://doi.org/10.1002/2013JD021374>, 2014.
- Merchant, C. J., Harris, A. R., Maturi, E., and Maccallum, S.: Probabilistic physically based cloud screening of satellite infrared imagery for operational sea surface temperature retrieval, *Q. J. Roy. Meteorol. Soc.*, 131, 2735–2755, <https://doi.org/10.1256/qj.05.15>, 2005.
- National Satellite Meteorological Center: Fengyun Satellite Remote Sensing Data Service Network, <http://satellite.nsmc.org.cn> (last access: 20 November 2024), 2024.
- Rossov, W. B. and Leonid, C. G.: Cloud detection using satellite measurements of infrared and visible radiances for IS-CCP, *J. Climate*, 12, 2341–2369, [https://doi.org/10.1175/1520-0442\(1993\)006<2341:CDUSMO>2.0.CO;2](https://doi.org/10.1175/1520-0442(1993)006<2341:CDUSMO>2.0.CO;2), 1993.
- Solvsteen, C.: Correlation based cloud-detection and an examination of the split-window method, *Proc. SPIE – The International Society for Optical Engineering*, 86–97, <https://doi.org/10.1117/12.228636>, 1995.
- Sassen, K., Wang, Z., and Liu, D.: Global distribution of cirrus clouds from CloudSat/Cloud-Aerosol Lidar and Infrared Pathfinder Satellite Observations (CALIPSO) measurements, *J. Geophys. Res.*, 113, D00A12, <https://doi.org/10.1029/2008JD009972>, 2008.
- Université de Lille: Cité Scientifique, AERIS/ICARE Data and Services Center – UAR 2877, <https://www.icare.univ-lille.fr/data-access/data-archive-access/> (last access: 20 November 2024), 2024.
- Xiang, S. P.: A Cloud Detection Algorithm for MODIS Images Combining Kmeans Clustering and Otsu Method, *IOP Conference Series: Materials Science and Engineering*, 392, 062199, <https://doi.org/10.1088/1757-899X/392/6/062199>, 2018.
- Yan, J., Guo, X., Qu, J., and Han, M.: An FY-4A/AGRI cloud detection model based on the naive Bayes algorithm, *Remote Sens.-Nat. Resour.*, 34, 33–42, <https://doi.org/10.6046/zrzyg.2021259>, 2022.
- Zhang, W., He, M., and Mak, M. W.: Cloud detection using probabilistic neural networks, *Geoscience and Remote Sensing Symposium, IEEE* 2373–2375, <https://doi.org/10.1109/IGARSS.2001.978006>, 2001.
- Zhang, Y., William, B. R., Andrew, A. L., Valdar, O. and Michael, I. M.: Calculation of radiative fluxes from the surface to the top of atmosphere based on ISCCP and other global data sets: Refinements of the radiative transfer model and the input data, *J. Geophys. Res.-Atmos.*, 109, 1–27, <https://doi.org/10.1029/2003JD004457>, 2004.