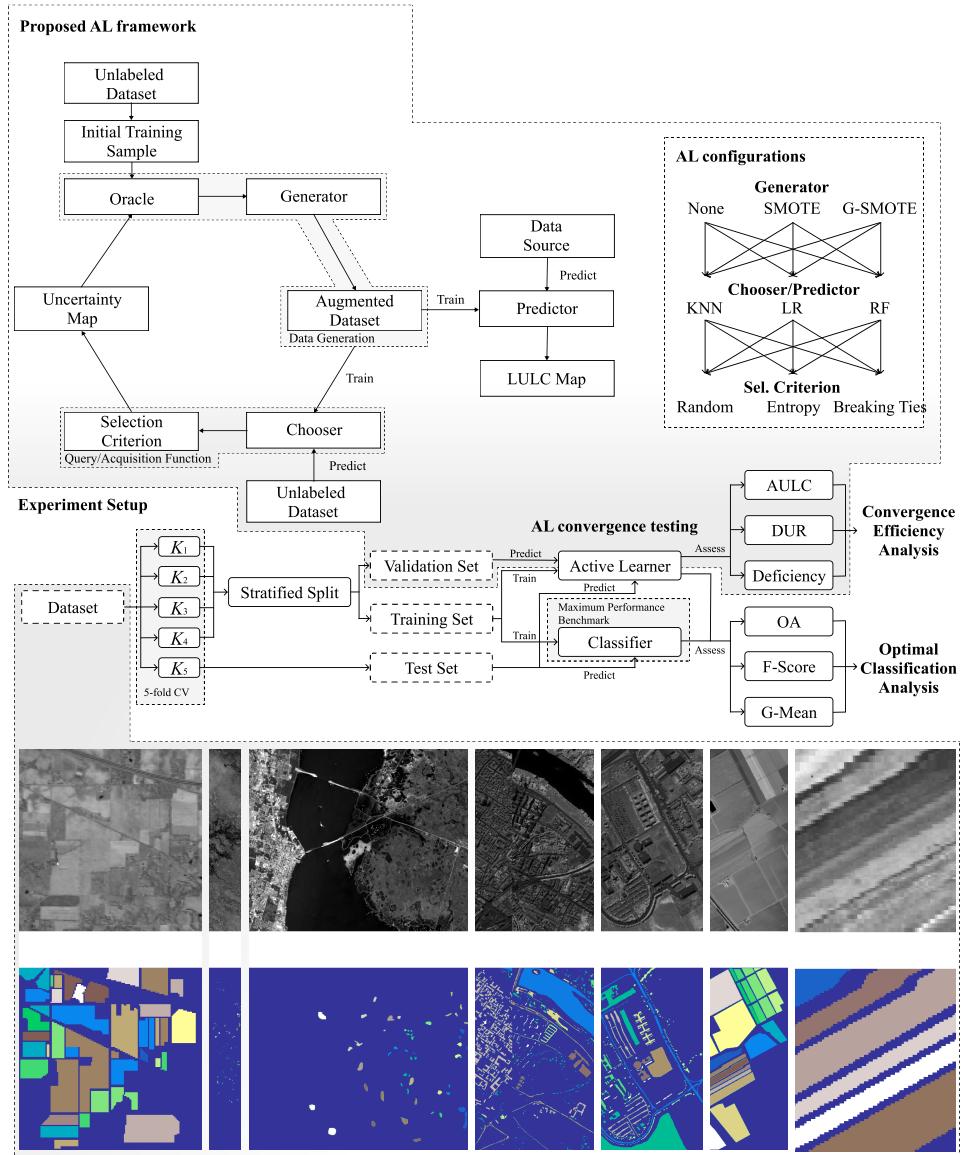


Graphical Abstract

Increasing the Effectiveness of Active Learning: Introducing Artificial Data Generation in Active Learning for Land Use/Land Cover Classification

Joao Fonseca, Georgios Douzas, Fernando Bacao



Highlights

**Increasing the Effectiveness of Active Learning:
Introducing Artificial Data Generation in Active Learning for Land Use/Land
Cover Classification**

Joao Fonseca, Georgios Douzas, Fernando Bacao

- We integrate artificial data generation into the Active Learning framework;
- The proposed modification significantly reduces the cost and time requirements for a successful AL implementation;
- The effectiveness of this framework was shown with simple, context-agnostic data generation heuristics;
- The amount of data required to reach the performance thresholds defined for the baseline AL methods were significantly reduced in all seven datasets used in this experiment;

Increasing the Effectiveness of Active Learning: Introducing Artificial Data Generation in Active Learning for Land Use/Land Cover Classification

Joao Fonseca^a, Georgios Douzas^a, Fernando Bacao^a

^a*NOVA Information Management School, Universidade Nova de Lisboa, Campus de Campolide, 1070-312 Lisboa, Portugal*

Abstract

In remote sensing, Active Learning (AL) has become an important technique to collect informative ground truth data “on-demand” for supervised classification tasks. In spite of its effectiveness, it is still significantly reliant on user interaction, which makes it both expensive and time consuming to implement. Most of the current literature focuses on the optimization of AL by modifying the selection criteria, the chooser and/or predictors used. Although improvements in these areas will result in more effective data collection, the use of artificial data sources to reduce human-computer interaction remains unexplored. In this paper, we introduce a new component to the typical AL framework, the data generator, a source of artificial data to reduce the amount of user-labeled data required in AL. The implementation of the proposed AL framework is done using SMOTE and Geometric SMOTE as data generators. We compare the new AL framework to the original one using similar acquisition functions and predictors over three AL-specific performance metrics in seven benchmark datasets. We show that this modification of the AL framework significantly reduces cost and time requirements for a successful AL implementation in the context of remote sensing.

Keywords: Active Learning, Artificial Data Generation, Land Use/Land Cover Classification, Oversampling, SMOTE

1. Introduction

The technological development of air and spaceborne sensors, as well as the increasing number of remote sensing missions have allowed the continuous collection of

large amounts of high quality remotely sensed data. This data is often composed of multi and hyper spectral satellite imagery, essential for numerous applications, such as Land Use/Land Cover (LULC) change detection, ecosystem management [1], agricultural management [2], water resource management [3], forest management, and urban monitoring [4]. Despite LULC maps being essential for most of these applications, their production is still a challenging task [5, 6]. They can be updated using one of the following strategies:

1. Photo-interpretation. This approach consists of evaluating a patch's LULC class by a human operator based on orthophoto and satellite image interpretation [7]. This method guarantees a decent level of accuracy, as it is dependent on the interpreter's expertise and human error. Typically, it is an expensive, time-consuming task that requires the expertise of a photo-interpreter. This task is also frequently applied to obtain ground-truth labels for training and/or validating Machine Learning (ML) algorithms for related tasks [8, 9].
2. Automated mapping. This approach is based on the usage of a ML method or a combination of methods in order to obtain an updated LULC map. The development of a reliable automated method is still a challenge among the ML and remote sensing community, since the effectiveness of existing methods varies across applications and geographical areas [5]. Typically, this method requires the existence of ground-truth data, which is frequently outdated or nonexistent for the required time frame [1]. On the other hand, employing a ML method provides readily available and relatively inexpensive LULC maps. The increasing quality of state-of-the-art classification methods have motivated the application and adaptation of these methods in this domain [10].
3. Hybrid approaches. These approaches employ photo-interpreted data to augment the training dataset and improve the quality of automated mapping [11]. It attempts to accelerate the photo-interpretation process by selecting a smaller sample of the study area to be interpreted. The goal is to minimize the inaccuracies found in the LULC map by supplying high-quality ground-truth data to the automated method. The final (photo-interpreted) dataset consists of only the most informative samples, *i.e.*, patches that are typically difficult to classify for a traditional automated mapping method [12].

The latter method is best known as AL. It is especially useful whenever there is a shortage or even absence of ground-truth data and/or the mapping region does not contain updated LULC maps [13]. In a context of limited sample-collection budget,

the collection of the most informative samples capable of optimally increasing the classification accuracy of a LULC map is of particular interest [13]. AL attempts to minimize the human-computer interaction involved in photo-interpretation by selecting the data points to include in the annotation process. These data points are selected based on an uncertainty measure and represent the points close to the decision borders. Afterwards, they are passed on for photo-interpretation and added to the training dataset, while the points with the lowest uncertainty values are ignored for photo-interpretation and classification. This process is iterated until a convergence criterion is reached [14].

The relevant work developed within AL is described in detail in Section 2. This paper attempts to address some of the challenges found in AL, mainly inherited from automated and photo-interpreted mapping: mapping inaccuracies and time consuming human-computer interactions. These challenges have different sources:

1. Human error. The involvement of photo-interpreters in the data labeling step carries an additional risk to the creation of LULC patches. The minimum mapping unit being considered, as well as the quality of the orthophotos and satellite images being used, are some of the factors that may lead to the overlooking of small-area LULC patches and label-noisy training data [15].
2. High-dimensional datasets. Although the amount of bands (*i.e.*, features) present in multi and hyper spectral images contain useful information for automated classification, they also introduce an increased level of complexity and redundancy in the classification step [16]. These datasets are often prone to the Hughes phenomenon, also known as the curse of dimensionality.
3. Class separability. Producing an LULC map considering classes with similar spectral signatures makes them difficult to separate [17]. A lower pixel resolution of the satellite images may also imply mixed-class pixels, which may lead to both lower class separability as well as higher risk of human error.
4. Existence of rare land cover classes. The varying morphologies of different geographical regions naturally implies an uneven distribution of land cover classes [18]. This is particularly relevant in the context of AL since the data selection method is based on a given uncertainty measure over data points whose class label is unknown. Consequently, AL’s iterative process of data selection may disregard wrongly classified land cover areas belonging to a minority class.

Research developed in the field of AL typically focus on the reduction of human error by minimizing the human interaction with the process through the development of more efficient choosers and selection criteria within the generally accepted AL framework. Concurrently, the problem of rare land cover classes is rarely addressed. This is a frequent problem in the ML community, known as the Imbalanced Learning problem. This problem exists whenever there is an uneven between-class distribution in the dataset [19]. Specifically, most classifiers are optimized and evaluated using accuracy-like metrics, which are designed to work primarily with balanced datasets. Consequently, these metrics tend to introduce a bias towards the majority class by attributing an importance to each class proportional to its relative frequency [10]. As an example, such a classifier could achieve an overall accuracy of 99% on a binary dataset where the minority class represents 1% of the overall dataset and still be deemed useless. A number of methods have been developed to deal with this problem. They can be categorized into three different types of approaches [20, 21]. Cost-sensitive solutions perform changes to the cost matrix in the learning phase. Algorithmic level solutions modify specific classifiers to reinforce learning on minority classes. Resampling solutions modify the dataset by removing majority samples and/or generating artificial minority samples. The latter is independent from the context and can be used alongside any classifier. Because of this we will focus on artificial data generation techniques, presented in Section 3.

In this paper, we propose a novel AL framework to address two limitations commonly found in the literature: minimize human-computer interaction and reduce the class imbalance bias. This is done with the introduction of an additional component in the iterative AL procedure (the generator), used to generate artificial data to both balance and augment the training dataset. The introduction of this component is expected to reduce the number of iterations required until convergence of the predictor's quality.

This paper is organized as follows: Section 1 explains the problem and its context, Sections 2 and 3 describe the state of the art in AL and Oversampling techniques, Section 4 explains the proposed method, Section 5 covers the datasets, evaluation metrics, ML classifiers and experimental procedure, Section 6 presents the experiment's results and discussion and Section 7 presents the conclusions drawn from our findings.

2. Active Learning Approaches

As the amount of unlabeled data increases, the interest and practical usefulness of AL follows that trend [22]. AL is used as the general definition of frameworks

aiming to train a learning system in multiple steps, where a set of new data points are chosen and added to the training dataset each time [11]. Typically, an AL framework is composed of 10 elements [23, 13, 11]:

1. Data source. In the context of LULC classification, the data source is usually a hyper/multi-spectral image, a Synthetic-aperture radar (SAR) image, or a composite image.
2. Unlabeled dataset. Consists of the original data source (or a sample thereof). It is used in combination with the chooser and the selection criterion to expand the training set in regions where the classification uncertainty is higher.
3. Initial training sample. It is a small sample of the unlabeled dataset, used to initiate the first AL iteration. The size of the initial training sample normally varies between no instances at all and 10% [24].
4. Augmented training dataset. This dataset is the concatenation of the labeled initial training sample along with the datasets labeled by the oracle in past iterations (discussed in point 6).
5. Uncertainty map. The dataset containing the highest uncertainty points/patches to be labeled by the oracle.
6. Oracle. An external entity to which the uncertainty map is presented to. The oracle is responsible for annotating unlabeled instances to be added to the augmented dataset. In remote sensing, the oracle is typically a photo-interpreter, as is the case in [25]. Some of the research also refers to the oracle as the *supervisor* [13, 26].
7. Chooser. Produces the class probabilities for each unlabeled instance. This is a classifier trained using the augmented dataset. It is used to estimate the class probabilities for each instance over the unlabeled dataset.
8. Selection criterion. It quantifies the chooser's uncertainty level for each instance belonging to the unlabeled dataset. It is typically based on the class probabilities assigned by the chooser. In some situations, the chooser and the selection criterion are grouped together under the concept *acquisition function* [11] or *query function* [13]. Some of the literature refers to the selection criterion by using the concept *sampling scheme* [12].

9. Predictor. The classifier used to infer the land cover classes for the final output map. Once a stopping criterion is met, the classifier is trained using the augmented dataset and the LULC classes are inferred from the data source.
10. Prediction output. In the context of LULC classification, the prediction output is the estimated LULC map raster.

Figure 1 schematizes the steps involved in a complete AL iteration. For a better context within the remote sensing domain, the prediction output is identified as the LULC map. This framework starts by collecting unlabeled data from the original data source. It is used to generate a random initial training sample and is labeled by the oracle. In practical applications, the oracle is frequently a group of photo-interpreters [22]. The chooser is trained on the resulting dataset and is used to predict the class probabilities on the unlabeled dataset. The class probabilities are fed into a selection criterion to estimate the prediction’s uncertainty, out of which the instances with the highest uncertainty will be selected. This calculation is motivated by the absence of labels in the uncertainty dataset. Therefore, it is impossible to estimate the prediction’s accuracy in a real case scenario. The iteration is completed when the selected points are tagged by the oracle and added to the training dataset (*i.e.*, the augmented dataset).

A common challenge found in AL tasks is ensuring the consistency of AL over different initializations [22]. There are two factors involved in this phenomenon. On one hand, the implementation of the same method over different initializations may result in significantly different initial training samples, amounts to varying accuracy curves. On the other hand, the lack of a robust selection criterion and/or chooser may also result in inconsistencies across AL experiments with different initializations. This phenomenon was observed and documented in a LULC classification context in [27].

Selecting an efficient selection criterion is particularly important to find the instances closest to the decision border (*i.e.*, instances difficult to classify) [26]. Therefore, most of AL related studies focus on the design of the query/acquisition function [13].

2.1. Non-informed selection criteria

Only one non-informed (*i.e.*, random) selection criterion was found in the literature. Random sampling selects unlabeled instances without considering any external information produced by the chooser. Since the method for selecting the unlabeled instances is random, this method disregards the usage of a chooser and is comparatively worse than any other selection criterion. However, random sampling is still a

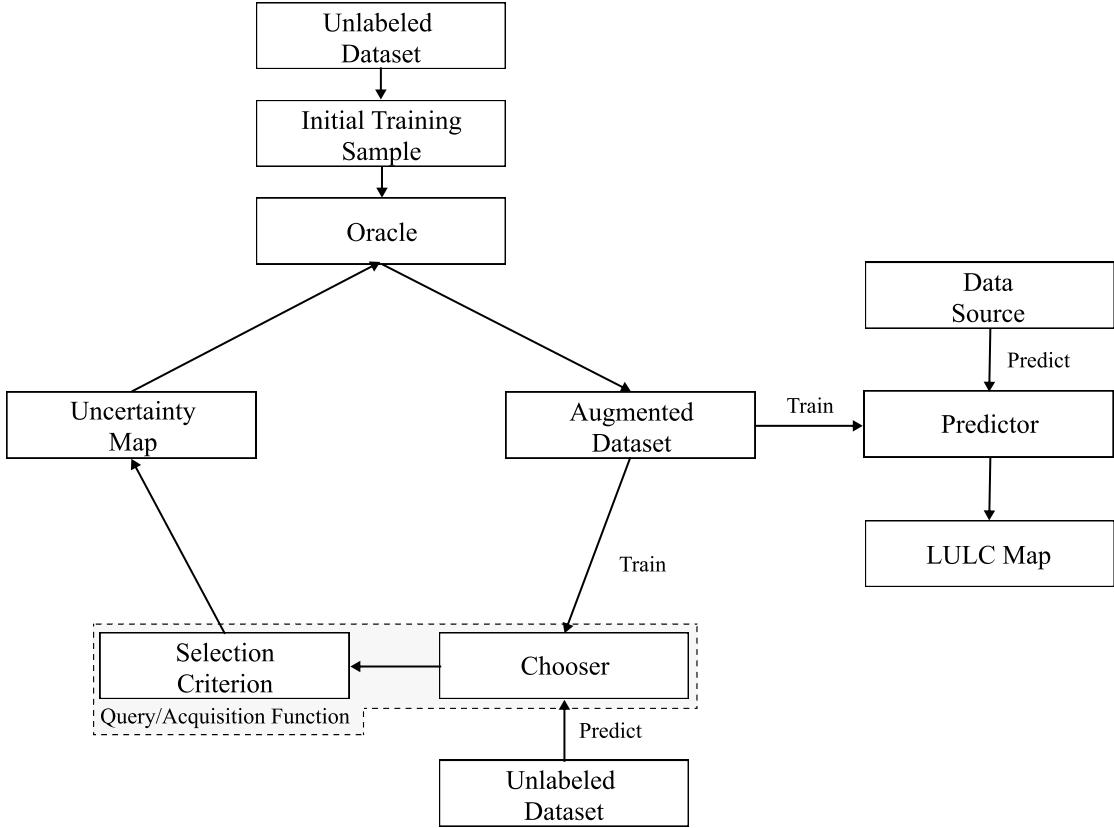


Figure 1: Diagram depicting the typical AL framework. Although the chooser and predictor are presented as separate entities, they are frequently the same classifier.

powerful baseline method [28]. Generally, different AL initializations result in high performance variability [22]. When this happens, the analysis of the mean performances over multiple repetitions is not of interest on its own. Instead, it is preferable to do pairwise comparison of different methods along with their corresponding variances.

2.2. Ensemble-based selection criteria

Ensemble disagreement is based on the class predictions of a set of classifiers. The disagreement between all the predictions for a given instance is a common measure for uncertainty, although computationally inefficient [11, 14]. It is calculated using the set of classifications over a single instance, given by the number of votes assigned to

the most frequent class [26]. This method was implemented successfully for complex applications such as deep active learning [11].

Multiview [29] consists on the training of multiple independent classifiers using different views, which correspond to the selection of subsets of features or instances in the dataset. Therefore, it can be seen as a bootstrap aggregation (bagging) ensemble disagreement method. It is represented by the maximum disagreement score out of set of disagreements calculated for each view [26]. A lower value for this metric means a higher classification uncertainty. Multiview-based maximum disagreement has been successfully applied to hyper-spectral image classification in [30] and [31].

An adapted disagreement criterion for an ensemble of k -nearest neighbors has been proposed in [14]. This method employs a k -nearest neighbors classifier and computes an instance's classification uncertainty based on the neighbors' class frequency using the maximum disagreement metric over varying values for k . As a result, this method is comparable to computing the dominant class' score over a weighted k -nearest neighbors classifier. This method was also used on a multimetric active learning framework [32].

Another relevant ensemble-based selection criterion is the binary random forest-based query model [13]. This method employs a one-versus-one ensemble method to demonstrate an efficient data selection method using the estimated probability of each binary random forest and determining the classification uncertainty based on the probabilities closest to 0.5 (*i.e.*, the least separable pair of classes are used to determine the uncertainty value). However, this study fails to compare the proposed method with other benchmark methods, such as random sampling.

2.3. Entropy-based criteria

A number of contributions have focused on entropy-based querying. The application of entropy is common among active deep learning applications [33], where the training of an ensemble of classifiers is often too expensive. The measure of entropy is formulated as follows:

$$H(x_i) = \sum_{\omega=1}^{N_i} p(y_i^* = \omega | x_i) \log_2[p(y_i^* = \omega | x_i)] \quad (1)$$

The measurement of entropy H is based on the observed probability $p(y_i^* = \omega | x_i)$ of obtaining class ω as the predicted class label y_i^* , where N_i is the number classes predicted for instance x_i .

Entropy query-by-bagging (EQB), also defined as maximum entropy [12], is an ensemble approach of the entropy selection criterion, originally proposed in [34]. This strategy uses the set of predictions produced by the ensemble classifier to calculate

those many entropy measurements. The estimated uncertainty measure for one instance is given by the maximum entropy within that set. EQB was observed to be an efficient selection criterion. Specifically, [26] applied EQB on hyper-spectral remote sensing imagery using Support Vector Machines (SVM) and Extreme Learning Machines (ELM) as choosers, achieving optimal results when combining EQB with ELM. Another study successfully implemented this method on an active deep learning application [12]. Another study improved over this method with a normalized EQB selection criterion [35].

2.4. Other relevant criteria

Margin Sampling is a SVM-specific criterion, based on the distance of a given point to the SVM’s decision boundary [26]. This method is less popular than the remaining methods because it is limited to one type of chooser (SVMs). One extension of this method is the multiclass level uncertainty [26], calculated by subtracting the instance’s distance to the decision boundaries of the two most probable classes [36].

The Mutual Information-based (MI) criterion selects the new training instances by maximizing the mutual information between the classifier and class labels in order to select instances from regions that are difficult to classify. Although this method is commonly used, it is frequently outperformed by the breaking ties selection criterion [37, 38].

The breaking ties (BT) selection criterion was originally introduced in [39]. It is formulated as follows:

$$BT(x_i) = \arg \min_{x_i, i \in S_u} \left\{ \max_{\omega \in N} p(y_i^* = \omega | x_i) - \max_{\omega \in N \setminus \{\omega^+\}} p(y_i^* = \omega | x_i) \right\} \quad (2)$$

Which is the subtraction of the probabilities of the two most likely classes. Another related method is Modified Breaking Ties scheme (MBT), which aims at finding the instances containing the largest probabilities for the dominant class [38, 40].

Another type of selection criteria identified is the loss prediction method [41]. This method replaces the selection criterion with a predictor whose goal is to estimate the chooser’s loss for a given prediction. This allows the new classifier to estimate the prediction loss on unlabeled instances and select the ones with the highest predicted loss.

Some of the literature fails to specify the strategy employed, although inferring it is generally intuitive. For example, [42] successfully used AL to address the imbalanced learning problem. They employed an ensemble of SVMs as the chooser and predictor, as well as an ensemble-based selection criterion. All of the research found related to this topic focused on the improvement of AL through modifications

on the selection criterion, chooser or predictor. None of these publications proposed significant variations to the original AL framework.

3. Artificial Data Generation Approaches

The generation of artificial data is a common approach to address imbalanced learning tasks [21], as well as improving the effectiveness of supervised learning tasks [43]. In recent years some sophisticated data generation approaches were developed. However, the scope of this work is to propose the integration of a generator within the AL framework. To do this, we will focus on heuristic data generation approaches, specifically, oversamplers.

Heuristic data resampling methods employ local and/or global information to generate new, relevant, non-duplicated instances. These methods are most commonly used to populate minority classes and balance the between-class distribution of a dataset. The Synthetic Minority Oversampling Technique (SMOTE) [44] is a popular heuristic oversampling algorithm, proposed in 2002. The simplicity and effectiveness of this method contributes to its prevailing popularity. It generates a new instance \vec{z} through a linear interpolation of a randomly selected minority-class instance \vec{x} and one of its randomly selected k -nearest neighbors \vec{y} such that $\vec{z} = \alpha \vec{x} + (1 - \alpha) \vec{y}$ where α is a random real number between 0 and 1, as shown in Figure 2.

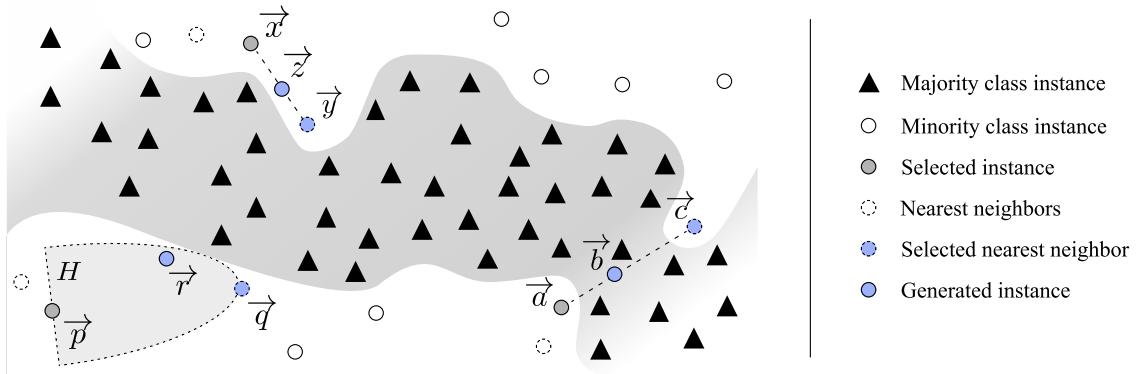


Figure 2: Examples of SMOTE and G-SMOTE generation process. SMOTE randomly selects instance \vec{x} and one of its nearest neighbors \vec{y} to produce instance \vec{z} . Noisy instance \vec{b} is generated with SMOTE using points \vec{d} and \vec{c} . Instance \vec{p} and one of its nearest neighbors \vec{q} are used by G-SMOTE to originate instance \vec{r} .

The implementation of SMOTE for LULC classification tasks has been found to improve the quality of the predictors used [45, 46]. Despite its popularity, its drawbacks motivated the development of other oversampling methods [47]. The following issues are addressed in more recent papers:

1. Generation of noisy instances due to the selection of k -nearest neighbors and initial instance. The selection of an instance and/or neighboring instance located inside a majority class region may produce artificial instances within that region and amplify noisy data. Borderline-SMOTE [48] is a modification of SMOTE in which only the minority examples near the borderline are oversampled. This method avoids the generation of noisy instances by disregarding minority class instances located in a majority class region as well as instances distant from the decision borders. The Adaptive Synthetic Sampling approach (ADASYN) [49] uses a density distribution ratio to address this limitation and focus the artificial data generation on minority class regions that are more difficult to classify.
2. Generation of noisy data due to the use of instances from two different minority class clusters. Choosing a minority class instance \vec{d} and one of its nearest neighbors \vec{b} belonging to a different minority cluster may lead to the generation of an instance \vec{c} located within the two classes, as shown in Figure 2. K-means SMOTE [50] and Self-Organizing map oversampling (SOMO) [51] reduce this effect by oversampling minority class instances within the same clusters.
3. Generation of nearly duplicated instances. The linear interpolation of parent instances that are close to each other produces an artificial instances with similar properties as its parents. Geometric SMOTE (G-SMOTE) [47] introduces a modification of the SMOTE algorithm in the data generation mechanism to produce artificial instances with higher variability.

The G-SMOTE algorithm is introduced as a generalization of the vanilla SMOTE. Instead of generating artificial data as a linear combination of the parent instances, it is done within a deformed, truncated hyper-spheroid. G-SMOTE generates an artificial instance \vec{r} within a hyper-spheroid H , formed by selecting a minority instance \vec{p} and one of its nearest neighbors \vec{q} , as shown in Figure 2. The truncation and deformation parameters define the shape of the spheroid's geometry. The method also modifies the selection strategy for the k -nearest neighbors, accepting the generation of artificial instances using instances from different classes. G-SMOTE has shown

superior performance when compared with other oversampling methods for LULC classification tasks, regardless of the classifier used [52].

4. Proposed method

Within the literature identified, most of the work developed in the AL domain revolved around improving the quality of the chooser, predictor and/or selection criterion. Although these methods allow earlier convergence of the AL iterative process, the impact of these methods are only observed between iterations. Consequently, none of these contributions focused on the definition of decision borders within iterations. The method proposed in this paper modifies the AL framework by introducing an artificial data generation step within AL’s iterative process. We define this component as the generator and is intended to be integrated into the AL framework as shown in Figure 3.

This method leverages the capability of artificial data to introduce more data variability into the augmented dataset and facilitate the chooser’s training phase with a more consistent definition of the decision boundaries at each iteration. Therefore, any algorithm capable of producing artificial data, be it agnostic or specific to the domain, can be employed. The artificial data is only used to train the classifiers involved in the process (chooser and predictor) and is discarded once the chooser’s training phase is completed. The remaining steps in the AL framework remain unchanged. This method addresses the limitations found in the previous sections:

1. The convergence of the predictor’s performance should be anticipated with the clearer definition of the decision boundaries across iterations.
2. Annotation cost is expected to reduce as the need for labeled instances reduces along with the early convergence of the classification performance.
3. The class imbalance bias observed in typical classification tasks, as well as in AL is mitigated by balancing the class frequencies at each iteration.

Although the performance of this method is shown within a LULC classification context, the proposed framework is independent from the domain. The high dimensionality of remotely sensed imagery make its classification particularly challenging when the availability of labeled data is scarce and/or comes at a high cost, being subjected to the curse of dimensionality. Consequently, it is a relevant and appropriate domain to test this method.

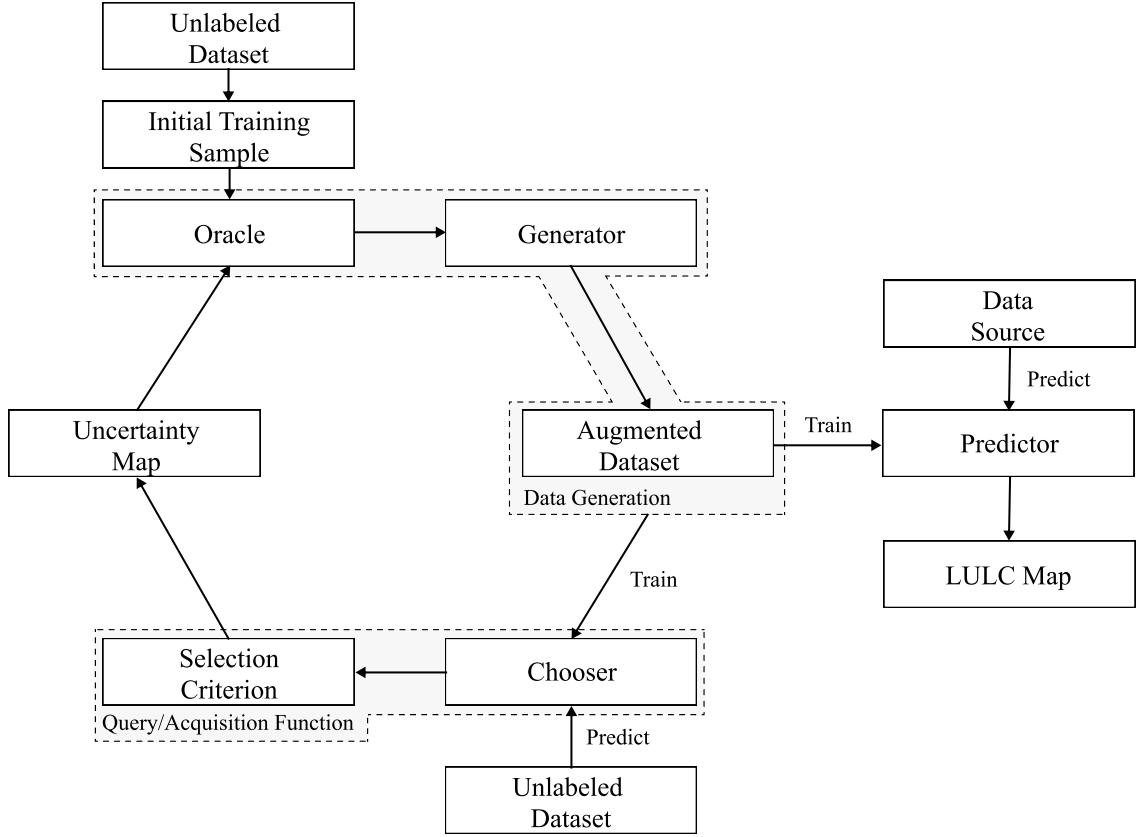


Figure 3: Proposed AL framework. The data generation mechanism is represented as the generator, which is used to complete the data generation phase. The remaining steps are left unchanged.

5. Methodology

In this section we describe the datasets, evaluation metrics, oversamplers, classifiers, software used and the procedure developed. We demonstrate the proposed method’s efficiency over 7 datasets, sampled from publicly available, well-known benchmark remote sensing landscapes frequently found in the literature. The datasets and sampling strategy are described in Subsection 5.1. On each of these datasets, we implement 3 different classifiers over the entire training set to estimate the optimal classification performance, the original AL framework as the baseline reference and the proposed method using two different generators, described in Subsection 5.2. The metrics used to estimate the performance of these algorithms are described in

Subsection 5.3. Finally, the experimental procedure is described in Subsection 5.4.

Our methodology focuses on two objectives: (1) Comparison of optimal classification performance among active learners and traditional supervised learning and (2) Comparison of classification convergence efficiency across AL frameworks.

5.1. Datasets

The datasets used were extracted from publicly available repositories with hyperspectral images. Additionally, all datasets were collected using the same sampling procedure. The description of the hyperspectral scenes used in this study is provided in Table 1. These scenes were chosen because of their popularity in the research community and their high baseline classification scores. Consequently, demonstrating an outperforming method in this context is particularly challenging and valuable.

The Indian Pines scene [53] is composed of agriculture fields in approximately two thirds of its coverage, low density buildup areas and natural perennial vegetation in the remainder of its area (see Figure 4a). The Pavia Centre and University scenes are hyperspectral, high-resolution images containing ground truth data composed of urban-related coverage (see Figures 4b and 4c). The Salinas and Salinas A scenes contain at-sensor radiance data. As subset of Salinas, the Salinas A scene contains the vegetables fields present in Salinas and the latter is also composed of bare soils and vineyard fields (see Figures 4d and 4e). The Botswana scene contains ground truth data composed of seasonal swamps, occasional swamps, and drier woodlands located in the distal portion of the Delta (see Figure 4f). The Kennedy Space Center scene contains a ground truth composed of both vegetation and urban-related coverage (see Figure 4g)

Dataset	Sensor	Location	Dimension	Bands	Res. (m)	Classes
Botswana	Hyperion	Okavango Delta	1476 x 256	145	30	14
Salinas A	AVIRIS	California, USA	86 x 83	224	3.7	6
Kennedy Space Center	AVIRIS	Florida, USA	512 x 614	176	18	16
Indian Pines	AVIRIS	NW Indiana, USA	145 x 145	220	20	16
Salinas	AVIRIS	California, USA	512 x 217	224	3.7	16
Pavia University	ROSIS	Pavia, Italy	610 x 610	103	1.3	9
Pavia Centre	ROSIS	Pavia, Italy	1096 x 1096	102	1.3	9

Table 1: Description of the hyperspectral scenes used in this experiment.

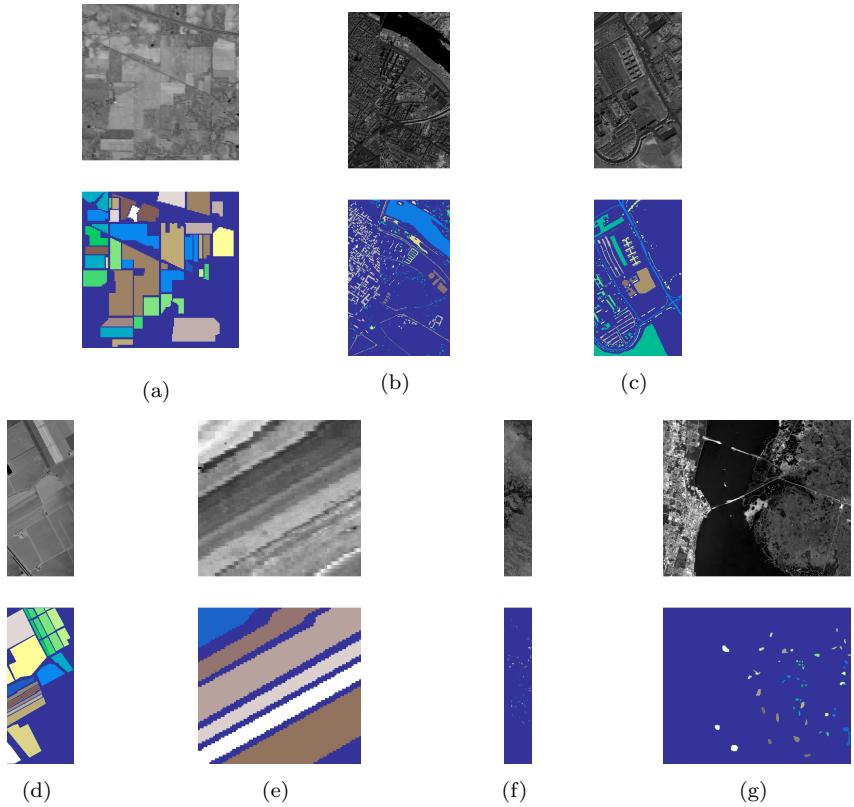


Figure 4: Gray scale visualization of a band (top row) and ground truth (bottom row) of each scene used in this study. (a) Indian Pines, (b) Pavia Centre, (c) Pavia University, (d) Salinas, (e) Salinas A, (f) Botswana, (g) Kennedy Space Center

The sampling strategy is similar to all datasets. The pixels without a ground truth label are first discarded. All the classes with cardinality lower than 150 are also discarded. This is done to maintain feasible Imbalance Ratios (IR) across datasets (where $IR = \frac{count(C_{maj})}{count(C_{min})}$). Finally, a stratified sample of 1500 instances are selected for the experiment. The resulting datasets are described in Table 2. The motivation for this strategy is three fold: (1) reduce the datasets to a manageable size and allow the experimental procedure to be completed within a feasible time frame, (2) ensure the relative class frequencies in the scenes are preserved and (3) ensure equivalent analyses across datasets and AL frameworks. In this context, a fixed number of instances per dataset is especially important to standardize the AL-related performance metrics.

Dataset	Features	Instances	Min. Instances	Maj. Instances	IR	Classes
Botswana	145	1500	89	154	1.73	12
Salinas A	224	1500	109	428	3.93	6
Kennedy Space Center	176	1500	47	272	5.79	12
Indian Pines	220	1500	31	366	11.81	12
Salinas	224	1500	25	312	12.48	16
Pavia University	103	1500	33	654	19.82	9
Pavia Centre	102	1500	27	668	24.74	9

Table 2: Description of the datasets collected from each corresponding scene. The sampling strategy is similar to all scenes.

5.2. Machine Learning Algorithms

We use two different types of ML algorithms. Data generation algorithms, used to form the generator, and classification algorithms, used to form the chooser and predictor. In order to maintain simplicity and a common approach to most of the literature in the topic, the classifiers used to play the chooser and predictor are the same.

Although any method capable of generating artificial data can be used as a generator, the ones used in this experiment are oversamplers, originally developed to deal with imbalanced learning problems. Specifically, we chose SMOTE for its popularity and simplicity. We also chose G-SMOTE as a better performing generalization of the former method.

Three classification algorithms are used as the chooser and predictor. We use different types of classifiers to test the framework’s performance under varying situations: neighbors-based, linear and ensemble models. The neighbors-based classifier chosen was K -nearest neighbors (KNN) [54], a logistic regression (LR) [55] is used as the linear model and a random forest classifier (RFC) [56] was used as the ensemble model.

The acquisition function is completed by testing three different selection criteria. Random selection is used as a baseline selection criterion, whereas entropy (see Formula 1) and breaking ties (see Formula 2) are used due to their popularity and classifier independence.

5.3. Evaluation Metrics

According to [5], nearly 80% of the satellite-based LULC studies employ the *Overall Accuracy* (OA) and *Kappa coefficient* performance metrics. However, these metrics are frequently insufficient to accurately depict classification performance [57, 58]. Metrics such as Producer's Accuracy (or *Recall*) and User's Accuracy (or *Precision*) are also commonly used. Since they consist of ratios based on True/False Positives (TP and FP) and Negatives (TN and FN), formulated as $Precision = \frac{TP}{TP+FP}$ and $Recall = \frac{TP}{TP+FN}$, they provide per class information regarding the classifier's classification performance. However, in this experiment, the meaning and number of classes available in each dataset varies, making these metrics difficult to synthesize.

While OA and Kappa tend to overestimate a classifier's performance on datasets with high IR, other metrics such as *F-score* and *Geometric mean* (G-mean) are less sensitive to the data imbalance bias [59, 60]. Therefore, we employ 3 performance metrics:

1. The G-mean scorer is the geometric mean of $Sensitivity = \frac{TN}{TN+FP}$ and *Sensitivity* (also known as *Recall*) [60]. Both metrics are calculated in a multiclass context considering a one-versus-all approach. For multiclass problems, the *G-mean* scorer is calculated as its average per class values:

$$G\text{-mean} = \sqrt{Sensitivity_i \times Specificity_i} \quad (3)$$

2. F-score is the harmonic mean of *Precision* and *Recall*. The two metrics are also calculated considering a one-versus-all approach. The *F-score* for the multiclass case can be calculated using its average per class values [61]:

$$F\text{-score} = 2 \frac{\overline{Precision} \times \overline{Recall}}{\overline{Precision} + \overline{Recall}} \quad (4)$$

3. OA consists of the ratio between the number of correctly classified instances and the total number of instances. This metric, because of its popularity and easy interpretability, is kept for discussion purposes. Considering C as the set of classes within a dataset, it is expressed as:

$$OA = \frac{\sum_i^C TP_i}{\sum_i^C (TP_i + FP_i)} \quad (5)$$

The comparison of classification convergence across AL frameworks and selection criteria is done using 3 AL-specific performance metrics. Particularly, we follow the

recommendations found in [22]. Each AL configuration is evaluated using the *Area Under the Learning Curve* (AULC) performance metric. It is the sum of the classification performance values of all iterations. To facilitate the analysis of the results, we fix the range of this metric between $[0, 1]$ by dividing it by the total amount of iterations (*i.e.*, the maximum performance area). The *Data Utilization Rate* (DUR) [62] metric consists of the ratio between the minimum number of instances necessary to reach a given performance threshold by an AL strategy and an equivalent baseline strategy. The deficiency score [63] is used to compare the performance between two active learners. The deficiency of algorithm A with respect to algorithm B is calculated with the areas between the respective learning curves and the maximum performance line MP :

$$\text{deficiency} = \frac{MP - AULC_A}{2MP - AULC_A - AULC_B} \quad (6)$$

This metric varies between $[0, 1]$, where values 0 and 1 are achieved by algorithms A and B , capable of achieving maximal performance from the first iteration onwards, respectively. A deficiency score of 0.5 means that active learners A and B are equivalent.

5.4. Experimental Procedure

A common practice in methodological evaluations is the implementation of an offline experiment [64]. It consists of using an existing set of labeled data as a proxy for the population of unlabeled instances. Because the dataset is already fully labeled, the oracle's typical annotation process involved in each iteration is done at zero cost. Each AL and classifier configuration is tested using a stratified 5-fold cross validation testing scheme. For each round, the larger partition is split in a stratified fashion to form a training and validation set (containing 20% of the original partition). The validation set is used to evaluate the convergence efficiency of active learners; the chooser's classification performance metrics and amount of data points used at each iteration are used to compute the AULC and DUR. Additionally, within the AL iterative process, the classifier with optimal performance on the validation set is evaluated using the test set. In order to further reduce possible initialization biases, this procedure is repeated 3 times with different seeds and the results of all runs are averaged (*i.e.*, each configuration is trained and evaluated 15 times). Finally, the maximum performance lines are calculated using the same approach. In those cases, the validation set is not used. The experimental procedure is depicted in Figure 5.

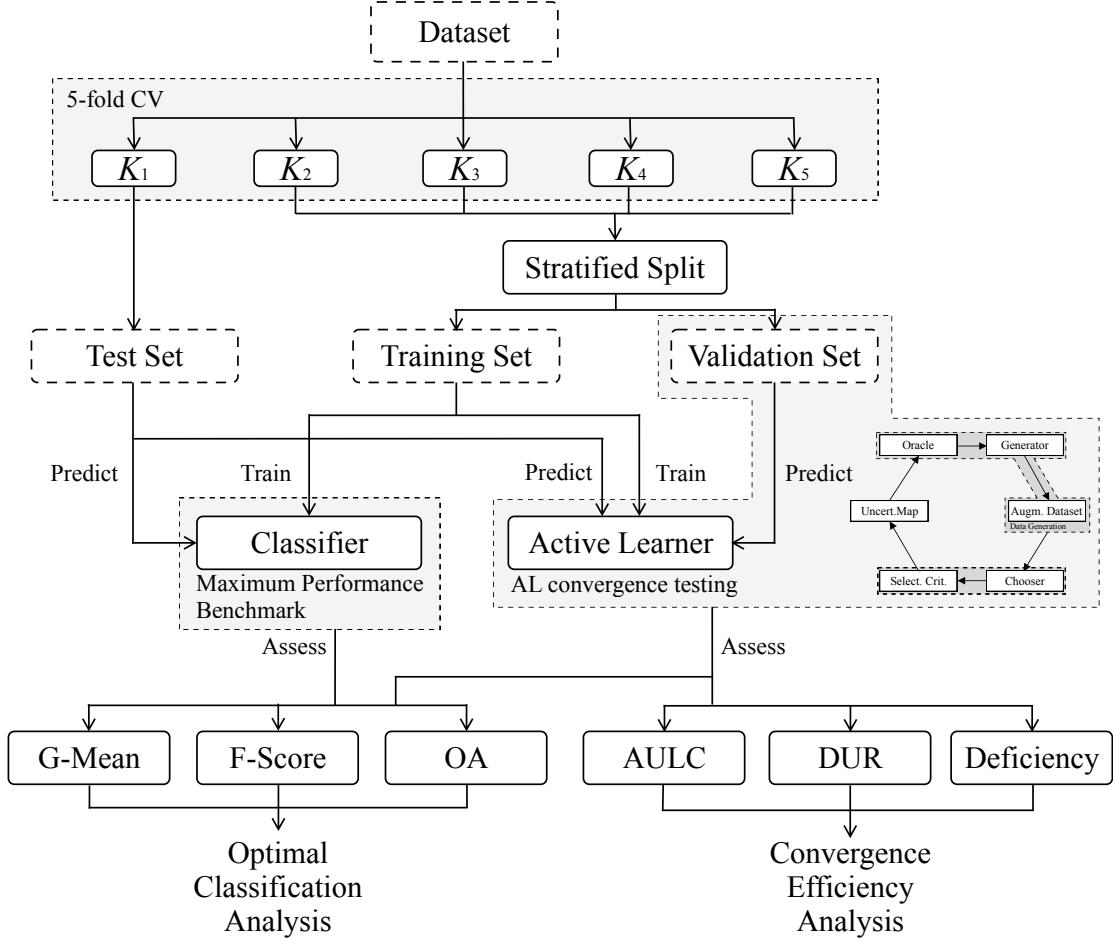


Figure 5: Experimental procedure. The performance metrics are averaged over the 5 folds across each of the 3 different initializations of this procedure for a given combination of generator, chooser/predictor and selection criterion.

To make the convergence metrics comparable across active learners, the configurations of the different frameworks must be similar. For each dataset, the number of instances is constant to facilitate the analysis of the same metrics.

In most practical AL applications it is assumed that the number of instances in the initial training sample is too small to perform hyperparameter tuning. Consequently, in order to ensure realistic results, our experimental procedure does not include hyperparameter optimization. The predefined hyperparameters are shown in Table 3. They were set up based on general recommendations and default settings for the

classifiers and generators used.

The AL iterative process is set up with a randomly selected initial training sample with 15 initial samples. At each iteration, an additional 15 samples are added to the training set. This process is stopped after 49 iterations, once 50% of the dataset is added to the augmented dataset.

Classifier	Hyperparameters	Values
LR	maximum iterations	10000
	solver	sag
	penalty	None
KNN	# neighbors	5
	weights	uniform
	metric	euclidean
RF	maximum tree depth	None
	# estimators	100
	criterion	gini
<hr/>		
Generator		
SMOTE	# neighbors	5
	# neighbors	5
G-SMOTE	deformation factor	0.5
	truncation factor	0.5

Table 3: Hyper-parameter definition for the classifiers and generators used in the experiment.

5.5. Software Implementation

The experiment was implemented using the Python programming language, along with the Python libraries Scikit-Learn [65], Imbalanced-Learn [66], Geometric-SMOTE, Cluster-Over-Sampling and Research-Learn libraries. All functions, algorithms, experiments and results are provided in the GitHub repository of the project.

6. Results & Discussion

The evaluation of the different AL frameworks in a multiple dataset context should not rely uniquely on the mean of the performance metrics across datasets. [67] recommends the use of mean ranking scores, since the performance levels of the different frameworks varies according to the data it is being used on. Consequently,

evaluating these performance metrics solely based on their mean values might lead to inaccurate analyses. Accordingly, the results of this experiment are analysed using both the mean ranking and absolute scores for each model. The rank values are assigned based on the mean scores resulting from three different initializations of 5-fold cross validation for each classifier and active learner.

6.1. Results

Table 4 shows the average rankings and standard deviations across datasets of the AULC scores for each active learner. A lower ranking value (*i.e.*, closer to 1) means a superior performance.

Classifier	Evaluation Metric	NONE	SMOTE	G-SMOTE
KNN	Accuracy	2.43 ± 0.9	2.29 ± 0.45	1.29 ± 0.45
KNN	F-score	3.00 ± 0.0	1.57 ± 0.49	1.43 ± 0.49
KNN	G-mean	3.00 ± 0.0	1.00 ± 0.0	2.00 ± 0.0
LR	Accuracy	1.71 ± 0.88	2.00 ± 0.76	2.29 ± 0.7
LR	F-score	2.43 ± 0.9	1.86 ± 0.64	1.71 ± 0.7
LR	G-mean	2.86 ± 0.35	2.00 ± 0.53	1.14 ± 0.35
RF	Accuracy	2.14 ± 0.64	2.14 ± 0.83	1.71 ± 0.88
RF	F-score	2.43 ± 0.73	2.29 ± 0.7	1.29 ± 0.45
RF	G-mean	2.57 ± 0.49	2.29 ± 0.7	1.14 ± 0.35

Table 4: Mean rankings of the AULC metric over the different datasets used in the experiment.

The mean AULC absolute scores are provided in Table 5. These values are computed as the mean of the sum of the scores of a specific performance metric over all iterations (for an AL configuration). In other words, these values correspond to the average AULC over *7 datasets* \times *5 folds* \times *3 initializations*.

Classifier	Evaluation Metric	NONE	SMOTE	G-SMOTE
KNN	Accuracy	0.811 ± 0.115	0.806 ± 0.141	0.820 ± 0.123
KNN	F-score	0.762 ± 0.131	0.796 ± 0.123	0.794 ± 0.123
KNN	G-mean	0.864 ± 0.079	0.892 ± 0.068	0.886 ± 0.073
LR	Accuracy	0.868 ± 0.114	0.868 ± 0.113	0.867 ± 0.115
LR	F-score	0.839 ± 0.119	0.843 ± 0.117	0.843 ± 0.116
LR	G-mean	0.907 ± 0.074	0.910 ± 0.071	0.911 ± 0.071
RF	Accuracy	0.851 ± 0.09	0.850 ± 0.09	0.851 ± 0.092
RF	F-score	0.810 ± 0.109	0.816 ± 0.097	0.819 ± 0.1
RF	G-mean	0.890 ± 0.068	0.896 ± 0.058	0.901 ± 0.059

Table 5: Average AULC of each AL configuration tested.

The mean deficiency scores of the proposed framework is shown in Table 6. This metric is calculated using equation 6 having the proposed framework using G-SMOTE and SMOTE as algorithms *A* and the baseline active learner as algorithm *B*. As mentioned previously, a deficiency score below 0.5 (represented in boldface) means that active learner *A* (the proposed AL framework) is more efficient than active learner *B*.

Classifier	Evaluation Metric	SMOTE	G-SMOTE
KNN	Accuracy	0.454 ± 0.129	0.373 ± 0.183
KNN	F-score	0.320 ± 0.11	0.318 ± 0.146
KNN	G-mean	0.004 ± 0.673	0.072 ± 0.651
LR	Accuracy	0.507 ± 0.038	0.509 ± 0.039
LR	F-score	0.484 ± 0.025	0.483 ± 0.025
LR	G-mean	0.462 ± 0.037	0.453 ± 0.039
RF	Accuracy	0.526 ± 0.063	0.487 ± 0.028
RF	F-score	0.490 ± 0.057	0.452 ± 0.041
RF	G-mean	0.471 ± 0.068	0.385 ± 0.095

Table 6: Mean deficiency scores. The scores were calculated by estimating the deficiency of the proposed framework with respect to the original AL framework.

The average DURs are shown in Figure 6. They were calculated for various threshold levels on each of the performance metrics used, varying at a step of 5% between 60% and 95% (represented in the vertical lines of the figures). The DURs shown in the figure use the typical AL framework as the baseline strategy. A DUR

above 1 means that the proposed framework requires that much percentage of extra data (when compared to the baseline framework) to reach that given performance threshold. Inversely, a DUR below 1 means that the proposed framework requires that much less data (as a percentage, relative to the amount of data required by the baseline framework). In other words, DUR is the percentage of data required by the proposed framework (relative to the baseline framework) to reach each of the given performance thresholds.

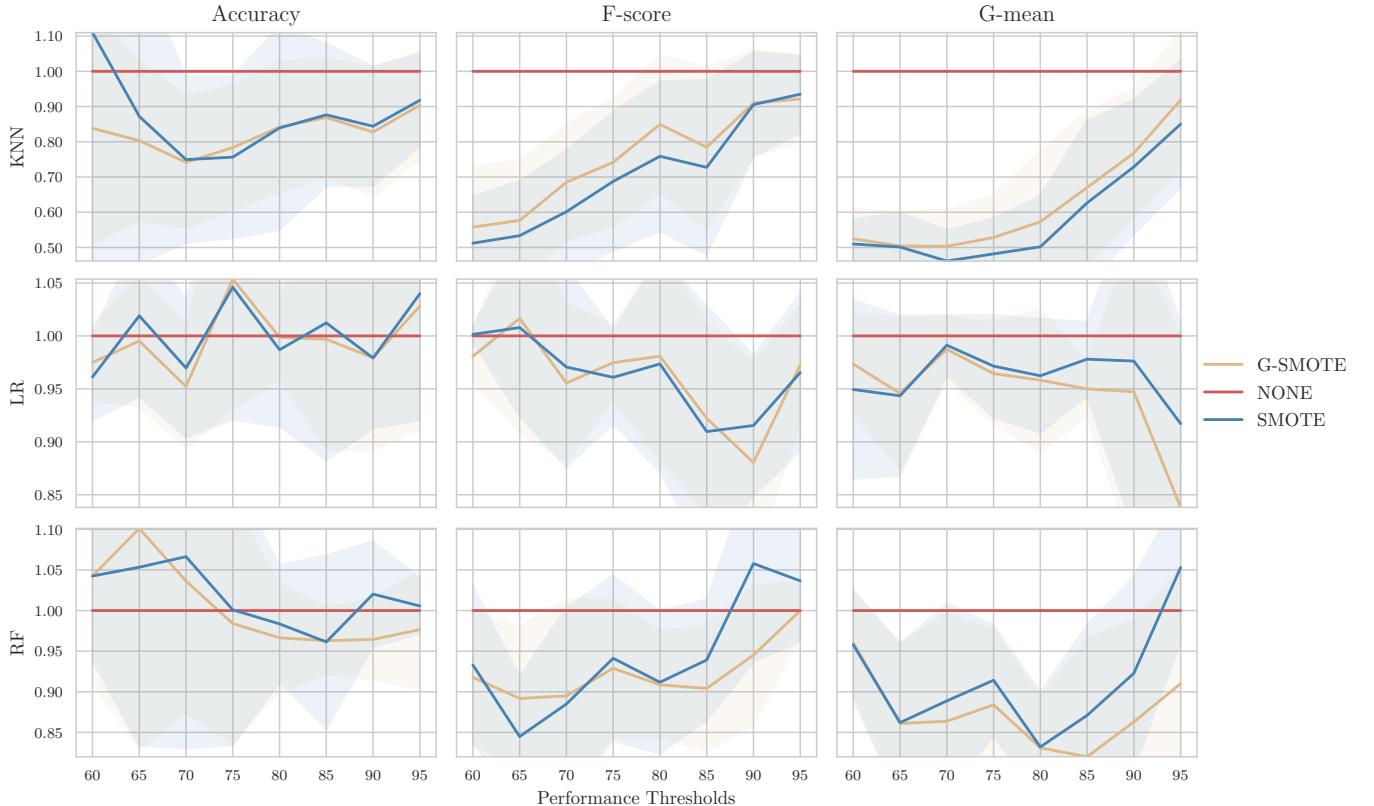


Figure 6: Mean data utilization rates. The y-axis shows the average amount of data necessary to reach the different performance thresholds (as a percentage of the baseline AL framework).

The mean optimal classification scores is shown in Table 7. The maximum performance (MP) classification scores are shown as a benchmark and represent the performance of the corresponding classifier using the entire training set. One of the goals of this study is ensuring that the classification performance of the predictors resulting from the proposed framework are not worse than the predictors produced

using the typical AL framework.

Classifier	Evaluation Metric	MP	NONE	SMOTE	G-SMOTE
KNN	Accuracy	0.864 ± 0.096	0.861 ± 0.101	0.847 ± 0.123	0.860 ± 0.102
KNN	F-score	0.838 ± 0.106	0.835 ± 0.115	0.841 ± 0.111	0.843 ± 0.105
KNN	G-mean	0.907 ± 0.063	0.904 ± 0.069	0.918 ± 0.056	0.912 ± 0.061
LR	Accuracy	0.902 ± 0.088	0.899 ± 0.094	0.897 ± 0.096	0.897 ± 0.097
LR	F-score	0.890 ± 0.084	0.883 ± 0.096	0.888 ± 0.095	0.887 ± 0.097
LR	G-mean	0.935 ± 0.052	0.931 ± 0.059	0.936 ± 0.057	0.938 ± 0.055
RF	Accuracy	0.884 ± 0.072	0.889 ± 0.07	0.890 ± 0.068	0.890 ± 0.071
RF	F-score	0.859 ± 0.083	0.866 ± 0.081	0.870 ± 0.074	0.869 ± 0.08
RF	G-mean	0.918 ± 0.051	0.921 ± 0.051	0.928 ± 0.042	0.930 ± 0.043

Table 7: Optimal classification scores. The Maximum Performance (MP) classification scores are calculated using classifiers trained using the entire training set.

The optimal AULC results of each method are reported in Table 8. These results depict in higher detail the findings drawn from this experiment.

Dataset	Classifier	Metric	NONE	SMOTE	G-SMOTE
Botswana	KNN	Accuracy	0.823 ± 0.027	0.851 ± 0.017	0.852 ± 0.018
Botswana	KNN	F-score	0.814 ± 0.02	0.856 ± 0.022	0.850 ± 0.022
Botswana	KNN	G-mean	0.898 ± 0.016	0.922 ± 0.011	0.921 ± 0.011
Botswana	LR	Accuracy	0.929 ± 0.01	0.935 ± 0.01	0.935 ± 0.011
Botswana	LR	F-score	0.929 ± 0.01	0.936 ± 0.008	0.935 ± 0.009
Botswana	LR	G-mean	0.962 ± 0.005	0.966 ± 0.005	0.966 ± 0.005
Botswana	RF	Accuracy	0.855 ± 0.02	0.856 ± 0.022	0.860 ± 0.02
Botswana	RF	F-score	0.856 ± 0.021	0.857 ± 0.023	0.861 ± 0.021
Botswana	RF	G-mean	0.920 ± 0.012	0.921 ± 0.013	0.924 ± 0.011
IP	KNN	Accuracy	0.595 ± 0.016	0.524 ± 0.026	0.586 ± 0.023
IP	KNN	F-score	0.506 ± 0.02	0.550 ± 0.024	0.552 ± 0.017
IP	KNN	G-mean	0.711 ± 0.012	0.755 ± 0.016	0.740 ± 0.015
IP	LR	Accuracy	0.620 ± 0.017	0.623 ± 0.016	0.620 ± 0.017
IP	LR	F-score	0.584 ± 0.025	0.590 ± 0.024	0.592 ± 0.025
IP	LR	G-mean	0.747 ± 0.015	0.754 ± 0.016	0.754 ± 0.015
IP	RF	Accuracy	0.682 ± 0.023	0.683 ± 0.023	0.679 ± 0.024
IP	RF	F-score	0.611 ± 0.03	0.641 ± 0.025	0.632 ± 0.029
IP	RF	G-mean	0.768 ± 0.017	0.791 ± 0.015	0.787 ± 0.015

Dataset	Classifier	Metric	NONE	SMOTE	G-SMOTE
KSC	KNN	Accuracy	0.829 ± 0.016	0.840 ± 0.014	0.846 ± 0.01
KSC	KNN	F-score	0.774 ± 0.013	0.789 ± 0.012	0.789 ± 0.014
KSC	KNN	G-mean	0.871 ± 0.008	0.885 ± 0.007	0.884 ± 0.007
KSC	LR	Accuracy	0.912 ± 0.017	0.912 ± 0.016	0.911 ± 0.015
KSC	LR	F-score	0.872 ± 0.018	0.871 ± 0.019	0.870 ± 0.016
KSC	LR	G-mean	0.929 ± 0.01	0.931 ± 0.01	0.931 ± 0.008
KSC	RF	Accuracy	0.861 ± 0.012	0.863 ± 0.009	0.866 ± 0.012
KSC	RF	F-score	0.803 ± 0.014	0.807 ± 0.009	0.811 ± 0.013
KSC	RF	G-mean	0.890 ± 0.008	0.893 ± 0.007	0.896 ± 0.008
PC	KNN	Accuracy	0.948 ± 0.012	0.954 ± 0.011	0.958 ± 0.011
PC	KNN	F-score	0.834 ± 0.038	0.857 ± 0.026	0.860 ± 0.029
PC	KNN	G-mean	0.912 ± 0.017	0.931 ± 0.011	0.931 ± 0.012
PC	LR	Accuracy	0.969 ± 0.011	0.968 ± 0.011	0.968 ± 0.01
PC	LR	F-score	0.885 ± 0.038	0.885 ± 0.03	0.886 ± 0.031
PC	LR	G-mean	0.939 ± 0.018	0.940 ± 0.015	0.941 ± 0.016
PC	RF	Accuracy	0.955 ± 0.011	0.952 ± 0.009	0.956 ± 0.01
PC	RF	F-score	0.860 ± 0.03	0.852 ± 0.024	0.865 ± 0.027
PC	RF	G-mean	0.925 ± 0.016	0.922 ± 0.014	0.932 ± 0.013
PU	KNN	Accuracy	0.749 ± 0.022	0.720 ± 0.023	0.737 ± 0.02
PU	KNN	F-score	0.654 ± 0.03	0.712 ± 0.012	0.702 ± 0.008
PU	KNN	G-mean	0.791 ± 0.02	0.847 ± 0.009	0.833 ± 0.006
PU	LR	Accuracy	0.850 ± 0.011	0.841 ± 0.019	0.840 ± 0.014
PU	LR	F-score	0.774 ± 0.023	0.784 ± 0.021	0.783 ± 0.024
PU	LR	G-mean	0.861 ± 0.013	0.876 ± 0.01	0.878 ± 0.01
PU	RF	Accuracy	0.794 ± 0.013	0.787 ± 0.014	0.791 ± 0.015
PU	RF	F-score	0.702 ± 0.038	0.726 ± 0.02	0.733 ± 0.014
PU	RF	G-mean	0.817 ± 0.021	0.839 ± 0.011	0.852 ± 0.009
Salinas	KNN	Accuracy	0.778 ± 0.021	0.787 ± 0.018	0.793 ± 0.02
Salinas	KNN	F-score	0.808 ± 0.011	0.843 ± 0.016	0.838 ± 0.018
Salinas	KNN	G-mean	0.896 ± 0.006	0.922 ± 0.007	0.916 ± 0.007
Salinas	LR	Accuracy	0.821 ± 0.014	0.822 ± 0.016	0.823 ± 0.016
Salinas	LR	F-score	0.860 ± 0.023	0.862 ± 0.014	0.863 ± 0.015
Salinas	LR	G-mean	0.924 ± 0.013	0.924 ± 0.01	0.927 ± 0.008
Salinas	RF	Accuracy	0.839 ± 0.023	0.841 ± 0.028	0.836 ± 0.029
Salinas	RF	F-score	0.871 ± 0.021	0.864 ± 0.021	0.868 ± 0.023
Salinas	RF	G-mean	0.932 ± 0.009	0.930 ± 0.009	0.934 ± 0.011

Dataset	Classifier	Metric	NONE	SMOTE	G-SMOTE
SA	KNN	Accuracy	0.956 \pm 0.006	0.968 \pm 0.011	0.970 \pm 0.012
SA	KNN	F-score	0.944 \pm 0.01	0.966 \pm 0.009	0.967 \pm 0.01
SA	KNN	G-mean	0.966 \pm 0.005	0.980 \pm 0.006	0.980 \pm 0.007
SA	LR	Accuracy	0.974 \pm 0.005	0.973 \pm 0.005	0.973 \pm 0.005
SA	LR	F-score	0.973 \pm 0.006	0.972 \pm 0.006	0.972 \pm 0.006
SA	LR	G-mean	0.983 \pm 0.005	0.983 \pm 0.005	0.983 \pm 0.005
SA	RF	Accuracy	0.969 \pm 0.012	0.968 \pm 0.012	0.969 \pm 0.01
SA	RF	F-score	0.964 \pm 0.021	0.963 \pm 0.022	0.965 \pm 0.018
SA	RF	G-mean	0.979 \pm 0.012	0.979 \pm 0.012	0.980 \pm 0.011

Table 8: Mean cross-validation AULC scores of AL algorithms for each dataset. Legend: IP Indian Pines, KSC Kennedy Space Center, PC Pavia Center, PU Pavia University, SA Salinas A.

6.2. Statistical Analysis

The methods used to test the experiment’s results must be appropriate for a multi-dataset context. Therefore the statistical analysis is performed using the Friedman test [68] and the Wilcoxon signed-rank test [69] for a post-hoc analysis. The variable used for this test is the AULC, considering the various performance metrics used.

Table 9 shows the results of the Friedman test. In most cases the null hypothesis is rejected, which indicates a statistically significant difference on the performance among AL frameworks and generators.

Classifier	Evaluation Metric	p-value	Significance
KNN	Accuracy	6.6e-02	True
KNN	F-score	5.1e-03	True
KNN	G-mean	9.1e-04	True
LR	Accuracy	5.6e-01	False
LR	F-score	3.7e-01	False
LR	G-mean	5.8e-03	True
RF	Accuracy	6.5e-01	False
RF	F-score	6.6e-02	True
RF	G-mean	1.8e-02	True

Table 9: Results for Friedman test. Statistical significance is tested at a level of $\alpha = 0.1$. The null hypothesis is that there is no difference in the classification across AL frameworks.

The Wilcoxon signed-rank test results are shown in Table 10. We test as null hypothesis that the performance of the proposed framework is the same as the original AL framework. The null hypothesis was rejected in all datasets.

Dataset	p-value	Significance
Botswana	3.9e-03	True
Indian Pines	3.9e-02	True
Kennedy Space Center	2.0e-02	True
Pavia Centre	1.2e-02	True
Pavia University	7.4e-02	True
Salinas	5.5e-02	True
Salinas A	7.4e-02	True

Table 10: Adjusted p-values using the Wilcoxon signed-rank method. Bold values are statistically significant at a level of $\alpha = 0.1$. The null hypothesis is that the performance of the proposed framework is similar to that of the original framework.

6.3. Discussion

The results found in Table 4 show that the convergence efficiency of the proposed method (*i.e.*, active learners using generators G-SMOTE and SMOTE) is consistently higher than the baseline AL framework (NONE). With the exception of two scenarios, the proposed framework using G-SMOTE was able to outperform the remaining methods. The only scenario where the baseline active learner was able to outperform the proposed AL framework was using the LR classifier and OA as the optimization goal.

We use the AULC metric to compare the convergence efficiency among the different AL architectures. The comparison of the mean AULC scores in Table 5 shows a significant performance superiority across active learners in some situations, particularly when the classifier KNN is used. The proposed AL framework using G-SMOTE as generator improved the baseline performance in all settings except one. When analysing the performance of the two active learners using a data generation mechanism jointly, we find that the mean performance of the proposed framework can be optimized to be better than or as good as the baseline framework. This optimization may be owed to either (or both):

1. The earlier convergence of the proposed AL framework, therefore requiring less data to achieve comparable performance levels. This increases the AULC scores before convergence is reached. This effect is studied with the various data utilization rates presented in Figure 6.

2. Higher optimal classification performance, therefore reaching higher performance levels overall. This increases the AULC scores once convergence is reached. This effect analyzed using the optimal classification scores shown in Table 7.

Through the analysis both the MP threshold and the optimal classification scores of AL algorithms shown in Table 7, we can confirm that in general the classification scores are high. This makes the variability of AULC among frameworks particularly meaningful, since this could limit the magnitude of the differences among frameworks. Although, the proposed framework was capable of outperforming the baseline framework and the maximum performance classification threshold in 7 out of 9 tests. Therefore, the increased classification quality of the proposed framework is one of the causes of higher AULC scores.

The results found in Figure 6 show a generalized decrease of data required to reach the performance thresholds in the various scenarios. For higher performance thresholds, the gap between the proposed and baseline methods tend to be reduced, since the amount of data required is larger and the benefits of data generation is less apparent. These results reflect the effect described previously, showing that AULC scores also improve due to the earlier convergence of the AL iterative process, since the amount of data necessary to reach a given performance threshold is generally reduced with the proposed AL framework.

The deficiency scores measured in terms of the baseline AL framework in Table 6 are used to measure the overall improvement or decline of the proposed AL framework. The results found with this metric show a consistent improvement over the baseline method with the exception of three situations, where the difference among active learners is marginal.

We used Table 8 to compare the performance of the different active learners across the different datasets, classifiers and framework settings. We found that in 54 out of 63 cases (approximately 86%) the proposed framework outperformed (or was as good as) the baseline. The datasets where the proposed method had more difficulties outperforming the baseline active learner were the Pavia University and Salinas A datasets, which is consistent with the p-values calculated in the Wilcoxon signed-rank test shown in Table 10. Although, the results were still statistically significant.

In most tests, our results were statistically significant. These tests were done considering two different perspectives:

1. Examining whether there is a difference across AL frameworks and generators used by conducting a Friedman test. The results of this statistical test is

shown in Table 9, where 6 out of 9 tests were statistically significant. The 3 tests whose null hypothesis was not rejected indicate the occasions in which the performance of the two configurations of the proposed framework was not significantly different both among themselves and the baseline method, which can be visually attested in Figure 6. Although, while Table 4 shows that in these situations the proposed framework performs better than the baseline method, the difference is less apparent in terms of DUR, as shown in Figure 6. Regardless, the quality the classification outputs are not affected by this difference, as in 2 out of those 3 situations the proposed framework was still capable of producing better predictions than the baseline. This proves that in our experiment, the implementation of the proposed framework did not negatively affect the quality of the predictions.

2. Examining whether the proposed framework outperforms the baseline framework with statistical significance on each dataset across different classifiers and performance metrics. This is done by comparing the optimal mean AULC scores of both the baseline method and the proposed method for each combination of classifiers and performance metrics, where the null hypothesis was rejected in all datasets. This indicates that regardless of the context under which an AL algorithm is used, the proposed framework outperforms the baseline framework with statistical significance.

This paper introduces the concept of applying data generation algorithm in the AL framework. This was done with the implementation of a popular data generation algorithm and a recent, SOTA generalization of this same method. Although, since these algorithms are based on heuristics, future work should focus on improving these results through the implementation of more sophisticated data generation mechanisms, at the cost of additional computational power. In addition, we also noticed significant standard errors in our experimental results. This indicates that AL procedures seem to be particularly sensitive to the initialization method, which is still a limitation of AL, regardless of the framework used. This is consistent with the findings in [22], which future work should attempt to address. Although the Generator component marginally reduced this standard error found in AULC and optimal classification scores found in Tables 5 and 7, it is not sufficient to address this specific limitation.

7. Conclusion

The aim of this experiment was to test the effectiveness of a new AL framework, where a new element is introduced to improve the convergence rate of AL through the use of artificial data generation. The experiment was designed to test the proposed method under particularly challenging conditions, where the maximum performance line is naturally high in most datasets (with exception to the Indian Pines dataset). In order to test basic setups for this new framework, the elements that constitute the Generator component were set up in a plug-and-play scheme, without significant tuning of the data generators (SMOTE and G-SMOTE). The tests showed that this new framework is able to consistently outperform the original AL framework in most scenarios, as shown in Table 8. These results could be further improved through the modification and more intense tuning of the data generation strategy. In our experiment, during each iteration, the new artificial data is generated only to match each non-majority class frequency with the majority class frequency, thus strictly balancing the class distribution. Generating a larger amount of data for all classes (especially in early iterations) can further improve these results.

We also consider the fast convergence of AL on these datasets. The high performance scores for the baseline AL framework made the achievement of significant improvements over the traditional AL framework under these conditions particularly meaningful. The advantage of the proposed AL framework is shown in Figure 6. In most of the presented scenarios there is a substantial reduction of necessary data to reach each of the tested performance metric thresholds.

The results from this experiment show that the inclusion of a data generator in the AL framework will yield significant improvements in the convergence of the method. The proposed method successfully anticipated the predictor’s optimal performance, as shown in Tables 4, 5 and 6. This means the annotation cost would have been reduced in a real application since the number of iterations and labeled instances necessary to reach near optimal classification performance is reduced, as shown in Figure 6. The class imbalance bias observed in AL tasks also is reduced, as shown in Table 6, where data imbalance appropriate metrics are always improved over the baseline scores.

8. Acknowledgements

The authors would like to thank Professor Victor Lobo (NOVA IMS, Universidade Nova de Lisboa, and CINAV, Escola Naval, CIDIUM) for reviewing this paper and providing important feedback throughout its development.

9. Funding

This research was funded by “Fundação para a Ciência e Tecnologia” (Portugal) [grant numbers PCIF/SSI/0102/2017 - foRESTER, DSAIPA/AI/0100/2018 - IPSTERS].

References

- [1] S. Nagai, K. N. Nasahara, T. K. Akitsu, T. M. Saitoh, H. Muraoka, Importance of the Collection of Abundant Ground-Truth Data for Accurate Detection of Spatial and Temporal Variability of Vegetation by Satellite Remote Sensing, in: Biogeochemical Cycles: Ecological Drivers and Environmental Impact, American Geophysical Union (AGU), 2020, pp. 223–244. doi:10.1002/9781119413332.ch11.
URL <https://onlinelibrary.wiley.com/doi/abs/10.1002/9781119413332.ch11>
- [2] Y. Huang, Z. xin CHEN, T. YU, X. zhi HUANG, X. fa GU, Agricultural remote sensing big data: Management and applications, Journal of Integrative Agriculture 17 (9) (2018) 1915–1931. doi:10.1016/S2095-3119(17)61859-8.
- [3] X. Wang, H. Xie, A review on applications of remote sensing and geographic information systems (GIS) in water resources and flood risk management, Water (Switzerland) 10 (5) (2018) 608. doi:10.3390/w10050608.
URL <http://www.mdpi.com/2073-4441/10/5/608>
- [4] R. Khatami, G. Mountarakis, S. V. Stehman, A meta-analysis of remote sensing research on supervised pixel-based land-cover image classification processes: General guidelines for practitioners and future research, Remote Sensing of Environment 177 (2016) 89–100. doi:10.1016/J.RSE.2016.02.028.
URL <https://www.sciencedirect.com/science/article/pii/S0034425716300578>
- [5] A. B. Gavade, V. S. Rajpurohit, Systematic analysis of satellite image-based land cover classification techniques: literature review and challenges, International Journal of Computers and Applications (2019) 1–10 doi:10.1080/1206212x.2019.1573946.
- [6] M. A. Wulder, N. C. Coops, D. P. Roy, J. C. White, T. Hermosilla, Land cover 2.0, International Journal of Remote Sensing 39 (12) (2018) 4254–4284. doi:10.1080/01431161.2018.1452075.

- [7] H. Costa, P. Benevides, F. Marcelino, M. Caetano, Introducing automatic satellite image processing into land cover mapping by photo-interpretation of airborne data, *The International Archives of Photogrammetry, Remote Sensing and Spatial Information Sciences* 42 (2020) 29–34.
- [8] E. F. Vermote, S. Skakun, I. Becker-Reshef, K. Saito, Remote sensing of coconut trees in tonga using very high spatial resolution worldview-3 data, *Remote Sensing* 12 (19) (2020) 3113.
- [9] D. Costantino, M. Pepe, G. Dardanelli, V. Baiocchi, USING OPTICAL SATELLITE AND AERIAL IMAGERY FOR AUTOMATIC COASTLINE MAPPING, *Geographia Technica* (2020) 171–190 doi:10.21163/gt_2020.152.17.
- [10] A. E. Maxwell, T. A. Warner, F. Fang, Implementation of machine-learning classification in remote sensing: An applied review, *International Journal of Remote Sensing* 39 (9) (2018) 2784–2817. doi:10.1080/01431161.2018.1433343. URL <https://doi.org/10.1080/01431161.2018.1433343>
- [11] V. Růžička, S. D'Aronco, J. D. Wegner, K. Schindler, Deep active learning in remote sensing for data efficient change detection, *arXiv preprint arXiv:2008.11201* (2020).
- [12] S.-J. Liu, H. Luo, Q. Shi, Active Ensemble Deep Learning for Polarimetric Synthetic Aperture Radar Image Classification, *IEEE Geoscience and Remote Sensing Letters* (2020) 1–5 arXiv:2006.15771, doi:10.1109/lgrs.2020.3005076.
- [13] T. Su, S. Zhang, T. Liu, Multi-spectral image classification based on an object-based active learning approach, *Remote Sensing* 12 (3) (2020) 504. doi:10.3390/rs12030504. URL <https://www.mdpi.com/2072-4292/12/3/504>
- [14] E. Pasolli, H. L. Yang, M. M. Crawford, Active-metric learning for classification of remotely sensed hyperspectral images, *IEEE Transactions on Geoscience and Remote Sensing* 54 (4) (2016) 1925–1939. doi:10.1109/TGRS.2015.2490482.
- [15] C. Pelletier, S. Valero, J. Inglada, N. Champion, C. M. Sicre, G. Dedieu, Effect of training class label noise on classification performances for land cover mapping with satellite image time series, *Remote Sensing* 9 (2) (2017) 173. doi:10.3390/rs9020173. URL <http://www.mdpi.com/2072-4292/9/2/173>

- [16] O. Stromann, A. Nascetti, O. Yousif, Y. Ban, Dimensionality Reduction and Feature Selection for Object-Based Land Cover Classification based on Sentinel-1 and Sentinel-2 Time Series Using Google Earth Engine, *Remote Sensing* 12 (1) (2020) 76. doi:10.3390/RS12010076.
- [17] F. Alonso-Sarria, C. Valdivieso-Ros, F. Gomariz-Castillo, Isolation forests to evaluate class separability and the representativeness of training and validation areas in land cover classification, *Remote Sensing* 11 (24) (2019) 3000. doi: 10.3390/rs11243000.
URL <https://www.mdpi.com/2072-4292/11/24/3000>
- [18] W. Feng, W. Huang, H. Ye, L. Zhao, Synthetic minority over-sampling technique based rotation forest for the classification of unbalanced hyperspectral data, in: International Geoscience and Remote Sensing Symposium (IGARSS), Vol. 2018-July, Institute of Electrical and Electronics Engineers Inc., 2018, pp. 2651–2654. doi:10.1109/IGARSS.2018.8518242.
- [19] N. V. Chawla, N. Japkowicz, A. Kotcz, Editorial: Special Issue on Learning from Imbalanced Data Sets, *ACM SIGKDD Explorations Newsletter* 6 (1) (2004) 1–6. doi:10.1145/1007730.1007733.
URL <http://portal.acm.org/citation.cfm?doid=1007730.1007733>
- [20] A. Fernández, V. López, M. Galar, M. J. del Jesus, F. Herrera, Analysing the classification of imbalanced data-sets with multiple classes: Binarization techniques and ad-hoc approaches, *Knowledge-Based Systems* 42 (2013) 97–110. doi:10.1016/J.KNOSYS.2013.01.018.
URL <https://www.sciencedirect.com/science/article/abs/pii/S0950705113000300>
- [21] H. Kaur, H. S. Pannu, A. K. Malhi, A systematic review on imbalanced data challenges in machine learning: Applications and solutions, *ACM Computing Surveys* 52 (4) (2019) 1–36. doi:10.1145/3343440.
URL <http://dl.acm.org/citation.cfm?doid=3359984.3343440>
- [22] D. Kottke, A. Calma, D. Huseljic, G. Kreml, B. Sick, Challenges of reliable, realistic and comparable active learning evaluation, in: CEUR Workshop Proceedings, Vol. 1924, 2017, pp. 2–14.
URL <http://dspace.library.uu.nl/handle/1874/359528>
- [23] Y. Sverchkov, M. Craven, A review of active learning approaches to experimental design for uncovering biological networks, *PLOS Computational Biology* 13 (6)

- (2017) e1005466. doi:10.1371/journal.pcbi.1005466.
URL <https://dx.plos.org/10.1371/journal.pcbi.1005466>
- [24] X. Li, Y. Guo, Active learning with multi-label svm classification, in: In IJCAI, 2013, pp. 1479–1485.
- [25] J. Li, X. Huang, X. Chang, A label-noise robust active learning sample collection method for multi-temporal urban land-cover classification and change analysis, ISPRS Journal of Photogrammetry and Remote Sensing 163 (January) (2020) 1–17. doi:10.1016/j.isprsjprs.2020.02.022.
URL <https://doi.org/10.1016/j.isprsjprs.2020.02.022>
- [26] V. K. Shrivastava, M. K. Pradhan, Hyperspectral Remote Sensing Image Classification Using Active Learning, in: Studies in Computational Intelligence, Vol. 907, Springer, 2021, pp. 133–152. doi:10.1007/978-3-030-50641-4_8.
URL https://link.springer.com/chapter/10.1007/978-3-030-50641-4_8
- [27] D. Tuia, E. Pasolli, W. J. Emery, Using active learning to adapt remote sensing image classifiers, Remote Sensing of Environment 115 (9) (2011) 2232–2242.
- [28] G. Cawley, Baseline Methods for Active Learning, Proceedings of Active Learning and Experimental Design workshop In conjunction with AISTATS 16 (2011) 47–57.
URL <http://proceedings.mlr.press/v16/cawley11a.html>
- [29] I. Muslea, S. Minton, C. A. Knoblock, Active learning with multiple views, Journal of Artificial Intelligence Research 27 (2006) 203–233. arXiv:1110.1073, doi:10.1613/jair.2005.
URL <https://www.jair.org/index.php/jair/article/view/10470>
- [30] W. Di, M. M. Crawford, View generation for multiview maximum disagreement based active learning for hyperspectral image classification, IEEE Transactions on Geoscience and Remote Sensing 50 (5 PART 2) (2012) 1942–1954. doi:10.1109/TGRS.2011.2168566.
- [31] X. Zhou, S. Prasad, M. Crawford, Wavelet domain multi-view active learning for hyperspectral image analysis, in: Workshop on Hyperspectral Image and Signal Processing, Evolution in Remote Sensing, Vol. 2014-June, IEEE Computer Society, 2014. doi:10.1109/WHISPERS.2014.8077528.

- [32] Z. Zhang, E. Pasolli, H. L. Yang, M. M. Crawford, Multimetric Active Learning for Classification of Remote Sensing Data, *IEEE Geoscience and Remote Sensing Letters* 13 (7) (2016) 1007–1011. doi:10.1109/LGRS.2016.2560623.
- [33] H. H. Aghdam, A. Gonzalez-Garcia, A. Lopez, J. Weijer, Active learning for deep detection neural networks, in: Proceedings of the IEEE International Conference on Computer Vision, Vol. 2019-Octob, 2019, pp. 3671–3679. arXiv:1911.09168, doi:10.1109/ICCV.2019.00377.
URL www.gitlab.com/haghdam/deep{ }active{ }learning
- [34] D. Tuia, F. Ratle, F. Pacifici, M. F. Kanevski, W. J. Emery, Active learning methods for remote sensing image classification, *IEEE Transactions on Geoscience and Remote Sensing* 47 (7) (2009) 2218–2232. doi:10.1109/TGRS.2008.2010404.
- [35] L. Copa, D. Tuia, M. Volpi, M. Kanevski, Unbiased query-by-bagging active learning for VHR image classification, in: L. Bruzzone (Ed.), *Image and Signal Processing for Remote Sensing XVI*, Vol. 7830, SPIE, 2010, p. 78300K. doi:10.1117/12.864861.
URL <http://proceedings.spiedigitallibrary.org/proceeding.aspx?doi=10.1117/12.864861>
- [36] B. Demir, C. Persello, L. Bruzzone, Batch-mode active-learning methods for the interactive classification of remote sensing images, *IEEE Transactions on Geoscience and Remote Sensing* 49 (3) (2011) 1014–1031. doi:10.1109/TGRS.2010.2072929.
- [37] J. Li, J. M. Bioucas-Dias, A. Plaza, Hyperspectral image segmentation using a new bayesian approach with active learning, *IEEE Transactions on Geoscience and Remote Sensing* 49 (10 PART 2) (2011) 3947–3960. doi:10.1109/TGRS.2011.2128330.
- [38] W. Liu, J. Yang, P. Li, Y. Han, J. Zhao, H. Shi, A novel object-based supervised classification method with active learning and random forest for PolSAR imagery, *Remote Sensing* 10 (7) (2018). doi:10.3390/rs10071092.
- [39] T. Luo, K. Kramer, D. Goldgof, L. O. Hall, S. Samson, A. Remsen, T. Hopkins, Learning to recognize plankton, in: Proceedings of the IEEE International Conference on Systems, Man and Cybernetics, Vol. 1, 2003, pp. 888–893. doi:10.1109/icsmc.2003.1243927.

- [40] J. Li, J. M. Bioucas-Dias, A. Plaza, Spectral-spatial classification of hyperspectral data using loopy belief propagation and active learning, *IEEE Transactions on Geoscience and Remote Sensing* 51 (2) (2013) 844–856. doi: 10.1109/TGRS.2012.2205263.
- [41] D. Yoo, I. S. Kweon, Learning loss for active learning, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2019.
- [42] S. Ertekin, J. Huang, C. L. Giles, Active learning for class imbalance problem, in: Proceedings of the 30th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR’07, ACM Press, New York, New York, USA, 2007, pp. 823–824. doi:10.1145/1277741.1277927. URL <http://portal.acm.org/citation.cfm?doid=1277741.1277927>
- [43] T. DeVries, G. W. Taylor, Dataset augmentation in feature space, in: 5th International Conference on Learning Representations, ICLR 2017 - Workshop Track Proceedings, International Conference on Learning Representations, ICLR, 2017. arXiv:1702.05538. URL <https://arxiv.org/abs/1702.05538v1>
- [44] N. V. Chawla, K. W. Bowyer, L. O. Hall, W. P. Kegelmeyer, SMOTE: Synthetic minority over-sampling technique, *Journal of Artificial Intelligence Research* 16 (2002) 321–357. arXiv:1106.1813, doi:10.1613/jair.953. URL <https://jair.org/index.php/jair/article/view/10302>
- [45] S. E. Jozdani, B. A. Johnson, D. Chen, Comparing Deep Neural Networks, Ensemble Classifiers, and Support Vector Machine Algorithms for Object-Based Urban Land Use/Land Cover Classification, *Remote Sensing* 11 (14) (2019) 1713. doi:10.3390/rs11141713. URL <https://www.mdpi.com/2072-4292/11/14/1713>
- [46] C. Bogner, B. Seo, D. Rohner, B. Reineking, Classification of rare land cover types: Distinguishing annual and perennial crops in an agricultural catchment in South Korea, *PLoS ONE* 13 (1) (jan 2018). doi:10.1371/journal.pone.0190476.
- [47] G. Douzas, F. Bacao, Geometric SMOTE a geometrically enhanced drop-in replacement for SMOTE, *Information Sciences* 501 (2019) 118–135. doi:10.1016/j.ins.2019.06.007.

- URL <https://www.sciencedirect.com/science/article/pii/S0020025519305353?via%}3Dihub>
- [48] H. Han, W. Y. Wang, B. H. Mao, Borderline-SMOTE: A new over-sampling method in imbalanced data sets learning, Lecture Notes in Computer Science 3644 (PART I) (2005) 878–887. doi:10.1007/11538059_91.
- [49] H. He, Y. Bai, E. A. Garcia, S. Li, ADASYN: Adaptive synthetic sampling approach for imbalanced learning, in: Proceedings of the International Joint Conference on Neural Networks, IEEE, 2008, pp. 1322–1328. doi:10.1109/IJCNN.2008.4633969.
URL <http://ieeexplore.ieee.org/document/4633969/>
- [50] G. Douzas, F. Bacao, F. Last, Improving imbalanced learning through a heuristic oversampling method based on k-means and SMOTE, Information Sciences 465 (2018) 1–20. doi:10.1016/j.ins.2018.06.056.
- [51] G. Douzas, F. Bacao, Self-Organizing Map Oversampling (SOMO) for imbalanced data set learning, Expert Systems with Applications 82 (2017) 40–52. doi:10.1016/j.eswa.2017.03.073.
- [52] G. Douzas, F. Bacao, J. Fonseca, M. Khudinyan, Imbalanced learning in land cover classification: Improving minority classes' prediction accuracy using the geometric SMOTE algorithm, Remote Sensing 11 (24) (2019) 3040. doi:10.3390/rs11243040.
URL <https://www.mdpi.com/2072-4292/11/24/3040>
- [53] M. F. Baumgardner, L. L. Biehl, D. A. Landgrebe, 220 band aviris hyperspectral image data set: June 12, 1992 indian pine test site 3 (Sep 2015). doi:doi: /10.4231/R7RX991C.
URL <https://purrr.purdue.edu/publications/1947/1>
- [54] T. Cover, P. Hart, Nearest neighbor pattern classification, IEEE Transactions on Information Theory 13 (1) (1967) 21–27. doi:10.1109/TIT.1967.1053964.
URL <http://ieeexplore.ieee.org/lpdocs/epic03/wrapper.htm?arnumber=1053964>
- [55] J. A. Nelder, R. W. Wedderburn, Generalized linear models, Journal of the Royal Statistical Society: Series A (General) 135 (3) (1972) 370–384.

- [56] T. K. Ho, Random decision forests, in: Proceedings of the Third International Conference on Document Analysis and Recognition (Volume 1) - Volume 1, ICDAR '95, IEEE Computer Society, USA, 1995, p. 278.
- [57] P. Olofsson, G. M. Foody, S. V. Stehman, C. E. Woodcock, Making better use of accuracy data in land change studies: Estimating accuracy and area and quantifying uncertainty using stratified estimation, *Remote Sensing of Environment* 129 (2013) 122–131. doi:10.1016/j.rse.2012.10.031.
- [58] R. G. Pontius, M. Millones, Death to Kappa: Birth of quantity disagreement and allocation disagreement for accuracy assessment, *International Journal of Remote Sensing* 32 (15) (2011) 4407–4429. doi:10.1080/01431161.2011.552923.
- [59] L. A. Jeni, J. F. Cohn, F. De La Torre, Facing imbalanced data - Recommendations for the use of performance metrics, in: Proceedings - 2013 Humaine Association Conference on Affective Computing and Intelligent Interaction, ACII 2013, 2013, pp. 245–251. doi:10.1109/ACII.2013.47.
- [60] M. Kubat, S. Matwin, et al., Addressing the curse of imbalanced training sets: one-sided selection, in: Icml, Vol. 97, Citeseer, 1997, pp. 179–186.
- [61] H. He, E. A. Garcia, Learning from Imbalanced Data, *IEEE Transactions on Knowledge and Data Engineering* 21 (9) (2009) 1263–1284. arXiv:arXiv:1011.1669v3, doi:10.1109/TKDE.2008.239.
- [62] T. Reitmaier, B. Sick, Let us know your decision: Pool-based active training of a generative classifier with the selection strategy 4ds, *Information Sciences* 230 (2013) 106–131.
- [63] E. Yanik, T. M. Sezgin, Active learning for sketch recognition, *Computers & Graphics* 52 (2015) 93–105.
- [64] J.-F. Kagy, T. Kayadelen, J. Ma, A. Rostamizadeh, J. Strnadova, The practical challenges of active learning: Lessons learned from live experimentation (2019). arXiv:1907.00038.
- [65] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, É. Duchesnay, Scikit-learn: Machine Learning in Python, *Journal of Machine Learning Research* 12 (Oct) (2011) 2825–2830.
URL <http://jmlr.csail.mit.edu/papers/v12/pedregosa11a.html>

- [66] G. Lemaître, F. Nogueira, C. K. Aridas, Imbalanced-learn: A python toolbox to tackle the curse of imbalanced datasets in machine learning, *Journal of Machine Learning Research* 18 (17) (2017) 1–5.
URL <http://jmlr.org/papers/v18/16-365.html>
- [67] J. Demšar, Statistical comparisons of classifiers over multiple data sets, *Journal of Machine learning research* 7 (Jan) (2006) 1–30.
- [68] M. Friedman, The use of ranks to avoid the assumption of normality implicit in the analysis of variance, *Journal of the american statistical association* 32 (200) (1937) 675–701.
- [69] F. Wilcoxon, Individual Comparisons by Ranking Methods, *Biometrics Bulletin* 1 (6) (1945) 80. doi:10.2307/3001968.