

# MACHINE LEARNING TECHNIQUES FOR KNOWLEDGE EXTRACTION FROM SATELLITE IMAGES: APPLICATION TO SPECIFIC AREA TYPES

C.O. Dumitru<sup>1\*</sup>, G. Schwarz<sup>1</sup>, M. Datcu<sup>1</sup>

<sup>1</sup>German Aerospace Center (DLR), Remote Sensing Technology Institute, Münchener Str. 20, 82234 Weßling, Germany  
(corneliu.dumitru, gottfried.schwarz, mihai.datcu)@dlr.de

**KEY WORDS:** Image classification, Image compression, Active learning, Latent Dirichlet Allocation (LDA), Convolutional Neural Networks (CNNs), Sentinel-1, Sentinel-2.

## ABSTRACT:

When we want to extract knowledge from satellite images, several well-known image classification and analysis techniques can be concatenated or combined to gain a more detailed target understanding. In our case, we concentrated on specific extended target areas such as polar ice-covered surfaces, forests shrouded by fire plumes, flooded areas, and shorelines. These image types can be described by characteristic features and statistical relationships. Here, we demonstrate that both multispectral (optical) as well as SAR (Synthetic Aperture Radar) images can be used for knowledge extraction. The free availability of image data provided by the European Sentinel-1 and Sentinel-2 satellites allowed us to conduct a series of experiments that verified our classification approaches. This could already be verified in our recent work by quantitative quality tests.

## 1. INTRODUCTION

During the past years, one could witness a wide range of innovative machine learning applications in the field of remote sensing. Typical applications often being dealt with are routine classifications of satellite images ranging exemplarily from cloud patterns in the atmosphere, to wave patterns and icebergs on the oceans, and to the analysis of time series of land cover and land use images taken in different spectral bands.

Currently, one can find quite a number of applications, where “deep” (i.e., multi-layer) machine learning algorithms based on well-known convolutional neural network (CNN) or auto-encoder (AE) techniques, yield rapid and sufficiently accurate image classification results – despite the fact that these algorithms are not yet adapting themselves in an optimal way to the varying local characteristics of typical satellite images nor their detailed semantic meaning and its representation. This led to the situation that a number of deep learning software packages can be downloaded as public domain software allowing every user to assign about ten to twenty land cover classes to typical satellite images. In most cases, however, these software packages have to be trained by hopefully representative examples. Therefore, their actual performance depends on the volume and selection of test cases, and any variation of these cases already provides a good measure of their classification robustness.

At the same time, a more detailed analysis of deep learning results shows that by small adaptations of the processing parameters and options users can still improve and optimize the classification performance of the general approaches. Hence, we were very much interested in a straightforward and robust approach to more detailed parameter settings for small-scale (i.e., high-resolution) semantic image content classifications.

This prompted us to investigate the potential of efficiently combining knowledge from different sources, such as from low-resolution and reliably classified bigger images, derived interactively by qualified image analysts with information derived in parallel from high-resolution smaller image patches processed routinely with fully-automated modern deep learning algorithms. However, this knowledge combination step had to be preceded by an investigation demonstrating on which level any type of knowledge being extracted from our satellite images by different techniques can be efficiently combined, for instance, knowledge derived from selected image patches, extracted feature vectors, semantic (multi-)labels, etc. Then, a combination of knowledge generated by different analysis techniques can be a way to improve local and overall classification accuracies.

We recognized that our aim should be to finally assign reliable and very detailed labels to very small locally confined image patches, thus avoiding the need for conventional feature vector extraction. As a consequence, we concentrated on purely semantic classifications and knowledge extraction techniques that – after being combined – are very useful for a wide range of remote sensing applications such as shipping route safety in arctic waters, coastal deltas, forest fires, and flood monitoring (see the selected locations in Figure 1).

This paper comprises six sections, one appendix, and one reference section. Section 2 contains a description of the selected remote sensing applications, followed in Section 3 by the descriptions of our test areas and the used datasets. Section 4 briefly summarizes the probed algorithms (with references to details of the applied methods), while Section 5 details our findings sorted by application, data, and the proposed methods. Conclusions and future research directions are described in Section 6 that completes this paper.

---

\* Corresponding author



**Figure 1.** Location of our target areas marked on Google Maps: a) Belgica Bank in Greenland, b) Danube Delta in Romania, c) Sydney in Australia, d) Amazon rainforest between Brazil, Bolivia, and Paraguay, and e) Montevideo in Uruguay.

## 2. SEA-ICE AND OTHER APPLICATIONS

In this section we give a short introduction to our typical satellite image applications that call for up-to-date data together with some dedicated and automated software to classify and understand these images. Therefore, we concentrate on four main application areas, namely ice-covered polar areas, coastlines and river deltas, burnt areas and fires, and flooding events.

These four application areas are typical for monitoring the Earth's surface by satellite images, however, users do not only need a single image, but time series of repeated and time-keeping image acquisitions, coverage of extended areas by big images, automated measurement and detection of expected and unexpected events, and their interpretation together with quality assessments.

A successful monitoring and image understanding is needed by many authorities such as civil protection agencies as well scientific institutions that care for the state of our planet. The data that we need can be provided by several airborne or spaceborne instruments, however, the original instrument data have to undergo accurate calibration and reprojection steps before becoming useful for general use. In addition, one has to be aware of the fact that the different types of instruments yield different types of results that may have to be combined and analysed prior to a final understanding.

The following four-part table (Table 1) summarizes the most important requirements and parameters of satellite imaging and image understanding.

Application field	Monitoring and analysis of polar areas	Coastlines and river deltas	Burnt areas and fires	Flooding events
User community	shipping agencies, charting services, climate change	Reservation Biosphere Danube Delta, UNESCO World Heritage	emergency operations centres, local authorities	emergency operations centres, local authorities
Monitoring instruments	SAR instruments	SAR and multispectral/optical instruments	Multispectral/optical instruments	Multispectral/optical instruments
Observation conditions	day and night (for SAR)	day and night (for SAR) day with low cloud cover (for optical imaging)	dayside imaging with low cloud cover (for optical imaging)	dayside imaging with low cloud cover (for optical imaging)
Observed targets	ice coverage, calving and floating icebergs, ice floes, snow-covered surfaces	sediments transported by water, surface water dynamics, snow coverage, occurrences of fire	legal or illegal deforestation, natural disasters after the fires	destroyed buildings and agricultural areas landslides
Pixel spacing	20 m for Sentinel-1	20 m for Sentinel-1 10-20-60 m for Sentinel-2	10-20-60 m for Sentinel-2	10-20-60 m for Sentinel-2
Spectral bands	dual-band for Sentinel-1	dual-band for Sentinel-1 13 bands for Sentinel-2	13 bands for Sentinel-2	13 bands for Sentinel-2
Repeat cycle	6 days for Sentinel-1A/B	6 days for Sentinel-1A/B 5 days for Sentinel-2A/B	5 days for Sentinel-2A/B	5 days for Sentinel-2A/B
Special algorithms	surface classification by active learning or LDA technique and physical scattering representation	surface classification by active learning (with expert users) and by LDA technique (fully automated)	surface classification by LDA technique	surface classification by LDA technique
Obtained results	semantic classification maps benchmarking datasets	semantic classification maps, benchmarking datasets	semantic classification maps, maps of burnt areas	semantic classification maps, maps of flooded areas

**Table 1:** Important requirements and parameters of our application using Sentinel-1/Sentinel-2 images.

## 3. SATELLITE DATA AND TARGET AREAS

### 3.1 Sea-Ice

For a first sea-ice application, we selected a target area around Belgica Bank in the north-east of Greenland which is an area of extensive fast land-locked ice (ExtremeEarth, 2021).

The Sentinel-1 C-band SAR (Synthetic Aperture Radar) spacecraft is one of the satellites that constantly monitor this area. From its available product types, we selected level-1 Ground Range Detected amplitude data in Interferometric Wide swath mode with dual polarization (HH and HV) and a pixel spacing of 10×10 meters. Our analysis period covers the time between Jan. 2018 and Dec. 2019, and from that we selected 24 images (*i.e.*, one image/month). See an example in Figure 2.

### 3.2 River Deltas and Coastlines

The Danube Delta is the second largest river delta in Europe and one of the best-preserved deltas on the continent (Dumitru, et al., 2019).

Here, both the Sentinel-1 SAR and Sentinel-2 multispectral instruments are covering our area of interest. For Sentinel-1 data products, we selected the same product types as in the previous Section 3.1 with a single modification: now the polarizations for this area are VV and VH. The acquisition period overlaps with the one for Sentinel-2 (see below). From the available Sentinel-2 data products, we selected level-1C data with radiometrical and geometrical corrections.

The pixel spacing of the images varies depending on the given band and lies between 10 and 60 m. Our analysis period covers the time between November 20<sup>th</sup> 2015 and May 18<sup>th</sup> 2016, and after removing some cloud-covered images we selected 16 clear-view images. A typical example is shown in Figure 4.

### 3.3 Fires

One of the most devastating fires in Australia occurred in the year 2019, when several million hectares of land were burnt and almost half a billion animals perished (Australia, 2019). In August 2019, many fires affected the Amazon rainforest. Here a second study application is one of the Amazon forest located between Brazil, Bolivia, and Paraguay (Amazon, 2019).

Again, the Sentinel-2 multispectral instrument is covering our area of interest, and we used the same configuration of product types as in the previous Section 3.2. For both areas, we selected a set of three images acquired before, during, and after the fires. As a typical example, Figure 5 shows the results of the image acquired during the Australian fires on December 31<sup>st</sup>, 2019, while Figure 6 shows an example of the Amazonian fires on August 25<sup>th</sup>, 2019.

### 3.4 Floods

A prominent flooding example is Montevideo in Uruguay, where the confluence of the two rivers Paraná and Uruguay was affected on October 18<sup>th</sup>, 2019 (Uruguay, 2019).

Here again, we consider the Sentinel-2 satellite with the same product types as in the Section 3.2. Our selected images were acquired before, during, and after October 18<sup>th</sup>, 2019. In Figure 7, we show the results of the image acquired on October 18<sup>th</sup>, 2019.

Due to space limitations for this paper, we show only one image per application from our full set of available images, while the numerical results are presented for all proposed methods. The results for the remaining images are similar.

## 4. ALGORITHMS

In this section we present the algorithms being used for the applications described in Section 3 with their related data.

In order to create semantic classification maps, topic representations or physical scattering representations, a sequence of (manual or automated) steps has to be executed. These are specific for each algorithm.

For all algorithms a *pre-processing step* is required, in which a given Sentinel image is split into patches of instrument-dependent sizes. For Sentinel-1, the patch size is 256×256 pixels, while for Sentinel-2 the patch size is 120×120 pixels.

### 4.1 Semantic Classification based on Gabor Filtering and Support Vector Machines (SVMs)

This algorithm is an active learning algorithm described in detail in (EOLib, 2018). Its specific steps are the following ones:

*Feature extraction step:* From each patch, a 60-element feature vector is automatically created that is based on Gabor filtering (we compute the means and variances of 5 scales and 6 orientations).

*Clustering step:* The extracted features are clusters using a Support Vector Machine (with a Chi-squared kernel) with relevance feedback. This is a manual step, where a GUI is operated by an expert user in order to create classes that make sense. Here, the number of classes is not fixed, and they are defined by the user and the content of the data.

*Semantic annotation step:* The retrieved classes are semantically annotated by the user (by choosing an existing semantic label or by defining a new one) based on his/her experience and using the existing ground-truth data.

*Map generation step:* After all patches have been labelled, we automatically generate a semantic classification map (see an example in Figure 2-b).

*This method can generate its results very fast, through only a few positive and negative examples given by the expert user, and the class he/she is searching can be found quickly.*

### 4.2 Semantic Classification based on Compression Rates

Here, the first part of the algorithms is the same. The part that differs is the classification step where two algorithms are compared (Dumitru et al., 2021). The steps of both methods are automated, except for the label assignment. These clustering algorithms are chosen as state-of-the-art methods, but other algorithms can be used, too.

#### 4.2.1 Compression based on Dictionaries and *k*-means Clustering

*Dictionary computation step:* To each patch tiled from the image, a Lempel-Ziv-Welch (LZW) compression algorithm is applied. LZW is a lossless data compression algorithm (Smith, 2004). Then a dictionary with compression rate results is created (Dumitru et al., 2021).

*Clustering step:* The dictionaries are string vectors containing the information from each patch in a compact form. For applying a clustering algorithm, all strings must have the same size; as a consequence, the strings are reduced to a short uniform length (Vaduva et al., 2015). In our case, an unsupervised *k*-means technique is used, where *k* is equal to the number of classes retrieved by the user in the first method.

*Semantic annotation step:* The output of the classification are classes (clusters) that are labelled manually by an expert user.

#### 4.2.2 Compression based on Dictionaries and Gaussian Mixture Models (GMMs)

*Dictionary computation step:* This step is identical with the one from Section 4.2.1.

*Clustering step:* The procedure is the same as the one from Section 4.2.1, but the clustering algorithm is Unsupervised Gaussian Mixture Models (GMMs).

*Semantic annotation step:* This step is the same as the one from Section 4.2, where the classes are labelled by an expert user.

*The LZW method compresses the data into dictionaries and reduces the storage space. Thus, it can be an alternative to other established feature extraction methods. Of course, the*

*selected clustering algorithm also has some impact on the classification results.*

#### 4.3 Topic Representations based on Latent Dirichlet Allocation (LDA)

After the *pre-processing step*, where the macro-patches of  $256 \times 256$  pixels are created, the next steps to be applied are the following ones (Karmakar et al., 2021):

*Macro-patch tiling step:* Once the macro-patches are generated, each patch is tiled again into smaller patches of  $4 \times 4$  pixels. A macro-patch of  $256 \times 256$  pixels will create 4,096 micro-patches of  $4 \times 4$  pixels.

*Clustering step:* A  $k$ -means algorithm is applied to the local descriptors; these descriptors are the linearized pixel brightness values of the micro-patches. The number of clusters is experimentally set to 50. Each cluster is then considered as a visual word (Karmakar et al., 2021).

*Bag-of-words modelling step:* The pixel values of each image patch are then assigned to the words of the dictionary. The image patches are then modelled as a bag of words based on the occurrence of the words.

*LDA step:* We apply LDA (Blei et al., 2003) in order to discover the latent semantics of the images as a set of topics. These topics are distributions over the words of the dictionary. The images can then be represented as distributions over the topics. The number of topics is set to match the number of semantic classes discovered by the user using the active learning algorithm.

*LDA can be applied for learning the high-level semantic structures in areas with no or poor existing prior knowledge (disaster areas, polar areas, etc.).*

#### 4.4 Physical Scattering Representations based on LDA and Convolutional Neural Networks (CNNs)

This hybrid method is specific for SAR data, for which two polarizations are required. The method is trying to make full use of the physical scattering mechanisms and spatial information of dual-pol SAR images (in our case, the HH and HV polarization data provided by Sentinel-1 tiled into patches of  $256 \times 256$  pixels). A summary of the steps is presented below but for more details, see (Huang et al., 2021).

*Scattering step:* Here we apply Cloude and Pottier's polarimetric decomposition based on scattering entropy (Cloude and Pottier, 1997) to our dual-polarized SAR data (with HH and HV polarizations). This is followed by a Wishart classifier (Lee et al., 1999) that generates nine types of scattering mechanisms.

*Topic modelling step:* A mixture of topics is generated for each patch based on the scattering labels generated by a Latent Dirichlet Allocation (LDA) model. Then, for each patch a bag of scattering topics is created.

*Unsupervised learning step:* The topic description can also be represented by image features extracted from a pre-trained convolutional neural network (CNN) model for SAR images; here, a soft constraint function is used for the learning step. For our applications, we use a pre-trained deep model of a

TerraSAR-X dataset transferred via a ResNet-18 backbone (Huang et al., 2020) to simulate high-quality Sentinel-1 data.

*Supervised label prediction step:* This final step uses the annotated (i.e., labelled) data to train the classification layer with given constraints. The subsequent classification is performed with a limited number of labelled data, and the entire network is fine-tuned with a constraint loss function.

*Integrating artificial intelligence (AI) learning with physics models produces results with higher generalization power and robustness and is increasing prediction performance. The physics-driven AI for SAR is regularising existing AI models by physical rules for the SAR image formation and target scattering, thus implementing hybrid paradigms where machine learning models substitute the unknowns or computationally expensive physics-based models.*

## 5. EXPERIMENTAL RESULTS

In this section, we start with the first application (i.e., sea-ice), for which all the methods from Section 4 are verified. Based on the obtained results, we identified the two best methods that are subsequently applied to the next applications, namely river deltas and coastlines.

For the last two applications, namely fires and floods, we selected a method for which we wanted to see the influence of the selected multispectral bands and the retrieval behaviour of different classes (e.g., the separation between smoke and clouds, the different types of water – such as muddy or ocean water).

### 5.1 Sea-Ice

The first method (active learning) was applied by an expert user, and the resulting classes were obtained after the necessary user interaction with the system (knowledge transfer from the user to the system). This method is based on Gabor filtering and was applied to a sequence of patches tiled from the full image, followed by an SVM-based feature classification step. The number of classes are defined again by an expert user. From the selected image acquired on April 17<sup>th</sup>, 2021 we obtained 8 semantic classes (excluding the black border class that appear in each Sentinel-1 image). The results of the semantic classification are shown in Figure 2-a.

The next method tests were conducted without user interaction and by using the number of classes of the first method.

For the second (compression with  $k$ -means) and third method (compression with GMM), an LZW compression was applied to each patch in order to generate a common dictionary.

For the second method, we first performed an unsupervised  $k$ -means clustering step keeping the number of classes to the value found by the expert user ( $k=8$ ) for the first method. Applying the algorithm, we noticed that the patches were grouped into 7 semantic classes (see Figure 2-b). Analysing the result from Figure 2-c by comparing it with the given ground-truth, we observed that one class is a mixture of two classes and the class *Mountains* is missing in contrast with the Figure 2-b results.

For the third method, the procedure is the same but this time we used some Gaussian Mixture Models (GMMs). The number of retrieved semantic classes is 6, two less than for the SVM method and one less than for the  $k$ -means method.

The fourth method (LDA) is a topic representation based on Latent Dirichlet Allocation (LDA). Here, we noticed a better grouping, and the generated map comes closer to the quick-look image (see Figure 2-a). These topics can be used to represent the corresponding semantic classes (e.g., from Figure 2-b) through a combination of topics (Karmakar et al., 2021). For example, the semantic class *Water bodies* from Figure 2-b has as dominant topics the categories *topic 6* (38%) and *topic 12* (36%), while the remainder is split into almost equal proportions of *topic 4* and *topic 11*. Another example is the class *Icebergs* which is mostly composed of *topic 12* (33%), *topic 1* (22%), *topic 5* (16%), and *topic 11* (13%).

An appropriate selection of the number of LDA topics is very important. We observed that there are some topics that do not fall into the scope of any of the semantic classes. Therefore, a more detailed study still needs to be made in future to find the optimum number of topics.

The last method (LDA with CCN) is a hybrid approach which is specific to dual-pol SAR data (in our case, to Sentinel-1). The method extracts the physical scattering phenomena into 9 classes. Unfortunately, these physical classes cannot be compared with the previous ones because there are no one-to-one correspondences.

As a conclusion of this first Section 5.1, we can say that the methods that are providing results close to our ground-truth data are the one based on SVM, and the one based on LDA.

## 5.2 River Deltas and Coastlines

Following the results from the first Section 5.1, we applied these two methods to our data of the Danube Delta acquired by

Sentinel-1 and Sentinel-2. Due to the different revisit periods of each satellite and the partly visible cloud coverage of the area by Sentinel-2, the closest acquisitions of both satellites were May 18<sup>th</sup>, 2016 for Sentinel-1, and April 28<sup>th</sup>, 2016 for Sentinel-2 (Dumitru, et al., 2019).

Based on the results from the previous Section 5.1, we applied for this application the two selected methods. The results for Sentinel-1 are depicted in Figure 3, while the results for Sentinel-2 are shown in Figure 4. From the two images, we noticed that the area covered by Sentinel-1 is larger than the one covered by Sentinel-2.

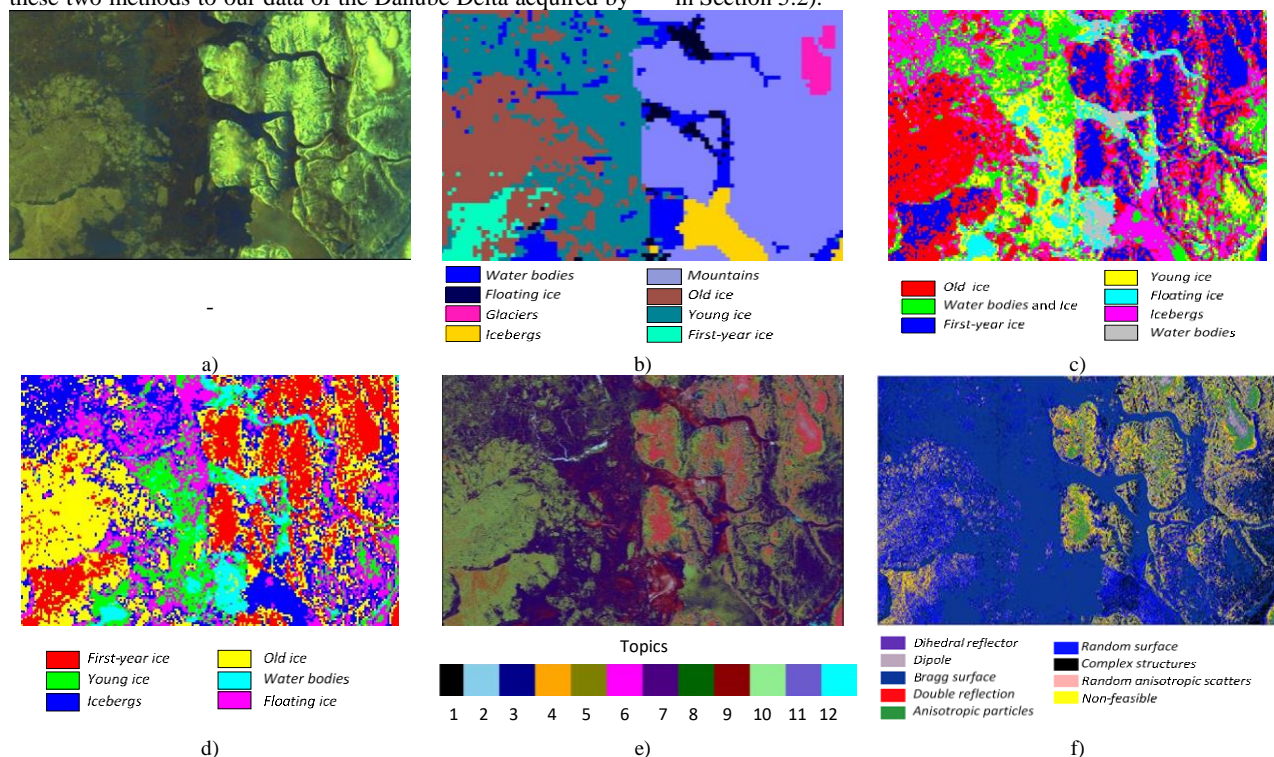
In the case of the first method, due to the differences in resolution between Sentinel-1 and Sentinel-2, the number of categories and the retrieved semantic classes may differ, and some categories can be mixed together or be missing. Here, there are two semantic class differences between the sensors.

In the case of the second method, the topic representation gives more details compared to the first method (e.g., the currents or the waves in the Black Sea). Here, the LDA method is applied only to the RGB 10 m bands (as a rule, B4, B3, B2) provided by the Sentinel-2 satellite.

As a conclusion for this Section 5.2, we propose to use Sentinel-2 data for more structural details, and LDA as an analysis method.

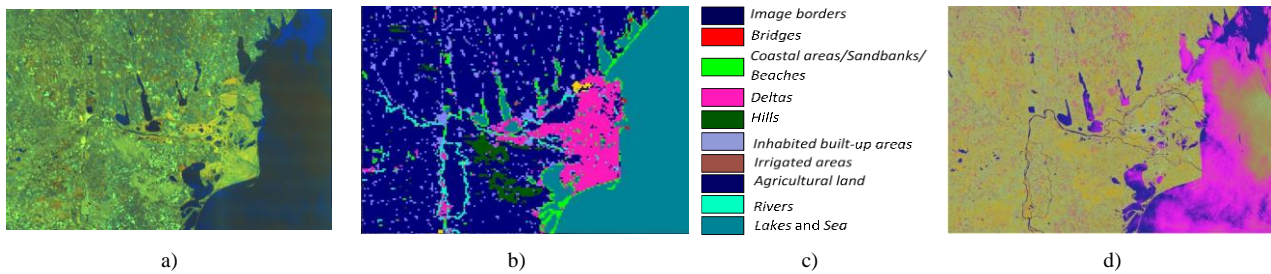
## 5.3 Fires

In this section, we first start by presenting the results obtained based on SVM and those with LDA using Sentinel-2 data (like in Section 5.2).

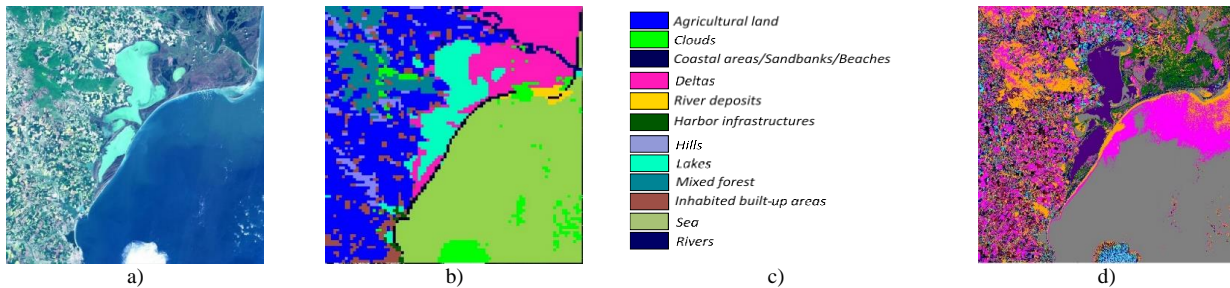


**Figure 2.** An image example of Belgica Bank, Greenland acquired on April 17<sup>th</sup>, 2018. (From left to right, upper part): The Sentinel-1 quick-look image, our semantic classification based on active learning (with Gabor filtering as feature vector and an SVM as classifier), our semantic classification based on compression (with dictionaries as feature vector and *k*-means as classifier); (From left to right, lower part): Our semantic classification based on compression (with dictionaries as feature vector and a GMM as classifier), the topic representation based on LDA, and the physical scattering representation based on a hybrid approach (LDA and CNNs). Each case includes a colour legend, except figure a) which is the quick-look of the analysed image.





**Figure 3.** An image example of the Danube Delta, Romania acquired on May 18<sup>th</sup>, 2016. a) The Sentinel-1 quick-look image, b) Our semantic classification based on active learning (with Gabor filtering as feature vector and an SVM as classifier), and c) Our topic representation based on LDA. The active learning process is supported by a colour legend shown in c); the colour legend of LDA is similar to the one from Figure 2.



**Figure 4.** An image example of the Danube Delta, Romania acquired on April 28<sup>th</sup>, 2016. a) The Sentinel-2 RGB (B4, B3, and B2) quick-look image, b) Our semantic classification based on active learning (with Gabor filtering as feature vector applied to each band and an SVM as classifier), and d) Our topic representation based on LDA. The active learning process is supported by a colour legend shown in c); the colour legend of LDA is similar to the one from Figure 2.

These results are followed by an impact analysis of Sentinel-2 spectral band selections, and how this combination can help find or easily separate classes (e.g., to separate clouds from smoke or muddy water from clean water).

Figure 5 presents the comparative results between the first method based on SVM and the fourth method based on LDA. These results are obtained using only the RGB bands (B4, B3, and B2). The conclusion we can draw from this figure is the same one as the previous case illustrated in Figure 4.

The Sentinel-2 instrument has 13 spectral bands, with a swath width of 290 km and a spatial resolution of 10 m (for the four visible and near-infrared bands), 20 m (for the six red-edge and shortwave infrared bands), and 60 m (for the three atmospheric correction bands). Table 2 summarizes the details for each band.

We prepared four combinations of bands and analysed the impact of these combinations. This study was made to provide a better separation of the different categories (e.g., smoke from clouds), and possibly to increase the number of classes. The selected combination of bands is: visible false-colour bands (see the green colour highlighting in Table 2), false-colour visible/infrared bands (see the orange colour in Table 2), false-colour infrared bands (see the blue colour in Table 2), and all 13 bands (see the pink colour in Table 2).

As a first example, we show the impact of different band combinations of Sentinel-2 channels. Figure 6 illustrates the band-dependent appearance of *Clouds*, *Smoke*, and *Fires* in the area of Sydney, Australia affected by fires in the end of 2019.

By carefully analysing each image from Figure 6, we can say that:

- The class *Clouds* and *Smoke* can be identified and classified (based on topics) from the two types of band combinations. A better separation between them is

obtained by the second combination. For these images, the *Forest* area that has been decimated by fire can be better separated and determined as a different topic in the second combination.

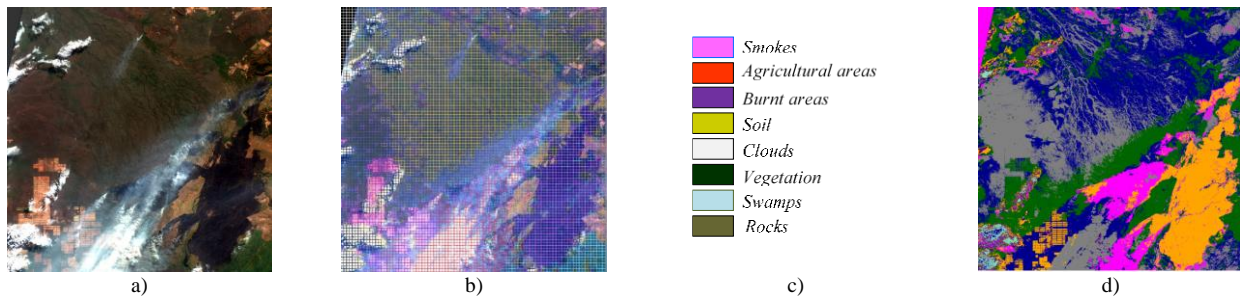
- When we analyse the third combination of bands, the two classes of *Smoke* and *Clouds* no longer appear, but instead we can clearly see the class *Fires* and that the area that was separated by the second combination is still aflame.
- From the last combination of bands, which comprises all bands, we can extract the classes *Smoke* and *Clouds*, however, the devastated forest area is mixed with the unaffected forest area.

		Band number	Central Wavelength	Combination	
10 m	2 - Blue		492.4 nm		
	3 - Green		559.8 nm		
	4 - Red		664.6 nm		
	8 - Near Infrared (NIR)		832.8 nm		
20 m	5 - Vegetation red edge		704.1 nm		
	6 - Vegetation red edge		740.5 nm		
	7 - Vegetation red edge		782.8 nm		
	8a - Narrow NIR		864.7 nm		
	11 - Shortwave IR (SWIR)		1613.7 nm		
60 m	12 - SWIR		2202.4 nm		
	1 - Coastal aerosol		442.7 nm		
	9 - Water vapour		945.1 nm		
	10 - SWIT - Cirrus		1373.5 nm		

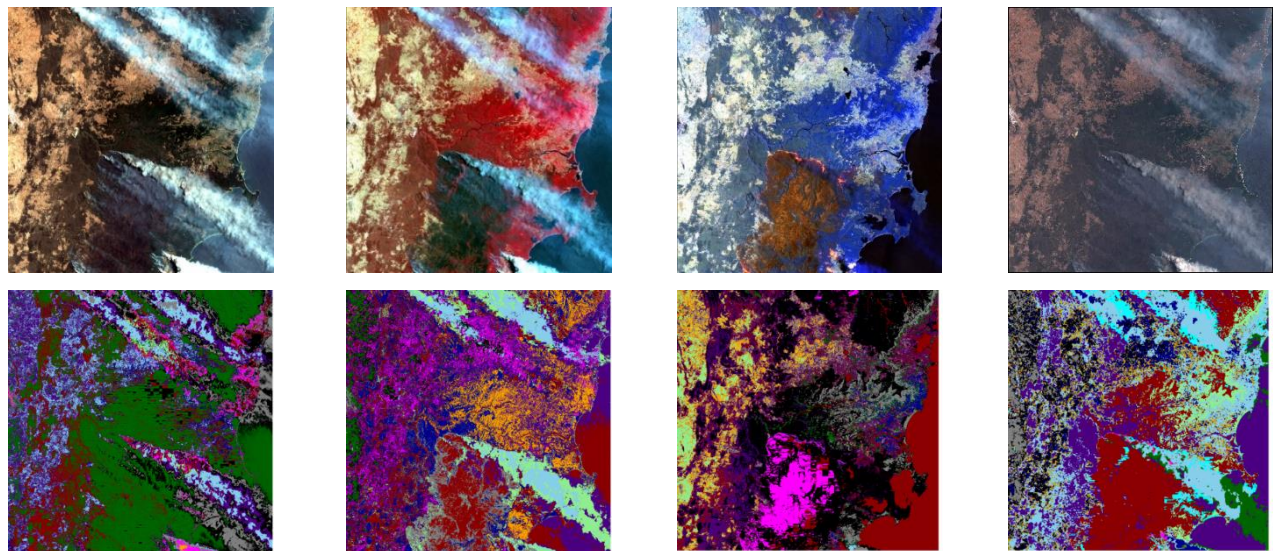
**Table 2:** Spectral bands of Sentinel-2 given by the European Space Agency (ESA).

As a conclusion for this application, to separate *Clouds* from *Smoke*, we recommend to use the second combination which can also delimit the area destroyed by *Fires*.

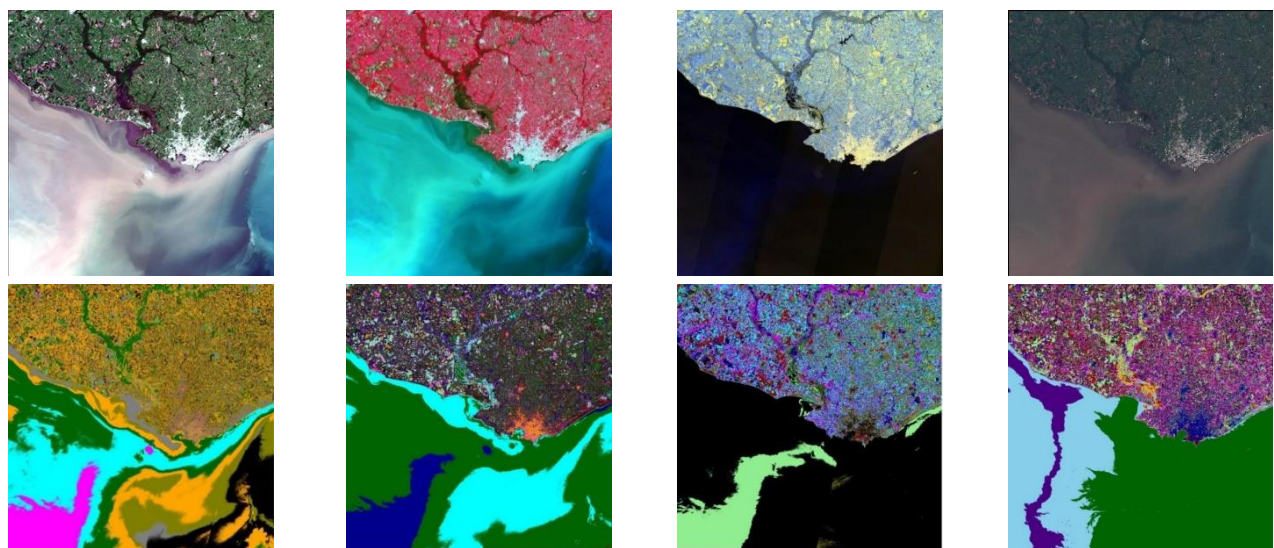




**Figure 5.** An image example of the Amazon rainforest, between Brazil, Bolivia, and Paraguay on August 25th, 2019. a) The Sentinel-2 RGB (B4, B3, and B2) quick-look image, b) Our semantic classification based on active learning (with Gabor filtering as feature vector applied to each band and an SVM as classifier), and d) Our topic representation based on LDA. The active learning process is supported by a colour legend shown in c); the colour legend of LDA is similar to the one from Figure 2.



**Figure 6.** The area of Sydney, Australia affected by forest fires on December 31<sup>st</sup>, 2019. The visibility of the different classes depends on the selection of the Sentinel-2 colour bands. (From left to right, upper part): A quick-look image of visible false-colour bands (B4, B3, and B2), false-colour visible/infrared bands (B8, B4, and B3), false-colour infrared bands (B12, B11, and B8A), and all 13 bands. (From left to right, lower part): The information that can be gained from the topics extracted by LDA for the four band combinations. The topic colours are similar to Figure 2.



**Figure 7.** The area of Montevideo, Uruguay affected by floods on October 18<sup>th</sup>, 2019. The visibility of the different classes depends on the selection of the Sentinel-2 colour bands. (From left to right, upper part): A quick-look image of visible false-colour bands, false-colour visible/infrared bands, false-colour infrared bands, and all 13 bands. (From left to right, lower part): The topic colours are similar to Figure 2.

## 5.4 Floods

Here, like in Section 5.3, we are analysing the impact of the selected Sentinel-2 spectral bands, and how different combinations can separate different classes for other applications. These combinations are the same as those from the *Fires* application depicted in Table 2.

As a second example, we show the impact of different band combinations of Sentinel-2 channels. Figure 7 illustrates the results for the area of Montevideo, Uruguay affected by floods.

Making the same type of analysis as in the previous example, we can state:

- In the first two types of band combinations, we can see the estuary formed by the Plata river at the confluence of the Paraná and Uruguay rivers, and also the river delta created by the Lucía river before entering the Plata river. Here, we can extract more topics that probably correspond to the *Ocean/River* currents, the mud brought by the floods, the alluvium silt-laden waters, etc.
- In the third combination only one extra class/topic can be retrieved in the water.
- For the last combination, we see the confluence of two waters and their separation as well as another class of *Water* that comes from the Lucia Sana River with alluvia.

As a conclusion for this last application, to better separate different water topics/classes, we recommend to use the first combination.

## 6. CONCLUSION AND FUTURE WORK

In conclusion, this paper presents a number of algorithms applied to Sentinel-1 and Sentinel-2 images for the analysis of four different applications. The output of these algorithms are semantic classification maps, topic representations or physical scattering representations. For one algorithm (active learning), a sequence of steps has to be executed by an expert user. Our proposed algorithms can be applied to both SAR and multispectral images except for the hybrid algorithm which is specific for SAR images.

In future, we plan to extend the area of applications and to include: volcanic eruptions, tsunamis/tornadoes, cyclones, landslides, and industrial accidents (Charter, 2021).

## ACKNOWLEDGEMENTS

Our approaches were tested within the framework of the European H2020 research project ExtremeEarth and the German HGF project Automated Scientific Discovery. We would like to thank our colleagues Chandra Karmakar and Reza Bahmanyar for their first insights into the domain of LDA.

## REFERENCES

Amazon– ESA Sentinel-2 applications, 2019. Available online: [https://www.esa.int/ESA\\_Multimedia/Images/2019/08/Wildfires\\_on\\_the\\_border\\_between\\_Bolivia\\_Paraguay\\_and\\_Brazil\\_from\\_Copernicus\\_Sentinel-2](https://www.esa.int/ESA_Multimedia/Images/2019/08/Wildfires_on_the_border_between_Bolivia_Paraguay_and_Brazil_from_Copernicus_Sentinel-2).

Australia - ESA Sentinel-2 applications, 2019. Available online: [http://www.esa.int/Applications/Observing\\_the\\_Earth/Copernicus/Australia\\_like\\_a\\_furnace](http://www.esa.int/Applications/Observing_the_Earth/Copernicus/Australia_like_a_furnace).

Blei, D.M., Ng, A.Y., and Jordan, M.I., 2003. Latent Dirichlet Allocation, *Journal of Machine Learning Research*, 3, pp. 993-1022.

Charter – The International Charter Space and Major Disasters, 2021. Available online: <https://disasterscharter.org/web/guest/about-the-charter>.

Cloude, S.R., and Pottier, E., 1997. An Entropy Based Classification Scheme for Land Applications of Polarimetric SAR, *IEEE TGRS*, 35(1), pp. 68-78.

Copernicus Open Access Hub, April 2021. Available online: <https://scihub.copernicus.eu/>.

Dumitru, C.O., Dax, G., Schwarz, G., Cazacu, C., Adamescu, M.C., and Datcu, M., 2019. Accurate Monitoring of the Danube Delta Dynamics using Copernicus Data, *Proc. of SPIE Remote Sensing*, Strasbourg, France, paper 1115002.

Dumitru, C.O., Schwarz, G., Datcu, M., Stillman, M., Karmakar, C., Ao, D., and Huang, Z., 2021. Improved Training for Machine Learning: The Additional Potential of Innovative Algorithmic Approaches, *Proc. of EGU General Assembly 2021*, Virtual meeting, paper EGU21-4683.

EOLib project, February 2018. Available online: <http://wiki.services.eoportal.org/tiki-index.php?page=EOLib>.

ExtremeEarth project, April 2021. Available online: <http://earthanalytics.eu/>.

Huang, Z., Dumitru, C.O., Pan, Z., Lei, B., and Datcu, M., 2020. Classification of Large-Scale High-Resolution SAR Images with Deep Transfer Learning, *IEEE GRSL*, 18(1), pp. 107-111.

Huang, Z., Dumitru, C.O., and Ren, J., 2021. Physics-Aware Feature Learning of SAR Images with Deep Neural Networks: A Case Study, *Proc. of IGARSS*, Brussels, Belgium, pp. 1-4.

Karmakar, C., Dumitru, C.O., Schwarz, G., and Datcu, M., 2021. Feature-Free Explainable Data Mining in SAR Images Using LDA, *IEEE JSTARS*, 14, pp. 676-689.

Lee, J.S., Grunes, M.R., Ainsworth, T.L., Du, J.L., Schuler, D.L., and Cloude, S.R., 1999. Unsupervised Classification Using Polarimetric Decomposition and the Complex Wishart Classifier, *IEEE TGRS*, 37(5), pp. 2249-2258.

Smith, S., 2004. Digital Signal Processing: A Practical Guide for Engineers and Scientists, *Elsevier Science & Technology Publisher*, 3rd Revised ed., Oxford, UK, 664 pages.

Uruguay - ESA Sentinel-2 applications, 2019. Available online: [http://www.esa.int/ESA\\_Multimedia/Images/2020/04/Montevideo\\_Uruguay](http://www.esa.int/ESA_Multimedia/Images/2020/04/Montevideo_Uruguay).

Vaduva, C., Georgescu, F.A., and Datcu, M., 2015. Dictionary-Based Compact Data Representation for Very High-Resolution Earth Observation Image Classification, *Lecture Notes in Computer Science*, 9386, pp. 816-825.