






Article

Improving Land Cover Classification Using Genetic Programming for Feature Construction

João E. Batista ^{1,*} , Ana I. R. Cabral ² , Maria J. P. Vasconcelos ² , Leonardo Vanneschi ³ , Sara Silva ¹ 

¹ LASIGE, Faculty of Sciences, University of Lisbon, Campo Grande, 1749-016 Lisbon, Portugal

² Forest Research Centre, School of Agriculture, University of Lisbon, Tapada da Ajuda, 1349-017, Lisbon, Portugal

³ NOVA Information Management School (NOVA IMS), Universidade Nova de Lisboa, Campus de Campolide, 1070-312 Lisbon, Portugal

* Correspondence: jebatista@fc.ul.pt

Abstract: Genetic Programming (GP) is a powerful Machine Learning (ML) algorithm that can produce readable white-box models. Although successfully used for solving an array of problems in different scientific areas, GP is still not well known in Remote Sensing. The M3GP algorithm, a variant of the standard GP algorithm, performs Feature Construction by evolving hyper-features from the original ones. In this work, we use the M3GP algorithm on several sets of satellite images over different countries to create hyper-feature from satellite bands to improve the classification of land cover types. We add the evolved hyper-features to the reference datasets and observe a significant improvement of the performance of three state-of-the-art ML algorithms (Decision Trees, Random Forests and XGBoost) on multiclass classifications and no significant effect on the binary classifications. We show that adding the M3GP hyper-features to the reference datasets brings better results than adding the well-known spectral indices NDVI, NDWI and NBR. We also compare the performance of the M3GP hyper-features in the binary classification problems with those created by other Feature Construction methods like FFX and EFS.

Keywords: Genetic Programming; Evolutionary Computation; Machine Learning; Classification; Multiclass Classification; Feature Construction; Hyper-features; Spectral Indices

1. Introduction

Since the establishment of the Warsaw Framework in 2013, Remote Sensing (RS) is recommended as an appropriate technology for monitoring and Measuring, Reporting and Verification (MRV) for countries reporting forest land cover and land cover change to the UNFCCC ¹. However, many difficulties, from the availability of adequate in-situ reference data to the spatial and temporal resolution of freely available satellite imagery and data processing power, have been hindering the operational use of this technology for MRV. Now, with the evolution of Earth Observation systems (with provision of higher spatial and temporal resolution images) and with novel open-data distribution policies, there is an opportunity of applying Machine Learning (ML) to induce models that automatically identify land cover types in satellite images and improve the capacity for producing frequent and accurate land cover maps.

Previous ML work in classification of satellite imagery for land cover mapping has been very successful. One simple practice that helps obtain good results is the inclusion of spectral indices as additional independent variables² in the reference dataset. Spectral indices are combinations of reflectance values from different wavelengths that represent the relative abundance of certain terrain elements. They have been used by the RS community for a long time to enhance the identification of vegetation (e.g., NDVI [1]), water (e.g., NDWI [2]), burnt areas (e.g., NBR [3]) and many other elements. Over the years, many indices were created and adapted to accommodate the particularities of different images. In the case of vegetation indices, this number is so vast that over one hundred of them were reviewed in [4].

Like indices, hyper-features are mathematical expressions that combine the original features of the data (the independent variables) with the goal of representing data properties that facilitate the learning of ML models. Spectral indices are, in fact, particular cases of hyper-features. Ideally, the hyper-features should be simple and meaningful, allowing the RS experts to easily understand the ML models that are based on them, or to directly use them in image analysis software to visualise what they represent.

Notwithstanding the success of ML methods when performing classification of satellite imagery, the reported results are often obtained by applying a model in the same images where it was trained (e.g., [5–7]), or in an image time series from the same location (e.g., [8–10]). Training models to be ready to be used outside their training images is not a trivial task due to the radiometric variations between different images. These variations can arise from multiple sources, such as the difference in the angle of the solar incidence on the ground; the weather; the conditions of the terrain; the type of terrain; or the growth stage of the vegetation. Spectral indices are also sensitive to these variations, despite the efforts to increase their robustness.

Our goal is to improve satellite imagery classification, by creating hyper-features that increase the performance of ML algorithms. In previous work [11], we used a Genetic Programming (GP) [12] classifier called M3GP [13] to evolve hyper-features that, when used instead of the original ones, were able to improve the accuracy of different ML algorithms in binary classification of images different from the ones used in training (although, for unseen data of the same images, there was no significant effect). GP is a powerful ML evolutionary algorithm that can produce readable white-box models. Successfully used for solving an array of problems in different scientific areas, GP is, however, still not well known in RS. The M3GP algorithm is a variant of standard GP that was originally developed as a multiclass classifier, but later used as a Feature Construction method for other algorithms, both for classification and for regression [11,14,15]. Creating hyper-features from one image and using them for classifying a different image falls under the area of Transfer Learning [15], which attempts to use knowledge from one problem to solve another similar problem.

The area of ML is divided into several fields, one of which is named Evolutionary Computation (EC). This field deals with the creation (or evolution) of models using an evolutionary cycle that was inspired by the evolution theories of Charles Darwin. Using this cycle, different flavors of EC (from which GP is one of the youngest) use a fitness function to guide a population of evolving models through a search space, until one of its individuals reaches a certain fitness and is returned as the best model. This cycle will be further explained in Section 3.3.1, in the context of the M3GP algorithm.

In this work, we perform a thorough study of the effects of adding M3GP-evolved hyper-features to the reference datasets. We test our approach on several datasets from different images in two types of problems that have been tackled several times over the last decades, the binary classification of burnt areas [16–22] and the multiclass classification of land cover types [23–27]. The images used in our study cover several different regions over developing countries: Angola, Brazil, Democratic

Republic of the Congo, Guinea-Bissau and Mozambique. We add the evolved hyper-features to the reference datasets and analyse the differences in the generalisation ability of different ML algorithms when tested on unseen data from the same images. Three common state-of-the-art algorithms are tested, namely Decision Trees [28], Random Forests [29] and XGBoost [30]. We also perform the same experiments when adding spectral indices instead of the hyper-features, comparing the results. The selected indices are the popular NDVI, NDWI and NBR. For the binary classification problems, we also compare our results with the ones obtained when adding hyper-features created by two different Feature Construction methods, EFS [31] and FFX [32].

It is important to emphasise the differences between the current work and the previous one [11]. On the previous work, a manual selection of evolved hyper-features completely replaced the original features of the reference datasets, while here all the hyper-features resulting from each run are automatically added to the reference datasets. The goal of the previous work was to explore feature spaces in order to explain the variable degrees of success of Transfer Learning to different images. Here, we concentrate on the performance inside each image, and compare our approach to alternative ones that use indices and other types of hyper-features. Finally, while the previous work only used binary classification datasets, this one greatly extends its reach by tackling also multiclass classification problems.

2. Related Work

Feature Engineering is an essential step in the knowledge discovery process and one of the keys to success in applied ML. The features used to induce a data model can directly influence the quality of the model itself and the results that it can achieve. Feature Engineering can be broadly partitioned into Feature Selection and Feature Construction. According to [33], Feature Selection is a process that chooses a subset of features from the original data variables, so that the feature space is optimally reduced according to a certain criterion, while Feature Construction/Extraction (also called Feature Generation, Feature Learning, or Constructive Induction) is a process that creates a new set of hyper-features from the original data. Feature Construction typically combines existing variables into more informative hyper-features. Both Feature Selection and Feature Construction attempt to improve model performance and can be used in isolation or in combination.

Feature Construction, the focus of this work, has been widely studied in the last two decades. Recent surveys can be found in [34–36], while the book [37] gives an in-depth presentation of the area. In all these references, the importance of EC as an effective method for Feature Construction is asserted, together with other Feature Construction methods such as the ones based on Decision Trees, Inductive Logic Programming and Clustering. A very recent survey of EC techniques for Feature Construction can be found in [38]. Among the different EC flavours, GP is probably the one that has been used more often and more successfully. Indeed, GP is particularly suited for Feature Construction because it naturally evolves functions of the original variables. The versatility offered by the user-defined fitness function of GP allows the user to choose among several possible criteria for evolving new hyper-features. Additionally, the fact that the evolved hyper-features are, in principle, readable and understandable, can play an important role in model interpretability. Several existing GP-based methods for Feature Construction are discussed in [39,40], and a deep analysis of previous work can be found in [15], where GP-constructed features are used for Transfer Learning.

Among the large set of Feature Construction methods available, in this paper we use M3GP [13] as our method of choice, and two others for comparison purposes: the non-EC method FFX [32] and the EC method EFS [31]. The remainder of this section will focus on GP-based Feature Construction, including applications, and on Feature Construction and GP in the context of RS.

2.1. Feature Construction with Genetic Programming

Among the several previous contributions in which GP was used for Feature Construction, Krawiec has shown that classifiers induced using the representation enriched by GP-constructed

hyper-features provide better accuracy on a set of benchmark classification problems [41]. Krawiec and colleagues have also used GP in a co-evolutionary system for Feature Construction [42,43].

The use of GP for Feature Construction was later deeply investigated by Zhang and colleagues. For instance, in [44], a GP approach was proposed that, instead of wrapping a particular classifier for single Feature Construction as in most of the existing methods, it used GP to construct multiple features from the original variables. The proposed method used a fitness function based on class dispersion and entropy, and thus was independent of any particular classification algorithm. The approach was tested using Decision Trees on the new obtained dataset and experimentally compared with the standard Decision Tree method, using the original features. The results showed that the proposed approach outperforms standard Decision Trees on the studied test problems in terms of the classification performance, dimension reduction and the learned Decision Tree size. Several years later, in [45], GP was used for both Feature Construction and implicit Feature Selection. The work presented a comprehensive study, investigating the use of GP for Feature Construction and Feature Selection on high-dimensional classification problems. Different combinations of the constructed and/or selected features were tested and compared on seven high-dimensional gene expression problems, and different classification algorithms were used to evaluate their performance. The results indicated that the constructed and/or selected feature sets can significantly reduce the dimensionality and maintain or even increase the classification accuracy in most cases. In [46], previous GP-based approaches for Feature Construction were extended to deal with incomplete data. The results indicated that the proposed approach can, at the same time, improve the accuracy and reduce the complexity of the learnt classifiers. While until a few years ago GP-based Feature Construction had been applied mainly to classification, in [47] it was applied with success to symbolic regression, thus giving a demonstration of the generality of the approach. In [48], different approaches based on GP to constructing multiple features were investigated. One of the most interesting results showed that multiple-feature construction achieves significantly better performance than single-feature construction. Consistently with that result, also the method presented in this paper uses GP to construct multiple features.

It should be noted that the use of GP for Feature Construction has been explored for some time, as surveyed in [40]. Although the most common approach to multiclass classification problems used to be splitting a classification problem with n classes into n binary classification problems, and evolving one hyper-feature for each class [49,50], some methods create several hyper-features to separate the classes within the feature space. In this category, the survey includes works that converted the datasets into hyper-datasets using exclusively the evolved hyper-features [41], and works that include the original hyper-features in the hyper-dataset [51] (similarly to our work).

GP-based Feature Construction methods have been used with success in several real-life applications. For instance, [52] proposed a novel method for breast cancer diagnosis using the features generated by GP. A few years later, in [53], GP-based Feature Construction was used for improving the accuracy of several classification algorithms for biomarker identification. In [54], a method to find smaller solutions of equally high quality compared to other state-of-the-art GP approaches was coupled with a GP-based Feature Construction method and applied to cancer radiotherapy dose reconstruction. One year later, in [55], GP-based Feature Construction was successfully applied to the classification of ten different categories of skin cancer from lesion images. Interestingly, while the application tackled in [54] is a symbolic regression problem, the one in [55] is a multiclass classification problem, thus confirming that the GP-based Feature Construction approach can be successfully applied to both types of problems. Finally, in [56], GP-based Feature Construction was extended for the first time to experimental physics. In particular, to be applicable to physics, dimensional consistency was enforced using grammars. The presented results showed that the constructed hyper-features can both significantly improve classification accuracy and be easily interpretable.

2.2. Feature Construction and Genetic Programming in Remote Sensing

In the RS domain, many techniques have been used to extract features from satellite images. These features include statistical descriptors, obtained by the Gray Level Co-occurrence Matrix (GLCM) and other methods [57]; features of interest, such as known structures (e.g., buildings, roads), using deep learning [58]; sets of generic features, using the Principal Component Analysis (PCA) [59]; and even temporal features, using the Continuous Change Detection and Classification (CCDC) algorithm [60].

GP-based algorithms, mainly the standard GP algorithm, have been previously used in the area of RS in tasks such as the creation of vegetation indices [61], the detection of riparian zones [62] and the estimation of soil moisture [62,63], the estimation of canopy nitrogen content at the beginning of the tasselling stage [64], the estimation of chlorophyll levels to monitor the water quality in reservoirs [65], the prediction of soil salinity by estimating the electrical conductivity on the ground [66] and also in geoscience projects reviewed in [67].

The expressions obtained by the GP-based algorithms can be used in Transfer Learning by exporting them to datasets under the form of hyper-features, in the attempt to improve the performance of ML algorithms. Our work continues to develop this kind of application, which was already explored in the area of RS using EC-based algorithms [68,69] and specifically GP-based algorithms [11,62].

3. Materials and Methods

In this section, we describe the used datasets and the respective climate and type of vegetation in each of their geographic locations. This section also includes a description of the Feature Construction and Classification algorithms used, with particular emphasis on the M3GP algorithm, whose algorithm is explained step by step.

3.1. Datasets and Study Areas

The datasets used in this work are meant to train ML models to classify burnt areas and land cover types on a pixel-level. We use a total of nine datasets, obtained from Landsat-7, Landsat-8 and Sentinel-2A satellite images. The characteristics of these images and datasets are summarised in Tables 1 and 2, and their associated geographic locations are highlighted in Figure 1.

3.1.1. Datasets

From the Landsat-7 images, we have one binary classification dataset (Gw2) and two multiclass classification datasets (IM-10 and IM-3). The IM-3 dataset was built, in previous work, from IM-10 by extracting only the pixels classified in-situ from the three forest land cover types that ML models failed to correctly discriminate. These images were both obtained over Guinea-Bissau.

From the Landsat-8 images, we have three binary classification datasets and two multiclass classification datasets. The binary classification datasets have the objective of training models to identify burnt areas, by classifying each pixel as "burnt" or "non-burnt". These three datasets were obtained from satellite images over Angola (Ao2), Brazil (Br2) and Democratic Republic of the Congo (Cd2). The multiclass classification datasets have the objective of training models to correctly classify each pixel as one of several different land cover types. These two datasets were extracted from satellite images over Angola (Ao8) and Guinea-Bissau (Gw10).

Lastly, from the Sentinel-2A satellite images, we have one multiclass classification dataset that was extracted from several satellite images from the entire country of Mozambique (Mz6). These images were obtained through 2016, between February 19th and October 6th [70].

3.1.2. Study Areas

In terms of size of the classified areas, the pixels used in the Landsat satellite images consist of 900m² areas and those used in the Sentinel-2A satellite images consist of 400m² areas. As such,

Table 1. Summary of the datasets used.

Dataset	Ref.	Country	Scene Identifier Path / Row	Acq. Date DD/MM/YYYY	No. Images	Satellite	KGCS
Ao2	^a	Angola	177 / 67	09/07/2013	1	LS-8	Cwa
Br2	[71]	Brazil	225 / 64	28/02/2015	1	LS-8	Af, Am
Cd2	[71]	DR Congo	175 / 62	08/06/2013	1	LS-8	Aw
Gw2	[71]	Guinea-Bissau	204 / 52	13/05/2002	1	LS-7	Am, Aw
IM-3	^b	Guinea-Bissau	203 / 51, 52	From: 02/01/2010 To: 01/04/2010	17	LS-7	Am, Aw
IM-10	[72]		204 / 51, 52				
			205 / 51				
Ao8	[73]	Angola	182 / 64, 65	18/06/2016	2	LS-8	Aw
Gw10	^c	Guinea-Bissau	204 / 51, 52	01/03/2019	3	LS-8	Am, Aw
			205 / 51	24/03/2019			
Mz6	[70]	Mozambique	Entire Country (122 S-2A tiles)	From: 19/02/2016 To: 06/10/2016	2806 ^d	S-2A	Am, Aw, BSh, Cwa, Cwb, Cfa

^a There is no reference paper for this dataset.

^b This is a sub-dataset, obtained by extracting three forest classes from the IM-10 dataset.

^c The reference paper for this dataset is under review.

^d An approximation obtained by considering that the S-2A mapped every tile of Mozambique once every 10 days for 230 days.

the classified areas can be calculated from Table 2. Next, we describe the climate (according to the Köppen–Geiger Classification System (KGCS) [74]) and vegetation in each of the study areas:

Brazil: The study area of the Br2 dataset is located in eastern Amazonia, in southeastern Pará, Brazil. According to the KGCS, the climate in this image is classified as Equatorial Monsoon (Am) and Equatorial rainforest, fully humid (Af) in the north and south sections, respectively. This area is drier than central and western Amazonia, with annual rainfall between 1500mm and 2000mm and average temperatures ranging from 23°C to 30°C. The vegetation in this image ranges from lowland Amazon forest in the north through submontane dense and open forests in the south [71].

Guinea-Bissau: The study area of the Gw2, IM-3, IM-10 and Gw10 datasets, is located in Guinea-Bissau, West Africa. According to the KGCS, the climate in this area is classified as *Am* and Equatorial savanna with dry winter (*Aw*) within the coastal and interior areas, respectively. This area is characterised by having a marshy coastal plain with a dry to moist (North to South) tropical climate. There are two marked seasons, a dry season between November and May, and a wet season between June and October. Total annual rain values vary from 1200 to 1400mm in the Northeast region, and from 2400 to 2600mm in the Southwest region. The monthly average temperature ranges from 25.9°C and 27.1°C. The vegetation consists of mangroves on the coast and gradually becomes composed of mainly dry forest and savanna inland [71].

Northern Angola: The study area of the Ao8 dataset is located in the Zaire province, northern Angola. According to the KGCS, the climate in this region is classified as *Aw* with a mean annual rainfall near 1300mm, distributed in two periods separated by a short dry spell. The monthly average temperature ranges from 20.5°C and 24.9°C. The vegetation is mainly savanna scrublands and some dense humid forests mostly located along rivers, creeks, and gullies. There are anthropic forests composed by native species and mango, cola, safou, avocado, citrus, and guava trees in ancestral settlements, abandoned due to forced relocation along the main roads by the colonial administration [75].

Eastern Angola: The study area of the Ao2 dataset is located in Lunda Sul, Eastern Angola. According to the KGCS, the climate in this area is classified as Warm temperate climate with dry winter and hot summer (*Cwa*) and a mean annual rainfall near 1300mm, distributed between October and April and a dry season from May to September. The monthly average temperature ranges from 20.0°C and 24.4°C. The vegetation is mainly dominated by woody and shrub savannas and gallery forests essentially located along the valleys of the great rivers [76].

Democratic Republic of Congo: The study area of the Cg2 dataset is located in the central-eastern Democratic Republic of Congo. According to the KGCS, the climate in this area is classified as *Aw* with

Table 2. Summary of the datasets used.

Dataset	Classes (No. Pixels)			No. Classes	No. Bands No. Features	Total Pixels
Ao2	Burnt (1573)	Non-Burnt (2309)		2	7	3882
Br2	Burnt (2033)	Non-Burnt (2839)		2	7	4872
Cd2	Burnt (877)	Non-Burnt (1972)		2	7	2849
Gw2	Burnt (1101)	Non-Burnt (3430)		2	7	4531
IM-3	Savanna Woodland (114)	Dense Forest (68)	Open Forest (140)	3	6	322
IM-10	Agriculture/Bare Soil (950)	Burnt (77)	Dense Forest (524)	10	6	6798
	Grassland (75)	Mangrove (1240)	Open Forest (723)			
	Savanna Woodland (1626)	Sand (166)	Mud (509)			
	Water (908)					
Ao8	Agriculture/Bare Soil (73)	Burnt (301)	Clouds (332)	8	10	2183
	Forest (662)	Grassland (12)	Urban (53)			
	Savanna Woodland (598)	Water (152)				
Gw10	Agriculture/Bare Soil (449)	Burnt (157)	Dense Forest (62)	10	7	5080
	Grassland (16)	Mangrove (1383)	Open Forest (646)			
	Savanna Woodland (1308)	Sand (50)	Water (620)			
	Wetland (389)					
Mz6	Agriculture/Bare Soil (33611)	Forest (63190)	Grassland (28406)	6	10	190202
	Urban (4194)	Wetland (35673)	Other (25128)			

a mean annual rainfall near 1600mm. There are two distinct seasons, a dry season (with temperatures ranging between 18°C and 27°C) from June to August, and a rainy season (with temperatures ranging between 22°C and 33°C) from September to May. The vegetation is characterised by a congolian lowland forest in the north to miombo woodlands in the south. In the southwestern region, the population pressure had conducted to the degradation of the miombo woodlands [71].

Mozambique: The study area of the Mz6 dataset includes the entire country of Mozambique. According to the KGCS, the climate is classified as *Aw* in the coastal area and near the Zambezi river; as *Cwa* in the interior, at the north and the west of the Zambezi river; as Warm temperate climate with dry winter and warm summer (*Cwb*) near Lichinga and west of Chibabava; as Hot semi-arid (*BSh*) in the interior in south Mozambique and east of Mungári and Derre, and as Hot desert (*Bwh*) in the area between Dindiza and the frontier between Mozambique and Zimbabwe, south of the Save river. It has a wet season from October to March and a dry season from April to September. The lowest average rainfall (300-800mm/year) occur in the interior southern regions and the highest average rainfall (over 1200 mm/year) occur in the area around Espungabera. The average temperatures are the highest along the coast in the northern regions (with temperatures ranging between 25°C and 27°C in summer and between 20°C and 23°C in winter) and in the southern regions (with temperatures ranging between 24°C and 26°C in summer and between 20°C and 22°C in winter), while the high inland regions have cooler temperatures (with temperatures ranging between 18°C and 20°C in summer and between 13°C and 16°C in winter). The northern areas are predominantly occupied by miombo woodlands and the western and southern borders by Zambezian and Mopane woodland. The most widespread vegetation in the north coast is the Zanzibar-Inhambane forest mosaic, followed by the African mangroves and the Maputaland forest mosaic in the south-east coast [70,77,78].

3.2. Methodology

The core of this work is to expand the reference datasets with hyper-features that improve the performance of different ML methods. Figure 2 illustrates the process of obtaining and using such hyper-features. As usual, the reference dataset is split in two datasets, one for training the classifiers, called the training set, and one for testing the classifiers on unseen data, called the test set. Based only on the training set, the Feature Construction algorithm creates a set of hyper-features that are used to expand the reference dataset, in both training and test sets. The expanded training set is used by the classification algorithm to obtain a trained classifier, that is applied to both (expanded) training and test sets in order to report the performance in terms of learning and generalisation, respectively.

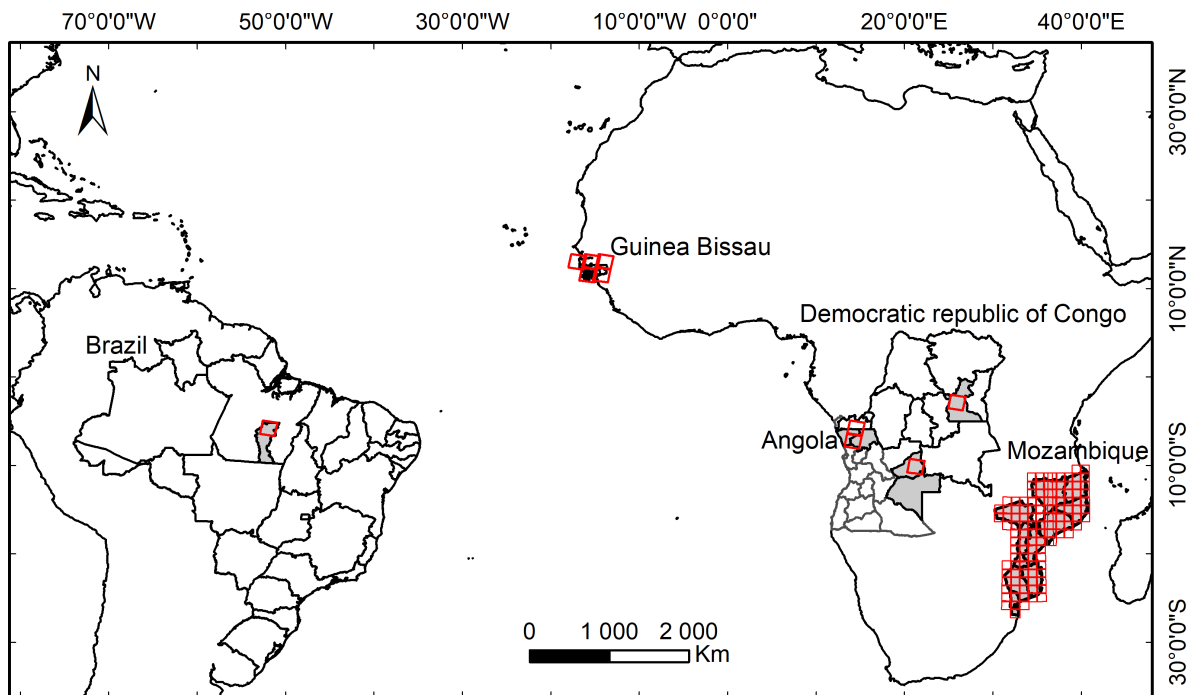


Figure 1. Location of the study areas, in red, in Brazil (Br2), Guinea-Bissau (Gw2, Gw10, IM-3, IM-10), Democratic Republic of Congo (Cd2), Angola (Ao2, Ao8) and Mozambique (Mz6) in South America and Africa continent.

A small deviation to this process has been made for the datasets IM-3 and IM-10, where the hyper-features used to expand the IM-10 datasets were obtained in the training data of its subset IM-3, and not the training data of the complete IM-10. The goal of this deviation was to check whether a larger dataset could also benefit from hyper-features obtained in a much limited context (results reported in Section 4.3).

As Feature Construction algorithms, we use M3GP and compare it with FFX and EFS, all described below. As classification algorithms, we use Decision Trees (DT), Random Forests (RF) and XGBoost (XGB), also briefly described below. The number and complexity of the created hyper-features is not predefined, but automatically determined by the Feature Construction algorithm.

We also experiment with expanding the reference datasets with the NDVI, NDWI and NBR indices, instead of doing Feature Construction. These indices were selected from the RS literature as being helpful to the ML algorithms for separating vegetation, water and burnt classes, since these elements are present among the pixels used in the datasets.

Each experiment is performed 30 times for each possible trio of reference dataset, Feature Construction algorithm and Classification algorithm (with the exception of the EFS and FFX algorithms, which are only used in binary classification datasets), each time with a different random split of the reference dataset in training and test sets. In other words, and limiting the explanation to Figure 2, our experimental process follows these steps:

Splitting the Dataset: The reference dataset is split randomly into training (70% of the pixels) and test sets (remaining 30%), stratified by class;

Creating and Adding Hyper-Features / Indices: The training set is used by a Feature Construction algorithm to create a new set of hyper-features, and the training and test sets are then extended using these hyper-features, or the indices;

Training and Testing a Classifier: A classifier is trained using the extended training set and tested on the extended test set, providing the final results.

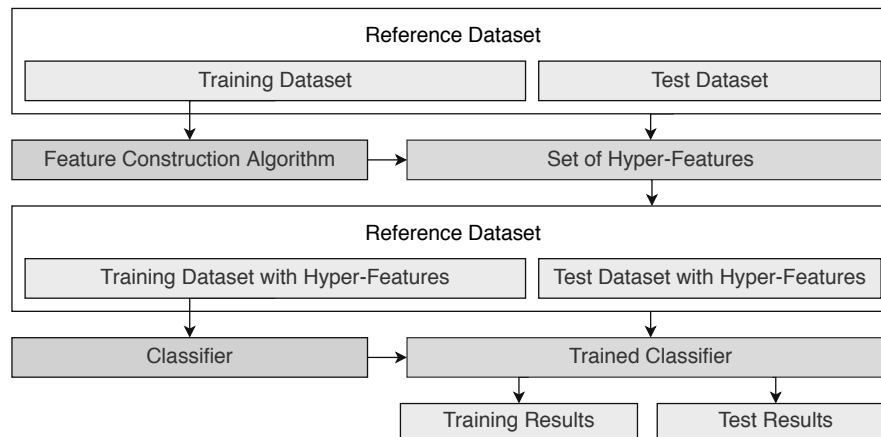


Figure 2. Representation of the methodology adopted to obtain and use hyper-features.

3.3. Feature Construction Algorithms

We use three different methods for Feature Construction. Our method of choice is the M3GP algorithm, because of the interpretability of the hyper-features it creates, and because it can evolve hyper-features for multiclass classification problems. For comparing our M3GP results with the results of other evolutionary and non-evolutionary methods, we selected the EFS and FFX algorithms, due to their running speed, availability of the authors' implementations and number of citations. However, EFS and FFX are focused on regression problems, rather than classification problems. They are easily adapted to binary classification, by defining a threshold separating the two classes, but there is no easy adaptation for multiclass problems, the reason why we test them only on the binary classification datasets.

M3GP algorithm: Multidimensional Multiclass GP with Multidimensional Populations (M3GP) is a GP-based algorithm that evolves a set of hyper-features that convert the original feature space into a new feature space, guided by a fitness function that measures the performance of a classifier in the new feature space. The M3GP is an all-in-one algorithm that both creates the hyper-features and uses them for solving both regression and classification problems. The inner workings of this algorithm are explained below, in section 3.3.1.

EFS algorithm: Evolutionary Feature Synthesis (EFS) is an evolutionary algorithm that uses pathwise LASSO [79] regression to optimise multiple linear regression models that are extended for nonlinear relationships between features. This extension is made using functions such as *cos*, *sin* and *log*, as well as functions with several inputs, e.g., multiplication of variables. This regression tool can produce a set of interpretable hyper-features in seconds.

FFX algorithm: Fast Function Extraction (FFX) is a deterministic algorithm that applies Pathwise Regularised Learning [80] to a large set of generated nonlinear functions to search for a set of hyper-features with minimal error. Although the hyper-features observed are simple, this algorithm generates hundreds of hyper-features, which leads us to consider the final model non-interpretable.

3.3.1. The M3GP Algorithm

The M3GP [13] is a GP-based algorithm that evolves models similar in structure to the models of standard GP [12]. Standard GP represents each model as a parse tree, to be read depth-first, where the terminal nodes are features and the non-terminal nodes are operators that combine features (e.g., arithmetic operators like multiplication, subtraction, etc). The main difference between the models evolved by M3GP and the ones of standard GP is that, in M3GP, each model is not a single tree, but a set of trees, as exemplified in Figure 3. These trees are what we call hyper-features and are evolved using the steps shown in Figure 4 and explained here:

Initialisation: The M3GP algorithm initialises its population of models as single, randomly generated, trees. Therefore, in the beginning of the evolutionary cycle, each individual is a simple model consisting of a single hyper-feature.

Evaluation: Each individual of the population is evaluated by the following procedure. The n hyper-features are used to convert the original features into a new n -dimensional dataset. The fitness of the model is then calculated by applying a fitness function to this new (hyper) feature space. This fitness function rewards the individuals whose set of hyper-features creates a space where the different classes are more easily separable. In the original M3GP algorithm, the fitness was the overall accuracy of the Mahalanobis Distance Classifier (described below), but in the current implementation we use the Weighted Average of F-measures (WAF) instead of the overall accuracy³, for its robustness to class imbalance, especially in multiclass classification.

Stopping Criteria: After the population is evaluated, the algorithm decides whether to stop the evolution or not. The most common stopping criteria are related to the number of generations/iterations already done, and to the quality of the best model achieved (in terms of accuracy or any other metric). In the current implementation, the evolution stops when 50 generations are completed or when one individual achieves 100% accuracy on the training set, whichever occurs first. If the evolution does not stop, a new generation of models is created, following the steps described next.

Selection: The parents of the next generation are selected using the *tournament* method. To select each parent, the algorithm randomly selects a small group of models and, out of these models, chooses the best. The tournament method is able to maintain enough selective pressure to choose mostly the best individuals, thus promoting the propagation of their good traits in the next generation, while allowing also the bad ones to become parents, thus avoiding the loss of genetic diversity that would stagnate the evolution.

Breeding: After selecting models to act as parents, each new model is created either through a mutation of one parent or through a crossover of two parents. When using a mutation genetic operator, the parent can either: create a new, randomly generated, tree and add it to its set of hyper-features; randomly select one of its hyper-features and remove it (if it contains more than one hyper-feature); or modify one of its hyper-features by replacing one of its branches with a new, randomly generated, tree. When using a crossover genetic operator, the parent models can swap either branches or entire hyper-features between each other. Unlike the mutation genetic operator, the crossover results in two offspring.

After a new population has been created, the algorithm returns to the evaluation step.

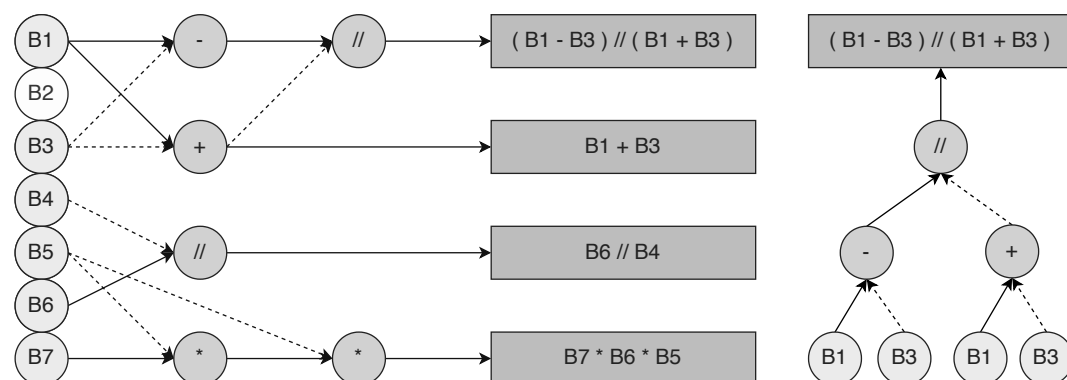


Figure 3. Example of an M3GP model that uses six of the seven available features to build four hyper-features (left) and a single hyper-feature (right). The solid and dashed lines indicate the first and second variables used by the operators. $//$ is a division operator protected against division by zero.

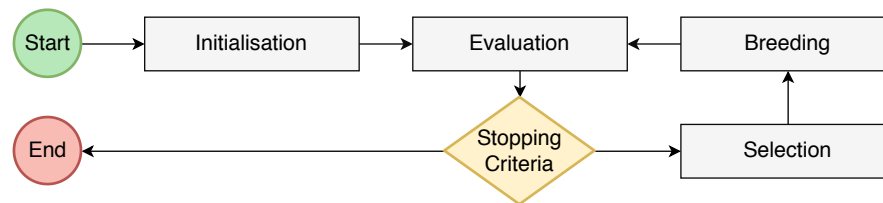


Figure 4. Evolutionary cycle used by M3GP.

3.4. Classification Algorithms

Four different classifiers are used in this work: MD, DT, RF and XGB. The MD classifier is used only as part of M3GP, but the other three are used to independently test the effectiveness of the indices and hyper-features added to the reference datasets.

Mahalanobis Distance classifier: The MD classifier is a non-parametric supervised cluster-based algorithm that classifies a data point by associating it with the closest cluster centroid using the Mahalanobis distance, where a cluster is defined as a set of pixels belonging to the same class.

Decision Tree classifier: The DT algorithm is a non-parametric supervised algorithm that infers simple decision rules from the training data. This algorithm can be used in both classification and regression problems.

Random Forest classifier: The RF algorithm is an ensemble algorithm that uses a set of DT to solve both classification and regression problems, by assigning each data point to the majority vote of all DT in classification problems, or to the average of the prediction of all DT in regression problems.

XGBoost classifier: The XGB algorithm is DT-based ensemble algorithm that uses an optimised gradient boosting to minimise errors. This algorithm can be used in both classification and regression problems.

3.5. Tools and Parameters

All the experiments involving M3GP are performed using our own implementation of the M3GP algorithm⁴, which includes the implementation of the MD classifier. The DT and RF implementations belong to the sklearn python library [81] and the XGB implementation belongs to the xgboost python library [30]. The EFS implementation⁵ is provided by the authors in their paper [31] and the FFX implementation⁶ belongs to the ffx python library.

The parameter settings used in this work are the standard within the ML community, with the main parameters and variations specified in Table 3. The EFS, FFX and M3GP algorithms used the same parameters as those used by the authors in their respective papers. The variations in this work include our implementation of the M3GP using the WAF of the MD classifier (untied with the number of hyper-features and then with the total size of the model) as fitness, rather than the accuracy, and only pruning the final individual, for consistency with previous work that had this variation [11]. Every run using the DT, RF and XGB classifiers used the default parameters of their respective implementations, except for the XGB runs in the Mz6 dataset. In this dataset, the XGB was unable to obtain perfect training accuracy with the default maximum depth for its models. As such, the maximum depth was increased from 6 to 20.

4. Results and Discussion

We start this section by presenting the results and hyper-features obtained by running M3GP by itself on all the datasets. Then, we discuss the interpretability of the hyper-features and the popularity

Table 3. Main parameters and variations used in the experiments.

General:	
Runs	30
Training Set	70% of the samples of each class
Statistical Significance	p -value < 0.01 (Kruskal-Wallis H-test)
M3GP:	
Stopping Criteria	50 generations or 100% training accuracy
Fitness	WAF (Weighted Average of F-measures)
Pruning	Final individual
XGBoost:	
Maximum Depth	20 in the Mz6 dataset and 6 (default) in the other datasets

of the different satellite bands in the solutions proposed by M3GP. Next, we compare the overall accuracy, and class accuracy, obtained when running the DT, RF and XGB algorithms on the original datasets and on the datasets expanded with indices or hyper-features. Regarding the class accuracy, we present the results only for XGB because the results for DT and RF were very similar in terms of the relationship between the classes, and therefore would not bring any new information.

The presentation of the results is split into three categories: binary classification datasets (Ao2, Br2, Cg2 and Gw2), regarding the detection of burnt areas; discrimination of similar classes (IM-3 and IM-10), regarding the separation of forest types; and discrimination of all classes (Ao8, Gw10 and Mz6), regarding the separation of different land cover types.

We present the results in tables, boxplots and confusion matrices. On the tables, each overall accuracy value is the median obtained in the 30 runs. The confusion matrices, rather than showing the class accuracy, show the difference (in percentage of pixels) between using hyper-features (or indices) and using the original dataset, to facilitate the identification of the effect produced by the hyper-features. Statistical significance is determined with the non-parametric Kruskal-Wallis H-test (from the scipy Python library) at $p < 0.01$.

4.1. M3GP Performance and Hyper-feature Analysis

Although we use M3GP as a Feature Construction method for other ML algorithms, M3GP can perform binary and multiclass classification by itself, as described in Section 3.4. While using M3GP to evolve the hyper-features, we have registered the accuracy values it achieved in each dataset, presented in Table 4. Although the accuracy is high, it is generally worse than the accuracy achieved by the other ML algorithms we used, and therefore we will not refer to the results of standalone M3GP again.

In terms of interpretability of the evolved hyper-features, in Table 4 we report the number of hyper-features and their median size (with minimum and maximum values, between parenthesis). While the number of evolved hyper-features seems to depend heavily on the number of classes of the problem, the median average size of each hyper-feature tends to be higher for the binary datasets, where a very large dispersion of values is observed.

To exemplify the variety of different sets of evolved hyper-features, we picked three examples. On the first two examples, a single hyper-feature was evolved, but with very different sizes. Both were evolved for the Gw2 dataset and obtained perfect test accuracy on the respective runs. The third example is a set of 16 hyper-features that were evolved in a run for the Ao8 dataset and obtained median test accuracy. This variety of hyper-features can be seen in Eqs. 1 through 7. Note that Bn refers to the n^{th} band of the satellite. As we can see, the M3GP algorithm can generate hyper-features that are as simple as (and perfectly equal to) the original features themselves (Eqs. 3), hyper-features that are simple enough to be interpreted (Eqs. 1, 4 and 5), and hyper-features which need to be decomposed for a proper analysis of the expression (Eqs. 2, 6 and 7).

Looking at Table 4 and the examples of hyper-features in Eqs. 1 through 7, we can state that, although the M3GP sometimes produces complex hyper-features, the general case seems to be the production of interpretable hyper-features. While this work focuses exclusively on datasets from the

Table 4. The median training and test overall accuracy, size, number of hyper-features and average size of the hyper-features obtained by the M3GP models in 30 runs in each dataset.

	Ao2	Br2	Cg2	Gw2	IM-3	IM-10	Ao8	Gw10	Mz6
Accuracy									
Training	1.000	0.992	0.993	1.000	0.996	0.932	1.000	0.988	0.620
Test	0.999	0.990	0.993	1.000	0.948	0.916	0.983	0.971	0.620
Hyper-Features									
Number	5(3-8)	8(1-15)	4(3-8)	2(1-3)	8.5(5-13)	23(17-29)	18(14-21)	21(15-23)	14.5(12-17)
Avg. Size	14(9-28)	14(6-39)	23(6-40)	11(2-43)	11(5-22)	10(8-13)	10(5-14)	8(6-12)	11(5-17)

RS domain, the same tendency regarding interpretability was already observed in the original M3GP paper [13], which used datasets from a much wider range of domains.

Gw2, Run#11, 1 Hyper-feature:

$$\frac{B5 (B3 + B5)}{B7 + B4 + 1} \quad (1)$$

Gw2, Run#26, 1 Hyper-feature:

$$\frac{B2^2 B3 B4 B5 B6 - B2^2 B4^2 B5 + B2 B3^2 B5^2 B6 - B2 B3 B4 B5^2 - B3 B4^3 B7}{B4^2 B7 (B2 B4 + B3 B5)} \quad (2)$$

Ao8, Run#18, 16 Hyper-features :

$$B3 \quad B5 \quad B6 \quad B10 \quad B11 \quad (3)$$

$$B9 - B2 \quad \frac{B3}{B5} \quad \frac{B1}{B2 B7^2} \quad B5 - B6 + B9 - 2 B10 \quad \frac{B2 B9}{B2 B11 - B10 - B2 B5} \quad (4)$$

$$\frac{(B1 + B4 - B10) (B3 + B9)}{B6} \quad B6 B9 - B1 B2 - B1 + B3 + B6 - B9 \quad (5)$$

$$\frac{B9 (B9 - B11)}{B7 (B3 B6 + B9 - B11)} \quad (B4 + B10 - \frac{B1^2}{B5 - B9}) (2 B2 + \frac{B3}{B4} - B4 + B5 + B10) \quad (6)$$

$$\frac{B7 B9^2 (B2 + B7 - B11)}{B5 B6 B11^2 (B2 + B3 - B9)} \quad \frac{B1 B2 + B1 B3 B6 + B1 B3 B9 + B3 B4 B5}{B11} \quad (7)$$

Regarding the popularity of the different satellite bands in the evolved hyper-features, Table 5 and Figure 5 show, for each band and each dataset, the fraction of hyper-features generated for that dataset (in 30 runs) that use the band. We only check whether a band appears in a hyper-feature. Measuring its importance inside the hyper-feature would be a complex exercise that we do not perform here. For each dataset, we subjectively identify a group of most popular bands as the ones ranked higher and at a larger distance from the rest (Figure 5). We do not identify any popular bands for Ao8, since on this dataset all the bands are ranked low, with very small distances between them.

In Binary Classification Datasets: The most popular band in all four datasets was the SWIR2 (B7 in both LS-7 and LS-8), which appears in 62.4% to 81.4% of the hyper-features across all datasets. This preference for the SWIR2 band is expected due to its usefulness when searching for dry earth, which may indicate a recent fire [82].

When Discriminating Similar Forest Classes: The most popular band in the IM-3 and IM-10 datasets was the NIR (B4), which appeared in 69.6% and 67.5% of the hyper-features, respectively. The popularity of the NIR band in the creation of hyper-features can be justified by its importance on the

Table 5. Fraction of hyper-features generated for each dataset (in 30 runs) that use a given band. The bands identified as popular are highlighted.

Band LS-8	Ao2	Br2	Cg2	Ao8	Gw10	Band LS-7	Gw2	IM-3	IM-10	Band S2-A	Mz6
B1	0.590	0.600	0.625	0.381	0.378	B1	0.197	0.563	0.542	B2	0.374
B2	0.565	0.564	0.708	0.359	0.427	B2	0.246	0.470	0.478	B3	0.505
B3	0.602	0.612	0.567	0.370	0.403	B3	0.492	0.466	0.551	B4	0.486
B4	0.553	0.528	0.642	0.376	0.562	B4	0.377	0.696	0.675	B5	0.422
B5	0.665	0.536	0.575	0.396	0.543	B5	0.639	0.551	0.473	B6	0.427
B6	0.528	0.552	0.717	0.419	0.483	B6	0.492	0.530	0.542	B7	0.390
B7	0.814	0.624	0.742	0.402	0.438	B7	0.705	—	—	B8	0.427
B9	—	—	—	0.325	—	—	—	—	—	B8A	0.516
B10	—	—	—	0.374	—	—	—	—	—	B11	0.397
B11	—	—	—	0.351	—	—	—	—	—	B12	0.349

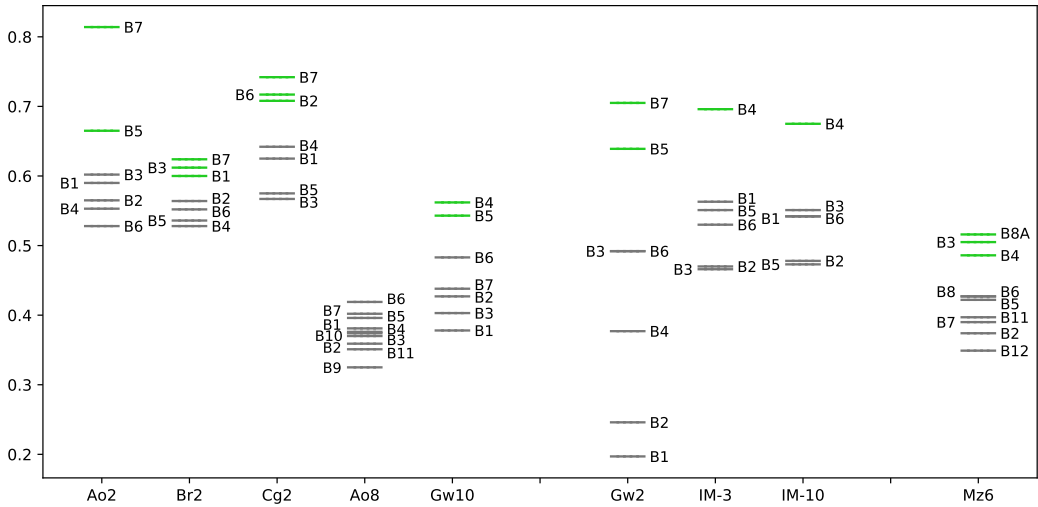


Figure 5. Fraction of hyper-features generated for each dataset (in 30 runs) that use a given band. The bands identified as popular are highlighted.

visualisation of healthy vegetation in the discrimination between *Dense Forest* and *Open Forest* pixels. The importance of this band has also led to the creation of indices, such as the NDVI [1,82].

When Discriminating All Classes: Since the target classes are not similar to each other, here we observe which are the most popular bands on each dataset, and attempt to explain their popularity based on which classes benefited the most from the hyper-features. We do not discuss the Ao8 dataset, not only because it lacks popular bands, but also because it did not benefit from the hyper-features, as we will see below in Section 4.4.

In the Gw10 dataset, the Red (B4) and NIR (B5) bands were the most popular, appearing in 56.2% and 54.3% of the hyper-features, respectively. We will see in the next section that on this dataset the hyper-features improved the classification of the *Mangrove*, *Savanna Woodland* and *Wetland* pixels. This suggests that the better discrimination of these land cover types took into account the amount of healthy vegetation and the composition of the soil.

Regarding the Mz6 dataset, the Vegetation Red Edge (B8A), Green (B3) and Red (B4) were identified as the most popular bands, appearing in 51.6%, 50.5% and 48.6% of the hyper-features, respectively. However, their effect is not so clear, since the improvement brought by the hyper-features affected several different classes. Taking into consideration only the classes with the highest improvement, which were *Agriculture / Bare Soil*, *Forest* and *Wetlands*, we suggest that the better discrimination of these land cover types considered the health and age of the vegetation, as well as the composition of the soil.

Table 6. Comparison of the median overall test accuracy obtained by the three ML algorithms in the original datasets, when adding indices, and when using hyper-features evolved by EFS, FFX and M3GP. The coloured *p*-values indicate significantly *better*/*worse* results.

Dataset	Decision Trees				Random Forests				XGBoost			
	Ao2	Br2	Cd2	Gw2	Ao2	Br2	Cd2	Gw2	Ao2	Br2	Cd2	Gw2
Orig. Dataset												
Test Accuracy	0.999	0.989	0.996	0.999	0.999	0.992	0.999	1.000	0.999	0.993	0.998	0.999
Indices												
Test Accuracy	0.999	0.990	0.996	0.999	0.999	0.993	0.999	1.000	0.999	0.993	0.998	0.999
<i>p</i> -value vs Orig.	0.694	0.678	0.909	0.669	0.675	0.917	0.436	0.871	1.000	1.000	1.000	1.000
EFS												
Test Accuracy	0.998	0.989	0.996	0.999	1.000	0.992	0.999	1.000	0.999	0.993	0.998	0.999
<i>p</i> -value vs Orig.	0.012	0.226	0.143	0.735	0.091	0.777	0.137	0.619	0.256	0.682	0.489	0.127
FFX												
Test Accuracy	0.999	0.990	0.996	1.000	0.999	0.993	0.999	1.000	0.999	0.993	0.998	1.000
<i>p</i> -value vs Orig.	0.224	0.941	0.294	0.000	0.363	0.988	0.148	0.730	0.739	0.794	0.886	0.000
M3GP												
Test Accuracy	0.999	0.990	0.996	0.999	0.999	0.992	0.999	1.000	0.999	0.993	0.998	0.999
<i>p</i> -value vs Orig.	0.908	0.947	0.672	0.813	0.500	0.846	0.688	0.871	1.000	1.000	1.000	1.000
<i>p</i> -value vs Ind.	0.782	0.761	0.598	0.849	0.780	0.982	0.738	1.000	1.000	1.000	1.000	1.000
<i>p</i> -value vs FFX	0.276	0.830	0.525	0.001	0.095	0.905	0.319	0.868	0.739	0.794	0.886	0.000
<i>p</i> -value vs EFS	0.017	0.272	0.291	0.572	0.286	0.682	0.309	0.757	0.256	0.682	0.489	0.127

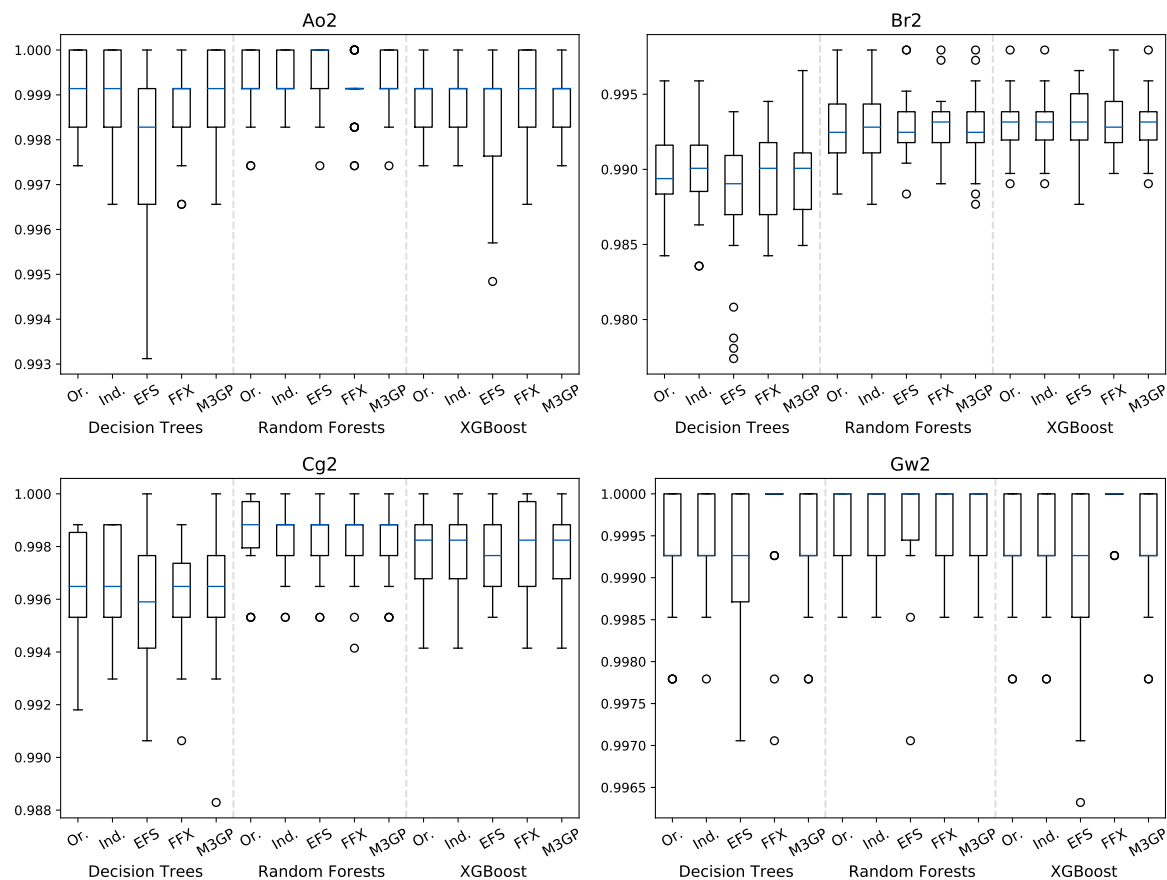
In some of the related work regarding the use of GP to build hyper-features in RS, the authors reveal what were considered the best hyper-features obtained. For comparison with our own, here we also comment on those hyper-features. It is important to say that these works address regression problems (rather than classification problems), which sometimes require more complex models in order to be solved. In comparison with our own, these works use an extensive list of mathematical operators to combine the original features (which also tends to cause the creation of larger models). The authors also include indices in the datasets, similarly to what is done in part of our work.

In [65], the authors use GP to monitor the quality of the water in reservoirs, by predicting the amount of chlorophyll in the water. The final model is quite simple (having a size of 8), according to our criteria, and only uses the Green, Red and NIR bands. While in this case the hyper-feature used is simple, that is not the case in the other two works. In [66], the authors attempt to predict the soil salinity. Their final model uses one band and five indices, and its size is near 50 (making it larger than any of our hyper-features). In [62], the authors attempt to predict the soil moisture and, although their final model only uses four terminals (SAR backscatter coefficient, slope, soil permeability (in/hr) and the NDVI), this model is the most complex out of these three.

4.2. Hyper-features in Binary Classification Datasets

The results obtained on the binary classification datasets (Ao2, Br2, Cd2, Gw2) are reported in Table 6 and Figure 6. In terms of training accuracy, the three classification algorithms managed to obtain perfect results in nearly every run, and therefore those results are not included in the table. In terms of test accuracy, the induced models achieved very high values, nearly all above 99% (both in terms of overall accuracy and class accuracy, as seen in Table 7), also on the original (non-expanded) datasets. The lowest results belong to DT that, when applied to the Br2 dataset, achieved a median overall test accuracy of 98.9%. Without much room for improvement, FFX was still able to create hyper-features that improved the test accuracy in two cases (DT and XGB in the Gw2 dataset), surpassing also the M3GP hyper-features, while neither the indices nor the M3GP or EFS hyper-features caused any significant difference in the results. The boxplots show a very low dispersion of accuracy values (the ranges of the y-axes are very limited), which seems to be marginally larger for the EFS results.

These seemingly uninteresting results agree with the findings of our previous work [11]. There, using a different method of selecting and using the hyper-features also had no effect in the cases where

Figure 6. Boxplots of the test accuracy obtained in the binary classification datasets in each test case.

the training and test datasets came from the same image. However, the hyper-features revealed to be beneficial when the induced models were applied to datasets that came from images not seen during training. This suggests that the current method of obtaining and using the hyper-features may also prove beneficial in a similar training and test setting.

4.3. Hyper-features to Discriminate Similar Classes in a Multiclass Classification Dataset

Before looking at these results, it is worth recalling that the IM-3 dataset was built from three similar classes within the IM-10 dataset. As such, even though it has a reduced number of classes, it is not unexpected to see a lower accuracy in this dataset. It is also worth specifying that the hyper-features used in the IM-10 dataset were obtained only in the IM-3 dataset, in an attempt to help discriminate these similar classes. Finally, we also recall that EFS and FFX are not used in the multiclass datasets.

The results for IM-10 and IM-3 are reported in Table 8 and Figure 7. Once again, the training results are omitted from the table because all three algorithms achieved perfect results in nearly every run. In terms of test accuracy, we observe that, although the values are high, they have a larger margin for improvement when compared to the binary classification results reported above. When adding indices to the original dataset, the test accuracy on the IM-10 dataset increased with two algorithms (RF and XGB). When adding the hyper-features evolved by M3GP, the test accuracy in the IM-10 dataset increased with all three algorithms, and in the IM-3 dataset it increased with the XGB algorithm. Neither the indices nor the M3GP hyper-features degraded the test accuracy. When comparing the performance of indices versus M3GP hyper-features, M3GP is better with two algorithms (DT and XGB). In the boxplots, we observe that IM-3 has a larger dispersion of values than IM-10 (notice the different y-axes ranges). On IM-10, the DT algorithm visibly falls behind RF and XGB.

Table 7. Average test accuracy in each class when using the XGBoost algorithm in the original datasets.

XGB - Original	Ao2	Br2	Cg2	Gw2	IM-3	IM-10	Gw10	Ao8	Mz6
Agriculture / Bare Soil	—	—	—	—	—	98.96%	96.34%	83.02%	72.59%
Burnt	99.81%	99.51%	99.63%	99.80%	—	93.19%	98.72%	99.22%	—
Clouds	—	—	—	—	—	—	—	100.00%	—
Forest	—	—	—	—	—	—	—	99.87%	88.05%
- Dense Forest	—	—	—	—	92.67%	87.67%	79.81%	—	—
- Open Forest	—	—	—	—	93.02%	93.65%	98.41%	—	—
Grassland	—	—	—	—	—	92.73%	81.67%	85.56%	64.02%
Mangrove	—	—	—	—	—	99.12%	98.41%	—	—
Mud	—	—	—	—	—	96.07%	—	—	—
Sand	—	—	—	—	—	95.78%	90.22%	—	—
Savanna Woodland	—	—	—	—	99.71%	84.41%	98.66%	98.66%	—
Urban	—	—	—	—	—	—	—	87.78%	56.82%
Water	—	—	—	—	—	97.33%	99.48%	99.48%	—
Wetland	—	—	—	—	—	—	95.29%	—	80.34%
Other	99.97%	99.16%	99.88%	99.98%	—	—	—	—	76.07%

Table 8. Comparison of the median overall test accuracy obtained by the three ML algorithms in the original datasets, when adding indices, and when adding hyper-features evolved by the M3GP algorithm. The coloured *p*-values indicate significantly **better** results.

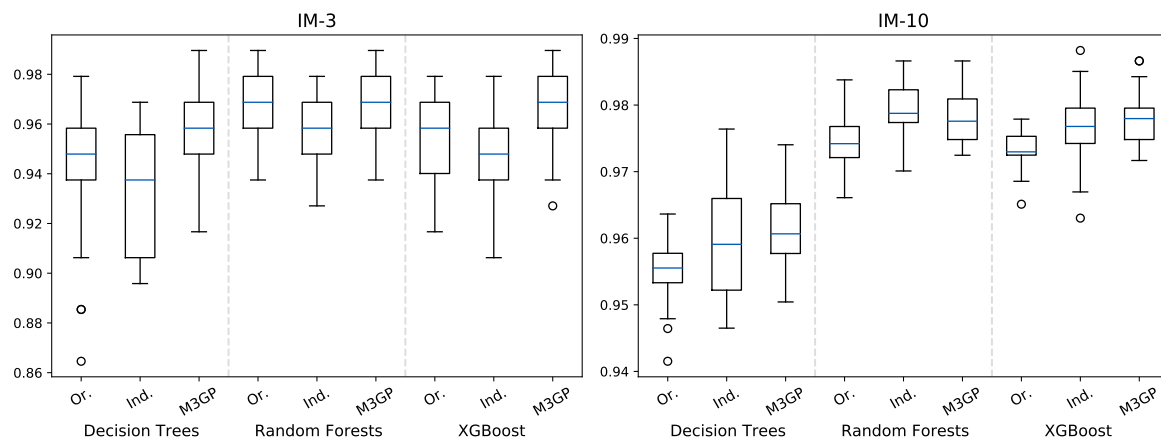
Dataset	Decision Trees		Random Forests		XGBoost	
	IM-3	IM-10	IM-3	IM-10	IM-3	IM-10
Original Dataset						
Test Accuracy	0.948	0.956	0.969	0.974	0.958	0.973
Indices						
Test Accuracy	0.938	0.959	0.958	0.979	0.948	0.977
<i>p</i> -value vs Original	0.151	0.051	0.062	0.000	0.178	0.001
M3GP						
Test Accuracy	0.958	0.961	0.969	0.978	0.969	0.978
<i>p</i> -value vs Original	0.020	0.000	0.844	0.000	0.009	0.000
<i>p</i> -value vs Indices	0.000	0.407	0.055	0.218	0.000	0.711

In terms of class accuracy in the IM-3 dataset, we can see in Table 9 that the hyper-features improved the discrimination between the *Dense Forest* and the *Open Forest* pixels, at the cost of reducing the accuracy on the *Savanna Woodland* class, by increasing its confusion with *Open Forest*. Looking at Table 7, we see that *Savanna Woodland* had almost perfect accuracy, and therefore any changes on this class would certainly be for the worse. In the end, the three classes became more balanced in terms of accuracy. Although these hyper-features do not seem to be helpful in the classification of *Savanna Woodland* on the IM-3 dataset, when applied to the IM-10 dataset (See Table 10) their largest impact is precisely in this class, by correcting pixels that were previously misclassified as *Agriculture/Bare soil*, *Grassland* and *Mangrove*. Their second biggest impact is in the classification of *Dense Forest*, by improving its discrimination from *Open Forest* and by correcting pixels that were previously misclassified as *Mangrove*.

In these two datasets, although the M3GP hyper-features performed better than the indices, these were also clearly beneficial when added to the original datasets. This behaviour was similar to all three classification algorithms and, as such, we only displayed the confusion matrices related to the XGBoost algorithm, which had the best results. It is worth noticing that the IM-3 and IM-10 datasets were extracted from a set of satellite images with different acquisition dates. Next, we will observe additional evidence that indices and hyper-features seem to be more useful in datasets coming from sets of images with different acquisition dates.

4.4. Hyper-features to Discriminate All Classes in Multiclass Classification Datasets

The results obtained on the three unrelated multiclass classification datasets (Ao8, Gw10, Mz6) are reported in Table 11 and Figure 8. Once again, the training results were omitted from the table, as

Figure 7. Boxplots of the test accuracy obtained in the IM-3 and IM-10 datasets in each test case.**Table 9.** Confusion matrix comparing the average test accuracy obtained by the XGBoost algorithm with and without hyper-features in the IM-3 dataset. This table shows the difference in the percentage of pixels in each line. Improvements are in **green** and deteriorations in **red**. The rightmost column indicates the accuracy obtained in each class without using hyper-features.

XGB IM-3	Dense Forest	Open Forest	Savanna Woodland	Original Accuracy
Dense Forest	2.83%	-2.83%	0.00%	92.67%
Open Forest	-3.33%	2.46%	0.87%	93.02%
Savanna Woodland	0.00%	1.57%	-1.57%	99.71%

perfect accuracy was achieved in almost every run. However, for the Mz6 dataset, XGB required a maximum tree depth larger than the implementation default in order to achieve it (see Section 3.5).

In terms of test accuracy, the indices improved the accuracy in two test cases (DT on Gw10, XGB on Mz6) and reduced the accuracy in one test case (RF on Mz6). On the other hand, the hyper-features evolved by M3GP improved the test accuracy in five test cases (Mz6 with all algorithms, and Gw10 with RF and XGB), when comparing the results with those on the original dataset, and in four test cases, when comparing with the results obtained with the indices (Mz6 with all algorithms, Gw10 with XGB). Once again, the hyper-features evolved by M3GP did not lead to a degradation of the test accuracy in any of the cases.

Both the indices and the M3GP-evolved hyper-features had an impact on the Gw10 and Mz6 datasets, which were obtained from a set of satellite images with different acquisition dates. Neither the indices nor the hyper-features had an impact on the Ao8 dataset, which was obtained from two images with the same acquisition date. These results, together with those displayed previously, seem to indicate that both the indices and the hyper-features are particularly useful in datasets obtained by mixing satellite images with different acquisition dates.

On the boxplots, once again we observe that DT falls behind RF and XGB, and completely struggles on the Mz6 problem.

In terms of class accuracy in the Gw10 dataset (Tables 12 and 13), when using the hyper-features with the DT algorithm, the hyper-features are particularly useful in the classification of *Grassland*, by correcting pixels that were misclassified as *Savanna Woodland* (although some of those are now misclassified as *Mangrove*); in the classification of *Dense Forest*, by correcting pixels that were misclassified as *Open Forest* (although some of those are now misclassified also as *Mangrove*); and in the classification of *Wetlands*, by correcting pixels previously misclassified as *Mangrove*. When using the XGBoost algorithm, the improvements are more general across the classes, with the exception of

Table 10. Confusion matrix comparing the average test accuracy obtained by the XGBoost algorithm with and without hyper-features in the IM-10 dataset. This table shows the difference in the percentage of pixels in each line. Improvements are in **green** and deteriorations in **red**. Only the 20 cells with the highest impact are coloured. The rightmost column indicates the accuracy obtained in each class without using hyper-features.

XGB IM-10	Water	Burnt	Sand	Agriculture / Bare soil	Open Forest	Dense Forest	Grassland	Mangrove	Savanna Woodland	Mud	Original Accuracy
Water	0.00%	-0.05%	0.00%	0.02%	0.00%	0.00%	0.00%	-0.01%	0.00%	0.04%	97.33%
Burnt	0.87%	0.29%	0.00%	-0.29%	0.00%	0.00%	0.14%	-0.87%	-0.58%	0.43%	93.19%
Sand	0.00%	0.00%	0.95%	-0.95%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	95.78%
Agriculture / Bare soil	0.00%	0.00%	-0.15%	0.19%	0.00%	0.00%	-0.01%	-0.18%	0.13%	0.02%	98.96%
Open Forest	0.00%	0.00%	0.00%	-0.08%	1.03%	-1.83%	-0.08%	0.79%	0.16%	0.00%	93.65%
Dense Forest	0.00%	0.00%	0.00%	0.00%	-2.50%	4.33%	0.00%	-1.83%	0.00%	0.00%	87.67%
Grassland	-0.15%	0.15%	0.00%	-0.15%	-0.45%	0.00%	0.15%	0.61%	-0.15%	0.00%	92.73%
Mangrove	-0.15%	0.00%	0.00%	-0.03%	-0.03%	0.00%	0.00%	0.23%	-0.04%	0.01%	99.12%
Savanna Woodland	0.00%	-0.29%	0.00%	-2.55%	0.10%	0.00%	-1.27%	-1.18%	5.49%	-0.29%	84.41%
Mud	-1.29%	0.09%	0.00%	-0.02%	0.00%	0.00%	0.00%	0.07%	0.00%	1.16%	96.07%

Table 11. Comparison of the overall test accuracy obtained by the three ML algorithms in the original datasets, when adding indices, and when adding hyper-features evolved by the M3GP algorithm. The coloured *p*-values indicate significantly **better**/**worse** results.

Dataset	Decision Trees			Random Forests			XGBoost		
	Ao8	Gw10	Mz6	Ao8	Gw10	Mz6	Ao8	Gw10	Mz6
Original Dataset									
Test Accuracy	0.977	0.964	0.662	0.988	0.981	0.773	0.985	0.979	0.780
Indices									
Test Accuracy	0.978	0.968	0.662	0.989	0.980	0.769	0.986	0.980	0.781
<i>p</i> -value vs Original	0.291	0.000	0.371	0.213	0.824	0.000	0.645	0.335	0.003
M3GP									
Test Accuracy	0.980	0.970	0.665	0.988	0.982	0.775	0.987	0.983	0.786
<i>p</i> -value vs Original	0.125	0.000	0.000	0.923	0.054	0.006	0.389	0.000	0.000
<i>p</i> -value vs Indices	0.693	0.847	0.000	0.228	0.038	0.000	0.650	0.002	0.000

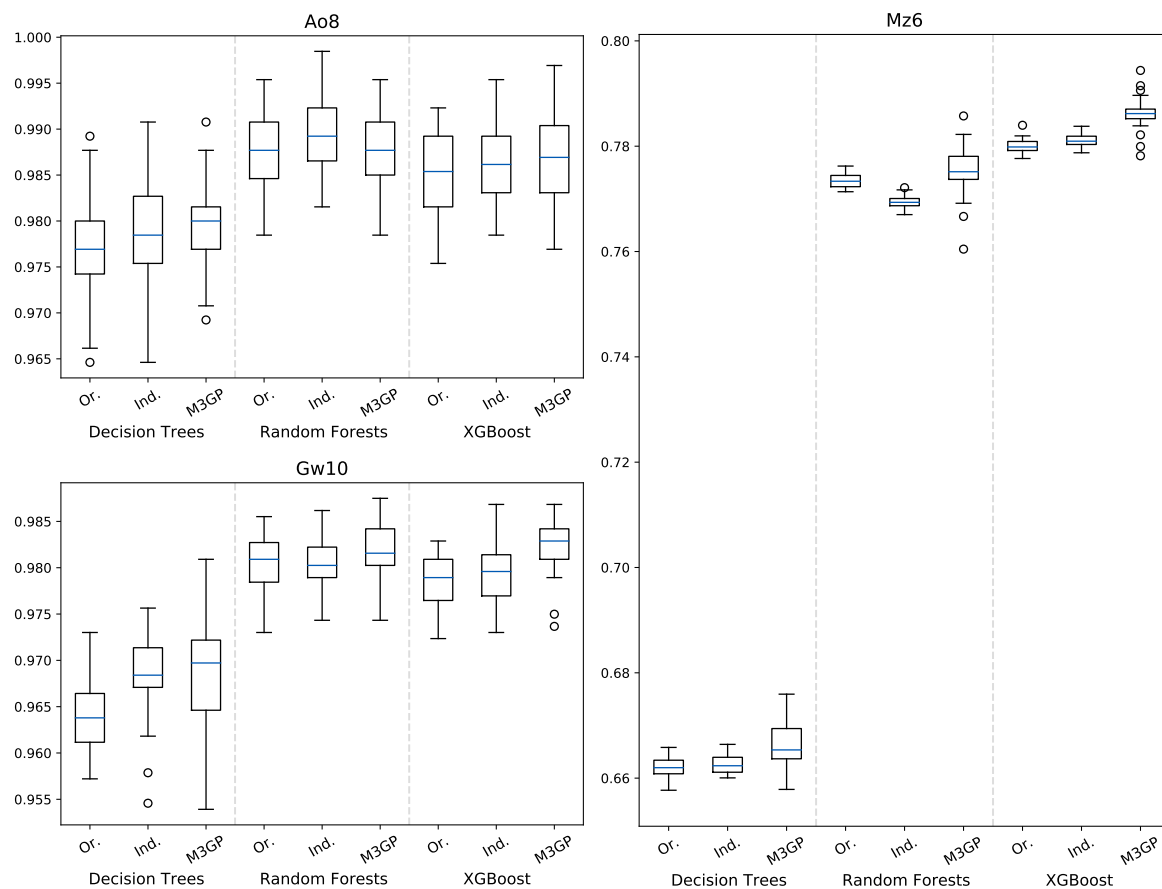
the *Grassland* pixels, which are now misclassified as *Savanna Woodland*, and the *Dense Forest* pixels, which were previously misclassified as *Open Forest*.

In this case, we omitted the results regarding the RF classifier since there was no statistically significant difference between the runs in the original and the extended datasets. We present the results of both the DT and XGBoost classifier to show that the same hyper-features can have different effects on two classifiers. In this case, they improved the DT accuracy primarily in three classes, while the improvements on the XGBoost accuracy were general.

Regarding class accuracy in the Mz6 case (Tables 14 and 15), in addition to showing the improvements when using hyper-features with the XGB algorithm, we show the degradation obtained by using indices in the RF algorithm. The idea is not to say that the indices are bad, but to show that adding more hyper-features will not necessarily bring an improvement. When using the hyper-features in the XGBoost algorithm, the improvement was general among all classes, with a higher impact on the *Urban* pixels that were previously misclassified as *Other*. The improvements in this class's pixels can be easily justified by being the class with the lowest accuracy in the original dataset, followed by *Grassland*, which was also improved. When using indices in the RF algorithm, the class accuracy was degraded, in particular in the classification of *Agriculture / Bare Soil*, *Grassland* and *Wetland*. When misclassified, these pixels tend to be classified as either *Agriculture / Bare Soil*, *Forest* or *Grassland*.

4.5. Impact on the MRV Performance

When training hyper-features to discriminate multiple classes, the results indicate an overall improvement, particularly in the classes that previously had a lower accuracy. The improvements are

Figure 8. Boxplots of the test accuracy obtained in the Ao8, Gw10, and Mz6 datasets in each test case.

significant in the IM-3, IM-10, Gw10 and Mz6 datasets. These datasets have one thing in common: they consist of mosaics derived from several acquisition dates. On the contrary, in the Ao8 dataset, where the two images of the mosaic are from the same day, there are no significant improvements. These results, together with those previously obtained in [11], indicate that both indices and hyper-features are more useful when training models in images with more than one acquisition date (or from different locations).

Monitoring forest land cover at country level, in compliance with UNFCCC standards, is a challenging endeavor, especially for vast countries covering various ecosystem types with distinct seasonality. The production of wall-to-wall maps derived from satellite imagery is especially attractive in these cases because remote sensing can cover large extents, greatly reducing costs, improving consistency, and increasing the periodicity of observations. However, in such cases, the image mosaics required to produce good quality maps are likely to include many different acquisition dates, maximizing image quality and observation date adequacy regarding vegetation cycles and climatic conditions. Thus, considering the results obtained, with hyper-features improving both the discrimination of analogous wooded vegetation classes and the classification accuracy of large image mosaics, it can be ascertained that the methods presented merit further development to exploit improvements in remote sensing based MRV performance.

Table 12. Confusion matrix comparing the average test accuracy obtained by the DT algorithm with and without hyper-features in the Gw10 dataset. This table shows the difference in the percentage of pixels in each line. Improvements are in **green** and deteriorations in **red**. Only the 20 cells with the highest impact are coloured. The rightmost column indicates the accuracy obtained in each class without using hyper-features.

DT Gw10	Agriculture /Bare soil	Burnt	Dense Forest	Grassland	Mangrove	Open Forest	Sand	Savanna Woodland	Water	Wetland	Original Accuracy
Agriculture /Bare soil	0.25%	0.00%	0.00%	0.02%	0.05%	0.00%	-0.22%	-0.27%	0.00%	0.17%	95.32%
Burnt	0.07%	0.64%	0.00%	0.00%	0.00%	0.00%	0.00%	-0.14%	0.00%	-0.57%	98.65%
Dense Forest	0.00%	0.00%	2.22%	0.00%	1.85%	-4.07%	0.00%	0.00%	0.00%	0.00%	71.85%
Grassland	0.00%	0.00%	0.00%	5.83%	1.67%	0.00%	0.00%	-7.50%	0.00%	0.00%	70.83%
Mangrove	0.02%	-0.06%	0.00%	-0.02%	0.20%	0.04%	0.00%	-0.02%	0.02%	-0.18%	97.28%
Open Forest	0.00%	0.00%	0.10%	0.02%	-0.03%	0.14%	0.00%	-0.17%	-0.02%	-0.03%	96.60%
Sand	0.22%	0.00%	0.00%	0.00%	0.00%	0.00%	-0.22%	0.00%	0.00%	0.00%	87.78%
Savanna Woodland	-0.26%	0.02%	-0.03%	0.11%	-0.31%	0.02%	0.00%	0.68%	0.00%	-0.22%	97.16%
Water	0.00%	-0.04%	0.00%	0.00%	0.11%	-0.02%	0.00%	0.00%	-0.02%	-0.04%	99.46%
Wetland	0.23%	-0.14%	0.00%	0.11%	-1.49%	-0.03%	0.00%	-0.34%	0.06%	1.61%	92.30%

Table 13. Confusion matrix comparing the average test accuracy obtained by the XGBoost algorithm with and without hyper-features in the Gw10 dataset. This table shows the difference in the percentage of pixels in each line. Improvements are in **green** and deteriorations in **red**. Only the 20 cells with the highest impact are coloured. The rightmost column indicates the accuracy obtained in each class without using hyper-features.

XGB Gw10	Agriculture /Bare soil	Burnt	Dense Forest	Grassland	Mangrove	Open Forest	Sand	Savanna Woodland	Water	Wetland	Original Accuracy
Agriculture /Bare soil	0.77%	0.00%	0.00%	0.00%	0.00%	0.00%	-0.62%	-0.17%	0.00%	0.02%	96.34%
Burnt	-0.14%	0.28%	0.00%	0.00%	-0.43%	-0.07%	0.00%	0.21%	0.00%	0.14%	98.72%
Dense Forest	0.00%	0.00%	0.93%	0.00%	0.93%	-1.85%	0.00%	0.00%	0.00%	0.00%	79.81%
Grassland	0.00%	0.00%	0.00%	-2.50%	0.00%	0.00%	0.00%	2.50%	0.00%	0.00%	81.67%
Mangrove	0.00%	0.00%	-0.03%	0.00%	0.40%	-0.02%	0.00%	-0.08%	-0.08%	-0.19%	98.41%
Open Forest	0.00%	0.00%	-0.19%	0.00%	-0.03%	0.31%	0.00%	-0.09%	0.00%	0.00%	98.41%
Sand	0.67%	0.00%	0.00%	0.00%	0.00%	0.00%	-0.67%	0.00%	0.00%	0.00%	90.22%
Savanna Woodland	-0.09%	0.00%	0.00%	0.00%	-0.02%	-0.14%	0.00%	0.29%	0.00%	-0.04%	98.66%
Water	0.00%	0.00%	0.00%	0.00%	-0.05%	0.00%	0.00%	0.00%	0.16%	-0.11%	99.48%
Wetland	-0.11%	0.06%	0.00%	0.00%	-0.34%	0.00%	0.00%	0.00%	0.06%	0.34%	95.29%

Table 14. Confusion matrix comparing the average test accuracy obtained by the RF algorithm with and without indices in the Mz6 dataset. This table shows the difference in the percentage of pixels in each line. Deteriorations are shown in **red**. The 10 cells with the highest impact are coloured. The rightmost column indicates the accuracy obtained in each class without using indices.

RF - Mz6	Agriculture /Bare soil	Forest	Grassland	Urban	Wetland	Other	Original Accuracy
Agriculture /Bare soil	-0.92%	0.44%	0.39%	-0.01%	0.02%	0.07%	71.79%
Forest	-0.07%	-0.01%	0.03%	0.00%	-0.01%	0.05%	89.06%
Grassland	0.33%	0.31%	-0.90%	-0.01%	0.31%	-0.04%	60.58%
Urban	-0.29%	0.03%	-0.16%	-0.14%	0.06%	0.50%	50.91%
Wetland	0.14%	0.25%	0.15%	0.02%	-0.45%	-0.09%	79.72%
Other	0.19%	-0.09%	0.32%	-0.04%	-0.02%	-0.36%	75.35%

Table 15. Confusion matrix comparing the average test accuracy obtained by the XGBoost algorithm with and without hyper-features in the Mz6 dataset. This table shows the difference in the percentage of pixels in each line. Improvements are shown in **green**. The 10 cells with the highest impact are coloured. The rightmost column indicates the accuracy obtained in each class without using hyper-features.

XGB - Mz6	Agriculture /Bare soil	Forest	Grassland	Urban	Wetland	Other	Original Accuracy
Agriculture /Bare soil	0.87%	-0.24%	-0.28%	-0.01%	-0.01%	-0.32%	72.59%
Forest	-0.00%	0.34%	-0.29%	0.00%	-0.04%	-0.01%	88.05%
Grassland	-0.26%	-0.32%	0.51%	0.00%	0.14%	-0.08%	64.02%
Urban	-0.03%	-0.02%	-0.25%	1.82%	-0.06%	-1.48%	56.82%
Wetland	-0.15%	-0.19%	-0.39%	0.02%	0.84%	-0.12%	80.34%
Other	-0.26%	-0.04%	-0.18%	-0.04%	-0.08%	0.60%	76.07%

5. Conclusions

We performed Feature Construction using M3GP, a variant of the standard Genetic Programming algorithm, with the goal of improving the performance of several Machine Learning algorithms by adding the new hyper-features to the reference datasets. We tested the approach in the tasks of binary classification of burnt areas and multiclass classification of land cover types. The datasets used were obtained from Landsat-7, Landsat-8 and Sentinel-2A satellite images over the countries of Angola, Brazil, Democratic Republic of Congo, Guinea-Bissau, and Mozambique.

The hyper-features produced by the M3GP algorithm, although variable in number and size, were generally not very complex, and considered to be quite interpretable. While a larger number of hyper-features were created on the multiclass classification problems, a higher dispersion of sizes was observed on the binary problems. Regarding the popularity of each satellite band in the binary and multiclass classification problems, the models frequently used the SWIR2 band when trying to detect burnt areas in the binary datasets. On the multiclass classification datasets, the models seemed to have a preference for the Vegetation Red Edge, NIR, Red, and Green bands when training hyper-features to discriminate different forest classes or when the hardest classes included vegetation (e.g., *Agriculture / Bare Soil* and *Forest*), and in some cases, also water (e.g., *Mangroves* and *Wetlands*).

The performance of Decision Trees, Random Forests and XGBoost was assessed on the original datasets and on the datasets expanded with the evolved hyper-features, and the results compared for statistical significance. For comparison purposes, we also assessed the performance of the same algorithms on all datasets expanded with the well-known spectral indices NDVI, NDWI and NBR, and on the binary datasets expanded with hyper-features created by the FFX and EFS Feature Construction algorithms. On the binary classification problems, we conclude that neither of the four alternatives (M3GP, indices, FFX, EFS) leads to substantial improvements. Only FFX was able to improve the results in 2 out of 12 test cases (both on the same dataset). On the multiclass classification problems, the hyper-features evolved by the M3GP caused significant improvements in 9 out of 15 test cases, with no degradation of results in any test case, while the indices caused significant improvements in 4 out of 15 test cases and significant degradation of results in one test case. The approach appears to be equally beneficial to all three Machine Learning algorithms.

Overall, both hyper-features and indices displayed the capability of improving the robustness of the machine learning models in multiclass classification datasets. However, this improvement seems to exist only in datasets built from collections of images with several acquisition dates, which indicates that both hyper-features and indices can be robust to the radiometric variations across images and can be used to improve the MRV performance of mechanisms such as REDD+.

Although the hyper-features have the advantage of being created automatically with specific goals, such as the discrimination of specific classes, there is a computational cost associated with this task. Taking this into consideration, one of our objectives for future work is to continue the validation of the efficacy of the hyper-features in the discrimination of similar classes and their robustness to the radiometric variations across different satellite images. We hope to be able to create reusable hyper-features, thus reducing the computational cost of generating them frequently. Besides this validation, we also want to expand this work into regression problems, such as the estimation of biomass from satellite images.

Author Contributions: Conceptualization, J.B. and S.S.; methodology, J.B.; software, J.B.; validation, S.S.; formal analysis, J.B.; investigation, J.B.; resources, S.S.; data curation, A.C.; writing—original draft preparation, J.B. and L.V.; writing—review and editing, A.C., J.B., M.V. and S.S.; visualization, A.C. and J.B.; supervision, S.S.; project administration, S.S.; funding acquisition, L.V., M.V. and S.S. All authors have read and agreed to the published version of the manuscript.

Funding: This work was partially supported by FCT through funding of LASIGE Research Unit (UIDB/00408/2020 and UIDP/00408/2020) and CEF (UIDB/00239/2020); projects BINDER (PTDC/CCI-INF/29168/2017), OPTOX (PTDC/CTA-AMB/30056/2017), PREDICT (PTDC/CCI-CIF/29877/2017), INTERPHENO (PTDC/ASP-PLA/28726/2017), GADgET (DSAIPA/DS/0022/2018), AICE (DSAIPA/DS/0113/2019); PhD Grant (SFRH/BD/143972/2019).

Conflicts of Interest: The authors declare no conflict of interest. The funders had no role in the design of the study; in the collection, analyses, or interpretation of data; in the writing of the manuscript, or in the decision to publish the results.

Abbreviations

The following abbreviations are used in this manuscript:

Af	Equatorial rainforest, fully humid
Am	Equatorial Monsoon
Ao	Angola
Aw	Equatorial savanna with dry winter
Bx	Band x
Br	Brazil
BSh	Hot semi-arid
Bwh	Hot desert
CCDC	Continuous Change Detection and Classification
Cd	Democratic Republic of the Congo
Cwa	Warm temperate climate with dry winter and hot summer
Cwb	Warm temperate climate with dry winter and warm summer
DT	Decision Tree
EC	Evolutionary Computation
EFS	Evolutionary Feature Synthesis (algorithm)
FFX	Fast Function Extraction (algorithm)
GLCM	Gray Level Co-occurrence Matrix
Gw	Guinea-Bissau
GP	Genetic Programming
KGCS	Köpper-Geiger Classification System
LS-7	Landsat 7
LS-8	Landsat 8
M3GP	Multidimensional Multiclass GP with Multidimensional Populations (algorithm or classifier)
MD	Mahalanobis Distance (classifier)
ML	Machine Learning
MRV	Measure, Report and Verify
Mz	Mozambique
NBR	Normalized Burn Ratio
NDVI	Normalized Difference Vegetation Index
NDWI	Normalized Difference Water Index
PCA	Principal Component Analysis
REDD+	Reducing Emissions from Deforestation and forest Degradation
RF	Random Forest
RS	Remote Sensing
S-2A	Sentinel-2A
UNFCCC	United Nations Framework Convention on Climate Change
XGB	XGBoost
WAF	Weighted Average of F-measures

References

1. Herring, J.A. Measuring Vegetation (NDVI EVI) : Feature Articles. 2000.
2. McFEETERS, S.K. The use of the Normalized Difference Water Index (NDWI) in the delineation of open water features. *International Journal of Remote Sensing* **1996**, *17*, 1425–1432. doi:10.1080/01431169608948714.
3. Key, C.; Benson, N., Landscape Assessment: Ground measure of severity, the Composite Burn Index; and Remote sensing of severity, the Normalized Burn Ratio.; 2006; pp. LA 1–51.
4. Jinru, X.; Su, B. Significant Remote Sensing Vegetation Indices: A Review of Developments and Applications. *Journal of Sensors* **2017**, *2017*, 1–17.
5. Maggiori, E.; Tarabalka, Y.; Charpiat, G.; Alliez, P. Convolutional Neural Networks for Large-Scale Remote Sensing Image Classification. *IEEE T. Geosci. Remote Sens.* **2016**, *55*. doi:10.1109/TGRS.2016.2612821.
6. Ribeiro, F.; Roberts, D.; Hess, L.; Davis, F.; Caylor, K.; Daldegan, G. Geographic Object-Based Image Analysis Framework for Mapping Vegetation Physiognomic Types at Fine Scales in Neotropical Savannas. *Remote Sensing* **2020**, *12*, 1721. doi:10.3390/rs12111721.

7. Dragozi, E.; Gitas, I.; Stavrakoudis, D.; Theocharis, J. Burned area mapping using support vector machines and the FuzCoC feature selection method on VHR IKONOS imagery. *Remote Sensing MDPI* **2014**, *Remote Sens.* **2014**, 12005–12036. doi:10.3390/rs61212005.
8. Solano Correa, Y.; Bovolo, F.; Bruzzone, L. A Semi-Supervised Crop-Type Classification Based on Sentinel-2 NDVI Satellite Image Time Series And Phenological Parameters. 2019. doi:10.1109/IGARSS.2019.8897922.
9. Orynbaikyzy, A.; Gessner, U.; Mack, B.; Conrad, C. Crop Type Classification Using Fusion of Sentinel-1 and Sentinel-2 Data: Assessing the Impact of Feature Selection, Optical Data Availability, and Parcel Sizes on the Accuracies. *Remote Sensing* **2020**, *12*, 2779. doi:10.3390/rs12172779.
10. Carrao, H.; Gonçalves, P.; Caetano, M. Contribution of multispectral and multitemporal information from MODIS images to land cover classification. *Remote Sensing of Environment* **2008**, *112*, 986–997. doi:10.1016/j.rse.2007.07.002.
11. Batista, J.E.; Silva, S. Improving the Detection of Burnt Areas in Remote Sensing using Hyper-features Evolved by M3GP. 2020 IEEE Congress on Evolutionary Computation (CEC). IEEE, 2020. doi:10.1109/cec48606.2020.9185630.
12. Poli, R.; B. Langdon, W.; McPhee, N. *A Field Guide to Genetic Programming*; 2008.
13. Muñoz, L.; Silva, S.; Trujillo, L. M3GP – Multiclass Classification with GP. 2015, Vol. 9025, pp. 78–91.
14. Muñoz, L.; Trujillo, L.; Silva, S.; Castelli, M.; Vanneschi, L. Evolving multidimensional transformations for symbolic regression with M3GP. *Memetic Comput.* **2019**, *11*, 111–126. doi:10.1007/s12293-018-0274-5.
15. Muñoz, L.; Trujillo, L.; Silva, S. Transfer learning in constructive induction with Genetic Programming. *Genetic Programming and Evolvable Machines* **2019**, pp. 1–41.
16. Bastarrika, A.; Chuvieco, E.; Martín, M.P. Mapping burned areas from Landsat TM/ETM+ data with a two-phase algorithm: Balancing omission and commission errors. *Remote Sensing of Environment* **2011**, *115*, 1003 – 1012. doi:https://doi.org/10.1016/j.rse.2010.12.005.
17. Chen, W.; Moriya, K.; Sakai, T.; Koyama, L.; Cao, C. Mapping a burned forest area from Landsat TM data by multiple methods. *Geomatics, Natural Hazards and Risk* **2016**, *7*, 384–402. doi:10.1080/19475705.2014.925982.
18. Daldegan, G.; de Carvalho, O.; Guimarães, R.; Gomes, R.; Ribeiro, F.; McManus, C. Spatial Patterns of Fire Recurrence Using Remote Sensing and GIS in the Brazilian Savanna: Serra do Tombador Nature Reserve, Brazil. *Remote Sensing* **2014**, *6*, 9873–9894. doi:10.3390/rs6109873.
19. Liu, J.; Heiskanen, J.; Maeda, E.E.; Pellikka, P.K. Burned area detection based on Landsat time series in savannas of southern Burkina Faso. *International Journal of Applied Earth Observation and Geoinformation* **2018**, *64*, 210 – 220. doi:https://doi.org/10.1016/j.jag.2017.09.011.
20. Silva, J.M.N.; Pereira, J.M.C.; Cabral, A.I.; Sá, A.C.L.; Vasconcelos, M.J.P.; Mota, B.; Grégoire, J.M. An estimate of the area burned in southern Africa during the 2000 dry season using SPOT-VEGETATION satellite data. *Journal of Geophysical Research: Atmospheres* **2003**, *108*, n/a–n/a. doi:10.1029/2002jd002320.
21. Stroppiana, D.; Bordogna, G.; Carrara, P.; Boschetti, M.; Boschetti, L.; Brivio, P. A method for extracting burned areas from Landsat TM/ETM images by soft aggregation of multiple Spectral Indices and a region growing algorithm. *ISPRS Journal of Photogrammetry and Remote Sensing* **2012**, *69*, 88–102. doi:10.1016/j.isprsjprs.2012.03.001.
22. Trisakti, B.; Nugroho, U.C.; Zubaidah, A. TECHNIQUE FOR IDENTIFYING BURNED VEGETATION AREA USING LANDSAT 8 DATA. *International Journal of Remote Sensing and Earth Sciences (IJReSES)* **2017**, *13*, 121. doi:10.30536/ijreses.2016.v13.a2447.
23. Cabral, A.I.R.; vasconcelos, M.J.P.; Pereira, J.M.C.; Martins, E.; Bartholomé, É. A land cover map of southern hemisphere Africa based on SPOT-4 Vegetation data. *International Journal of Remote Sensing* **2006**, *27*, 1053–1074. doi:10.1080/01431160500307409.
24. Cabral, A.; Vasconcelos, M.; Oom, D.; Sardinha, R. Spatial dynamics and quantification of deforestation in the central-plateau woodlands of Angola (1990–2009). *Applied Geography* **2011**, *31*, 1185 – 1193. doi:https://doi.org/10.1016/j.apgeog.2010.09.003.
25. Ceccarelli, T.; Smiraglia, D.; Bajocco, S.; Rinaldo, S.; Angelis, A.D.; Salvati, L.; Perini, L. Land cover data from Landsat single-date imagery: an approach integrating pixel-based and object-based classifiers. *European Journal of Remote Sensing* **2013**, *46*, 699–717. doi:10.5721/eujrs20134641.
26. Midekisa, A.; Holl, F.; Savory, D.J.; Andrade-Pacheco, R.; Gething, P.W.; Bennett, A.; Sturrock, H.J.W. Mapping land cover change over continental Africa using Landsat and Google Earth Engine cloud computing. *PLOS ONE* **2017**, *12*, e0184926. doi:10.1371/journal.pone.0184926.

27. Phiri, D.; Morgenroth, J. Developments in Landsat Land Cover Classification Methods: A Review. *Remote Sensing* **2017**, *9*, 967. doi:10.3390/rs9090967.
28. Quinlan, J.R. Induction of Decision Trees. *Mach. Learn.* **1986**, *1*, 81–106. doi:10.1023/A:1022643204877.
29. Breiman, L. Random Forests. *Mach. Learn.* **2001**, *45*, 5–32. doi:10.1023/A:1010933404324.
30. Chen, T.; Guestrin, C. XGBoost: A Scalable Tree Boosting System. *ArXiv* **2016**, *abs/1603.02754*.
31. Arnaldo, I.; O'Reilly, U.M.; Veeramachaneni, K. Building Predictive Models via Feature Synthesis. 2015, pp. 983–990. doi:10.1145/2739480.2754693.
32. Mcconaghy, T., FFX: Fast, Scalable, Deterministic Symbolic Regression Technology; 2011; pp. 235–260. doi:10.1007/978-1-4614-1770-5_13.
33. Liu, H.; Motoda, H. *Feature Extraction, Construction and Selection: A Data Mining Perspective*; Kluwer Academic Publishers: USA, 1998.
34. Sondhi, P. Feature construction methods: a survey. *Sifaka. cs. uiuc. edu* **2009**, *69*, 70–71.
35. Khalid, S.; Khalil, T.; Nasreen, S. A survey of feature selection and feature extraction techniques in machine learning. 2014 Science and Information Conference, 2014, pp. 372–378.
36. Rasan, N.; Mani, D. A Survey on Feature Extraction Techniques. *International Journal of Innovative Research in Computer and Communication Engineering* **2015**, *03*, 52–55. doi:10.15680/ijirce.2015.0301009.
37. Dong, G.; Liu, H. *Feature Engineering for Machine Learning and Data Analytics*, 1st ed.; CRC Press, Inc.: USA, 2018.
38. Swesi, I.M.A.O.; Bakar, A.A., Recent Developments on Evolutionary Computation Techniques to Feature Construction. In *Intelligent Information and Database Systems: Recent Developments*; Huk, M.; Maleszka, M.; Szczerbicki, E., Eds.; Springer International Publishing: Cham, 2020; pp. 109–122. doi:10.1007/978-3-030-14132-5_9.
39. Xue, B.; Zhang, M. Evolutionary computation for feature manipulation: key challenges and future directions. 2016 IEEE Congress on Evolutionary Computation (CEC), 2016, pp. 3061–3067.
40. Espejo, P.G.; Ventura, S.; Herrera, F. A Survey on the Application of Genetic Programming to Classification. *IEEE Transactions on Systems, Man, and Cybernetics, Part C (Applications and Reviews)* **2010**, *40*, 121–144. doi:10.1109/TSMCC.2009.2033566.
41. Krawiec, K. Genetic Programming-based Construction of Features for Machine Learning and Knowledge Discovery Tasks. *Genetic Programming and Evolvable Machines* **2002**, *3*, 329–343.
42. Krawiec, K.; Bhanu, B. Coevolutionary Feature Learning for Object Recognition. *Machine Learning and Data Mining in Pattern Recognition*; Perner, P.; Rosenfeld, A., Eds.; Springer Berlin Heidelberg: Berlin, Heidelberg, 2003; pp. 224–238.
43. Krawiec, K.; Włodarski, L. Coevolutionary feature construction for transformation of representation of machine learners. *Intelligent Information Processing and Web Mining*; Kłopotek, M.A.; Wierzchoń, S.T.; Trojanowski, K., Eds.; Springer Berlin Heidelberg: Berlin, Heidelberg, 2004; pp. 139–150.
44. Neshatian, K.; Zhang, M.; Johnston, M. Feature Construction and Dimension Reduction Using Genetic Programming. *Proceedings of the 20th Australian Joint Conference on Advances in Artificial Intelligence*; Springer-Verlag: Berlin, Heidelberg, 2007; AI'07, p. 160–170.
45. Tran, B.; Xue, B.; Zhang, M. Genetic programming for feature construction and selection in classification on high-dimensional data. *Memetic Computing* **2016**, *8*, 3–15.
46. Tran, C.T.; Zhang, M.; Andreae, P.; Xue, B. Genetic Programming Based Feature Construction for Classification with Incomplete Data. *Proceedings of the Genetic and Evolutionary Computation Conference*; Association for Computing Machinery: New York, NY, USA, 2017; GECCO '17, p. 1033–1040. doi:10.1145/3071178.3071183.
47. Chen, Q.; Zhang, M.; Xue, B. Genetic Programming with Embedded Feature Construction for High-Dimensional Symbolic Regression. *Intelligent and Evolutionary Systems*; Leu, G.; Singh, H.K.; Elsayed, S., Eds.; Springer International Publishing: Cham, 2017; pp. 87–102.
48. Tran, B.; Xue, B.; Zhang, M. Genetic programming for multiple-feature construction on high-dimensional classification. *Pattern Recognition* **2019**, *93*, 404 – 417. doi:https://doi.org/10.1016/j.patcog.2019.05.006.
49. Lin, J.Y.; Ke, H.R.; Chien, B.C.; Yang, W.P. Designing a classifier by a layered multi-population genetic programming approach. *Pattern Recognition* **2007**, *40*, 2211–2225. doi:10.1016/j.patcog.2007.01.003.

50. Kishore, J.K.; Patnaik, L.M.; Mani, V.; Agrawal, V.K. Application of genetic programming for multicategory pattern classification. *IEEE Transactions on Evolutionary Computation* **2000**, *4*, 242–258. doi:10.1109/4235.873235.
51. Smith, M.; Bull, L. Feature Construction and Selection Using Genetic Programming and a Genetic Algorithm. 2003. doi:10.1007/3-540-36599-0_21.
52. Guo, H.; Nandi, A.K. Breast cancer diagnosis using genetic programming generated feature. *Pattern Recognition* **2006**, *39*, 980 – 987. doi:https://doi.org/10.1016/j.patcog.2005.10.001.
53. Ahmed, S.; Zhang, M.; Peng, L.; Xue, B. Multiple Feature Construction for Effective Biomarker Identification and Classification Using Genetic Programming. Proceedings of the 2014 Annual Conference on Genetic and Evolutionary Computation; Association for Computing Machinery: New York, NY, USA, 2014; GECCO '14, p. 249–256. doi:10.1145/2576768.2598292.
54. Virgolin, M.; Alderliesten, T.; Bel, A.; Witteveen, C.; Bosman, P.A.N. Symbolic Regression and Feature Construction with GP-GOMEA Applied to Radiotherapy Dose Reconstruction of Childhood Cancer Survivors. Proceedings of the Genetic and Evolutionary Computation Conference; Association for Computing Machinery: New York, NY, USA, 2018; GECCO '18, p. 1395–1402. doi:10.1145/3205455.3205604.
55. Ain, Q.U.; Xue, B.; Al-Sahaf, H.; Zhang, M. Genetic Programming for Multiple Feature Construction in Skin Cancer Image Classification. 2019 International Conference on Image and Vision Computing New Zealand (IVCNZ), 2019, pp. 1–6.
56. Cherrier, N.; Poli, J.; Defurne, M.; Sabatié, F. Consistent Feature Construction with Constrained Genetic Programming for Experimental Physics. IEEE Congress on Evolutionary Computation, CEC 2019, Wellington, New Zealand, June 10-13, 2019. IEEE, 2019, pp. 1650–1658. doi:10.1109/CEC.2019.8789937.
57. Gong, P.; Marceau, D.J.; Howarth, P.J. A comparison of spatial feature extraction algorithms for land-use classification with SPOT HRV data. *Remote Sensing of Environment* **1992**, *40*, 137 – 151. doi:https://doi.org/10.1016/0034-4257(92)90011-8.
58. Romero, A.; Gatta, C.; Camps-Valls, G. Unsupervised Deep Feature Extraction for Remote Sensing Image Classification. *IEEE Transactions on Geoscience and Remote Sensing* **2016**, *54*, 1349–1362. doi:10.1109/TGRS.2015.2478379.
59. Ren, J.; Zabalza, J.; Marshall, S.; Zheng, J. Effective Feature Extraction and Data Reduction in Remote Sensing Using Hyperspectral Imaging [Applications Corner]. *IEEE Signal Processing Magazine* **2014**, *31*, 149–154.
60. Pasquarella, V.J.; Holden, C.E.; Woodcock, C.E. Improved mapping of forest type using spectral-temporal Landsat features. *Remote Sensing of Environment* **2018**, *210*, 193 – 207. doi:https://doi.org/10.1016/j.rse.2018.02.064.
61. Puente, C.; Olague, G.; Smith, S.; Bullock, S.; González-Botello, M.; Hinojosa, A. A Genetic Programming Approach to Estimate Vegetation Cover in the Context of Soil Erosion Assessment. *Photogrammetric Engineering and Remote Sensing* **2011**, *77*, 363–375. doi:10.14358/PERS.77.4.363.
62. Makkeasorn, A.; Chang, N.B.; Li, J. Seasonal change detection of riparian zones with remote sensing images and genetic programming in a semi-arid watershed. *Journal of Environmental Management* **2009**, *90*, 1069 – 1080. doi:https://doi.org/10.1016/j.jenvman.2008.04.004.
63. Makkeasorn, A.; Chang, N.B.; Beaman, M.; Wyatt, C.; Slater, C. Soil moisture estimation in a semiarid watershed using RADARSAT-1 satellite imagery and genetic programming. *Water Resources Research* **2006**, *42*.
64. Chion, C.; Landry, J.A.; Costa, L. A Genetic-Programming-Based Method for Hyperspectral Data Information Extraction: Agricultural Applications. *Geoscience and Remote Sensing, IEEE Transactions on* **2008**, *46*, 2446 – 2457. doi:10.1109/TGRS.2008.922061.
65. Chen, L. A study of applying genetic programming to reservoir trophic state evaluation using remote sensor data. *International Journal of Remote Sensing* **2003**, *24*, 2265–2275.
66. Taghizadeh-Mehrjardi, R.; Ayoubi, S.; Namazi, Z.; Malone, B.; Zolfaghari, A.; Roustaei-Sadrabadi, F. Prediction of soil surface salinity in arid region of central Iran using auxiliary variables and genetic programming. *Arid Land Research and Management* **2016**, *30*. doi:10.1080/15324982.2015.1046092.
67. Lary, D.J.; Alavi, A.H.; Gandomi, A.H.; Walker, A.L. Machine learning in geosciences and remote sensing. *Geoscience Frontiers* **2016**, *7*, 3 – 10. Special Issue: Progress of Machine Learning in Geosciences, doi:https://doi.org/10.1016/j.gsf.2015.07.003.

68. Costa, L.; Nunes, L.; Ampatzidis, Y. A new visible band index (vNDVI) for estimating NDVI values on RGB images utilizing genetic algorithms. *Computers and Electronics in Agriculture* **2020**, *172*, 105334. doi:<https://doi.org/10.1016/j.compag.2020.105334>.
69. Kabiri, P.; Pandi, M.H.; Nejat, S.K.; Ghaderi, H. NDVI Optimization Using Genetic Algorithm. 2011 7th Iranian Conference on Machine Vision and Image Processing, 2011, pp. 1–5.
70. Lopes, C.; Leite, A.; Vasconcelos, M.J. Open-access cloud resources contribute to mainstream REDD+: The case of Mozambique. *Land Use Policy* **2019**, *82*, 48 – 60. doi:<https://doi.org/10.1016/j.landusepol.2018.11.049>.
71. Cabral, A.I.; Silva, S.; Silva, P.C.; Vanneschi, L.; Vasconcelos, M.J. Burned area estimations derived from Landsat ETM+ and OLI data: Comparing Genetic Programming with Maximum Likelihood and Classification and Regression Trees. *ISPRS Journal of Photogrammetry and Remote Sensing* **2018**, *142*, 94 – 105. doi:<https://doi.org/10.1016/j.isprsjprs.2018.05.007>.
72. Vasconcelos, M.; Cabral, A.; B. Melo, J.; Pearson, T.; Pereira, H.; Cassamá, V.; Yudelman, T. Can blue carbon contribute to clean development in West Africa? The case of Guinea-Bissau. *Mitigation and Adaptation Strategies for Global Change* **2014**, *20*. doi:10.1007/s11027-014-9551-x.
73. Temudo, M.; Cabral, A.; Talhinhos, P. Petro-Landscapes: Urban Expansion and Energy Consumption in Mbanza Kongo City, Northern Angola. *Human Ecology* **2019**, *47*, 565–575. doi:10.1007/s10745-019-00088-6.
74. Kotteck, M.; Grieser, J.; Beck, C.; Rudolf, B.; Rubel, F. World Map of the Köppen-Geiger climate classification updated. *Meteorologische Zeitschrift* **2006**, *15*, 259–263. doi:10.1127/0941-2948/2006/0130.
75. Temudo, M.P.; Cabral, A.I.; Talhinhos, P. Urban and rural household energy consumption and deforestation patterns in Zaire province, Northern Angola: A landscape approach. *Applied Geography* **2020**, *119*, 102207. doi:<https://doi.org/10.1016/j.apgeog.2020.102207>.
76. Dinis, A.C. *Características mesológicas de Angola: descrição e correlação dos aspectos fisiográficos, dos solos e da vegetação das zonas agrícolas angolanas*; IPAD - Instituto Português de Apoio ao Desenvolvimento: Lisboa, 2006.
77. *Climate Risk and Adaptation Country Profile: Mozambique*, (accessed November 17, 2020). <https://www.gfdr.org/en/publication/climate-risk-and-adaptation-country-profile-mozambique>.
78. *Climate Analysis Mozambique*, (accessed November 17, 2020). <https://fscluster.org/mozambique/document/climate-analysis-mozambique>.
79. Tibshirani, R. Regression Shrinkage and Selection via the Lasso. *Journal of the Royal Statistical Society. Series B (Methodological)* **1996**, *58*, 267–288.
80. Friedman, J.; Hastie, T.; Tibshirani, R. Regularization Paths for Generalized Linear Models via Coordinate Descent. *Journal of Statistical Software* **2010**, *33*. doi:10.18637/jss.v033.i01.
81. Pedregosa et al, F. Scikit-learn: Machine Learning in Python. *Journal of Machine Learning Research* **2011**, *12*, 2825–2830.
82. *Landsat 8 Bands*, (accessed November 17, 2020). <https://landsat.gsfc.nasa.gov/landsat-8/landsat-8-bands>.